

Jan-Cornelius Molnar

Nichtparametrische Statistik und statistische Lerntheorie

Aufbauend auf einer Vorlesung von Prof. Dr. I. Steinwart

Stuttgart, Sommersemester 2010

Version vom 1. März 2012

Diese Rohfassung enthält sicher noch zahlreiche Druckfehler. Für Hinweise auf solche und Kommentare jeder Art bin ich stets dankbar.¹ Viel Spaß!

¹Jan-Cornelius Molnar jan.molnar@studentpartners.de

Inhaltsverzeichnis

0 Grundlagen	5
o-1 Statistische Lerntheorie (SLT)	5
o-1-1 Inhalt der Vorlesung	5
o-2 Lernszenarien und -ziele	5
o-2-1 Überwachtes Lernen	5
o-2-2 Halbüberwachtes Lernen	7
o-2-3 Unüberwachtes Lernen	8
o-2-4 Weitere Lernszenarien	10
o-2-5 Datengenerierung	10
o-3 Klassische Verfahren zur Klassifikation	11
1 Grundkonzepte der statistischen Lerntheorie	15
1-1 Grundannahmen	15
1-2 Verlustfunktionen und Risiken	17
1-3 Universelle Konsistenz	23
1-3-1 Plug-in-rules für Klassifikation	24
1-4 Lernraten	36
2 Verlustfunktionen	47
2-1 Eigenschaften von Verlustfunktionen	47
2-2 Margin basierte Verlustfunktionen	52
2-3 Distanzbasierte Verlustfunktionen	56
3 Konzentrationsungleichungen	59
4 Empirische Risikominimierung (ERM)	65
4-1 ERM über endliche Funktionenmengen	65
4-2 ERM über unendlichen Funktionenmengen	67

5	Reproduzierende Kernhilberträume (RKHS)	73
5-1	Kerne und Beispiele	73
5-2	RKHS (Reproduzierende Kern-Hilbert-Räume)	80
5-3	Eigenschaften von Kernen und RKHS	86
5-4	Große RKHS	93
6	Support Vector Machines (SVMs)	99
6-1	Definition, einfache Eigenschaften und Beispiele	99
6-1-1	Ausflug in das Reich der konvexen Analysis	105
6-2	Orakelungleichungen für SVMs	111
6-3	Die Funktion $A(\lambda)$	117
6-4	Konsistenz und Lernraten für SVMs	123
7	Verlustfunktionen II	131
7-1	Marginbasierte Verlustfunktionen	131
7-2	Distanzbasierte Verlustfunktionen	133

o Grundlagen

o-1 Statistische Lerntheorie (SLT)

Die statistische Lerntheorie ist Teil der nichtparametrischen Statistik sowie der Informatik. Sie bildet den mathematischen Arm des maschinellen Lernens und befasst sich sowohl mit den theoretischen Grundlagen wie auch Fragen der Implementierung und besitzt zahlreiche Anwendungen.

o-1-1 Inhalt der Vorlesung

- Was ist Lernen?
- Klassische Verfahren
- Empirische Risikominimierung (ERM)
- Regularisierte ERMs
- Kernbasierte Verfahren (support vector machines, SVMs)
- Verlustfunktionen

o-2 Lernszenarien und -ziele

o-2-1 Überwachtes Lernen

Beim überwachten Lernen wird der Lernprozess durch einen **Lehrer** geführt. Den **Eingaberaum** bezeichnen wir mit X , den **Ausgaberaum** mit Y . Weiterhin sind **Trainingsdaten** in der Form

$$((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$$

gegeben. Das Ziel ist es eine Funktion

$$f_D : X \rightarrow Y$$

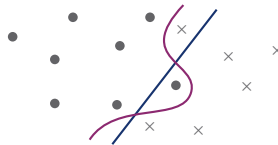
zu so finden, dass $f_D(x) \approx y$ für neue Paare (x, y) . Es geht also nicht darum, die Daten einfach auswendig zu lernen, sondern y für zukünftiges x vorherzusagen. Was “ \approx ” genau bedeutet werden wir später erklären.

Beim überwachten Lernen, erhält der Lehrer ein x_i und ordnet diesem ein **Label** y_i zu. Der Lernende erhält dann vom Lehrer das Tupel (x_i, y_i) und soll darauf f_D ableiten.

Klassifikation

Gegeben sei der Ausgaberaum $Y = \{-1, 1\}$ und Ziel ist f_D so zu bestimmen, dass $f_D(x) = y$ „möglichst häufig“.

Es ist nun grundsätzlich zu entscheiden, ob man den Trainingsdaten traut oder eher ein Rauschen vermutet.



o.1 Zwei Möglichkeiten zur Klassifikation von Trainingsdaten.

Geht man davon aus, dass auf den Trainingsdaten keine Fehler sind, besteht die Gefahr des sogenannten **overfitting**, d.h. man traut den Daten zu sehr. Vermutet man andererseits ein Rauschen besteht die Gefahr des **underfitting**, d.h. die Annahmen an die Zielfunktion sind unzureichend.

Over- und underfitting konkurrieren miteinander. Gute Lernverfahren finden einen guten bzw. den “besten” Kompromiss.

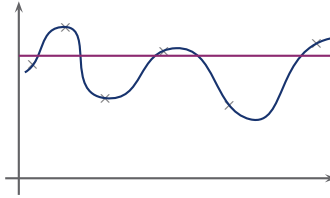
Anwendungen der Klassifikation

- Spamfilter
- Fraud Detection (Betrugsvorfälle bei Kreditkartenunternehmen)

- Schriftenerkennung
- Diagnoseverfahren

Regression

Sei $Y = [-M, M]$ oder $Y = \mathbb{R}$. Auch hier ist das Ziel $f_D(x) \approx y$ möglichst häufig und es besteht ebenfalls die Gefahr des over- bzw. underfittings.



0.2 Zwei Möglichkeiten zur Regression (datentreu oder linear).

Was man genau unter “ \approx ” versteht, fällt unter den Terminus **Verlustfunktion**. Je nachdem wie man “ \approx ” mathematisch präzisiert, erreicht man unterschiedliche Lernziele.

Anwendung der Regression

- Finanzindustrie: Kursvorhersage von Aktien, Ausfallrisiken von Krediten, ...
- Data Mining
- Klassifikation $P(Y = y | x) = ?$
- Schätzungen von Restlebenszeiten von Teilen

0-2-2 Halbüberwachtes Lernen

Hier sind X und Y wie beim überwachten Lernen definiert. Die Trainingsdaten sind jedoch nur noch zum Teil mit einem Label versehen. Sie bestehen also aus

einer „kleinen“ Menge an

$$((x_1, y_1), \dots, (x_n, y_n)) \in (X, Y)^n,$$

wobei $n = 2$ den Extremfall bildet, sowie einer „großen“ Menge an

$$(x_{n+1}, \dots, x_m) \in X^{m-n}, \quad m \gg 1.$$

Nun kann man ebenfalls Klassifikation oder Regression betreiben.

Die Motivation hierbei ist, dass die Beschaffung der y_i entweder “teuer” oder gar unmöglich ist. Beispiele dafür sind medizinische Studien oder Diagnoseverfahren zur Überwachung von Brücken.

Die Methoden für das halbüberwachte Lernen sind noch weitgehend unausgereift. Es gibt noch sehr viele Fragestellungen an denen man arbeiten kann.

0-2-3 Unüberwachtes Lernen

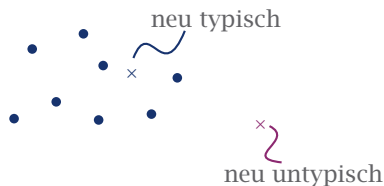
Beim unüberwachten Lernen sind nur Trainingsdaten ohne Label gegeben

$$x_1, \dots, x_n \in X.$$

Es lassen sich nun verschiedene Lernziele definieren. Wir wollen hier nur einen kurzen Überblick geben; ausführlich werden wir uns lediglich mit der Ausreißeridentifikation beschäftigen.

Ausreißeridentifikation

Es geht darum “untypische” Daten zu identifizieren und von „typischen“ zu unterscheiden.

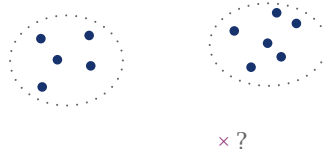


0.3 Zwei Möglichkeiten zur Klassifikation von Trainingsdaten.

Die Ausreißeridentifikation hat viele Ähnlichkeiten mit der Klassifikation und bereits gut verstanden.

Clusteranalyse

Hier ist zunächst festzustellen, ob die Daten Cluster bilden, und diese gegebenenfalls zu ermitteln.



0.4 Clusterbildung.

Im Gegensatz zur Ausreißeridentifikation bietet die Clusteranalyse noch sehr viele offene Fragen.

Dimensionsreduzierung

Viele Modelle erzeugen sehr hochdimensionale Daten. Hochdimensionale Daten sind schlecht numerisch zu handhaben und außerdem schlecht zu visualisieren.



0.5 Dimensionsreduktion durch Approximation der Daten durch eine Funktion.

BSP Eine Kaffeekanne wird unter Rotation fotografiert, um so ein dreidimensionales Bild zu erhalten. Jedes Foto besteht aus n Pixeln, die man als Punkte im \mathbb{R}^n ansehen kann. Diese gewaltige Datenmenge lässt sich jedoch oft auf den \mathbb{R}^3 reduzieren. ■

Auch die Dimensionsreduzierung ist noch weitgehend unerforscht.

Dichteschätzung

Die Dichteschätzung ist mit klassischen Verfahren gut handhabbar und bereits sehr ausgereift.



o.6 Zur Dichteschätzung.

o-2-4 Weitere Lernszenarien

■ *Reinforcement Learning*

Reinforcement Learning hat das Lernen von Handlungen zum Ziel. Beispielsweise misst ein Roboter Umwelteinflüsse und entscheidet sich danach für Handlungen. Welche Entscheidung “gut” war ist jedoch nur mit einer gewissen Verzögerung feststellbar.

■ *Multi-armed Bandit Learning*

Man hat n Optionen, wovon jede einen Gewinn ausschüttet. Aufgrund beschränkter Ressourcen muss man eine Priorisierung vornehmen, um den größtmöglichen Gewinn zu erzielen.

o-2-5 Datengenerierung

■ *Batch Learning*

Hier liegen alle Trainingsdaten vor Trainingsbeginn vor.

■ *Online Learning*

Das Online Learning ist bei den Informatikern sehr beliebt, da sich damit sehr gut “worst case”-Szenarien realisieren lassen.

Ein möglicher Algorithmus wäre

- Habe bereites $(x_1, y_1), \dots, (x_i, y_i)$
- Finde $f_i : X \rightarrow Y$
- Erhalte x_{i+1}
- Mache Vorhersage $f_i(x_{i+1})$
- Bekomme y_{i+1}
- Zurück zu a)

Eine "billige" Generierung von f_i ist wesentlich für den Erfolg des Online Learning.

- *Active Learning*

Bestimme x_1, \dots, x_n selbst und erhalte anschließend die zugehörigen Labels y_1, \dots, y_n .

Natur der Daten

- *Zufällige Daten.* Dies ist typisch für batch learning. Wir werden uns später damit ausführlich beschäftigen.
- *Pseudozufällige Daten.* Die Daten werden eigentlich von einem deterministischen System generiert, dieses ist aber chaotisch.
- *"worst case" Daten.* Typisch für online learning.

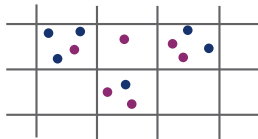
0-3 Klassische Verfahren zur Klassifikation

Sei $X = \mathbb{R}^d$ und $Y = \{-1, 1\}$.

- *Histogrammregel.*

Sei $(Q_j)_{j \geq 1}$ eine Zerlegung von X in Würfel. Es gilt dann offenbar

$$\forall x \in X \exists j \geq 1 : x \in Q_j.$$



0.7 Zerlegung von X .

Wir definieren nun

$$f_D(x) := \text{sign} \left(\sum_{x_i \in Q_i(x)} y_i \right).$$

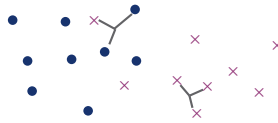
■ *Nearest Neighbor.*

Sei (X, d) ein metrischer Raum (z. B. $X = \mathbb{R}^d$). Fixiere die Anzahl k der betrachteten Nachbarn. Für $x \in X$ suche die k Beispiele

$$(x_{i_1}, y_{i_1}), \dots, (x_{i_k}, y_{i_k})$$

aus D , die am nächsten liegen.

$$f_D(x) = \text{sign} \left(\sum_{j=1}^k y_{i_j} \right).$$



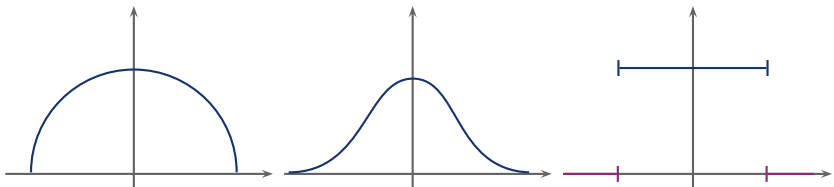
o.8 **Nearest Neighbor Identifikation.**

■ *Moving window / kernel rules.*

Sei $X = \mathbb{R}^d$ ein euklidischer Raum und

$$K : \mathbb{R} \rightarrow [0, \infty)$$

symmetrisch und auf $[0, \infty)$ monoton fallend gegeben.



o.9 **Beispiele für Funktionen $K : \mathbb{R} \rightarrow [0, \infty)$.**

Fixiere $h > 0$ und setze

$$f_D(x) = \text{sign} \left(\sum_{i=1}^k y_i K(h^{-1} \|x - x_i\|) \right)$$



0.10 Identifikation mittels moving window.

1 Grundkonzepte der statistischen Lerntheorie

1-1 Grundannahmen

Folgende Notationen werden durchgehend verwendet.

- (X, \mathcal{A}) sei ein metrischer Raum.
- $Y \subset \mathbb{R}$ sei abgeschlossen mit $Y = \{-1, 1\}$, $Y = [-M, M]$ oder $Y = \mathbb{R}$. Y sei stets mit der Borel- σ -Algebra versehen.
- P sei ein W -Maß auf $X \times Y$.
- $D = ((x_1, \mathcal{Y}_1), \dots, (x_n, \mathcal{Y}_n)) \in (X \times Y)^n$.
- $\mathcal{L}_0(X) := \{f : X \rightarrow \mathbb{R} : f \text{ messbar}\}$.

Wir benötigen folgenden Satz zu regulären bedingten Wahrscheinlichkeiten.

Satz *Es existiert eine Abbildung*

$$P(\cdot | \cdot) : \mathcal{B} \times X \rightarrow [0, 1]$$

mit folgenden Eigenschaften.

- 1) $P(\cdot | x) : \mathcal{B} \rightarrow [0, 1]$ ist ein W -Maß auf \mathcal{B} für alle $x \in X$.
- 2) $x \mapsto P(B | x)$ ist messbar für alle $B \in \mathcal{B}$.

3) Für alle messbaren $A \in \mathcal{A}$ und $B \in \mathcal{B}$ gilt

$$P(A \times B) = \int_A P(B \mid x) dP_X(x) \quad (*)$$

wobei P_X ein W -Maß auf (X, \mathcal{A}) ist, das durch $P_X(A) := P(A \times Y)$ für $A \in \mathcal{A}$ definiert ist.

4) $P(\cdot \mid \cdot)$ ist P_X -f.s. eindeutig. \times

Interpretation. Falls $P = P_Y \otimes P_X$, so ist $P(B \mid x) = P_Y(B)$ und daher

$$(*) = P(A \times B) = P_Y(B)P_X(A).$$

$P(\cdot \mid \cdot)$ ermöglicht es, P im Sinne von (*) aufzuspalten ("desintegrate"). \rightarrow

Als Folgerung erhalten wir eine Verallgemeinerung des Satzes von Fubini von Produktmaßen auf solche P .

Korollar Sei $f : X \times Y \rightarrow \mathbb{R}$ P -integrierbar, so gilt

$$\mathbb{E}_P f := \int_{X \times Y} f dP = \int_X \int_Y f(x, y) P(dy \mid x) dP_X(x). \quad \times$$

Bemerkung. Für allgemeine messbare Räume (Y, \mathcal{A}') ist dies *nicht* möglich. (Wesentlich ist hier, dass $Y \subset \mathbb{R}$ abgeschlossen). \rightarrow

Weitere Grundannahmen sind.

- P ist uns völlig unbekannt. Wir wissen lediglich, dass P existiert.
- Die Trainingsdaten D sind Beobachtungen einer Folge $Z = (Z_i)_{i=1}^n$ mit $Z_i = (X_i, Y_i)$ von Zufallsvariablen, die $X \times Y$ -wertig und unabhängig sind und alle die Verteilung P haben.

Interpretation.

$X_i \sim P_X$, dies gibt eine Beobachtung $x_i \in X$.

$Y_i \sim P(\cdot \mid x_i)$, dies gibt eine Beobachtung $y_i \in Y$.

Datengenerierungsalgorithmus.

- $i = 1$.
 - “Würfle” x gemäß der Verteilung P_X und unabhängig von $(x_1, y_1), \dots, (x_{i-1}, y_{i-1})$.
 - “Würfle” y gemäß der Verteilung $P(\cdot \mid x_i)$.
 - Wiederhole bis $i = n$.
- Alle zukünftigen Daten (x, y) sind Beobachtungen von einer Zufallsvariable $\tilde{Z} = (\tilde{X}, \tilde{Y})$, die $X \times Y$ -wertig ist, Verteilung P hat und unabhängig von der Folge $(Z_i)_{i=1}^n$ ist.
 - Für eine große Anzahl zukünftiger Beobachtungen, betrachte die Folge $(Z_i)_{i=n+1}^\infty$ von $(X \times Y)$ -wertigen Zufallsvariablen, die untereinander und von Z_1, \dots, Z_n unabhängig sind und Verteilung P haben.

1-1.1 **Definition** Eine *Lernmethode* oder auch *Lernverfahren* ist eine Folge $(L_n)_{n \geq 1}$ von Abbildungen

$$L_n : (X \times Y)^n \rightarrow \mathcal{L}_0(X), \quad D \mapsto f_D. \quad \times$$

- BSP** a) Histogrammregel.
- b) Nearest Neighbor.
- c) Moving windows, kernel rules. ■

1-1.2 **Bemerkung.** Wir werden später benötigen, dass die Abbildungen

$$(X \times Y)^n \times X \rightarrow \mathbb{R}, \quad (D, x) \mapsto f_D(x)$$

messbar sind. \rightarrow

1-2 Verlustfunktionen und Risiken

Ziel dieses Abschnitts ist es, genauer zu beschreiben, was $f_D(x) \approx y$ bedeutet.

1-2.1 **Definition** Eine messbare Abbildung

$$L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$$

heißt **Verlustfunktion**. Selbiges gilt für Funktionen

$$L : X \times \mathbb{R} \rightarrow [0, \infty), \quad \text{oder} \quad L : Y \times \mathbb{R} \rightarrow [0, \infty). \quad \times$$

Interpretation. $L(x, y, f(x)) \in [0, \infty)$ beschreibt den “Verlust” bei der Entscheidung $f_D(x)$, falls (x, y) beobachtet wird. Generell sind kleine Verluste besser als große! \rightarrow

1-2.2 **Definition** Sei L eine Verlustfunktion. Dann ist das **Risiko** einer messbaren Funktion $f : X \rightarrow \mathbb{R}$ durch

$$\begin{aligned} \mathcal{R}_{L,P}(f) &= \int_{X \times Y} L(x, y, f(y)) \, dP(x, y) \\ &= \int_X \int_Y L(x, y, f(x)) P(dy | x) \, dP_X(x) \end{aligned}$$

definiert. \times

Interpretation. Ist $Z = (Z_i)_{i=1}^n$ eine Folge von $X \times Y$ -wertigen, unabhängigen und P -verteilten Zufallsvariablen und $\mathcal{R}_{L,P}(f) < \infty$, so besagt das starke Gesetz der großen Zahlen, dass

$$\mathcal{R}_{L,P}(f) = \lim_{m \rightarrow \infty} \frac{1}{m - n - 1} \sum_{i=n+1}^m L(x_i, y_i, f(x_i)) \quad P\text{-f.s.}$$

$\mathcal{R}_{L,P}(f)$ beschreibt daher den mittleren zukünftigen Verlust. \rightarrow

Wir wollen jetzt genauer spezifizieren, was wir unter “kleinen” Risiken zu verstehen ist.

1-2.3 **Definition** Sei L eine Verlustfunktion, dann heißt das **kleinstmögliche Risiko**

$$\mathcal{R}_{L,P}^* := \inf \{ \mathcal{R}_{L,P}(f) : f : X \rightarrow \mathbb{R} \text{ messbar} \}$$

Bayes-Risiko. Eine messbare Funktion $f_{L,P}^* : X \rightarrow \mathbb{R}$ für die gilt

$$\mathcal{R}_{L,P}(f_{L,P}^*) = \mathcal{R}_{L,P}^*$$

heißt **Bayes-Entscheidungsfunktion**. \times

$f_{L,P}^*$ ist im Allgemeinen nicht P -f.s. eindeutig.

Unser informelles Lernziel besteht nun darin, ein Lernverfahren \mathcal{L} zu finden, für das das **Überschussrisiko**

$$\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^*, \quad f_D \in \mathcal{L}$$

mit hoher Wahrscheinlichkeit klein ist. In Abschnitt 1-3 werden wir präzisieren was wir unter "hoher Wahrscheinlichkeit" und "klein" verstehen.

BSP KLASSIFIKATION Hier ist unser Ausgaberaum $Y = \{-1, 1\}$ und unser Ziel, y "richtig vorausszusagen". Mithilfe der vorangegangenen Definitionen lässt sich dies nun mathematisch formulieren.

Wir definieren eine Verlustfunktion wie folgt

$$L_{\text{class}} : Y \times \mathbb{R} \rightarrow [0, \infty), \quad (y, t) \mapsto \mathbf{I}_{(-\infty, 0]}(y \text{ sign } t),$$

$$L_{\text{class}}(y, t) = \begin{cases} 0, & y = -1, t < 0, \\ 1, & y = -1, t \geq 0, \\ 1, & y = 1, t < 0, \\ 0, & y = 1, t \geq 0. \end{cases}$$

L_{class} bestraft somit Vorhersagen, für die $y \neq \text{sign } t$. Wir können nun das Risiko angeben,

$$\begin{aligned} \mathcal{R}_{L_{\text{class}},P}(f) &= \int_{X \times Y} \mathbf{I}_{(-\infty, 0]}(y \text{ sign } f(x)) \, dP(x, y) \\ &= P \{ (x, y) : y \neq \text{sign } f(x) \} \\ &= \int_X \int_Y L_{\text{class}}(y, f(x)) P(dy | x) \, dP_X(x). \end{aligned}$$

Setzen wir nun $\eta(x) = P(y = 1 | x)$ so können wir das Integral schreiben als

$$\int_X \int_Y \eta(x) \mathbf{I}_{(-\infty, 0)}(f(x)) + (1 - \eta(x)) \mathbf{I}_{[0, \infty)}(f(x)) \, dP_X(x).$$

Damit der Integrand minimal wird muss gelten,

$$\begin{aligned} \eta(x) > \frac{1}{2} &\Rightarrow f(x) \geq 0, \\ \eta(x) < \frac{1}{2} &\Rightarrow f(x) < 0. \end{aligned}$$

Daher minimiert f das Risiko genau dann, wenn

$$\begin{aligned} f(x) &\geq 0, && \text{auf } [\eta > \frac{1}{2}], \\ f(x) &< 0, && \text{auf } [\eta < \frac{1}{2}]. \end{aligned}$$

In diesem Fall ist $f = f_{L_{\text{class}},P}^*$ und folglich

$$\mathcal{R}_{L,P}^* = \int_X \min\{\eta, 1 - \eta\} dP_X.$$

Wir wollen noch zeigen, dass

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* = \int_X |2\eta - 1| \mathbf{I}_{(-\infty,0]}((2\eta - 1) \text{sign } f(x)) dP_X.$$

» *Beweisskizze.* Dazu betrachtet man einfach die 6 möglichen Kombinationen von

$$\begin{array}{ccc} \eta(x) > \frac{1}{2} & \eta(x) < \frac{1}{2} & \eta(x) = \frac{1}{2} \\ f(x) \geq 0 & f(x) < 0. & \end{array}$$

z.B. $\eta(x) > \frac{1}{2}$ und $f(x) \geq 0$, so ist

$$|2\eta(x) - 1| \mathbf{I}_{(-\infty,0]}(\underbrace{(2\eta - 1) \text{sign } f(x)}_{>0}) = 0$$

und

$$\eta(x) \underbrace{\mathbf{I}_{(-\infty,0]}(f(x))}_{=0} + (1 - \eta(x)) \underbrace{\mathbf{I}_{[0,\infty)}(f(x))}_{=1} - \underbrace{\min\{\eta(x), 1 - \eta(x)\}}_{1 - \eta(x)} = 0.$$

oder $\eta(x) > \frac{1}{2}$ und $f(x) < 0$, so ist

$$|2\eta(x) - 1| \mathbf{I}_{(-\infty,0]}(\underbrace{(2\eta - 1) \text{sign } f(x)}_{<0}) = 2\eta(x) - 1$$

und

$$\begin{aligned} & \eta(x) \underbrace{\mathbf{I}_{(-\infty,0)}(f(x))}_{=1} + (1 - \eta(x)) \underbrace{\mathbf{I}_{[0,\infty)}(f(x))}_{=0} - \underbrace{\min\{\eta(x), 1 - \eta(x)\}}_{1-\eta(x)} \\ & = 2\eta(x) - 1. \end{aligned}$$

Die übrigen Fälle folgen analog. «

In unserer Definition von L_{class} haben wir alle Fehler gleich gewichtet. Oft nimmt man eine unterschiedliche Gewichtung vor (Übungsaufgabe). ■

REGRESSION MIT KLEINSTEN FEHLERQUADRATEN Sei $Y \subset \mathbb{R}$ ein Intervall oder ganz \mathbb{R} . Unser Ziel ist es f zu finden mit $f(x) \approx y$. Dies wollen wir nun mathematisch fassen.

Dazu definieren wir folgende Verlustfunktion

$$L_{\text{LS}} : Y \times \mathbb{R} \rightarrow [0, \infty), \quad (y, t) \mapsto (y - t)^2.$$

LS steht hier für least squares. Das Risiko ist dann gegeben durch

$$\begin{aligned} \mathcal{R}_{L,P}(f) &= \int_X \int_Y (y - f(x))^2 P(\mathrm{d}y \mid x) \mathrm{d}P_X(x) \\ &= \int_X \int_Y y^2 - 2yf(x) + f^2(x) P(\mathrm{d}y \mid x) \mathrm{d}P_X(x) \\ &= \int_X \int_Y y^2 P(\mathrm{d}y \mid x) - 2f(x) \int_Y y P(\mathrm{d}y \mid x) \\ &\quad + f^2(x) \mathrm{d}P_X(x). \end{aligned}$$

Setzen wir $t = f(x)$ und $h(t) = -2t \int_Y y P(\mathrm{d}y \mid x) + t^2$, so erhalten wir ein Minimum durch

$$0 = h'(t) = -2 \int_Y y P(\mathrm{d}y \mid x) + 2t \Rightarrow t = \int_Y y P(\mathrm{d}y \mid x).$$

Somit ist $f : X \rightarrow \mathbb{R}$ eine Bayes-Entscheidungsfunktion, wenn

$$f(x) = \int_Y y P(\mathrm{d}y \mid x) = \mathbf{E}_P(Y \mid x) \text{ P-f.s.}$$

f ist in diesem Fall tatsächlich eindeutig.

Das alles funktioniert jedoch nur, wenn $\mathcal{R}_{L,S}^* < \infty$, denn sonst ist jede Funktion eine Bayes-Entscheidungsfunktion.

Setzen wir nun in $\mathcal{R}_{L,S}(\cdot)$ ein, erhalten wir

$$\begin{aligned} \mathcal{R}_{L,S}(f_{L,S}^*) &= \int_X \int_Y y^2 P(dy | x) - 2(\mathbf{E}_P(Y | x))^2 + (\mathbf{E}_P(Y | x))^2 dP_X(x) \\ &= \int_X \int_Y y^2 P(dy | x) - (\mathbf{E}_P(Y | x))^2 dP_X(x) \\ &= \text{mittlere Varianz der Label } y. \end{aligned}$$

Somit ist das Überschussrisiko

$$\begin{aligned} \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* &= \int_X \int_Y y^2 - 2f(x)y + f^2 - y^2 \\ &\quad - (\mathbf{E}_P(Y | x))^2 P(dy | x) dP_X(x) \\ &= \int_X (f(x) - f_{L,P}^*(x))^2 dP_X(x) = \|f - f_{L,S}^*\|_{L^2(P_X)}^2. \end{aligned}$$

Minimierung des Risikos ist also äquivalent zur L^2 -Approximation von $f_{L,S}^*$. ■

Bis jetzt waren alle konkreten Verlustfunktionen von x unabhängig.

1-2.4 **Definition** Eine messbare Funktion $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt *strikt überwachte Verlustfunktion*. ✕

1-2.5 **Definition** Eine messbare Funktion $L : X \times \mathbb{R} \rightarrow [0, \infty)$ heißt *unüberwachte Verlustfunktion*. ✕

Es liegt die Vermutung nahe, dass unüberwachte Verlustfunktionen nutzlos für überwachtes Lernen ist. Dem ist aber nicht so, wie folgendes Beispiel zeigen wird.

BSP MEDIANREGRESSION $Y \subset \mathbb{R}$ ist hier typischerweise ein Intervall oder \mathbb{R} .

Für $x \in X$ ist der **bedingte Median** (Median von der bedingten Erwartung in x $P(\cdot | x)$) gegeben durch

$$F_{1/2}^* := \left\{ t^* \in \mathbb{R} : \frac{1}{2} \leq P((-\infty, t^*] | x) \text{ und } \frac{1}{2} \leq P([t^*, \infty) | x) \right\}.$$

Im Allgemeinen ist dieser nicht eindeutig, wir nehmen nun aber an, dass dem so ist.

Das Ziel ist es nun f^* bei bekanntem $D = ((x_1, y_1), \dots, (x_n, y_n))$ abzuschätzen. Wie genau das gemacht wird behandeln wir später, jetzt soll es nur darum gehen, die Verluste zu bewerten.

Betrachte die Verlustfunktion

$$L : X \times \mathbb{R} \rightarrow [0, \infty), \quad (x, t) \mapsto |f^*(x) - t|.$$

Das Risiko ist dann

$$\mathcal{R}_{L,P}(f) := \int_X |f^*(x) - f(x)| \, dP_X(x),$$

$$\mathcal{R}_{L,P}^* = \mathcal{R}_{L,P}(f^*) = 0.$$

Somit ist das Überschussrisiko

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* = \|f^* - f\|_{L^1(P_X)}.$$

Die Verlustfunktion ist also tatsächlich unabhängig von den Daten, obwohl wir überwachte Labels haben.

Problematisch hierbei ist, dass L auf dem bedingten Median beruht, den wir nicht kennen und so L überhaupt nicht berechnen können. Man muss also einen anderen Ansatz wählen.

Dies ist auch ein gutes Beispiel dafür, dass ein Unterschied darin besteht, eine Verlustfunktion zu finden, die das Lernziel gut beschreibt, und eine, die in der Praxis gut verwendbar ist. ■

1-3 Universelle Konsistenz

Ziel dieses Abschnittes ist es die Formulierung

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^*$$

ist mit hoher Wahrscheinlichkeit klein, mathematisch zu präzisieren.

1-3.1 **Definition** Sei $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ eine Verlustfunktion. Eine Lernmethode \mathcal{L} heißt ***P-konsistent*** bezüglich L , falls für alle $\varepsilon > 0$ gilt

$$\lim_{n \rightarrow \infty} P^n \{D \in (X \times Y)^n : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^* < \varepsilon\} = 1.$$

\mathcal{L} heißt ***universell konsistent*** bezüglich L , falls \mathcal{L} *P-konsistent* bezüglich L für *W-Maße* P auf $X \times Y$ für die $\mathcal{R}_{L,P}^* < \infty$. \times

Die universelle Konsistenz ist einerseits eine schwache Forderung, da sie eine rein asymptotische Aussage macht. Sie ist andererseits aber auch eine sehr starke Forderung wie das folgende Beispiel zeigen wird.

BSP Wir betrachten Klassifikation auf $X = [0, 1]$. Seien dazu

$$\begin{aligned} X_1 &= C \subset [0, 1], && \text{“Cantor Menge”,} \\ X_{-1} &= (\mathbb{Q} \cap [0, 1]) \setminus C \end{aligned}$$

Nun existiert ein *W-Maß* μ auf X_1 (z.B. das Hausdorffmaß) und ein *W-Maß* ν auf X_{-1} (z.B. Zählmaß+Dichte), so dass für $P_X := \frac{1}{2}(\mu + \nu)$ gilt

$$P(Y = 1 \mid x) = \begin{cases} 0.51, & x \in X_1, \\ 0.49 & x \in X_{-1}. \quad \blacksquare \end{cases}$$

1-3-1 Plug-in-rules für Klassifikation

Wir wollen nun für ein spezielles Verfahren nachweisen, dass es universell konsistent ist.

Sei $Y = \{-1, 1\}$ und $\eta(x) = P(Y = 1 \mid x)$ die Wahrscheinlichkeit dafür, dass x ein positives Label besitzt. Gegeben sei eine Lernmethode \mathcal{L} , die die Entscheidungsfunktion f_D konstruiert und $f_D \approx \eta$.

Wir haben bereits gezeigt, dass für die Bayes-Klassifikationsentscheidungsfunktion gilt

$$f_{L_{\text{class}},P}^*(x) = \begin{cases} 1, & \eta(x) > \frac{1}{2}, \\ -1, & \eta(x) < \frac{1}{2}. \end{cases}$$

1-3.2 **Definition** Das zu \mathcal{L} gehörige *plug-in-Verfahren* ist durch

$$\hat{f}_D(x) := \begin{cases} 1, & f_D(x) > \frac{1}{2}, \\ -1, & f_D(x) < \frac{1}{2}, \end{cases}$$

gegeben. \times

Wie gut das Verfahren arbeitet, beschreibt das folgende

1-3.3 **Lemma** Sei P ein W -Maß auf $X \times Y$, $Y = \{-1, 1\}$, $\eta(x) = P(Y = 1 \mid x)$ und $x \in X$. Für $h : X \rightarrow \mathbb{R}$ gilt dann

$$\mathcal{R}_{L_{\text{class}}}(2h - 1) - \mathcal{R}_{L_{\text{class}}, P}^* \leq 2 \int_X |\eta - h| dP_X. \quad \times$$

Interpretation. Ist h eine "gute" Schätzung von η , dann ist $2h - 1$ (bzw. $\text{sign}(2h - 1)$) eine "gute" Klassifikationsentscheidungsfunktion.

» *Beweis des Lemmas.* Sei $f = 2h - 1$. Wir haben bereits gezeigt, dass

$$\mathcal{R}_{L_{\text{class}}}(f) - \mathcal{R}_{L_{\text{class}}, P}^* = \int_X \underbrace{|2\eta - 1| \mathbf{I}_{(-\infty, 0]}((2\eta - 1) \text{sign} f)}_{(1)} dP_X.$$

Falls $(2\eta - 1) \text{sign} f > 0$, so ist der Integrand Null und daher

$$(1) \leq 2 |\eta - h|.$$

Falls $(2\eta - 1) \text{sign} f < 0$, so wertet die Indikatorfunktion zu 1 aus. Ist nun $\eta > \frac{1}{2}$, so ist $f < 0$ und somit $-f - 1 > -1$, also

$$\begin{aligned} |2\eta - 1| &= 2\eta - 1 \leq 2\eta - 1 - f = 2\eta - 1 - (2h - 1) = 2(\eta - h) \\ &\leq 2 |\eta - h|. \end{aligned}$$

Ist dagegen $\eta < \frac{1}{2}$, so ist $f \geq 0$ und somit $f + 1 \geq 1$, also

$$|2\eta - 1| = 1 - 2\eta \leq f + 1 - 2\eta = 2 |\eta - h|.$$

Ist $\eta = \frac{1}{2}$, so ist (1) = 0. \ll

Betrachten wir die Histogrammregel mit $Y = \{-1, 1\}$, $X = \mathbb{R}^d$ und

$$D = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n.$$

Sei nun $\mathcal{A} = (A_i)_{i \geq 1}$ mit A_i messbar eine Partition von X , d.h.

$$X = \bigcup_{i \geq 1} A_i, \quad A_i \cap A_j = \emptyset, \quad i \neq j.$$

Für jedes $x \in X$ existiert ein eindeutiger Index $j \geq 1$, so dass $x \in A_j$. Setzen wir

$$N_D(x) := \sum_{i=1}^n \mathbf{1}_{A_i(x)}(x_i) = \# \{ \text{Samples } (x_i, y_i) \text{ für die } x_i \in A(x) \},$$

so erhalten wir

$$\hat{n}_D(x) = \begin{cases} \frac{1}{N_D(x)} \sum \{i: y_i=1\} \mathbf{1}_{A(x)}(x_i), & N_D(x) \neq 0, \\ 0, & \text{sonst.} \end{cases}$$

Damit können wir nun eine plug-in-Regel definieren.

$$f_D^{\mathcal{A}}(x) := \begin{cases} -1, & \hat{n}_D(x) < \frac{1}{2}, \\ 1 & \hat{n}_D(x) \geq \frac{1}{2}. \end{cases}$$

Für $A \subset \mathbb{R}^d$ ist der **Durchmesser** definiert als

$$\text{diam}A := \sup_{x, x' \in A} \|x - x'\|.$$

Für $n \geq 1$ betrachten wir nun eine weitere Partition $\mathcal{A}_n = (A_{n,j})_{j \geq 1}$. Das Lernverfahren \mathcal{L} sei nun durch

$$(X \times Y) \ni D \mapsto f_D^{\mathcal{A}_n} =: f_D \quad (*)$$

gegeben.

1-3.4 **Satz** Sei P ein W -Maß auf $X \times Y$, $Y = \{-1, 1\}$ und $X = \mathbb{R}^d$. Ferner gelten

$$1) \lim_{n \rightarrow \infty} \sup_j \text{diam}A_{n,j} = 0.$$

$$2) \lim_{n \rightarrow \infty} P^n \otimes P_X(\{(D, x) : N_D(x) \leq l\}) = 0 \text{ f.s. für } l \geq 1.$$

Dann ist die Lernmethode (*) P -konsistent bezüglich L_{class} . \times

» Nach Lemma 1-3.3 gilt für $h = \hat{\eta}_D$ und $f_D = \text{sign}(2\hat{\eta}_D - 1)$,

$$\mathcal{R}_{L_{\text{class}}}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \leq 2 \int_X |\eta - \hat{\eta}_D| dP_X.$$

Definiere nun $\bar{\eta} : X \rightarrow [0, 1]$,

$$\bar{\eta}(x) = \begin{cases} \frac{1}{P_X(A(x))} \int_{A(x)} \eta dP_X, & P_X(A(x)) > 0, \\ 0, & \text{sonst,} \end{cases}$$

als "mittleres η auf $A(x)$ ". $\bar{\eta}$ ist somit auf jedem $A(x)$ konstant. Dann gilt

$$\int |\eta - \hat{\eta}_D| dP_X \leq \underbrace{\int |\eta - \bar{\eta}| dP_X}_{(1)} + \underbrace{\int |\bar{\eta} - \hat{\eta}_D| dP_X}_{(2)}.$$

Wir betrachten zunächst den von D abhängigen Term (2),

$$\begin{aligned} & \int_{X \times Y} \int_X |\bar{\eta}(x) - \hat{\eta}_D(x)| dP_X(x) dP^n(D) \\ &= \int_X \int_{X \times Y} |\bar{\eta}(x) - \hat{\eta}_D(x)| dP^n(D) dP_X(x) \\ &= \int_X \sum_{k=0}^n \int_{\{D: N_D(x)=k\}} |\bar{\eta}(x) - \hat{\eta}_D(x)| dP^n(D) dP_X(x) \end{aligned}$$

Sei $N_D(x) = k$, dann existieren eindeutige $x_{i_1}, \dots, x_{i_k} \in A(x)$. Die zugehörigen y_{i_1}, \dots, y_{i_k} sind unabhängige Zufallsvariablen, wovon jede einzelne mit der Wahrscheinlichkeit $\bar{\eta}(x)$ positiv ist.

Somit ist $\sum_{\{j: Y_{i_j}=1\}} y_{i_j}$ binomialverteilt, genauer $B(k, \bar{\eta}(x))$ -verteilt und

$$\sum_{\{i: y_i=1\}} \mathbf{1}_{A(x)}(x_i) = N_D(x) \hat{\eta}_D(x) = k \cdot \hat{\eta}_D(x).$$

Sei $\xi : \Omega \rightarrow \mathbb{R}$ eine $B(k, \bar{\eta}(x))$ -verteilte Zufallsvariable. Dann gilt $\xi/k \sim \hat{\eta}_D(x)$ für $k \geq 1$ und folglich

$$\begin{aligned} & \frac{1}{P^n(\{D : N_D(x) = k\})} \int_{\{D : N_D(x) = k\}} |\bar{\eta}(x) - \hat{\eta}_D(x)| \, dP^n(D) \\ &= \mathbf{E} \left| \bar{\eta}(x) - \frac{1}{k} \xi \right| \\ &\leq \frac{1}{k} \left(\mathbf{E} \left((k\bar{\eta}(x) - \xi)^2 \right) \right)^{1/2} \\ &= \frac{1}{k} \left(\mathbf{E}(\mathbf{E}\xi - \xi)^2 \right)^{1/2}, \end{aligned}$$

nach Hölders Ungleichung. Die Varianz einer $B(k, \bar{\eta}(x))$ -verteilten Zufallsvariablen ist $k\bar{\eta}(x)(1 - \bar{\eta}(x)) \leq k/4$ und damit erhalten wir für obigen Ausdruck die Abschätzung

$$\frac{1}{k} \left(\mathbf{E}(\mathbf{E}\xi - \xi)^2 \right)^{1/2} = \frac{1}{k} (k\bar{\eta}(x)(1 - \bar{\eta}(x)))^{1/2} \leq \frac{1}{2\sqrt{k}}.$$

Für $k = 0$ gilt auf $\{D : N_D(x) = k\}$ nach Definition $\hat{\eta} \equiv 0$ sowie $|\bar{\eta}(x)| \leq 1$. Deshalb ist

$$\int_{\{D : N_D(x) = k\}} |\bar{\eta}(x) - \hat{\eta}_D(x)| \, dP^n(x) \leq P^n(\{D : N_D(x) = k\}).$$

Damit folgt

$$\begin{aligned} & \int_{(X \times Y)^n} \int_X |\bar{\eta}(x) - \hat{\eta}_D(x)| \, dP_X(x) \, dP^n(D) \\ &\leq \int_X \sum_{k=1}^n \frac{1}{2\sqrt{k}} P^n(\{D : N_D(x) = k\}) + P^n(\{D : N_D(x) = 0\}) \, dP_X(x) \\ &\leq \int_X \sum_{k=1}^n \int_{\{D : N_D(x) = k\}} \frac{1}{2\sqrt{k}} \, dP^n(D) + P^n(\{D : N_D(x) = 0\}) \, dP_X(x) \\ &\leq \frac{1}{2} \int_X \int_{\{D : N_D(x) > 0\}} \frac{dP^n(D)}{\sqrt{N_D(x)}} \, dP_X(x) + P^n \otimes P(\{(D, x) : N_D(x) = 0\}) \end{aligned} \quad (**)$$

Für $l \geq 1$ gilt dann

$$\begin{aligned} & \int_X \int_{\{D: N_D(x) > 0\}} \frac{dP^n(D)}{\sqrt{N_D(x)}} dP_X(x) \\ & \leq \int_X \int_{\{D: 0 < N_D(x) \leq l\}} \underbrace{\frac{dP^n(D)}{\sqrt{N_D(x)}}}_{\leq 1} dP_X(x) + \int_X \int_{\{D: l < N_D(x)\}} \underbrace{\frac{dP^n(D)}{\sqrt{N_D(x)}}}_{\leq \sqrt{l}^{-1}} dP_X(x) \\ & \leq P^n \otimes P(\{(D, x) : N_D(x) \leq l\}) + \frac{1}{\sqrt{l}}. \end{aligned} \quad (***)$$

Insgesamt gilt also

$$(**) \leq \frac{3}{2} P^n \otimes P(\{(D, x) : N_D(x) \leq l\}) + \frac{1}{2\sqrt{l}}.$$

Sei nun $\varepsilon > 0$ und $l > \varepsilon^{-2}$, so ist $\frac{1}{\sqrt{l}} < \varepsilon$ und nach Voraussetzung (ii) gilt für hinreichend große n ,

$$P^n \otimes P(\{(D, x) : N_D(x) \leq l\}) \leq \varepsilon$$

und somit (***) $\leq \frac{5}{2} \varepsilon$. «

1-3.5 **Satz** Sei μ ein endliches Borelmaß auf \mathbb{R}^d und

$$f : \mathbb{R}^d \rightarrow \mathbb{R}, \quad \mu\text{-integrierbar.}$$

Dann gibt es für jedes $\varepsilon > 0$ eine gleichmäßig stetige Abbildung $h : \mathbb{R}^d \rightarrow \mathbb{R}$, die μ -integrierbar ist und für die gilt

$$\int |f - h| d\mu \leq \varepsilon. \quad \times$$

» Der Beweis findet sich in [7], oder ergibt sich direkt, wenn das Lebesgue-Integral nach der Methode von Daniel Stone eingeführt wird. «

P_X ist als W-Maß insbesondere endlich und $\eta = P(Y = 1 \mid \cdot) \in [0, 1]$ ist beschränkt und damit P_X -integrierbar. Für $\varepsilon > 0$ gibt es nun eine gleichmäßig stetige Abbildung $\eta_\varepsilon : \mathbb{R}^d \rightarrow \mathbb{R}$, die ebenfalls P_X -integrierbar ist und

$$\int |\eta - \eta_\varepsilon| dP_X < \varepsilon$$

erfüllt. Setzen wir also

$$\bar{\eta}_\varepsilon(x) := \begin{cases} \frac{1}{P_X(A(x))} \int_{A(x)} \eta_\varepsilon \, dP_X, & P_X(A(x)) > 0, \\ 0, & \text{sonst,} \end{cases}$$

so gilt

$$\int |\eta - \bar{\eta}| \, dP_X \leq \underbrace{\int |\eta - \eta_\varepsilon| \, dP_X}_{(I)} + \underbrace{\int |\eta_\varepsilon - \bar{\eta}_\varepsilon| \, dP_X}_{(II)} + \underbrace{\int |\bar{\eta}_\varepsilon - \bar{\eta}| \, dP_X}_{(III)}$$

(I) ist nach Konstruktion $< \varepsilon$. Analog folgt

$$\begin{aligned} (III) &= \int |\bar{\eta}_\varepsilon - \bar{\eta}| \, dP_X \\ &\stackrel{0}{=} \int_X \frac{1}{P_X(A(x))} \int_{A(x)} |\eta_\varepsilon(x') - \eta(x')| \, dP_X(x') \, dP_X(x) \\ &= \int_X \int_X \frac{1}{P_X(A(x))} \mathbf{I}_{A(x)}(x) |\eta_\varepsilon(x') - \eta(x')| \, dP_X(x) \, dP_X(x') \\ &= \int_X |\eta_\varepsilon(x') - \eta(x')| \, dP_X(x') \leq \varepsilon. \end{aligned}$$

Da η_ε gleichmäßig stetig, gibt es ein $\delta > 0$, so dass

$$\|x - x'\| < \delta \Rightarrow |\eta_\varepsilon(x) - \eta_\varepsilon(x')| < \varepsilon$$

und nach Voraussetzung (i) des Satzes gilt für große Zahlen n ,

$$\sup_{j \geq 1} \text{diam} A_{n,j} \leq \delta.$$

Außerdem erhalten wir

$$\begin{aligned} (II) &= \int_X |\eta_\varepsilon - \bar{\eta}_\varepsilon| \, dP_X \\ &= \int_X \left| \eta_\varepsilon(x) - \frac{1}{P_X(A(x))} \int_{A(x)} \eta_\varepsilon(x') \, dP_X(x') \right| \, dP_X(x) \\ &\leq \int_X \frac{1}{P_X(A(x))} \int_{A(x)} \underbrace{|\eta_\varepsilon(x) - \eta_\varepsilon(x')|}_{< \varepsilon} \, dP_X(x') \, dP_X(x) \leq \varepsilon. \end{aligned}$$

Insgesamt haben wir somit gezeigt,

$$\forall \varepsilon > 0 \exists N_0 \in \mathbb{N} \forall n \geq N : \int |\eta - \bar{\eta}| dP_X < \varepsilon.$$

Damit haben wir gezeigt, dass

$$\mathbf{E}_{D \sim P^n} \int |\eta - \hat{\eta}_D| dP_X \rightarrow 0, \quad n \rightarrow \infty.$$

Jetzt müssen wir uns noch davon überzeugen, dass dies auch für die universelle Konsistenz genügt.

Schreibe $h_D := \mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \geq 0$, so gilt nach Lemma 1-3.3 und dem soeben gezeigten, dass

$$\mathbf{E}_{D \sim P^n} h_D \rightarrow 0, \quad n \rightarrow \infty.$$

Ferner gilt

$$\mathbf{E}_{D \sim P^n} h_D = \mathbf{E}_{D \sim P^n} |h_D| = \int_{[h_D < \varepsilon]} |h_D| dP^n(D) + \underbrace{\int_{[h_D \geq \varepsilon]} |h_D| dP^n(D)}_{\geq \varepsilon P^n(\{D : |h_D| \geq \varepsilon\})}.$$

Das heißt

$$\lim_{n \rightarrow \infty} P^n(\{D : \mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \geq \varepsilon\}) = 0,$$

da $\mathbf{E}_{D \sim P^n} h_D \rightarrow 0$.

Wir wollen jetzt zeigen, dass es Lernmethoden gibt, die universell konsistent sind.

1-3.6 **Universelle Konsistenz der Histogrammregel** Sei $\mathcal{A}_n := (A_{n,j})_{j \geq 1}$ eine Partition von \mathbb{R}^d , so dass jedes $A_{n,j}$ ein Hyperwürfel der Kantenlänge h_n ist, d.h. wir haben zu $x_{n,j} \in \mathbb{R}^d$,

$$A_{n,j} = x_{n,j} + [0, h_n]^d.$$

Es gelten

$$1) \lim_{n \rightarrow \infty} h_n = 0,$$

$$2) \lim_{n \rightarrow \infty} n \cdot h_n = \infty.$$

Dann ist die zugehörige Histogrammregel universell konsistent. \times

» Wir weisen die Voraussetzungen von Satz 1-3.4 nach.

“1)”: $\limsup_{n \rightarrow \infty} \max_{j \geq 1} \text{diam} A_{n,j} = 0$ ist erfüllt.

“2)”: Dazu führen wir $S > 0$ ein. Dann gibt es höchstens $c_1 + c_2 h_n^{-d}$ Hyperwürfel $A_{n,j}$ mit $A_{n,j} \cap B_S(0) \neq \emptyset$ und somit

$$\begin{aligned} P^n \otimes P_X (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l\}) &\leq P^n \otimes P_X (\{(D, \mathbf{x}) : \mathbf{x} \notin B_S(0)\}) \\ &\leq \sum_{j: A_{n,j} \cap B_S(0) \neq \emptyset} P^n \otimes P_X (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l\} \cap (X \times Y)^n \times A_{n,j}) \end{aligned}$$

Dann ist $P^n \otimes P_X (\{(D, \mathbf{x}) : \mathbf{x} \notin B_S(0)\}) = P_X(\mathbb{R}^d \setminus B_S(0))$ also

$$\begin{aligned} &\sum_{j: A_{n,j} \cap B_S(0) \neq \emptyset} P^n \otimes P_X (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l, \mathbf{x} \in A_{n,j}\}) \\ &\leq \sum_{\substack{j: A_{n,j} \cap B_S(0) \neq \emptyset \\ P_X(A_{n,j}) \leq \frac{2l}{n}}} P^n \otimes P_X (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l, \mathbf{x} \in A_{n,j}\}) \end{aligned} \quad (1)$$

$$+ \sum_{\substack{j: A_{n,j} \cap B_S(0) \neq \emptyset \\ P_X(A_{n,j}) > \frac{2l}{n}}} (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l, \mathbf{x} \in A_{n,j}\}) \quad (2)$$

Zu (1): $\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l, \mathbf{x} \in A_{n,j}\} \subset (X \times Y)^n \times A_{n,j}$ und daher

$$\begin{aligned} P^n \otimes P_X (\{(D, \mathbf{x}) : N_D(\mathbf{x}) \leq l, \mathbf{x} \in A_{n,j}\}) &\leq P^n \otimes P_X ((X \times Y)^n \times A_{n,j}) \\ &= P_X(A_{n,j}) \end{aligned}$$

Folglich ist

$$(1) \leq \sum_{\substack{j: A_{n,j} \cap B_S(0) \neq \emptyset \\ P_X(A_{n,j}) \leq \frac{2l}{n}}} P_X(A_{n,j}) \leq \frac{2l}{n} (c_1 + c_2 h_n^{-d}) \rightarrow 0,$$

nach Voraussetzung 2).

Zu (2): Setze $\mu_{D_X} := n^{-1} \sum_{i=1}^n \delta_{\{x_i\}}$, wobei

$$\delta_{\{x_i\}}(A) = \begin{cases} 1, & x_i \in A, \\ 0, & \text{sonst.} \end{cases}$$

μ_{D_X} heißt **empirisches Maß** bezüglich des Datensatzes $D_X = (x_1, \dots, x_n)$.

Für $x \in A_{n,j}$ gilt

$$N_D(x) = \sum_{i=1}^n \mathbf{I}_{A(x)}(x_i) = n \cdot \mu_{D_X}(A(x)).$$

Daher folgt

$$\begin{aligned} P^n \otimes P_X \left(\{ (D, x) : N_D(x) \leq l, x \in A_{n,j} \} \right) \\ \leq P_X(A_{n,j}) P^n \left(\{ D : n \cdot \mu_{D_X}(A_{n,j}) \leq l \} \right) \end{aligned}$$

Es folgt wegen $P_X(A_{n,j}) > \frac{2l}{n}$

$$\begin{aligned} P^n \left(\left\{ D : \mu_{D_X}(A_{n,j}) \leq \frac{l}{n} \right\} \right) \\ = P^n \left(\left\{ D : \mu_{D_X}(A_{n,j}) - P_X(A_{n,j}) \leq \frac{l}{n} - P_X(A_{n,j}) \right\} \right) \\ \leq P^n \left(\left\{ D : \mu_{D_X}(A_{n,j}) - P_X(A_{n,j}) \leq -\frac{1}{2} P_X(A_{n,j}) \right\} \right). \end{aligned}$$

Zur Fortsetzung des Beweises benötigen wir folgendes Lemma. «

1-3.7 **Lemma** Sei (Ω, \mathcal{A}, P) ein W -Raum und $h : \Omega \rightarrow \mathbb{R}$ messbar mit $\int h^2 dP < \infty$.
Dann gilt für alle $n \geq 1$ und $t > 0$,

$$P^n \left(\left\{ \omega \in \Omega : \frac{1}{n} \sum_{i=1}^n h(\omega_i) - \mathbf{E}_P h \geq t \right\} \right) \leq \frac{\mathbf{E}_P h^2}{t^2 \cdot n}. \quad \times$$

» Die Tschebyscheff-Ungleichung besagt für quadratintegrierbares f ,

$$P[|f| \geq t] \leq \frac{\mathbf{E}_P f^2}{t^2}.$$

Wenden wir dies auf $f = \frac{1}{n} \sum_{i=1}^n h \circ \pi_i - \mathbf{E}_P h$, wobei

$$\pi_i : \Omega^n \rightarrow \Omega$$

die i -te Projektion.

$$\begin{aligned} P^n \left(\left\{ \omega \in \Omega^n : \frac{1}{n} \sum_{i=1}^n h(\omega_i) - \mathbf{E}_P h \geq t \right\} \right) \\ \leq \frac{\mathbf{E}_{P^n} \left(\frac{1}{n} \sum_{i=1}^n h \circ \pi_i - \mathbf{E}_P h \right)^2}{t^2} \end{aligned}$$

und mit $\eta_i := h \circ \pi_i - \mathbf{E}_P h$ gilt

$$\begin{aligned} \mathbf{E}_{P^n} \left(\frac{1}{n} \sum_{i=1}^n (h \circ \pi_i - \mathbf{E}_P h) \right)^2 &= \frac{1}{n^2} \mathbf{E}_{P^n} \left(\sum_{i=1}^n \eta_i \right)^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{P^n} h^2 \circ \pi_i + \frac{1}{n^2} \sum_{i \neq j} \underbrace{\mathbf{E}_{P^n} \eta_i \eta_j}_{=0} \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{P^n} h^2 \circ \pi_i = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_P h^2 = \frac{1}{n} \mathbf{E}_P h^2. \quad \ll \end{aligned}$$

» Fortsetzung des Beweises von Satz 1-3.6. Setze nun in Lemma 1-3.5, $h = \mathbf{I}_{A_{n,j}}$ und $t = \frac{1}{2} P_X(A_{n,j})$, so gilt

$$\begin{aligned} P^n \left(\{D : \mu_{D_X}\} (A_{n,j}) - P(A_{n,j}) \leq -\frac{1}{2} P_X(A_{n,j}) \right) &\leq \frac{\mathbf{E}_P \mathbf{I}_{A_{n,j}}^2}{4P_X^2(A_{n,j})n} \\ &= \frac{1}{4n \cdot P_X(A_{n,j})}. \end{aligned}$$

Somit gilt schließlich

$$\begin{aligned} (2) &= \sum_{\substack{j: A_{n,j} \cap B_S(0) \neq \emptyset \\ P_X(A_{n,j}) > \frac{2l}{n}}} \left(\{ (D, x) : N_D(x) \leq l, x \in A_{n,j} \} \right) \\ &\leq \sum_{\substack{j: A_{n,j} \cap B_S(0) \neq \emptyset \\ P_X(A_{n,j}) > \frac{2l}{n}}} \frac{1}{4n} \leq \frac{c_1 + c_2 h_n^{-d}}{4n} \xrightarrow{n \rightarrow \infty} 0. \end{aligned}$$

Zum Abschluss des Beweises, müssen wir noch S passend wählen. Sei also $\varepsilon > 0$, dann existiert ein $\delta > 0$, so dass $P_X(\mathbb{R}^d \setminus B_S(0)) \leq \varepsilon$ und für hinreichend große n gilt ferner

$$\frac{2l}{n} (c_1 + c_2 h_n^{-d}) \leq \varepsilon,$$

$$(c_1 + c_2 h_n^{-d}) \frac{1}{4n} \leq \varepsilon,$$

und letztlich

$$P^n \otimes P_X(\{(D, x) : N_D(x) \leq \varepsilon\}) \leq 3\varepsilon. \quad \ll$$

Schlussbemerkungen A. Es gibt universell konsistente Klassifizierungsmethoden. Stone hat 1977 gezeigt, dass das Nearest-Neighbour-Verfahren universell konsistent ist, falls $K_n \rightarrow \infty$ und $\frac{K_n}{n} \rightarrow 0$.

B. Gibt es ein $S > 0$, so dass $P_X(\mathbb{R}^d \setminus B_S(0)) = 0$, dann ist die Konvergenz von

$$P^n \otimes P_X(\{(D, x) : N_D(x) \leq l\}) \quad (*)$$

für $n \rightarrow \infty$ von P unabhängig. Der wesentliche Grund dafür ist, dass Lemma ?? von P unabhängig ist.

C. Die Konvergenzgeschwindigkeit von

$$\int_{(X \times Y)^n} \int_X |\bar{\eta} - \hat{\eta}_D| dP_X dP^n \xrightarrow{n \rightarrow \infty} 0$$

hängt lediglich von (*) ab und ist daher von P unabhängig.

D. In Satz 1-3.4 benötigen wir eine gleichmäßig stetige Approximation von η und $\bar{\eta}$ und können daher keine Aussage zur Konvergenzrate von

$$\mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \leq \int |\eta - \bar{\eta}| dP_X \rightarrow 0$$

machen.

E. Ist η lipschitz stetig mit Konstante ≤ 1 , so können wir eine Konvergenzrate von

$$\mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \leq \int |\eta - \bar{\eta}| dP_X \rightarrow 0$$

berechnen. Diese ist unabhängig von P , in dem Sinn, dass man nur die P betrachtet für die η lipschitz ist. \rightarrow

1-4 Lernraten

Ziel dieses Abschnitts ist es, die Konvergenzraten von

$$\mathcal{R}_{L_{\text{class}},P}(f_D) - \mathcal{R}_{L_{\text{class}},P}^* \rightarrow 0$$

zu untersuchen.

1-4.1 **Definition** Sei $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ eine Verlustfunktion, P ein W -Maß auf $X \times Y$, \mathcal{L} eine Lernmethode und (a_n) eine Folge echt positiver Zahlen mit $a_n \downarrow 0$.

Wir sagen, dass \mathcal{L} **im Mittel mit Rate (a_n) lernt**, falls es ein $c > 0$ gibt, so dass

$$\mathbb{E}_{D \sim P^n} (\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^*) \leq c a_n, \quad n \geq 1.$$

Ferner sagen wir, dass \mathcal{L} **in Verteilung mit Rate (a_n) lernt**, falls es eine Familie $(c_\tau)_\tau \subset [1, \infty)$ und $\tau \in [0, 1]$ gibt, so dass gilt

$$P^n (\{D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^* \leq c_\tau a_n\}) \geq 1 - \tau, \quad n \geq 1, \tau \in [0, 1]. \quad \times$$

Bevor wir klären, ob es überhaupt Lernraten gibt, die für alle P gelten, wollen wir erstmal beide Begriffe vergleichen.

Bemerkungen. A. Ist \mathcal{L} konsistent, so kann man zeigen, dass \mathcal{L} im Mittel mit einer von P abhängigen Rate im Mittel lernt.

B. Im Beweis von Satz 1-3.6 haben wir benutzt,

$$\mathbb{E}_{D \sim P^n} (\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^*) \geq \varepsilon P^n (\{D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^* > \varepsilon\})$$

Ist nun $\mathbb{E}_{D \sim P^n} (\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^*) \leq c a_n$, dann ist

$$P^n (\{D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^* > \varepsilon\}) \leq \frac{c a_n}{\varepsilon}.$$

D.h. für $\tau = \frac{c a_n}{\varepsilon} \Leftrightarrow \varepsilon = \frac{c}{\tau} a_n$ und $c_\tau := c \tau^{-1}$ ist

$$P^n (\{D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P}^* \leq c_\tau a_n\}) \geq 1 - \tau.$$

Lernraten im Mittel implizieren somit Lernraten in Verteilung.

c. Später werden wir häufig Lernraten in Verteilung beweisen, die dann aufgrund der Struktur der c_τ auch im Mittel gelten. \rightarrow

Offene Frage Gibt es Verfahren \mathcal{L} , die mit Rate (a_n) für alle P lernen? \times

Wir werden im Folgenden sehen, dass dies “normalerweise” nicht so ist.

1-4.2 **Definition** Sei (X, \mathcal{A}, μ) ein endlicher Maßraum. Ein $A \in \mathcal{A}$ heißt **Atom**, wenn

$$\forall B \subset A, \quad B \in \mathcal{A} \text{ gilt entweder } \mu(B) = 0 \text{ oder } \mu(A \setminus B) = 0.$$

(X, \mathcal{A}, μ) heißt **atomfrei**, falls kein $A \in \mathcal{A}$ ein Atom ist. \times

BSP a) Versetzen wir die natürlichen Zahlen versehen mit dem Zählmaß ν mit $\nu(\{n\}) = 1$, so ist jede einelementige Menge $\{n\}$ ein Atom.

b) Das Intervall $[0, 1]$ versehen mit dem Lebesgue-Maß ist atomfrei. \blacksquare

1-4.3 **Satz von Lyapunov** Ist (X, \mathcal{A}, μ) ein endlicher atomfreier Maßraum, so ist

$$\{\mu(A) : A \in \mathcal{A}\} = [0, \mu(X)]. \quad \times$$

» Ein Beweis findet sich in [8]. \leftarrow

1-4.4 **No-free-lunch Theorem (Devroye 1982)** Sei $(a_n) \subset [0, 1/32]$ eine fallende Nullfolge, (X, \mathcal{A}, μ) ein atomfreier W -Raum und $Y := \{-1, 1\}$. Dann gibt es zu jeder Lernmethode \mathcal{L} ein W -Maß P auf $X \times Y$ mit

$$1) P_X = \mu.$$

$$2) \mathcal{R}_{L_{\text{class}}, P}^* = 0.$$

$$3) \mathbb{E}_{D \sim P^n} \mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \geq a_n \text{ für } n \geq 1. \quad \times$$

Bemerkungen. A. Das NFL-Theorem besagt, dass es im Allgemeinen keine universellen mittleren Raten gibt. Dies gilt auch für Raten in Verteilung.

B. Das NFL-Theorem gilt für *praktisch jede* Verlustfunktion.

C. Für $X = \mathbb{N}$ gilt das NFL-Theorem, falls $\mu = P_X$ nicht vorher festgelegt wird.

D. Für *endliche* X ist das NFL-Theorem falsch. \rightarrow

1-4.5 **Lemma** Sei $(a_n) \subset (0, 1/16]$ mit $a_n \downarrow 0$, dann existiert eine Folge $(p_n) \subset (0, 1)$ mit $p_1 \geq p_2 \geq \dots$ und $\sum_{n=1}^{\infty} p_n = 1$ und

$$\sum_{i=n+1}^{\infty} p_i \geq \max \{8a_n, 32np_{n+1}\}. \quad \times$$

» Da $p_n \geq p_{n+1}$ gezeigt wird, genügt es zu zeigen, dass

$$\sum_{i=n+1}^{\infty} p_i \geq \max \{8a_n, 32np_n\}. \quad \times$$

Für $l \leq m$ definiere

$$H(m, l) := \sum_{i=l}^{m-1} \frac{1}{i}.$$

Sei nun $n_1 := 1$, so gilt

$$8a_{n_1} = 8a_1 \leq 8 \frac{1}{16} \leq 2^{-1}.$$

Ferner gibt es ein $n_2 > n_1$ mit $H(n_2, n_1) \geq 32$, da $\sum_{i \geq 1} i^{-1} = \infty$. Außerdem ist $H(\cdot, n_1)$ wachsend und (a_n) fallend nach Voraussetzung, also können wir $8a_{n_2} \leq 2^{-2}$ für n_2 hinreichend groß garantieren. Sukzessive finden wir n_k für $k \geq 3$, so dass

$$n_k > n_{k-1}, \quad H(n_k, n_{k-1}) \geq H(n_{k-1}, n_{k-2}), \quad 8a_{n_k} \leq 2^{-k}.$$

Definiere nun

$$c_k := \frac{32}{2^k H(n_{k+1}, n_k)}, \quad k \geq 1,$$

so ist c_k fallend und

$$\frac{1}{32} \sum_{k \geq 1} c_k H(n_{k+1}, n_k) = \sum_{k \geq 1} \frac{1}{2^k} = 1.$$

Für $n \in [n_k, n_{k+1})$ definiere $p_n := \frac{c_k}{32n}$, dann fällt p_n monoton und

$$\sum_{n \geq 1} p_n = \sum_{k \geq 1} \sum_{i=n_k}^{n_{k+1}-1} \frac{c_k}{32i} = \sum_{k \geq 1} \frac{c_k}{32} \underbrace{\sum_{i=n_k}^{n_{k+1}-1} \frac{1}{i}}_{H(n_{k+1}, n_k)} = 1.$$

Ferner gilt für $n \in [n_k, n_{k+1})$

$$\begin{aligned} \sum_{i \geq n+1} p_i &\geq \sum_{i \geq n_{k+1}} p_i = \sum_{i \geq k+1} \sum_{j=n_i}^{n_{i+1}-1} \frac{c_i}{32j} = \sum_{i \geq k+1} \frac{c_i}{32} H(n_{i+1}, n_i) \\ &= \sum_{i \geq k+1} 2^{-i} = 2^{-k}. \end{aligned}$$

Da außerdem $H(n_{k+1}, n_k) \geq H(n_2, n_1) \geq 32$ folgt

$$\frac{1}{2^k} \geq \frac{32}{2^k H(n_{k+1}, n_k)} = c_k = 32np_n.$$

Zudem ist

$$\frac{1}{2^k} \geq 8a_{n_k} \geq 8a_n. \quad \ll$$

» *Beweis des No-free-lunch Theorems.* Ohne Einschränkung sei $f_D \in \{-1, 1\}$, ansonsten betrachte $\text{sign} f_D$. Fixiere eine Folge $(p_i) \downarrow$ gemäß dem vorangegangenen Lemma und zerlege X in Partitionen $(A_j)_{j \geq 1}$ mit $\mu(A_j) = p_j$ (Anwendung des Satzes von Lyapunov). Sei nun $\bar{\nu}$ ein W-Maß auf $\{0, 1\}$ mit $\bar{\nu}(\{0\}) = 1/2$, so ist

$$\nu := \bigotimes_{i=1}^{\infty} \bar{\nu},$$

ein W-Maß auf $\Omega = \{0, 1\}^{\infty}$. Für $\omega = (\omega_j) \in \Omega$ schreibe

$$\eta_{\omega}(x) := \sum_{j=1}^{\infty} \omega_j \mathbf{1}_{A_j}(x), \quad \in \{0, 1\}.$$

Nun sei P_{ω} das W-Maß auf $X \times Y$, das durch

$$\begin{aligned} P_X &= \mu, \\ P(Y = 1 \mid x) &= \eta_{\omega}(x), \quad P\text{-f.s.} \end{aligned}$$

charakterisiert wird, so gilt

$$\mathcal{R}_{L,P_\omega}^* = \int_X \underbrace{\min\{\eta_\omega, 1 - \eta_\omega\}}_{=0} dP_X = 0.$$

Wir zeigen nun, dass

$$\int_\Omega \inf_{n \geq 1} \frac{1}{a_n} \int_{(X \times Y)^n} \mathcal{R}_{L,P_\omega}(f_D) dP_\omega^n(D) d\nu(\omega) \geq \frac{1}{2}, \quad (*)$$

denn dann existiert ein ω , so dass auch

$$\inf_{n \geq 1} \frac{1}{a_n} \int_{(X \times Y)^n} \mathcal{R}_{L,P_\omega}(f_D) dP_\omega^n(D) \geq \frac{1}{2},$$

und für die Folge $(2a_n)_{n \geq 1}$ folgt die Behauptung.

Idee: Im besten Fall genügt ein Punkt, um alle Informationen auf dem zugehörigen A_j zu liefern, denn η_ω ist dort konstant. Haben wir beispielsweise 4 Datenpunkte mit je einem in A_1, \dots, A_4 , dann können wir auf diesen "richtig" entscheiden, jedoch haben wir keine Möglichkeit über A_4 hinaus noch "richtig" zu entscheiden, denn egal wie entschieden wird, existieren immer ω , für die die Entscheidung schlecht ist. Wir können dann also nur noch raten. . .

Der folgende Beweis ist jedoch sehr technisch und daher ist es schwer, die anschauliche Idee dort wieder zu finden.

Sei $D = (X \times Y)^\infty$ und schreibe $D_n = ((x_1, y_1), \dots, (x_n, y_n))$, so gilt

$$\begin{aligned} (*) &= \int_\Omega \inf_{n \geq 1} \frac{1}{a_n} \int_{(X \times Y)^n} \mathcal{R}_{L,P_\omega}(f_D) dP_\omega^n(D) d\nu(\omega) \\ &\geq \int_\Omega \int_{(X \times Y)^n} \inf_{n \geq 1} \frac{1}{a_n} \mathcal{R}_{L,P_\omega}(f_D) dP_\omega^n(D) d\nu(\omega). \end{aligned}$$

Ist $\mathcal{R}_{L,P_\omega}(f_D) \geq a_n$, so ist auch $\frac{1}{a_n} \mathcal{R}_{L,P_\omega}(f_D) \geq 1$ und daher

$$\dots \geq \int_\Omega \int_{(X \times Y)^n} \mathbf{I}_{\bigcap_{n=1}^\infty \{\mathcal{R}_{L,P_\omega}(f_D) \geq a_n\}} dP_\omega^n(D) d\nu(\omega).$$

Setzen wir $B_n = \{\mathcal{R}_{L,P_\omega}(f_D) \geq a_n\}$, so ist $\bigcap_{n \geq 1} B_n = X \setminus \left(\bigcup_{n \geq 1} X \setminus B_n\right)$ und für

das Maß $Q = P_\omega^\infty$ gilt folglich

$$\begin{aligned} Q\left(\bigcap_{n \geq 1} B_n\right) &= 1 - Q\left(\bigcup_{n \geq 1} X \setminus B_n\right) \geq 1 - \sum_{n \geq 1} Q(X \setminus B_n) \\ &\geq 1 - \underbrace{\sum_{n \geq 1} \int_{\Omega} \int_{(X \times Y)^\infty} \mathbf{I}_{\{\mathcal{R}_{L,P}(f_D) \leq a_n\}} dP_\omega^n(D) dv(\omega)}_{(**)}. \end{aligned}$$

Im Folgenden wollen wir (**) nach oben abschätzen. Für $n \geq 1$ definiere

$$\bar{f}_{D_n} := \begin{cases} 1, & \mu(\{f_D = 1\} \cap A_j) \geq \mu(\{f_D = -1\} \cap A_j), \\ 0, & \text{sonst,} \end{cases}$$

und schreibe ferner

$$E_{\omega,j}(f_{D_n}) := A_j \cap \{f_D \neq 2\eta_\omega - 1\}$$

für die Menge der Fehlentscheidungen von f_{D_n} auf A_j für das Maß P_ω . Auf A_j ist $\eta_\omega(x) \equiv \omega_j$ und daher gilt

$$\bar{f}_{D_n}(j) \neq 2\omega - 1 \Rightarrow \mu(E_{\omega,j}(f_{D_n})) \geq \frac{p_j}{2}.$$

Um dies einzusehen sei z.B. $\omega_j = 1$, dann ist $\bar{f}_{D_n}(j) = -1$ und daher nach Definition von \bar{f}_{D_n}

$$\mu(\{f_{D_n} = 1\} \cap A_j) = \mu(\{f_{D_n} = -1\} \cap A_j)$$

also

$$\begin{aligned} p_j = \mu(A_j) &= \mu(\{f_{D_n} = 1\} \cap A_j) + \mu(\{f_{D_n} = -1\} \cap A_j) \\ &\leq 2\mu\left(\underbrace{\{f_{D_n} = -1\} \cap A_j}_{E_{\omega,j}(f_{D_n})}\right) \end{aligned}$$

Damit folgt $\mathbf{I}_{\{\mu(E_{\omega,j}(f_{D_n})) \geq \frac{p_j}{2}\}} \geq \mathbf{I}_{\{\bar{f}_{D_n} \neq 2\omega_j - 1\}}$ und wegen

$$\mu(E_{\omega,j}(f_{D_n})) \geq \frac{p_j}{2} \mathbf{I}_{\{\mu(E_{\omega,j}(f_{D_n})) \geq \frac{p_j}{2}\}}$$

ist

$$\begin{aligned}
 \mathcal{R}_{L,p_\omega}(f_D) &= \mathcal{R}_{L,p_\omega}(f_D) - \underbrace{\mathcal{R}_{L,p_\omega}^*}_{=0} \\
 &= \int_X |2\eta_\omega - 1| \mathbf{I}_{(-\infty,0]}((2\eta_\omega - 1) \operatorname{sign} f_{D_n}) \, d\mu \\
 &= \mu(\{f_{D_n} \neq (2\eta_\omega - 1)\}) = \sum_{j=1}^{\infty} \mu(E_{\omega,j}(f_{D_n})) \\
 &\geq \frac{1}{2} \sum_{j=1}^{\infty} p_j \mathbf{I}_{\{f_{D_n} \neq 2\omega_{j-1}\}} \geq \frac{1}{2} \sum_{\substack{j:j \leq n \\ x_i \in A_j}} p_j \mathbf{I}_{\{\bar{f}_{D_n}(j) \neq 2\omega_{j-1}\}}
 \end{aligned}$$

Daraus folgt

$$\begin{aligned}
 &\int_{\Omega} \int_{(X \times Y)^\infty} \mathbf{I}_{\{\mathcal{R}_{L,p_\omega}(f_{D_n}) \leq a_n\}} \, dP_\omega^\infty(D) \, d\nu(\omega) \\
 &\leq \int_{\Omega} \int_{(X \times Y)^\infty} \mathbf{I}_{\left\{ \sum_{\substack{j:j \leq n \\ x_i \in A_j}} p_j \mathbf{I}_{\{\bar{f}_{D_n} \neq 2\omega_{j-1}\}} \leq 2a_n \right\}} \, dP_\omega^\infty(D) \, d\nu(\omega).
 \end{aligned}$$

Für alles Weitere fixieren wir $J := \{j : j \leq n, x_i \in A_j\}$.

Zu jedem $x \in X$ existiert genau ein $j_x \geq 1$ mit $x \in A_{j_x}$, wir schreiben dafür $j(x) := j_x$. Nach Konstruktion sind die x_i von ω unabhängig und es gilt

$$y_i = 2\omega_{j(x_i)} - 1.$$

Die Labels y_i lassen sich somit explizit berechnen und daher hängt \bar{f}_D nur von $D_X = (x_1, \dots, x_n)$ und $\omega_{D_X} = (\omega_{j(x_1)}, \dots, \omega_{j(x_n)})$ ab. Wir schreiben daher $\bar{f}_{D_X, \omega}$ für \bar{f}_D und erhalten,

$$\begin{aligned}
 &\int_{\Omega} \int_{(X \times Y)^\infty} \mathbf{I}_{\left\{ \sum_{j \in J} p_j \mathbf{I}_{\{\bar{f}_D(j) \neq 2\omega_{j-1}\}} \leq 2a_n \right\}} \, dP_\omega^n(D) \, d\nu(\omega) \\
 &= \int_{X^n} \int_{\Omega} \mathbf{I}_{\left\{ \sum_{j \in J} p_j \mathbf{I}_{\{\bar{f}_{D_X, \omega}(j) \neq 2\omega_{j-1}\}} \leq 2a_n \right\}} \, d\nu(\omega) \, d\mu^n(D_X).
 \end{aligned}$$

Sei weiterhin Ω_{D_X} das Kreuzprodukt, das durch die Koordinaten $j(x_1), \dots, j(x_n)$ beschrieben wird, und $\Omega_{\mathcal{P}_X}$ stehe für die "übrigen" Koordinaten. So ist $\Omega =$

$\Omega_{D_X} \times \Omega_{\mathcal{D}_X}$, in selber Weise definieren wir ν_{D_X} und $\nu_{\mathcal{D}_X}$. $\bar{f}_{D_X, \omega}$ ändert sich nur in den Koordinaten von ω , die zu Ω_{D_X} gehören, also ist

$$\begin{aligned} & \int_{\Omega} \mathbf{I} \left\{ \sum_{j \in J} p_j \mathbf{I} \left\{ \bar{f}_{D_X, \omega(j) \neq 2\omega_{j-1}} \leq 2a_n \right\} \right\} d\nu(\omega) \\ &= \int_{\Omega_{D_X}} \int_{\Omega_{\mathcal{D}_X}} \mathbf{I} \left\{ \sum_{j \in J} p_j \mathbf{I} \left\{ \bar{f}_{D_X, \omega(j) \neq 2\omega_{j-1}} \leq 2a_n \right\} \right\} d\nu_{\mathcal{D}_X}(\omega_{\mathcal{D}_X}) d\nu_{D_X}(\omega_{D_X}). \end{aligned}$$

Außer dem ist $\bar{f}_{D_X, \omega(j)} \equiv 1$ oder -1 für alle ω_{D_X} bei festem D_X und ω_{D_X} , wobei $\omega_j = 0$ oder 1 jeweils mit Wahrscheinlichkeit $\frac{1}{2}$, also

$$\begin{aligned} \dots &= \int_{\Omega_{D_X}} \int_{\Omega_{\mathcal{D}_X}} \mathbf{I} \left\{ \sum_{j \in J} p_j \underbrace{\mathbf{I} \left\{ \omega_j = 1 \right\}}_{=\omega_j} \leq 2a_n \right\} d\nu_{\mathcal{D}_X}(\omega_{\mathcal{D}_X}) d\nu_{D_X}(\omega_{D_X}) \\ &\geq \int_{\Omega} \mathbf{I} \left\{ \sum_{j=n+1}^{\infty} p_j \omega_j \leq 2a_n \right\} d\nu(\omega). \end{aligned} \quad (**)$$

Da p_i monoton fällt und (**) unabhängig von D_X , ist

$$\begin{aligned} (**) &= \nu \left(\left\{ \omega : - \sum_{j=n+1}^{\infty} p_j \omega_j > -2a_n \right\} \right) \\ &= \nu \left(\left\{ \omega : \exp \left(2sa_n - s \sum_{j=n+1}^{\infty} p_j \omega_j \right) > 1 \right\} \right) \end{aligned}$$

für ein $s > 0$, das wir später wählen. Anwendung der Markov-Ungleichung ergibt,

$$\begin{aligned} \dots &\leq \mathbf{E}_{\omega \sim \nu} \exp \left(2sa_n - s \sum_{j=n+1}^{\infty} p_j \omega_j \right) = \exp(2sa_n) \prod_{j=n+1}^{\infty} \mathbf{E}_{\omega \sim \nu} \exp(-sp_j \omega_j) \\ &= \exp(2sa_n) \prod_{j=n+1}^{\infty} \left(\frac{1}{2} + \frac{1}{2} \exp(-sp_j) \right) \end{aligned}$$

Verwenden wir nun, dass $e^{-t} \leq 1 - t + \frac{t^2}{2}$ für $t \geq 0$, so ist

$$\dots \leq \exp(2sa_n) \prod_{j=n+1}^{\infty} \frac{1}{2} \left(2 - sp_j + \frac{s^2 p_j^2}{2} \right)$$

und, da $1 - t \leq e^{-t}$ ist

$$\begin{aligned}
 \dots &\leq \exp(2sa_n) \prod_{j=n+1}^{\infty} \left(\exp\left(\frac{-sp_j}{2} + \frac{s^2 p_j^2}{4}\right) \right) \\
 &= \exp\left(2sa_n - s \sum_{j=n+1}^{\infty} p_j + \frac{s^2}{4} \sum_{j=n+1}^{\infty} p_j^2\right) \\
 &\leq \exp\left(2sa_n - s \sum_{j=n+1}^{\infty} p_j + \frac{s^2 p_{n+1}}{4} \sum_{j=n+1}^{\infty} p_j\right) \\
 &= \exp\left(2sa_n + \left(\frac{s^2 p_{n+1}}{4} - s\right) \sum_{j=n+1}^{\infty} p_j\right)
 \end{aligned}$$

Setzen wir nun $A := \sum_{j=n+1}^{\infty} p_j$ und

$$s = \frac{\sum_{j=n+1}^{\infty} p_j - 4a_n}{p_{n+1} \sum_{j=n+1}^{\infty} p_j} = \frac{A - 4a_n}{Ap_{n+1}},$$

so ist $s > 0$ nach Konstruktion der p_j . Eine längere aber elementare Rechnung zeigt,

$$\exp\left(2sa_n + \left(\frac{s^2 p_{n+1}}{4} - s\right) \sum_{j=n+1}^{\infty} p_j\right) = \exp\left(-\frac{1}{4} \frac{(A - 4a_n)^2}{Ap_{n+1}}\right)$$

Die Funktion $x \mapsto (A - 4x)^2$ hat ein globales Minimum bei $x = \frac{A}{4}$ und ist auf $(-\infty, \frac{A}{4}]$ monoton fallend. Folglich ist $(A - 4x)^2 \geq (A - \frac{A}{2})^2 = \frac{A^2}{4}$ für $x \leq \frac{A}{8}$ und daher

$$\dots \leq \exp\left(-\frac{A}{16p_{n+1}}\right) \leq \exp(-2n),$$

denn $A \geq 32p_{n+1}n$.

Schließlich folgt

$$\begin{aligned} & \int_{\Omega} \inf_{n \geq 1} \frac{1}{a_n} \int_{(X \times Y)^n} \mathcal{R}_{L,P}(f_D) \, dP_{\omega}^n(D) \, d\nu(\omega) \\ & \geq 1 - \sum_{n=1}^{\infty} \int_{\Omega} \int_{(X \times Y)^{\infty}} \mathbf{I}_{\{\mathcal{R}_{L,P}(f_{D_n}) \leq a_n\}} \, dP_{\omega}^{\infty}(D) \, d\nu(\omega) \\ & \geq 1 - \sum_{n=1}^{\infty} e^{-2n} = \frac{e^2 - 2}{e^2 - 1} \geq \frac{1}{2}. \quad \ll \end{aligned}$$

Abschlussbemerkungen. A. Das NFL gilt auch in den folgenden Situationen:

- $X = \mathbb{R}^d$ und $\eta \in C^{\infty}([0, 1])$.
- $X = \mathbb{R}^2$ und η ist **unimodal**, d.h. es gibt ein $x_0 \in X$ mit $\lambda \mapsto \eta(\lambda x_0)$ fallend für wachsendes $\lambda > 0$.
- $\eta \in \{0, 1\}$, $X \subset \mathbb{R}^2$ und $\{\eta = 1\}$ ist kompakt, konvex und $0 \in \{\eta = 1\}$.

B. Es gibt keine Super-Klassifikationsmethoden. Ist \mathcal{L} eine Klassifikationsmethode, so gibt es eine universell konsistente Klassifikationsmethode \mathcal{L}' und ein W -Maß P auf $X \times Y$ mit

$$\mathbf{E}_{D \sim P^n} \mathcal{R}_{L_{\text{class}}, P}(f_D) > \mathbf{E}_{D \sim P^n} \mathcal{R}_{L_{\text{class}}, P}(f'_D), \quad n \geq 1.$$

Dies erscheint auf den ersten Blick natürlich negativ, es hat jedoch auch positive Auswirkungen, denn dadurch wird das Feld stetig neu belebt. . .

C. Für jede Methode \mathcal{L} , die $\mathcal{R}_{L_{\text{class}}, P}^*$ abschätzt und jedes $n \geq 1$, $\varepsilon > 0$ gibt es ein P mit

$$\mathbf{E}_{D \sim P^n} \left| \mathcal{R}_{L_{\text{class}}, P}(f_D) - \mathcal{R}_{L_{\text{class}}, P}^* \right| \geq \frac{1}{4} - \varepsilon. \quad \rightarrow$$

Wir wollen dieses Kapitel mit einer offenen Frage beenden. Dazu die folgende

Definition Eine Klassifikationsmethode \mathcal{L} heißt **smart**, wenn $\mathbf{E}_{D \sim P^n} \mathcal{R}_{L_{\text{class}}, P}(f_D)$ monoton fallend für $n \rightarrow \infty$ und alle P . \times

Offene Frage Gibt es eine universell konsistente und smarte Klassifikationsmethode? \times

Bis jetzt ist noch keine solche Lernmethode bekannt. Dabei ist die universelle Konsistenz entscheidend. Lässt man diese Voraussetzung fallen, so gibt es natürlich zahllose smarte Lernmethoden.

2 Verlustfunktionen

Sei L eine Verlustfunktion, dann ist

$$\mathcal{R}_{L,D}(f) := \frac{1}{n} \sum_{i=1}^n L(x_i, y_i, f(x_i))$$

sinnvoll, wenn man D mit dem empirischen Maß

$$\mu = n^{-1} \sum_{i=1}^n \delta_{\{(x_i, y_i)\}}$$

identifiziert. Sei weiterhin \mathcal{F} eine Menge von Funktionen $f : X \rightarrow \mathbb{R}$, dann suchen wir $(f, D) \in \mathcal{F} \times (X \times Y)^n$, so dass das Infimum

$$\inf_{f \in \mathcal{F}} \mathcal{R}_{L,D}(f) + \gamma(f)$$

für einen **Strafterm** γ angenommen wird.

Um diese Frage zu klären, wollen wir nach Eigenschaften der Verlustfunktion suchen, die sich auf das Risiko übertragen.

2-1 Eigenschaften von Verlustfunktionen

Wir betrachten im Folgenden Verlustfunktionen

$$L : X \times Y \times \mathbb{R} \rightarrow [0, \infty).$$

Während X und Y relativ “strukturlos” sind, hat \mathbb{R} sehr viel “Struktur”. Wir suchen nun nach Eigenschaften von \mathbb{R} wie Topologie, Metrik, Konvexität, ... die sich auf L übertragen.

2-1.1 **Definition** Eine Verlustfunktion $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt (strikt) *konvex*, falls $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ (strikt) konvex. \times

Erinnerung. $\mathcal{L}_0(X) = \{f : X \rightarrow \mathbb{R} : f \text{ messbar}\}$. \rightarrow

2-1.2 **Lemma** Sei L eine (strikt) konvexe Verlustfunktion, dann ist auch

$$\mathcal{R}_{L,P}(\cdot) : \mathcal{L}_0(X) \rightarrow [0, \infty]$$

(strikt) konvex. \times

» Übung. «

2-1.3 **Definition** Eine Verlustfunktion $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt *stetig*, falls $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ stetig für alle $(x, y) \in (X \times Y)$. L heißt *lokal lipschitz-stetig*, falls

$$\forall a > 0 \exists c_a : |L(x, y, t) - L(x, y, t')| \leq c_a |t - t'|$$

für alle $(x, y) \in (X \times Y)$ und $t, t' \in [-a, a]$. Die kleinste Konstante c_a wird mit $|L|_{a,1}$ bezeichnet.

L heißt *lipschitz-stetig*, falls $|L|_1 := \sup_{a \geq 0} |L|_{a,1} < \infty$. \times

2-1.4 **Lemma** Sei L eine stetige Verlustfunktion und P ein W -Maß auf $X \times Y$, so gilt für alle Folgen (f_n) in $\mathcal{L}_0(X)$ und $f \in \mathcal{L}_0(X)$ mit $f_n \xrightarrow{P} f$,

$$\mathcal{R}_{L,P}(f) \leq \liminf_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n). \quad \times$$

Das Risiko ist also "halbstetig von unten".

» Da $f_n \xrightarrow{P} f$ existiert eine Teilfolge mit $f_{n_k} \rightarrow f$ P -f.s. und daher existiert auch eine weitere Teilfolge $(f_{n_{k_l}})$ mit

$$\lim_{l \rightarrow \infty} \mathcal{R}_{L,P}(f_{n_{k_l}}) = \liminf_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n).$$

Wir schreiben nun kürzer (f_{n_k}) für die Teilfolge mit beiden Eigenschaften. Da L stetig gilt $L(x, y, f_{n_k}(x)) \rightarrow L(x, y, f(x))$ P_X -f.s. und mit dem Lemma von Fatou folgt,

$$\begin{aligned} \mathcal{R}_{L,P}(f) &= \int_{X \times Y} L(x, y, f(x)) \, dP(x, y) = \int_{X \times Y} \lim_{k \rightarrow \infty} L(x, y, f_{n_k}(x)) \, dP(x, y) \\ &\leq \liminf_{k \rightarrow \infty} \int_{X \times Y} L(x, y, f_{n_k}(x)) \, dP(x, y) = \lim_{k \rightarrow \infty} \mathcal{R}_{L,P}(f_{n_k}) \\ &= \liminf_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n). \quad \ll \end{aligned}$$

Optimal für unsere Zwecke wäre “=” anstat “ \leq ” im Lemma 2-1.4. Es stellt sich jedoch heraus, dass dies nicht ohne zusätzliche Voraussetzungen an L möglich ist.

2-1.5 **Definition** Eine Verlustfunktion $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt **Nemitski-Verlustfunktion (NVF)**, falls eine messbare Funktion $b : X \times Y \rightarrow [0, \infty)$ und eine messbare und monoton wachsende Funktion $h : \mathbb{R} \rightarrow [0, \infty)$ existiert, so dass

$$L(x, y, t) \leq b(x, y) + h(|t|), \quad \forall x, y, t.$$

L heißt **NVF der Ordnung $p \in (0, \infty)$** , falls ein $c > 0$ existiert, so dass

$$L(x, y, t) \leq b(x, y) + c \cdot |t|^p, \quad \forall x, y, t.$$

Ist P ein W -Maß auf $(X \times Y)$ und L eine NVF, so heißt L **P -integrierbar**, falls b P -integrierbar. \times

2-1.6 **Lemma** Sei P ein W -Maß auf $X \times Y$ und L eine stetige, P -integrierbare NVF, dann gelten.

1) Sei (f_n) eine gleichmäßig beschränkte Folge in $\mathcal{L}_0(P_X)$ und $f \in \mathcal{L}_\infty(P_X)$ mit $f_n \rightarrow f$ P -f.s. Dann folgt

$$\mathcal{R}_{L,P}(f) = \lim_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n).$$

2) $\mathcal{R}_{L,P}(\cdot) : \mathcal{L}_\infty(P_X) \rightarrow [0, \infty)$ ist stetig.

3) Ist L außerdem von der Ordnung $p \in [1, \infty)$, so ist

$$\mathcal{R}_{L,P}(f) : \mathcal{L}_p(P_X) \rightarrow [0, \infty)$$

wohldefiniert und stetig. \times

» “1)”: f ist beschränkt mit $\|f\|_\infty \leq B$. Ferner gilt

$$\lim_{n \rightarrow \infty} \underbrace{|L(x, y, f_n(x)) - L(x, y, f(x))|}_{g_n(x)} = 0 \text{ } P\text{-f.s.}$$

und

$$\begin{aligned} g_n(x) &\leq L(x, y, f_n(x)) + L(x, y, f(x)) \\ &\leq b(x, y) + h(|f_n(x)|) + b(x, y) + h(|f(x)|) \\ &\leq 2(b(x, y) + h(B)) \text{ P-f.s.} \end{aligned}$$

wobei die rechte Seite als Funktion in (x, y) P -integrierbar ist. Mit dem Satz von Lebesgue folgt nun,

$$|\mathcal{R}_{L,P}(f_n) - \mathcal{R}_{L,P}(f)| \leq \int_{X \times Y} g_n(x) \, dP(x, y) \rightarrow 0, \quad n \rightarrow \infty.$$

“2)”: Sei (f_n) Folge in $\mathcal{L}_\infty(X)$ und $f \in \mathcal{L}_\infty(X)$ mit $\|f_n - f\|_\infty \rightarrow 0$. Dann ist f_n gleichmäßig beschränkt und $f_n \rightarrow f$ P_X -f.s., also $\mathcal{R}_{L,P}(f_n) \rightarrow \mathcal{R}_{L,P}(f)$.

“3)”: Es gilt für $f \in \mathcal{L}_p(P_X)$,

$$\mathcal{R}_{L,P}(f) = \int L(x, y, f(x)) \, dP(x, y) \leq \int b(x, y) + c |f(x)|^p \, dP(x, y) < \infty.$$

Somit ist $\mathcal{R}_{L,P}(\cdot) : \mathcal{L}_p(P_X) \rightarrow [0, \infty)$ wohldefiniert.

Sei nun $(f_n) \subset \mathcal{L}_p(P_X)$ und $f \in \mathcal{L}_p(P_X)$ mit $\|f_n - f\|_p \rightarrow 0$, dann gilt auch $f_n \xrightarrow{P} f$. Mit Lemma 2-1.4 folgt

$$\mathcal{R}_{L,P}(f) \leq \liminf_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n). \quad (*)$$

Definiere nun $\bar{L}(x, y, t) := b(x, y) + c |t|^p - L(x, y, t) \geq 0$, so ist $\bar{L} : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ eine stetige Verlustfunktion. Erneute Anwendung von 2-1.4 ergibt

$$\begin{aligned} \|b\|_1 + c \|f\|_p^p - \mathcal{R}_{L,P}(f) &= \mathcal{R}_{\bar{L},P}(f) \leq \liminf_{n \rightarrow \infty} \mathcal{R}_{\bar{L},P}(f_n) \\ &= \|b\|_1 + \liminf_{n \rightarrow \infty} c \|f_n\|_p^p - \mathcal{R}_{L,P}(f_n). \end{aligned}$$

Da $\|f_n\|_p \rightarrow \|f\|_p$ ist $-\mathcal{R}_{L,P}(f) \leq -\limsup_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n)$, d.h. mit (*) folgt

$$\limsup_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n) \leq \liminf_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n)$$

also $\lim_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n) = \mathcal{R}_{L,P}(f)$. «

2-1.7 **Lemma** Ist $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ strikt überwacht und konvex und Y endlich. Dann ist L lokal lipschitz. \times

» Wir benutzen, dass jede konvexe Abbildung $g : [-a, a] \rightarrow [0, \infty)$ lokal lipschitz ist. Nach dieser Aussage ist $L(y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ lokal lipschitz für alle $y \in Y$. Da Y endlich, folgt die Aussage. «

2-1.8 **Lemma** 1) Ist L lokal lipschitz, so ist L NVF.

2) Ist L lokal lipschitz und $\mathcal{R}_{L,p}(0) < \infty$, so ist L P -integrierbare NVF.

3) Ist L lipschitz stetig, so ist L NVF der Ordnung $p = 1$. \times

» 1) $|L(x, y, t) - L(x, y, 0)| \leq |L|_{|t|,1} |t|$. Somit ist

$$L(x, y, t) \leq \underbrace{|L|_{|t|,1} |t|}_{h(t)} + \underbrace{L(x, y, 0)}_{b(x,y)}.$$

2) Falls $\mathcal{R}_{L,p}(0) < \infty$ ist $b(x, y)$ P -integrierbar.

3) Falls L lipschitz, ist $|L|_{|t|,1} \leq |L|_1$ und damit

$$L(x, y, t) \leq L(x, y, 0) + |L|_1 |t|. \quad \llcorner$$

2-1.9 **Lemma** Sei L lokal lipschitz, $B \geq 0$ und $f, g \in \mathcal{L}_\infty(P_X)$ mit $\|f\|_\infty, \|g\|_\infty \leq B$. Dann gilt

$$|\mathcal{R}_{L,p}(f) - \mathcal{R}_{L,p}(g)| \leq |L|_{B,1} \|f - g\|_{\mathcal{L}_1(P_X)}. \quad \times$$

» Übung. «

2-1.10 **Definition** Eine Verlustfunktion L kann bei $M > 0$ abgeschnitten werden, wenn für alle $(x, y) \in X \times Y$ und $t \in \mathbb{R}$ gilt

$$L(x, y, \hat{t}) \leq L(x, y, t),$$

wobei

$$\hat{t} = \begin{cases} -M, & t \leq -M, \\ t, & t \in (-M, M) \\ M, & t \geq M. \end{cases} \quad \times$$

Man kann dies so interpretieren, dass Abschneiden den Verlust *nicht* erhöht.

2-1.11 **Lemma** Sei L eine konvexe Verlustfunktion und $M > 0$. Dann sind folgende Aussagen äquivalent:

(i) L kann bei M abgeschnitten werden.

(ii) Für alle $(x, y) \in X \times Y$ hat die Funktion $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ mindestens ein globales Minimum in $[-M, M]$. \times

» Schreibe $M_{x,y} := \{t^* \in \mathbb{R} : L(x, y, t^*) = \inf_{t \in \mathbb{R}} L(x, y, t)\}$. Da L konvex ist, ist $M_{x,y}$ ein Intervall.

“(i) \Rightarrow (ii)”: Angenommen es gibt ein $(x, y) \in X \times Y$ mit $M_{x,y} \cap [-M, M] = \emptyset$.

1. Fall $M_{x,y} = \emptyset$. L ist konvex, also ist $L(x, y, \cdot)$ strikt monoton, denn falls $L(x, y, \cdot) \in C^2(\mathbb{R})$, so ist aufgrund der Konvexität $L''(x, y, \cdot) \geq 0$. Da aber $M_{x,y} = \emptyset$ folgt $L'(x, y, t) \neq 0$ für alle t . Für allgemeines $L(x, y, \cdot)$ folgt die strikte Monotonie aus der Betrachtung von Subdifferenzialen.

Aber da $L(x, y, \cdot)$ strikt monoton, kann L nicht abgeschnitten werden. $\not\Leftarrow$

2. Fall $M_{x,y} \neq \emptyset$. Da $M_{x,y}$ ein abgeschlossenes Intervall, folgt ohne Einschränkung $t := \inf M_{x,y}$ erfüllt $M < t < \infty$, d.h. $M_{x,y}$ liegt rechts von $[-M, M]$.

Somit ist $L(x, y, t) = L(x, y, M) > L(x, y, t)$ da $t > M$ und $L(x, y, \cdot)$ aufgrund der Konvexität strikt fallend links von $M_{x,y}$. $\not\Leftarrow$

“(ii) \Rightarrow (i)”: Es gilt $M_{x,y} \cap [-M, M] \neq \emptyset$ und daher ist $\inf M_{x,y} \leq M$ und $\sup M_{x,y} \geq M$.

$L(x, y, \cdot)$ ist strikt konvex, also ist $L(x, y, \cdot)$ auf $[\sup M_{x,y}, \infty)$ wachsend und auf $(-\infty, \inf M_{x,y}]$ fallend. Somit kann L abgeschnitten werden. \llcorner

2-2 Margin basierte Verlustfunktionen

Sei $Y = \{-1, 1\}$ und $\eta(x) = P(Y = 1 \mid x)$.

2-2.1 **Definition** Eine strikt überwachte Verlustfunktion $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt *margin-basiert*, falls eine repräsentative Funktion $\varphi : \mathbb{R} \rightarrow [0, \infty)$ existiert, d.h. $L(y, t) = \varphi(y \cdot t)$. \times

2-2.2 **Lemma** Sei L margin-basiert und φ die repräsentative Funktion. Dann gelten:

- 1) L ist genau dann (strikt) konvex, wenn φ (strikt) konvex.
- 2) L ist genau dann stetig, wenn φ stetig.
- 3) L ist genau dann (lokal) lipschitz, wenn φ (lokal) lipschitz.
- 4) Ist L konvex, so ist L (lokal) lipschitz.
- 5) L ist P -integrierbare NVF.
- 6) Ist L lipschitz stetig, so ist L P -integrierbare NVF der Ordnung $p = 1$. \times

» Übung. «

2-2.3 **BSP** *Kleinste Quadrate*. Scharfes Hinsehen ergibt,

$$L_{\text{LS}}(\gamma, t) = (\gamma - t)^2 = (1 - \gamma t)^2, \quad \Rightarrow \quad \varphi(t) = (1 - t)^2.$$

Damit ist L_{LS} strikt konvex, da φ strikt konvex. $|L_{\text{LS}}|_{a,1} = 2a + 2$ für $a \geq 0$. L_{LS} kann bei $M = 1$ abgeschnitten werden (φ hat ein globales Minimum bei $t = 1$ und damit $L(1, \cdot)$ bei $t = 1$ und $L(-1, \cdot)$ bei $t = -1$). ■

2-2.4 **BSP** *Hinge loss*.

$$L_{\text{Hinge}}(\gamma, t) = \max\{0, 1 - \gamma t\}, \quad \Rightarrow \quad \varphi(t) = \max\{0, 1 - t\}.$$

L_{Hinge} ist konvex, lipschitz und $|L_{\text{Hinge}}|_1 = 1$. Sie ist *nicht* strikt konvex und kann bei $M = 1$ abgeschnitten werden. ■

2-2.5 **BSP** *Quadrierter Hinge-Loss*.

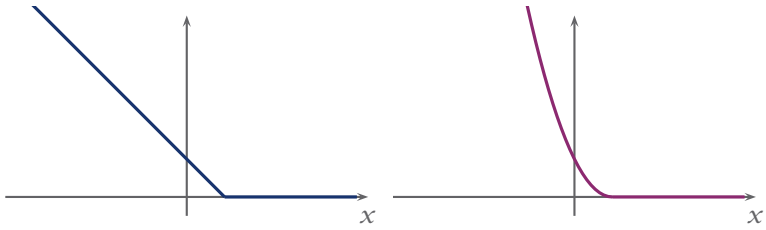
$$L(\gamma, t) = (\max\{0, 1 - \gamma t\})^2, \quad \Rightarrow \quad \varphi(t) = (\max\{0, 1 - t\})^2.$$

Sie stellt eine Mischung aus LS und HL dar. L ist konvex (nicht strikt), lokal lipschitz, $|L|_{a,1} = 2a + 2$ und kann bei $M = 1$ abgeschnitten werden. ■

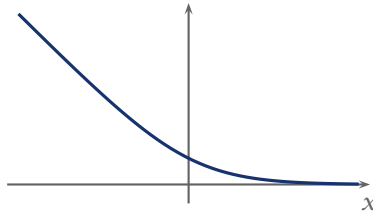
2-2.6 **BSP** *Logistische Verlustfunktion für Klassifikation*.

$$L_{\log}(\gamma, t) = \log(1 + \exp(-\gamma t)), \quad \Rightarrow \quad \varphi(t) = \log(1 + \exp(-\gamma t)).$$

L_{\log} ist strikt konvex, lipschitz $|L_{\log}|_1 = 1$, kann aber *nicht* abgeschnitten werden. ■



2.1 Hinge loss und quadrierte hinge loss Verlustfunktion



2.2 Logistische Verlustfunktion

In Übung 6 wurde gezeigt

$$\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^* \leq \sqrt{\mathcal{R}_{L_{\text{LS}},P}(f) - \mathcal{R}_{L_{\text{LS}},P}^*}$$

2-2.7 **Lemma** Für $\eta \in [0, 1]$ und $t \in [-1, 1]$ gilt

$$|2\eta - 1| \mathbf{I}_{(-\infty, 0]}((2\eta - 1) \text{sign } t) \leq |2\eta - 1| |t - \text{sign}(2\eta - 1)|. \quad \times$$

- » 1. Fall $\eta = \frac{1}{2}$ ist klar.
 2. Fall $\eta < \frac{1}{2}$. Für $t \in [-1, 0)$ ist

$$\underbrace{(2\eta - 1)}_{<0} \underbrace{\text{sign } t}_{<0} > 0,$$

somit verschwindet die linke Seite und die rechte ist ≥ 0 .

Für $t \in [0, 1)$ gilt umgekehrt,

$$\underbrace{(2\eta - 1)}_{<0} \underbrace{\text{sign } t}_{=1} < 0,$$

somit ist die linke Seite gleich $|2\eta - 1| < 1$ und $|t - \text{sign}(2\eta - 1)| = |t + 1| \geq 1$.
 «

2-2.8 **Zhang's (Un-)Gleichung** Sei $f_{L_{\text{class}}}^*(x) := \text{sign}(2\eta(x) - 1)$ für $x \in X$. Dann gilt für $f : X \rightarrow [-1, 1]$:

$$\mathcal{R}_{L_{\text{hinge}},P}(f) - \mathcal{R}_{L_{\text{hinge}},P}^* = \int_X |f(x) - f_{L_{\text{class}},P}^*(x)| |2\eta(x) - 1| dP_X(x)$$

„Überschuss-Lhinge-Risiko = gewichtete L^1 -Norm von $f - f_{L_{\text{class}},P}^*$ “.
 Und für $h : X \rightarrow \mathbb{R}$ gilt

$$\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^* \leq \mathcal{R}_{L_{\text{hinge}},P}(f) - \mathcal{R}_{L_{\text{hinge}},P}^* \quad \times$$

» $L_{\text{hinge}}(y, t) = \max\{0, 1 - yt\} \stackrel{t \in [-1, 1]}{=} 1 - yt$. Für $f : X \rightarrow [-1, 1]$ gilt daher

$$\begin{aligned} \mathcal{R}_{L_{\text{hinge}},P}(f) &= \int_X \eta(x)(1 - f(x)) + (1 - \eta(x))(1 + f(x)) dP_X \\ &= \int_X 1 + f(x)(1 - 2\eta(x)) dP_X. \end{aligned}$$

$f(x)(1 - 2\eta(x))$ ist minimal genau dann, wenn

$$\left. \begin{aligned} f &= 1 \text{ auf } \left\{ \eta > \frac{1}{2} \right\} \\ f &= -1 \text{ auf } \left\{ \eta < \frac{1}{2} \right\} \end{aligned} \right\} \text{d.h. } f = f_{L_{\text{class}},P}^* \text{ auf } \left\{ \eta \neq \frac{1}{2} \right\}$$

Da L_{hinge} bei $M = 1$ abgeschnitten werden kann folgt somit

$$\mathcal{R}_{L_{\text{hinge}},P}^* = \inf_{f: X \rightarrow \mathbb{R}} \mathcal{R}_{L_{\text{hinge}},P}(f) = \inf_{f: X \rightarrow [-1, 1]} \mathcal{R}_{L_{\text{hinge}},P}(f).$$

Damit ist

$$\begin{aligned} \mathcal{R}_{L_{\text{hinge}},P}(f) - \mathcal{R}_{L_{\text{hinge}},P}^* &= \int_X 1 + f \cdot (1 - 2\eta) - 1 - f_{L_{\text{class}},P}^* \cdot (1 - 2\eta) dP_X \\ &= \int_X \underbrace{(f - f_{L_{\text{class}},P}^*) \cdot (1 - 2\eta)}_{\geq 0} dP_X, \end{aligned}$$

wobei der Integrand positiv ist, da $f_{L_{\text{class},P}}^*$ punktweise das Überschussrisiko minimiert und dieses ist positiv, also

$$\mathcal{R}_{L_{\text{hinge},P}(f)} - \mathcal{R}_{L_{\text{hinge},P}}^* = \int_X |f - f_{L_{\text{class},P}}^*| |1 - 2\eta| \, dP_X.$$

Ungleichung. Da L_{hinge} bei $M = 1$ abgeschnitten werden kann, folgt

$$\mathcal{R}_{L_{\text{hinge},P}(\widehat{f})} - \mathcal{R}_{L_{\text{hinge},P}}^* \leq \mathcal{R}_{L_{\text{hinge},P}(f)} - \mathcal{R}_{L_{\text{hinge},P}}^*$$

und

$$\mathcal{R}_{L_{\text{class},P}(\widehat{f})} - \mathcal{R}_{L_{\text{class},P}}^* = \mathcal{R}_{L_{\text{class},P}(f)} - \mathcal{R}_{L_{\text{class},P}}^*.$$

Daher ist ohne Einschränkung $f : X \rightarrow [-1, 1]$. In Kapitel 1-2 haben wir gezeigt

$$\begin{aligned} \mathcal{R}_{L_{\text{class},P}(f)} - \mathcal{R}_{L_{\text{class},P}}^* &= \int_X |2\eta - 1| \mathbf{I}_{(-\infty, 0]}((2\eta - 1) \text{sign } f) \, dP_X \\ &\stackrel{2-2.7}{\leq} \int_X |2\eta - 1| \left| f - \underbrace{\text{sign}(2\eta - 1)}_{=f_{L_{\text{class},P}}^*} \right| \, dP_X \\ &\stackrel{\text{Zhang}}{=} \mathcal{R}_{L_{\text{hinge},P}(f)} - \mathcal{R}_{L_{\text{hinge},P}}^* \quad \ll \end{aligned}$$

2-3 Distanzbasierte Verlustfunktionen

2-3.1 **Definition** Eine Verlustfunktion $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ heißt *distanzbasiert*, falls eine Funktion $\psi : \mathbb{R} \rightarrow [0, \infty)$ existiert mit

$$L(y, t) = \psi(y - t), \quad y \in Y, \quad t \in \mathbb{R}.$$

Eine distanzbasierte Verlustfunktion heißt *symmetrisch*, falls $\psi(r) = \psi(-r)$ für alle $r \in \mathbb{R}$ heißt. \times

2-3.2 **Lemma** Sei L distanzbasiert, so gelten:

- 1) L ist genau dann (strikt) konvex, wenn ψ (strikt) konvex.
- 2) L ist genau dann stetig, wenn ψ stetig.

3) L ist genau dann lipschitz, wenn ψ lipschitz. (Gilt im Allgemeinen nicht für lokal lipschitz, siehe L_{LS})

4) Ist ferner $Y \subset [-M, M]$, so ist L genau dann lokal lipschitz, wenn ψ lokal lipschitz.

5) Ist L konvex, so ist ψ lokal lipschitz.

6) L ist eine P -integrierbare NVF. \times

» 4) „ \Leftarrow “: Sei $t \in [-a, a]$ und $a > 0$, so gilt

$$|L(y, t) - L(y, t')| = |\psi(y - t) - \psi(y - t')| \leq |\psi|_{1+M, 1} |t - t'|.$$

„ \Rightarrow “: Analog folgt $|L|_{a, 1} \leq |\psi|_{a+M, 1}$.

5) Sei L konvex, dann folgt für $y = 0$, dass auch

$$t \mapsto \psi(-t) = \psi(0 - t) = L(0, t)$$

konvex und daher ist ψ konvex und folglich auch lokal lipschitz. Mit 4) folgt nun, dass L lokal lipschitz.

6) folgt aus Lemma 2-1.9. \ll

2-3.3 BSP Kleinste Quadrate

$$L_{LS}(y, t) := (y - t)^2 \Rightarrow \psi(r) = r^2.$$

ψ ist strikt konvex, nicht lipschitz aber lokal lipschitz und symmetrisch. ■

2-3.4 BSP Betragsfunktion.

$$L_{\text{abs}}(y, t) := |y - t| \Rightarrow \psi(r) = |r|.$$

ψ ist konvex (nicht strikt), symmetrisch und Lipschitz mit $|L|_1 = 1$.

2-3.5 BSP Pinball für $\varepsilon \in (0, 1)$.

$$\psi(r) := \begin{cases} -1(1 - \varepsilon)r, & r < 0, \\ \varepsilon r, & r \geq 0 \end{cases}$$

ψ ist konvex (nicht strikt), symmetrisch genau dann, wenn $\varepsilon = \frac{1}{2}$ und lipschitz mit $|L|_1 = \min\{\varepsilon, 1 - \varepsilon\}$. ■

3 Konzentrationsungleichungen

Im Lemma 1-3.7 hatten wir gezeigt:

$$\mu^n \left(\left\{ \omega \in \Omega^n : \frac{1}{n} \sum_{i=1}^n h(\omega_i) - \mathbf{E}_\mu h \geq t \right\} \right) \leq \frac{\mathbf{E}_\mu h^2}{t^2 n}$$

für alle Maße μ , $n \geq 1$, $t > 0$ und $h : \Omega \rightarrow \mathbb{R}$.

3-0.6 **Satz** Sei $(\Omega, \mathcal{A}, \mu)$ ein Wahrscheinlichkeitsraum, dann gilt für alle $f : \Omega \rightarrow \mathbb{R}$ mit $\mathbf{E}_\mu |f| < \infty$ und $t > 0$,

$$\mu \left(\{ \omega \in \Omega : f(\omega) \geq t \} \right) \leq \frac{\mathbf{E}_\mu |f|}{t}. \quad \times$$

3-0.7 **Lemma** Für $x > -1$ gilt

$$(1+x) \ln(1+x) - x \geq \frac{3}{2} \frac{x^2}{x+3}. \quad \times$$

» Seien $f(x) = (1+x) \ln(1+x) - x$ und $g(x) = \frac{3}{2} \frac{x^2}{x+3}$. Dann folgt

$$\begin{aligned} f'(x) &= \ln(1+x), & g'(x) &= \frac{3}{2} \frac{x^2 + 6x}{(x+3)^2}, \\ f''(x) &= \frac{1}{1+x}, & g''(x) &= \frac{27}{(x+3)^3}, \end{aligned}$$

und es gelten

$$\begin{aligned} f'(0) &= 0, & f(0) &= 0, \\ f''(0) &= 0, & g(0) &= 0. \end{aligned}$$

Wir zeigen nun $f''(x) \geq g''(x)$ für alle $x > -1$. Es ist

$$\begin{aligned} 0 &\leq x^2(x+9) = x^3 + 9x^2 \\ \Rightarrow x^3 + 9x^2 + 27x + 27 &= (x+3)^3 \geq 27(1+x) = 27x + 27 \\ \Rightarrow \frac{27}{(x+3)^3} &\leq \frac{1}{1+x}. \end{aligned}$$

Für $x \geq 0$ folgt mit dem Hauptsatz der Differential- und Integralrechnung,

$$f'(x) = f'(x) - f'(0) = \int_0^x f''(t) dt \geq \int_0^x g''(t) dt = g'(x) - g'(0) = g'(x).$$

Das selbe Argument angewandt auf f liefert, $f(x) \geq g(x)$ und für $x \in (-1, 0]$ folgt $-f(x) \geq -g(x)$. \llcorner

3-0.8 **Bernsteins Ungleichung** Sei (Ω, \mathcal{A}, P) ein Wahrscheinlichkeitsraum und $B > 0$, $\sigma > 0$ und $n \geq 1$. Ferner seien $\xi_1, \dots, \xi_n : \Omega \rightarrow \mathbb{R}$ i.i.d. (unabhängige und identisch verteilte) Zufallsvariablen, so dass gilt

$$\left. \begin{aligned} \mathbf{E}_P \xi_i &= 0, \\ \|\xi_i\|_\infty &\leq B, \\ \mathbf{E}_P \xi_i^2 &\leq \sigma^2 \end{aligned} \right\} \text{für alle } i = 1, \dots, n.$$

Dann gilt

$$P\left(\frac{1}{n} \sum_{i=1}^n \xi_i \geq \sqrt{\frac{2\sigma^2\tau}{n}} + \frac{2B\tau}{3n}\right) \leq e^{-\tau}, \quad \text{für alle } \tau > 0. \quad \times$$

» Fixiere $t \geq 0$ und $\varepsilon > 0$, dann gilt

$$\begin{aligned} P\left(\sum_{i=1}^n \xi_i \geq \varepsilon n\right) &= P\left(\exp\left(t \sum_{i=1}^n \xi_i\right) \geq \exp(t\varepsilon n)\right) \\ &\stackrel{\text{Markov}}{\leq} e^{-t\varepsilon n} \mathbf{E}_P \exp\left(t \sum_{i=1}^n \xi_i\right) \\ &\stackrel{\text{Unabh.}}{=} e^{-t\varepsilon n} \prod_{i=1}^n \mathbf{E}_P e^{t\xi_i}. \end{aligned}$$

Folglich ist

$$\mathbf{E}_P e^{t\xi_i} = \mathbf{E}_P \sum_{k=0}^{\infty} \frac{t^k}{k!} \xi_i^k = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbf{E}_P \xi^k$$

und da $\mathbf{E}_P \xi_i = 0$ und $\mathbf{E}_P \xi^k \leq \sigma^2 B^{k-2}$ folgt

$$\begin{aligned} \mathbf{E}_P e^{t\xi} &= 1 + \sum_{k=2}^{\infty} \frac{t^k}{k!} \mathbf{E}_P \xi^k \leq 1 + \sum_{k=2}^{\infty} \frac{t^k}{k!} \sigma^k B^{k-2} = 1 + \frac{\sigma^2}{B^2} \sum_{k=2}^{\infty} \frac{t^k}{k!} B^k \\ &= 1 + \frac{\sigma^2}{B^2} (e^{tB} - tB - 1). \end{aligned}$$

Damit gilt

$$\begin{aligned} P \left(\sum_{i=1}^n \xi_i \geq \varepsilon n \right) &\leq e^{-tn\varepsilon} \left(1 + \frac{\sigma^2}{B^2} (e^{tB} - tB - 1) \right)^n \\ &\leq e^{-tn\varepsilon} \exp \left(n \frac{\sigma^2}{B^2} (e^{tB} - tB - 1) \right) \end{aligned}$$

Sei $h(t) := -tn\varepsilon + \frac{n\sigma^2}{B^2} (e^{tB} - tB - 1)$, dann suchen wir kritische Punkte

$$h'(t) = -n\varepsilon + \frac{n\sigma^2}{B^2} (Be^{tB} - B) = -n\varepsilon + \frac{n\sigma^2}{B} (e^{tB} - 1) \stackrel{!}{=} 0.$$

t ist genau dann kritisch, wenn

$$\frac{\sigma^2 n}{B} e^{tB} = \varepsilon n + \frac{\sigma^2 n}{B} \Leftrightarrow e^{tB} = \frac{\varepsilon B}{\sigma^2} + 1 \Leftrightarrow t^* = \frac{1}{B} \log \left(1 + \frac{\varepsilon B}{\sigma^2} \right).$$

t^* ist der einzige Kandidat für ein Optimum. Weiterhin ist

$$\lim_{t \rightarrow \pm\infty} h(t) = \infty$$

und folglich hat h ein globales Minimum bei t^* . Setzen wir $y = \frac{\varepsilon B}{\sigma^2}$, dann gilt

$t^* = \frac{1}{B} \log(1 + \gamma)$ und somit

$$\begin{aligned} & -t^* \varepsilon n + \frac{\sigma^2 n}{B^2} (e^{t^* B} - t^* B - 1) \\ &= -\frac{\varepsilon n}{B} \log(1 + \gamma) + \frac{\sigma^2 n}{B^2} (1 + \gamma - \log(1 + \gamma) - 1) \\ &= \frac{\sigma^2 n}{B^2} (-\gamma \log(1 + \gamma) + \gamma - \log(1 + \gamma)) \\ &= -\frac{\sigma^2 n}{B^2} ((1 + \gamma) \log(1 + \gamma) - \gamma) \\ &\stackrel{3-0.7}{\leq} -\frac{\sigma^2 n}{B^2} \frac{3}{2} \frac{\gamma^2}{\gamma + 3} \stackrel{\text{Rechnung}}{=} -\frac{3n\varepsilon^2}{2\varepsilon B + 6\sigma^2}. \end{aligned}$$

Damit ist

$$P\left(\frac{1}{n} \sum_{i=1}^n \xi_i \geq \varepsilon\right) \leq \exp\left(-\frac{3n\varepsilon^2}{2\varepsilon B + 6\sigma^2}\right).$$

Wähle nun

$$\tau := \frac{3n\varepsilon^2}{2\varepsilon B + 6\sigma^2},$$

dann folgt

$$\varepsilon = \sqrt{\frac{2\sigma^2\tau}{n} + \frac{B^2\tau^2}{9n^2}} + \frac{B\tau}{3n} \leq \sqrt{\frac{2\sigma^2\tau}{n} + \frac{2B\tau}{3n}}. \quad \ll$$

3-0.9 **Hoeffdings Ungleichung** Sei (Ω, \mathcal{A}, P) ein Wahrscheinlichkeitsraum, $a < b$, $1 \leq n \in \mathbb{N}$ und $\xi_1, \dots, \xi_n : \Omega \rightarrow [a, b]$ i.i.d. Dann gilt für $\tau > 0$,

$$P\left(\frac{1}{n} \sum_{i=1}^n (\xi_i - \mathbb{E}_P \xi_i) \geq (b - a) \sqrt{\frac{\tau}{2n}}\right) \leq e^{-\tau}. \quad \times$$

Beobachtung. Für $a = -b$ betrachte $\eta_i = \xi_i - \mathbb{E}_P \xi_i$.

Dann ist $\mathbb{E}_P \eta_i = 0$, $\mathbb{E}_P \eta_i^2 = \mathbb{E}_P \xi_i^2 - (\mathbb{E}_P \xi_i)^2 \leq b^2 =: \sigma^2$ und $\|\eta_i\|_\infty \leq 2b =: B$. Mit Bernstein folgt,

$$P\left(\frac{1}{n} \sum_{i=1}^n \eta_i \geq \sqrt{\frac{2b^2\tau}{n} + \frac{4b\tau}{3n}}\right) \leq e^{-\tau},$$

wobei $\sqrt{\frac{2b^2\tau}{n}} = 2b\sqrt{\frac{\tau}{2n}} = (b - a)\sqrt{\frac{\tau}{2n}}. \quad \rightarrow$

3-0.10 **Lemma (Union bound)** Sei (Ω, \mathcal{A}, P) ein Wahrscheinlichkeitsraum, $f_1, \dots, f_n : \Omega \rightarrow \mathbb{R}$ messbar. Dann gilt für alle $t \in \mathbb{R}$

$$\mu \left(\sup_{i=1, \dots, n} f_i \geq t \right) \leq \sum_{i=1}^n \mu(f_i \geq t). \quad \times$$

$$\gg \left\{ \sup_{i=1, \dots, n} f_i \geq t \right\} = \bigcup_{i=1}^n \left\{ f_i \geq t \right\}. \quad \ll$$

Bemerkung. Betrachte $\xi_1, \dots, \xi_n, -\xi_1, \dots, -\xi_n$ und das union bound für die Summe über 1. Block (= f_1) + 2. Block (= f_2). Somit ist

$$\left| \frac{1}{n} \sum_{i=1}^n \xi_i \right| = \sup \{f_1, f_2\}$$

und folglich gilt

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n \xi_i \right| \geq \sqrt{\frac{2\sigma^2\tau}{n}} + \frac{2B\tau}{3n} \right) \leq 2e^{-\tau}. \quad \rightarrow$$

„Zweiseitige Bernstein-Ungleichung“ (Hoeffdings analog!)

4 Empirische Risikominimierung (ERM)

Motivation. Sei D das empirische Wahrscheinlichkeitsmaß,

$$\mathcal{R}_{L,D}(f) := \frac{1}{n} \sum_{i=1}^n L(x_i, y_i, f(x_i)) \xrightarrow{n \rightarrow \infty} \mathcal{R}_{L,P}(f).$$

Ziel des Lernens: Finde f mit $\mathcal{R}_{L,P}(f)$ klein.

Heuristik von ERM: Minimiere Schätzung $\mathcal{R}_{L,D}(\cdot)$ von $\mathcal{R}_{L,P}(\cdot)$. \rightarrow

4-1 ERM über endliche Funktionenmengen

4-1.1 **Definition** Sei L eine Verlustfunktion, $\emptyset \neq \mathcal{F} \subset \mathcal{L}_0(X)$. Dann heißt eine Lernmethode \mathcal{L} **ERM** über \mathcal{F} , wenn

$$\mathcal{R}_{L,D}(f_D) = \inf_{f \in \mathcal{F}} \mathcal{R}_{L,D}(f), \quad \forall n \geq 1, D \in (X \times Y)^n \text{ und } f_D \in \mathcal{F}. \quad \times$$

Bemerkungen. A. Im Allgemeinen gibt es keine ERM.

B. Im Allgemeinen ist ERM nicht eindeutig.

C. Für $\mathcal{F} = \mathcal{L}_0(X)$ oder $\mathcal{F} = \mathcal{L}_\infty(X)$ erhält man im Allgemeinen overfitted ERM.

D. Für „kleine“ \mathcal{F} erhält man im Allgemeinen underfitted ERM.

Im Folgenden sei \mathcal{F} endlich.

4-1.2 **Orakelungleichung I** Sei L eine Verlustfunktion und $\mathcal{F} \subset \mathcal{L}_0(X)$ nicht leer und endlich. Ferner sei $B > 0$ mit

$$L(x, y, f(x)) \leq B, \quad \forall (x, y) \in X \times Y, \quad f \in \mathcal{F}.$$

Für den Wert des Orakels

$$\mathcal{R}_{L,P,\mathcal{F}}^* := \inf_{f \in \mathcal{F}} \mathcal{R}_{L,P}(f)$$

und ERM gilt für $n \geq 1$ und $\tau > 0$,

$$P^n \left(\left\{ D \in (X \times Y)^n : \mathcal{R}_{L,P}(f_D) \leq \mathcal{R}_{L,P,\mathcal{F}}^* + B \sqrt{\frac{2\tau + 2 \log(2|\mathcal{F}|)}{n}} \right\} \right) \geq 1 - e^{-\tau}. \quad \times$$

» Für $\delta > 0$ existiert ein $f_\delta \in \mathcal{F}$ mit $\mathcal{R}_{L,P}(f_\delta) \leq \mathcal{R}_{L,P,\mathcal{F}}^* + \delta$. Damit ist

$$\begin{aligned} \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P,\mathcal{F}}^* &\leq \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,D}(f_D) + \mathcal{R}_{L,D}(f_D) - \mathcal{R}_{L,P}(f_\delta) + \delta \\ &\leq \left| \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,D}(f_D) \right| + \left| \mathcal{R}_{L,P}(f_\delta) - \mathcal{R}_{L,D}(f_\delta) \right| + \delta \\ &\leq 2 \sup_{f \in \mathcal{F}} \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| + \delta. \end{aligned}$$

Für „ $\delta \rightarrow 0$ “ folgt

$$\mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P,\mathcal{F}}^* \leq 2 \sup_{f \in \mathcal{F}} \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right|$$

Damit folgt

$$\begin{aligned} &P^n \left(\left\{ D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P,\mathcal{F}}^* \geq B \sqrt{\frac{2\tau}{n}} \right\} \right) \\ &\leq P^n \left(\left\{ D : \sup_{f \in \mathcal{F}} \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| \geq B \sqrt{\frac{\tau}{2n}} \right\} \right) \\ &\stackrel{\text{union bound}}{\leq} \sum_{f \in \mathcal{F}} P^n \left(\left\{ D : \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| \geq B \sqrt{\frac{\tau}{2n}} \right\} \right). \end{aligned}$$

Setzen wir $\mathbf{E}_P \xi_1 = \mathcal{R}_{L,P}(f)$ und $\xi_i = L(x_i, y_i, f(x_i))$, so ist $\mathcal{R}_{L,D}(f) = \frac{1}{n} \sum_{i=1}^n \xi_i$ und mit der 2-seitigen Hoeffding Ungleichung folgt

$$\dots \leq \sum_{f \in \mathcal{F}} 2e^{-\tau} = 2|\mathcal{F}| e^{-\tau} = e^{\log(2|\mathcal{F}|) - \tau}.$$

Mit $-t' = \log(2|\mathcal{F}|) - \tau$ folgt die Behauptung. «

Idee: Falls $|\mathcal{F}| = \infty$ finde $\mathcal{F}_\varepsilon \subset \mathcal{F}$ endlich, das \mathcal{F} gut „approximiert“.

4-2 ERM über unendlichen Funktionenmengen

4-2.1 **Definition** Sei (T, d) ein metrischer Raum (MR) und $\varepsilon > 0$.

(a) $S \subset T$ heißt ε -Netz falls

$$\forall t \in T \exists s \in S : d(t, s) \leq \varepsilon.$$

(b) Die **Überdeckungszahl** (covering number) ist definiert als

$$\mathcal{N}(T, d, \varepsilon) := \inf \left\{ n \geq 1 : \exists s_1, \dots, s_n \in T : T \subset \bigcup_{i=1}^n B_d(s_i, \varepsilon) \right\}.$$

D.h. $S = \{s_1, \dots, s_n\}$ ist ε -Netz. Damit entspricht $\mathcal{N}(T, d, \varepsilon)$ der Größe des kleinsten ε -Netzes von T .

(c) Sei $T \subset E$, E normierter Raum. Dann ist

$$\mathcal{N}(T, \|\cdot\|_E, \varepsilon) := \mathcal{N}(T, d_E, \varepsilon), \quad d_E(x, x') = \|x - x'\|_E. \quad \times$$

Interpretation. Ist T kompakt, dann ist $\mathcal{N}(T, d, \varepsilon) < \infty$ für jedes $\varepsilon > 0$. Genauer ist T genau dann präkompakt, wenn $\mathcal{N}(T, d, \varepsilon) < \infty$ für jedes $\varepsilon > 0$.

Da für $\varepsilon < \varepsilon'$ gilt $\mathcal{N}(T, d, \varepsilon') < \mathcal{N}(T, d, \varepsilon)$, stellt das Wachstumsverhalten von $\mathcal{N}(T, d, \varepsilon)$ für $\varepsilon \downarrow 0$ ein quantitatives Maß für die Kompaktheit dar.

$\mathcal{N}(T, d, \varepsilon) \rightarrow \infty$ für $\varepsilon \rightarrow 0$ genau dann, wenn T unendlich. \rightarrow

4-2.2 **Definition** Sei (T, d) ein metrischer Raum und $n \geq 1$.

(a) Die n -te (b -a-dische) **Entropiezahl** von T ist

$$e_n(T, d) = \inf \left\{ \varepsilon > 0 : \exists s_1, \dots, s_{2^{n-1}} \in T : T \subset \bigcup_{i=1}^{2^{n-1}} B_d(s_i, \varepsilon) \right\}.$$

(b) $e_n(T, \|\cdot\|_E)$ für $T \subset (E, \|\cdot\|_E)$ ist wie üblich definiert.

(c) Sei $S : E \rightarrow F$ ein stetige linearer Operator zwischen den normierten Räumen E und F und es sei $B_E := \{x \in E : \|x\|_E \leq 1\}$, dann schreibe

$$e_n(S) = e_n(S : E \rightarrow F) := e_n(SB_E, \|\cdot\|_F). \quad \times$$

Bemerkungen (ohne Beweis). A. $e_1(S) = \|S\|$.

B. $e_{n+1}(S) \leq e_n(S)$ für $n \geq 1$.

C. $\lim_{n \rightarrow \infty} e_n(S) = 0$ genau dann, wenn S kompakt.

D. $e_{n+k-1}(S+T) \leq e_n(S) + e_k(T)$.

E. $e_{n+k-1}(S \circ T) \leq e_n(S)e_k(T)$.

F. Sei $d = \text{rank } S = \dim SE$. Dann ist $d < \infty$ genau dann, wenn

$$\exists c_1, c_2 : c_1 2^{-\frac{n-1}{d}} \leq e_n(S) \leq c_2 2^{-\frac{n-1}{d}}.$$

G. Sei $X \subset \mathbb{R}^d$ offen und beschränkt und

$$C_b^m(X) := \left\{ f \in C^m(X \rightarrow \mathbb{R}) : \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty < \infty \right\}$$

versehen mit der Norm

$$\|f\|_{C_b^m(X)} = \max_{|\alpha| \leq m} \|\partial^\alpha f\|_\infty.$$

Dann ist $e_n(\text{id} : C_b^m(X) \rightarrow C_b(X)) \leq C_{m,d}(X) n^{-\frac{m}{d}}$ und die Abschätzung ist scharf, d.h. es gibt ein $\tilde{C}_{m,d}(X)$ mit $e_n(\dots) \geq \tilde{C}_{m,d}(X) n^{-\frac{m}{d}}$. Weiterhin ist

$$C_{m,d}(r \cdot X) = r^m C_{m,d}(X), \quad r \geq 1. \quad \circ$$

4-2.3 **Lemma** Sei (T, d) MR und $a, q > 0$ mit

$$e_n(T, d) \leq a n^{-\frac{1}{q}}, \quad n \geq 1.$$

Dann gilt

$$\log(\mathcal{N}(T, d, \varepsilon)) \leq \log(4) \left(\frac{a}{\varepsilon}\right)^q, \quad \varepsilon > 0. \quad \times$$

» Für $\varepsilon > a$ ist $\mathcal{N}(T, d, \varepsilon) = 1$, d.h. $\log(\mathcal{N}(T, d, \varepsilon)) = 0$ und es ist nichts zu beweisen. Sei $\delta > 0$ und $\varepsilon \in (0, a]$, dann existiert ein $n \geq 1$ mit

$$a(1 + \delta)(n + 1)^{-\frac{1}{q}} \leq \varepsilon \leq a(1 + \delta)n^{-\frac{1}{q}}. \quad (*)$$

Da $e_n(T, d) < a(1 + \delta)n^{-\frac{1}{q}}$, gibt es ein $a(1 + \delta)n^{-\frac{1}{q}}$ -Netz der Größe 2^{n-1} , d.h.

$$\mathcal{N}(T, d, a(1 + \delta)n^{-\frac{1}{q}}) \leq 2^{n-1}.$$

Nun ist

$$\begin{aligned} a(1 + \delta)n^{-\frac{1}{q}} &= a(1 + \delta)(n + 1)^{-\frac{1}{q}} \underbrace{\left(\frac{n + 1}{n}\right)^{\frac{1}{q}}}_{\leq 2} \leq 2^{\frac{1}{q}} a(1 + \delta)(n + 1)^{-\frac{1}{q}} \\ &\stackrel{(*)}{\leq} 2^{\frac{1}{q}} \varepsilon. \end{aligned}$$

Und ferner folgt mit der rechten Seite von (*),

$$n \leq \left(\frac{a(1 + \delta)}{\varepsilon}\right)^q.$$

Somit ist

$$\begin{aligned} \log(\mathcal{N}(T, d, 2^{1/q}\varepsilon)) &\leq \log(\mathcal{N}(T, d, a(1 + \delta)n^{-1/q})) \leq \log(2^{n-1}) \\ &\leq n \log 2 \leq (\log 2) \left(\frac{a(1 + \delta)}{\varepsilon}\right)^q \end{aligned}$$

Für $\delta \rightarrow 0$ und $2^{1/q}\varepsilon = \tilde{\varepsilon}$ folgt die Behauptung. «

4-2.4 **Orakelungleichung II** Sei $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ eine lokal lipschitz-stetige Verlustfunktion und P ein W -Maß auf $(X \times Y)$. Ferner sei

$$\mathcal{F} \subset \mathcal{L}_\infty(X) = \{f : X \rightarrow \mathbb{R} : f \text{ messbar und beschränkt}\}$$

für die $B > 0$ und $M > 0$ existieren, so dass

$$\begin{aligned} \|f\|_\infty &< M, && \text{für alle } f \in \mathcal{F}, \\ L(x, y, f(x)) &\leq B, && \text{für alle } (x, y) \in (X \times Y), f \in \mathcal{F}. \end{aligned}$$

Dann gilt für ERM über f ,

$$P^n \left(\left\{ D : \mathcal{R}_{L,P}(f_D) < \mathcal{R}_{L,P,\mathcal{F}}^* + B \sqrt{\frac{2\tau + 2 \log(2\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon))}{n}} + 4\varepsilon |L|_{M,1} \right\} \right) \geq 1 - e^{-\tau}, \quad \text{für alle } n \geq 1, \varepsilon > 0 \text{ und } \tau > 0. \quad \times$$

Beachte, dass falls $\log(\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)) \leq \mu(\varepsilon)$ für $\varepsilon > 0$, der Ausdruck nach ε optimiert werden kann. Somit erhält man eine Abschätzung in $n \geq 1$ und $\tau > 0$.

» Fixiere $\varepsilon > 0$ und ein minimales ε -Netz \mathcal{F}_ε von \mathcal{F} bezüglich $\|\cdot\|_\infty$, d.h.

$$\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon) = |\mathcal{F}_\varepsilon|$$

und

$$\forall f \in \mathcal{F} \exists g \in \mathcal{F}_\varepsilon \text{ mit } \|f - g\|_\infty \leq \varepsilon.$$

Somit ist für solches f, g

$$\begin{aligned} & \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| \\ & \leq \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}(g) \right| + \left| \mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g) \right| + \left| \mathcal{R}_{L,D}(g) - \mathcal{R}_{L,D}(f) \right| \\ & \stackrel{??}{\leq} 2\varepsilon |L|_{M,1} + \left| \mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g) \right|. \end{aligned}$$

D.h.

$$\sup_{f \in \mathcal{F}} \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| \leq 2\varepsilon |L|_{M,1} + \sup_{g \in \mathcal{F}_\varepsilon} \left| \mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g) \right|.$$

Der Beweis von 4-1.2 zeigte außerdem, dass

$$\begin{aligned} \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P,\mathcal{F}}^* & \leq 2 \sup_{f \in \mathcal{F}} \left| \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f) \right| \\ & \leq 4\varepsilon |L|_{M,1} + 2 \sup_{g \in \mathcal{F}_\varepsilon} \left| \mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g) \right|. \end{aligned}$$

Anwendung der Hoeffding Ungleichung und es union bounds liefert,

$$\begin{aligned}
 & P^n \left(\left\{ D : \mathcal{R}_{L,P}(f_D) - \mathcal{R}_{L,P,F}^* \geq B \sqrt{\frac{2\tau}{n}} + 4\varepsilon |L|_{M,1} \right\} \right) \\
 & \leq P^n \left(\left\{ D : \sup_{g \in \mathcal{F}_\varepsilon} |\mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g)| \geq B \sqrt{\frac{\tau}{2n}} \right\} \right) \\
 & \leq 2\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon) e^{-\tau} = \exp(\log(2\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)) - \tau).
 \end{aligned}$$

Die Behauptung folgt nun mit einer Variablentransformation

$$-\tau' = \log(2\mathcal{N}(\mathcal{F}, \|\cdot\|_\infty, \varepsilon)) - \tau. \quad \ll$$

5 Reproduzierende Kernhilberträume (RKHS)

Es gibt zwei schöne Klassen von Banachräumen

- Funktionenräume.
- Hilberträume.

RKHS sind beides, der $L^2([0, 1])$ übrigens *kein* Funktionenraum!

5-1 Kerne und Beispiele

5-1.1 **Definition** Eine Funktion $k : X \times X \rightarrow \mathbb{R}$ für die ein Hilbertraum H und ein $\Phi : X \rightarrow H$ existieren, so dass

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle_H, \quad \forall x, x' \in X,$$

heißt **Kern**. \times

Bemerkungen. A. Im Allgemeinen sind H und Φ nicht eindeutig.

Betrachte z.B. $k(x, x') = x \cdot x'$ für $x, x' \in \mathbb{R}$, so ist k ein Kern, denn

$$H_1 = \mathbb{R}, \quad \Phi_1 = \text{id}_{\mathbb{R}}, \quad \langle \Phi_1(x), \Phi_1(x') \rangle = x \cdot x'.$$

$$H_2 = \mathbb{R}^2, \quad \Phi_2(x) = \frac{1}{\sqrt{2}} (x, x)$$

B. H heißt **Featurespace** (Merkmalraum), Φ heißt Featuremap. \rightarrow

5-1.2 **Lemma** Seien $f_n : X \rightarrow \mathbb{R}$, so dass $(f_n(x))_{n \geq 1} \in l_2$ für alle $x \in X$, so definiert

$$k(x, x') = \sum_{n \geq 1} f_n(x) f_n(x'), \quad x, x' \in X$$

einen Kern. \times

» k ist wohl definiert, denn nach Cauchy-Schwarz ist

$$\sum_{n \geq 1} |f_n(x) f_n(x')| \leq \|(f_n(x))\|_{l_2} \|(f_n(x'))\|_{l_2} < \infty.$$

Sei $H = l^2$ und $\Phi : X \rightarrow H$ gegeben durch

$$\Phi(x) = (f_n(x))_{n \geq 1}, \quad x \in X.$$

Dann ist $\langle \Phi(x), \Phi(x') \rangle = k(x, x')$ und daher k ein Kern. \llcorner

5-1.3 **Lemma** Ist k ein Kern auf X und $A : X' \rightarrow X$ eine Abbildung, dann definiert

$$k'(x, x') = k(A(x), A(x')), \quad x, x' \in X'$$

einen Kern auf X' . \times

» Sei H ein Featurespace von k und $\Phi : X \rightarrow H$ die zugehörige Featuremap. Setze

$$\Psi(x) = \Phi(A(x)),$$

so gilt

$$k'(x, x') = k(A(x), A(x')) = \langle \Phi(A(x)), \Phi(A(x')) \rangle_H = \langle \Psi(x), \Psi(x') \rangle_H,$$

d.h. Ψ ist Featuremap von k' . \llcorner

5-1.4 **Lemma** Sind k_1 und k_2 Kerne auf X und $\alpha \geq 0$. Dann sind auch $k_1 + k_2$ und $\alpha \cdot k_1$ Kerne auf X . \times

» Der Beweis ist eine leichte Übung. \llcorner

5-1.5 **Lemma** Sind k_1 und k_2 Kerne auf X , so ist auch $k_1 \cdot k_2$ ein Kern auf X . \times

» Der Beweis erfolgt unter Verwendung des Tensorprodukts von Hilberträumen. Wir überspringen diesen, da wir die Aussage nicht benötigen. «

5-1.6 **Lemma für Polynom-Kerne** Seien $m \geq 0$, $d \geq 1$ natürliche Zahlen und $0 \leq c \in \mathbb{R}$. Dann definiert

$$k(x, x') = (\langle x, x' \rangle_{\mathbb{R}^d} + c)^m, \quad x, x' \in \mathbb{R}^d$$

einen Kern auf \mathbb{R}^d . \times

» Nach Konstruktion existiert ein Polynom p der Ordnung m mit

$$k(x, x') = p(\langle x, x' \rangle), \quad x, x' \in \mathbb{R}^d,$$

wobei p nur nichtnegative Koeffizienten besitzt. Es genügt daher zu zeigen, dass die Monome

$$(x, x') \mapsto \langle x, x' \rangle^m$$

einen Kern definieren. Nun ist

$$\langle x, x' \rangle_{\mathbb{R}^d}^m = \left(\sum_{i=1}^d x_i x'_i \right)^m$$

also können wir Lemma 5-1.2 anwenden mit $f_i(x) = \pi_i(x)$ der Projektion auf die i -te Koordinate und deren Potenzen. «

5-1.7 **Lemma für Taylor-Kerne** Sei $f: (-r, r) \rightarrow \mathbb{R}$ mit $r \in (0, \infty]$,

$$f(x) = \sum_{n=0}^{\infty} a_n t^n, \quad t \in (-r, r).$$

Falls $a_n \geq 0$ für alle $n \geq 0$, definiert

$$k(x, x') = f(\langle x, x' \rangle_{\mathbb{R}^d}) = \sum_{n=0}^{\infty} a_n \langle x, x' \rangle_{\mathbb{R}^d}^n$$

einen Kern auf $\sqrt{r} \overset{\circ}{B}_{\mathbb{R}^d} = \{x \in \mathbb{R}^d : \|x\|_2 < \sqrt{r}\}$. \times

» Für $x, x' \in \sqrt{r} \mathring{B}_{\mathbb{R}^d}$ ist

$$|\langle x, x' \rangle_{\mathbb{R}^d}| \leq \|x\|_2 \|x'\|_2 < r$$

und daher die Konstruktion sinnvoll. Ferner seien

$$c_{j_1, \dots, j_d} = \frac{n!}{\prod_{i=1}^d j_i!}, \quad \sum_{j=1}^d j_i = n,$$

dann gilt

$$\begin{aligned} k(x, x') &= \sum_{n=0}^{\infty} a_n \left(\sum_{i=1}^d x_i x'_i \right)^n = \sum_{n=0}^{\infty} a_n \sum_{\substack{j_1, \dots, j_d \geq 0 \\ j_1 + \dots + j_d = n}} c_{j_1, \dots, j_d} \prod_{i=1}^d x_i^{j_i} (x'_i)^{j_i} \\ &= \sum_{n=0}^{\infty} a_n \sum_{\substack{j_1, \dots, j_d \geq 0 \\ j_1 + \dots + j_d = n}} c_{j_1, \dots, j_d} \left(\prod_{i=1}^d x_i^{j_i} \right) \left(\prod_{i=1}^d (x'_i)^{j_i} \right). \end{aligned}$$

Setze $\Phi : X \rightarrow l_2(\mathbb{N}_0^d)$, wobei $l_2(I) := \{(a_i)_{i \in I} : \|(a_i)\|_2 < \infty\}$, dann ist

$$\Phi(x) = \left(\sqrt{a_{j_1 + \dots + j_d}} c_{j_1, \dots, j_d} \prod_{i=1}^d x_i^{j_i} \right)_{(j_1, \dots, j_d) \in \mathbb{N}_0^d}$$

eine Featuremap, denn

$$\begin{aligned} &\sum_{n=0}^{\infty} a_n \sum_{\substack{j_1, \dots, j_d \geq 0 \\ j_1 + \dots + j_d = n}} c_{j_1, \dots, j_d} \left(\prod_{i=1}^d x_i^{j_i} \right) \left(\prod_{i=1}^d (x'_i)^{j_i} \right) \\ &= \sum_{n=0}^{\infty} \sum_{\substack{j_1, \dots, j_d \geq 0 \\ j_1 + \dots + j_d = n}} \sqrt{a_{j_1 + \dots + j_d}} c_{j_1, \dots, j_d} \left(\prod_{i=1}^d x_i^{j_i} \right) \sqrt{a_{j_1 + \dots + j_d}} c_{j_1, \dots, j_d} \left(\prod_{i=1}^d (x'_i)^{j_i} \right) \\ &= \langle \Phi(x), \Phi(x') \rangle_{l_2(\mathbb{N}_0^d)}. \quad \llcorner \end{aligned}$$

5-1.8 **BSP** *Exponential Kern*. Die Abbildung

$$k(x, x') = \exp(\langle x, x' \rangle_{\mathbb{R}^d}), \quad x, x' \in \mathbb{R}^d$$

ist ein Kern. ■

5-1.9 **BSP** *Gauß Kern* oder *Gauß'scher RBF-Kern*. Die Abbildung

$$k_\sigma(x, x') = \exp\left(-\sigma^2 \|x - x'\|_{\mathbb{R}^d}^2\right), \quad x, x' \in \mathbb{R}^d$$

ist ein Kern.

» Eine einfache Rechnung zeigt

$$\begin{aligned} \|x - x'\|^2 &= \langle x - x', x - x' \rangle = \langle x, x \rangle - 2 \langle x, x' \rangle + \langle x', x' \rangle \\ &= \|x\|^2 + \|x'\|^2 - 2 \langle x, x' \rangle. \end{aligned}$$

Setzen wir $h(x) = \exp(\sigma^2 \|x\|_{\mathbb{R}^d}^2)$, so können wir den Gauß-Kern schreiben als

$$\exp\left(-\sigma^2 \|x - x'\|_{\mathbb{R}^d}^2\right) = \frac{\exp(2\sigma^2 \langle x, x' \rangle_{\mathbb{R}^d})}{h(x)h(x')}.$$

Ferner ist $(x, x') \mapsto \exp(2\sigma^2 \langle x, x' \rangle)$ ein Kern nach BSP 5-1.8 und Lemma 5-1.3. Betrachten wir einen zugehörigen Featurespace H und eine Featuremap Φ und setzen

$$\Phi_\sigma = \frac{1}{h} \Phi : X \rightarrow H,$$

so gilt

$$\begin{aligned} \langle \Phi_\sigma(x), \Phi_\sigma(x') \rangle &= \left\langle \frac{\Phi(x)}{h(x)}, \frac{\Phi(x')}{h(x')} \right\rangle = \frac{1}{h(x)h(x')} \langle \Phi(x), \Phi(x') \rangle \\ &= \frac{\exp(2\sigma^2 \langle x, x' \rangle)}{h(x)h(x')}. \quad \ll \blacksquare \end{aligned}$$

5-1.10 **Definition** Eine Funktion $k : X \times X \rightarrow \mathbb{R}$ heißt *positiv definit*, falls

$$\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0$$

für alle $n \geq 1$ und $\alpha_1, \dots, \alpha_n \in \mathbb{R}$, d.h. die Grammatritzen $(k(x_i, x_j))_{i,j=1}^n$ sind positiv definit.

k heißt *strikt positiv definit*, falls für alle paarweise verschiedenen x_1, \dots, x_n und $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ gilt

$$\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) = 0 \Rightarrow \alpha_1 = \dots = \alpha_n = 0.$$

k heißt *symmetrisch* falls $k(x, x') = k(x', x)$ für alle $x, x' \in X$. \times

5-1.11 **Satz** $k : X \times X \rightarrow \mathbb{R}$ ist genau dann ein Kern, wenn k positiv definit und symmetrisch. \times

» „ \Rightarrow “: Sei $\Phi : X \rightarrow H$ eine Featuremap von k , dann gilt

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

und jedes Skalarprodukt ist positiv definit und symmetrisch

$$\begin{aligned} 0 \leq \langle h, h \rangle &= \left\langle \sum_{i=1}^n \alpha_i \Phi(x_i), \sum_{j=1}^n \alpha_j \Phi(x_j) \right\rangle \\ &= \sum_{i,j=1}^n \alpha_i \alpha_j \langle \Phi(x_i), \Phi(x_j) \rangle = \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j). \end{aligned}$$

„ \Leftarrow “: Wir müssen Featurespace und -map konstruieren. Dazu setzen wir

$$H_{\text{pre}} = \left\{ \sum_{i=1}^n \alpha_i k(x_i, \cdot) : n \geq 1, \alpha_1, \dots, \alpha_n \in \mathbb{R}, x_1, \dots, x_n \in X \right\}$$

und für $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$ und $g = \sum_{j=1}^n \beta_j k(x_j, \cdot)$ setze in H_{pre} ,

$$\langle f, g \rangle_{H_{\text{pre}}} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \beta_j k(x_i, x_j).$$

Weiterhin ist

$$\begin{aligned} \langle f, g \rangle_{H_{\text{pre}}} &= \sum_{j=1}^m \beta_j f(x'_j), \\ \langle f, g \rangle_{H_{\text{pre}}} &= \sum_{i=1}^m \alpha_i g(x_i), \end{aligned}$$

also ist $\langle \cdot, \cdot \rangle_{H_{\text{pre}}}$ von der speziellen Darstellung von f oder g unabhängig und daher wohldefiniert.

Offensichtlich ist $\langle \cdot, \cdot \rangle_{H_{\text{pre}}}$ bilinear und symmetrisch, da k symmetrisch.

Zu zeigen ist also noch, dass $\langle f, f \rangle = 0 \Rightarrow f = 0$. Wir zeigen dazu, dass $\langle \cdot, \cdot \rangle_{H_{\text{pre}}}$ die Cauchy-Schwarz'sche Ungleichung erfüllt.

1. Fall. Ist $\langle f, f \rangle = 0$ und $\langle g, g \rangle = 0$, dann ist

$$\begin{aligned} 0 &\leq \langle f + g, f + g \rangle = \langle f, f \rangle + 2 \langle f, g \rangle + \langle g, g \rangle = 2 \langle f, g \rangle, \\ 0 &\leq \langle f - g, f - g \rangle = \langle f, f \rangle - 2 \langle f, g \rangle + \langle g, g \rangle = -2 \langle f, g \rangle. \end{aligned}$$

Folglich ist $\langle f, g \rangle = 0$ und die Cauchy-Schwarz'sche Ungleichung folgt.

2. Fall. Ohne Einschränkung ist $\langle g, g \rangle > 0$, setzen wir also $\alpha = -\frac{\langle f, g \rangle}{\langle g, g \rangle}$. Dann folgt

$$0 \leq \langle f + \alpha g, f + \alpha g \rangle = \dots = \langle f, f \rangle - \frac{\langle f, g \rangle^2}{\langle g, g \rangle}.$$

Sei nun $f \in H_{\text{pre}}$ mit $\langle f, f \rangle = 0$ und $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$. Dann ist

$$\begin{aligned} |f(x)|^2 &= \left| \sum_{i=1}^n \alpha_i k(x_i, \cdot) \right|^2 = |\langle f, k(x, \cdot) \rangle|^2 \\ &\leq \langle f, f \rangle \langle k(x, \cdot), k(x, \cdot) \rangle = 0. \end{aligned}$$

Folglich ist $\langle \cdot, \cdot \rangle_{H_{\text{pre}}}$ ein Skalarprodukt und damit H_{pre} ein Prähilbertraum.

Zu H_{pre} existiert eine Vervollständigung H mit einer isometrischen Einbettung $J : H_{\text{pre}} \rightarrow H$. Setze nun $\Phi : X \rightarrow H, x \mapsto Jk(x, \cdot)$, dann gilt

$$\begin{aligned} \langle \Phi(x), \Phi(x') \rangle_H &= \langle Jk(x, \cdot), Jk(x', \cdot) \rangle_H = \langle k(x, \cdot), k(x', \cdot) \rangle_{H_{\text{pre}}} \\ &= k(x, x'). \quad \ll \end{aligned}$$

5-1.12 **Korollar** Sind $k_n, n \geq 1$ Kerne auf X und gibt es eine Funktion $k : X \times X \rightarrow \mathbb{R}$ mit $k_n(x, x') \rightarrow k(x, x')$ für alle $x, x' \in X$, so ist k ein Kern. \times

» Alle k_n sind symmetrisch und daher ist es die Grenzfunktion k auch. Weiterhin gilt

$$0 \leq \sum_{i,j=1}^d \alpha_i \alpha_j k_n(x_i, x_j) \rightarrow \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j)$$

und daher ist k positiv definit. \ll

5-2 RKHS (Reproduzierende Kern-Hilbert-Räume)

Bisher waren Featurespace und -map nicht eindeutig, es gibt aber eine kanonische Wahl, den RKHS.

5-2.1 **Definition** (a) Ein Hilbertraum H , der aus Funktionen $f : X \rightarrow \mathbb{R}$ besteht, heißt Hilbertfunktionenraum über X . (HFS)

(b) Ein HFS H über X heißt RKHS, falls die Dirac-Funktionale

$$\delta_x : H \rightarrow \mathbb{R}, \quad f \mapsto f(x)$$

für alle $x \in X$ stetig sind.

(c) Eine Funktion $k : X \times X \rightarrow \mathbb{R}$ heißt *reproduzierender Kern* eines HFS über X , falls für jedes $x \in X$ $k(x, \cdot) \in H$ und die *reproduzierende Eigenschaft* erfüllt ist,

$$f(x) = \langle f, k(x, \cdot) \rangle, \quad \text{für alle } f \in H, x \in X. \quad \times$$

Bemerkungen. A. $L_2([0, 1])$ ist kein RKHS (nicht mal ein HFS).

B. Norm-Konvergenz in einem RKHS H impliziert punktweise Konvergenz.

» Sei $x \in X$ und gelte $\|f_n - f\|_H \rightarrow 0$, dann gilt

$$f(x) = \delta_x(f) = \delta_x\left(\lim_{n \rightarrow \infty} f_n\right) = \lim_{n \rightarrow \infty} \delta_x(f_n) = \lim_{n \rightarrow \infty} f_n(x). \quad \leftarrow$$

C. (c) wurde für den im Beweis von Satz 5-1.11 konstruierten Hilbertraum schon „fast“ gezeigt. \rightarrow

5-2.2 **Lemma** Sei H ein HFS über X und $k : X \times X \rightarrow \mathbb{R}$ ein reproduzierender Kern von H . Dann gelten

1) H ist ein RKHS.

2) k ist ein Kern, H ist ein FS von k und

$$\Phi : X \rightarrow H, \quad x \mapsto k(x, \cdot)$$

die *kanonische Featuremap* von k . \times

» 1) Eine direkte Rechnung zeigt

$$|\delta_x(f)| = |f(x)| = |\langle f, k(x, \cdot) \rangle_H| \leq \|f\|_H \|k(x, \cdot)\|_H.$$

Somit ist δ_x stetig und $\|\delta_x\| \leq \|k(x, \cdot)\|$.

2) Für $x' \in X$ fest schreibe $f := k(x', \cdot)$. Dann ist $f \in H$ per definitionem und

$$\begin{aligned} \langle \Phi(x), \Phi(x') \rangle_H &= \langle k(x, \cdot), f \rangle_H = \langle f, k(x, \cdot) \rangle_H = f(x) = k(x', x) \\ &= \langle \Phi(x'), \Phi(x) \rangle. \quad \ll \end{aligned}$$

5-2.3 **Existenz und Eindeutigkeit des reproduzierenden Kerns** Sei H ein RKHS über X . Dann gilt

$$k : X \times X \rightarrow \mathbb{R}, \quad k(x, x') = \langle \delta_x, \delta_{x'} \rangle,$$

ist der einzige reproduzierende Kern von H . Ist ferner $(e_i)_{i \in \mathcal{I}}$ eine Orthonormalbasis (ONB), so gilt

$$k(x, x') = \sum_{i \in \mathcal{I}} e_i(x) e_i(x'), \quad x, x' \in X. \quad \times$$

» *Existenz.* Der Darstellungssatz von Fréchet-Riesz besagt, dass

$$I : H \rightarrow H', \quad f \mapsto \langle f, \cdot \rangle_H$$

ein isometrischer Isomorphismus ist, d.h. $H = H'$. Sei $J = I^{-1} : H' \rightarrow H$, so ist auch J ein isometrischer Isomorphismus und für $f \in H$ und $\omega \in H'$ gilt

$$\omega(f) = \langle J\omega, f \rangle_H. \quad (*)$$

Für x, x' folgt

$$k(x, x') = \langle \delta_x, \delta_{x'} \rangle_{H'} = \langle J\delta_x, J\delta_{x'} \rangle_H \stackrel{(*)}{=} \delta_x(J\delta_{x'}) = J\delta_{x'}(x).$$

Folglich ist $J\delta_{x'} = k(\cdot, x')$ und damit

$$f(x') - \delta_{x'}(f) = \langle f, J\delta_{x'} \rangle_H = \langle f, k(\cdot, x') \rangle_H = \langle f, k(x', \cdot) \rangle_H,$$

d.h. k ist reproduzierender Kern von H .

Eindeutigkeit. Sei \tilde{k} ein beliebiger reproduzierender Kern von H .

Für $x' \in X$ gilt $k(x', \cdot) \in H$. Für die ONB $(e_i)_{i \in I}$ folgt

$$k(x', \cdot) = \sum_{i \in I} \langle k(x', \cdot), e_i \rangle_H e_i = \sum_{i \in I} e_i(x') e_i,$$

wobei

$$k(x', x) = \sum_{i \in I} e_i(x') e_i(x)$$

punktweise existiert. Insbesondere ist \tilde{k} der einzige reproduzierende Kern von H und die Darstellung mithilfe der ONB $(e_i)_{i \in I}$ hängt nicht von der Wahl der ONB ab. «

5-2.4 **Satz** Sei k ein Kern über X mit Featurespace H_0 und Featuremap $\Phi_0 : X \rightarrow H_0$. Dann gelten

1) $H := \{f : X \rightarrow \mathbb{R} : \exists \omega \in H_0 : f(x) = \langle \omega, \Phi_0(x) \rangle_{H_0} \forall x \in X\}$ mit

$$\|f\|_H := \inf \{ \|\omega\|_{H_0} : \langle \omega, \Phi_0(\cdot) \rangle_{H_0} = f \}$$

ist der einzige RKHS von dem k ein reproduzierender Kern ist.

2) Der Operator $V : H_0 \rightarrow H$, $\omega \mapsto \langle \omega, \Phi_0(\cdot) \rangle_{H_0}$ ist eine metrische Surjektion, d.h. linear, beschränkt und $V \overset{\circ}{B}_{H_0} = \overset{\circ}{B}_H$.

3) $H_{\text{pre}} = \{ \sum_{i=1}^n \alpha_i k(x_i, \cdot) : n \geq 1, \alpha_i \in \mathbb{R}, x_i \in X \}$ ist dicht in H .

4) Für $f \in H_{\text{pre}}$ mit $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$ gilt

$$\|f\|_H^2 = \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j). \quad \times$$

Verleich. 5-1.11: Der FS ist eine Vervollständigung von H_{pre} .

5-2.4: Es existiert eine Vervollständigung, die aus Funktionen $X \rightarrow \mathbb{R}$ besteht und die RKHS ist. \sim

- » 1) Zeige H ist HFS über X . Offensichtlich besteht H aus Funktionen $X \rightarrow \mathbb{R}$ und $V : H_0 \rightarrow H$ ist linear. Für $f \in H$ gilt ferner

$$\|f\|_H = \inf \{ \|\omega\|_{H_0} : \omega \in V^{-1}(\{f\}) \}.$$

Zeige $\|\cdot\|_H$ ist Hilbertraumnorm auf H .

- a. $\ker V = \{ \omega \in H_0 : V\omega = 0 \} = V^{-1}(\{0\})$ ist abgeschlossen, da V stetig.
 b. Sei \tilde{H} das orthogonale Komplement von $\ker V$ in H_0 . Dann ist

$$V|_{\tilde{H}} : \tilde{H} \rightarrow H$$

injektiv nach Konstruktion und außerdem surjektiv, denn V ist surjektiv, d.h. zu $f \in H$ existiert ein $\omega \in H_0$ mit $V\omega = f$. Dann ist

$$\begin{aligned} \omega &= \tilde{\omega} + \omega_0, \quad \tilde{\omega} \in \tilde{H}, \quad \omega_0 \in H_0, \\ f &= V(\tilde{\omega} + \omega_0) = V(\tilde{\omega}) + \underbrace{V(\omega_0)}_{=0} = V|_{\tilde{H}}(\tilde{\omega}). \end{aligned}$$

- c. Zeige $\|f\|_H^2 = \|(V|_{\tilde{H}})f\|_H^2$.

$$\begin{aligned} \|f\|_H^2 &= \inf \{ \|\omega_0 + \tilde{\omega}\|_{H_0}^2 : \omega_0 \in \ker V, \tilde{\omega} \in \tilde{H} \} \\ &= \inf \{ \|\omega_0\|_{H_0}^2 + \|\tilde{\omega}\|_{H_0}^2 : \omega_0 \in \ker V, \tilde{\omega} \in \tilde{H} \} \\ &= \left\| V|_{\tilde{H}}^{-1}(f) \right\|_{\tilde{H}}^2. \end{aligned}$$

Damit ist H ein Hilbertraum, da \tilde{H} ein Hilbertraum ist, $V|_{\tilde{H}} : \tilde{H} \rightarrow H$ bijektiv linear und isometrisch.

Ferner folgt, dass V eine metrische Surjektion ist.

- 2) Zeige H ist ein RKHS und k ein reproduzierender Kern von H .

- a. Für $x \in X$ gilt

$$k(x, \cdot) = \langle \Phi_0(x), \Phi_0(\cdot) \rangle_{H_0} = V\Phi_0(x)$$

und folglich ist $k(x, \cdot) \in H$.

b. $\Phi_0(x) \in \tilde{H} = (\ker V)^\perp$, denn es gilt für $\omega \in \ker V$,

$$\langle \omega, \Phi_0(x) \rangle_{H_0} = V_\omega(x) = 0, \quad x \in X.$$

c. Zeige, dass $f(x) = \langle f, k(x, \cdot) \rangle_H$. Da $f = V(V|_{\tilde{H}})^{-1}f$ und $V|_{\tilde{H}}$ Isometrie, folgt

$$f(x) = \langle (V|_{\tilde{H}})^{-1}f, \Phi_0(x) \rangle_{H_0} = \langle f, V|_{\tilde{H}}\Phi_0(x) \rangle_H = \langle f, k(x, \cdot) \rangle_H.$$

Somit ist k ein reproduzierender Kern und nach Lemma 5-2.2 H ein RKHS.

3) Zeige 3) und 4) gilt für jeden RKHS \overline{H} für den k reproduzierender Kern ist. Wir zeigen dazu, dass

$$H_{\text{pre}} = \left\{ \sum_{i=1}^n \alpha_i k(x_i) : n \geq 1, \alpha_i \in \mathbb{R}, x_i \in X \right\}$$

dicht ist in \overline{H} und für $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$ gilt

$$\|f\|_{\overline{H}}^2 = \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j).$$

a. $k(x, \cdot) \in H_{\text{pre}}$ und folglich ist $H_{\text{pre}} \subset \overline{H}$ und $k(x, \cdot) \in \overline{H}$.

b. Zeige, dass $\overline{H_{\text{pre}}} = \overline{H}$. Angenommen $\overline{H_{\text{pre}}} \subsetneq \overline{H}$, dann ist $H_{\text{pre}}^\perp \neq (0)$, d.h. es existiert ein $f \in H_{\text{pre}}^\perp$ und ein $x \in X$ mit $f(x) \neq 0$ und $f \perp k(x, \cdot)$.

$$\Rightarrow 0 \neq f(x) = \langle f, k(x, \cdot) \rangle_{\overline{H}} = 0. \quad \text{↯}$$

D.h. $\overline{H_{\text{pre}}} = \overline{H}$, also ist H_{pre} dicht in \overline{H} .

Sei nun $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$, dann gilt

$$\begin{aligned} \|f\|_{\overline{H}}^2 &= \left\| \sum_{i=1}^n \alpha_i k(x_i, \cdot) \right\|_{\overline{H}}^2 = \left\langle \sum_{i=1}^n \alpha_i k(x_i, \cdot), \sum_{j=1}^n \alpha_j k(x_j, \cdot) \right\rangle \\ &= \sum_{i,j=1}^n \alpha_i \alpha_j \langle k(x_i, \cdot), k(x_j, \cdot) \rangle_{\overline{H}} = \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) \end{aligned}$$

aufgrund der repräsentierenden Eigenschaft angewandt auf $g = k(x_i, x \cdot)$.

- 4) Zeige schließlich H ist der einzige RKHS für den k ein reproduzierender Kern ist.

Seien dazu H_1 und H_2 RKHS für die k ein reproduzierender Kern ist. Nach dem bisher Gezeigten ist H_{pre} dicht in H_1 und H_2 und die Normen von H_1 und H_2 sind auf H_{pre} identisch. Wir zeigen nun, dass $H_1 \subset H_2$ und $\|f\|_{H_1} = \|f\|_{H_2}$ für alle $f \in H_1$. Dann folgt die Gleichheit aus der Symmetrie des Problems und H ist eindeutig.

Sei also $f \in H_1$, dann existiert eine Folge (f_n) in H_{pre} mit

$$\|f_n - f\|_{H_1} \rightarrow 0, \quad n \rightarrow \infty.$$

Folglich ist (f_n) Cauchyfolge in H_{pre} bezüglich $\|\cdot\|_{H_1}$ und damit auch bezüglich $\|\cdot\|_2$, also ist (f_n) Cauchyfolge in H_2 , d.h. es existiert ein $g \in H_2$, so dass

$$\|f_n - g\|_{H_2} \rightarrow 0, \quad n \rightarrow \infty.$$

Da die Konvergenz in $H_{1,2}$ der punktweisen Konvergenz entspricht gilt daher auch

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = g(x), \quad x \in X.$$

Folglich ist $f = g \in H_2$ und daher $H_1 \subset H_2$. Da außerdem

$$\|f_n - f\|_{H_2} \rightarrow 0, \quad n \rightarrow \infty$$

folgt $\|f_n\|_{H_2} \rightarrow \|f\|_{H_2}$, d.h.

$$\|f\|_{H_1} = \lim_{n \rightarrow \infty} \|f_n\|_{H_1} = \lim_{n \rightarrow \infty} \|f_n\|_{H_2} = \|f\|_{H_2}. \quad \ll$$

Interpretation. - Zu jedem Kern k existiert genau ein RKHS, so dass k ein reproduzierender Kern von H ist.

- Zu jedem RKHS H existiert genau ein reproduzierender Kern von H und dieser ist auch ein Kern.

Zwischen Kern und RKHS besteht also eine 1:1 Relation.

Der RKHS ist der „kleinste Featurespace“ im Sinne von Surjektionen. H besteht aus den Funktionen, die entstehen, wenn Daten in einen Featurespace H_0 abgebildet werden und dann ein linearer Ansatz gemacht wird. \rightarrow

5-3 Eigenschaften von Kernen und RKHS

Ziel dieses Abschnittes ist es, eine Beziehung zwischen den Eigenschaften von Kernen und den Eigenschaften der Funktionen eines RKHS herzustellen.

5-3.1 **Lemma** *Sei k ein Kern auf X und H der RKHS von k , dann sind äquivalent*

(i) k ist beschränkt.

(ii) Alle $f \in H$ sind beschränkt.

In diesem Fall gilt zudem

$$\|\text{id} : H \rightarrow \mathcal{L}_\infty(X)\| = \|k\|_\infty, \quad \|k\|_\infty^2 = \sup_{x \in X} k(x, x) < \infty. \quad \times$$

» „(i) \Rightarrow (ii)“: Für $f \in H$ und $x \in X$ gilt da $k(x, \cdot) \in H_{\text{pre}}$ nach Satz 5-2.4,

$$\begin{aligned} |f(x)| &= |\langle f, k(x, \cdot) \rangle_H| \leq \|f\|_H \|k(x, \cdot)\|_H = \|f\|_H \sqrt{k(x, x)} \\ &\leq \|f\|_H \underbrace{\|k\|_\infty}_{< \infty}. \end{aligned}$$

Somit ist $f \in \mathcal{L}_\infty(X)$ und $\|f\|_\infty \leq \|k\|_\infty \|f\|_H$. Also gilt insbesondere

$$\|\text{id} : H \rightarrow \mathcal{L}_\infty(X)\| \leq \|k\|_\infty.$$

„(ii) \Rightarrow (i)“: $\text{id} : H \rightarrow \mathcal{L}_\infty(X)$ ist wohldefiniert nach Voraussetzung. Wir zeigen, dass $\text{id} : H \rightarrow \mathcal{L}_\infty(X)$ stetig ist. Dazu benutzen wir den Satz vom abgeschlossenen Graphen. Sei dazu (f_n) eine Folge in H und $g \in \mathcal{L}_\infty(X)$ mit

$$\|f_n - f\|_H \rightarrow 0, \quad \|\text{id}(f_n) - g\|_{\mathcal{L}_\infty(X)} \rightarrow 0, \quad (*)$$

so ist zu zeigen, dass $f = g$.

Nach (*) gelten

$$f_n(x) \rightarrow f(x), \quad f_n(x) \rightarrow g(x), \quad \forall x \in X.$$

Somit ist $f(x) = g(x)$ für alle $x \in X$ und daher ist

$$\text{id} : H \rightarrow \mathcal{L}_\infty(X)$$

stetig. Weiterhin ist

$$\begin{aligned} k(x, x) &\leq \|k(x, \cdot)\|_\infty \leq \|\text{id} : H \rightarrow \mathcal{L}_\infty(x)\| \|k(x, \cdot)\|_H \\ &= \|\text{id} : H \rightarrow \mathcal{L}_\infty(x)\| \sqrt{k(x, x)} \end{aligned}$$

und daher

$$\begin{aligned} \sqrt{k(x, x)} &\leq \|\text{id} : H \rightarrow \mathcal{L}_\infty(x)\|, \\ \Rightarrow \|k\|_\infty &\leq \|\text{id} : H \rightarrow \mathcal{L}_\infty(x)\|. \end{aligned}$$

Insbesondere ist k auf der Diagonalen beschränkt. Schließlich ist

$$\begin{aligned} |k(x, x')| &= |\langle k(x, \cdot), k(x', \cdot) \rangle_H| \leq \|k(x, \cdot)\|_H \|k(x', \cdot)\|_H \\ &= \sqrt{k(x, x)} \sqrt{k(x', x')} \leq \|k\|_\infty^2. \quad \llcorner \end{aligned}$$

5-3.2 **Lemma** Sei k ein Kern auf X und H der RKHS von k , dann sind äquivalent

(i) $k(x, \cdot)$ ist messbar für alle $x \in X$.

(ii) Alle $f \in H$ sind messbar. \times

» „(ii) \Rightarrow (i)“: $k(x, \cdot) \in H$.

» „(i) \Rightarrow (ii)“: $k(x, \cdot)$ ist messbar und folglich besteht

$$H_{\text{pre}} = \{h(x, \cdot) : x \in X\}$$

aus messbaren Funktionen. Sei nun $f \in H$, dann existiert eine Folge (f_n) in H_{pre} mit $\|f_n - f\|_H \rightarrow 0$ und folglich ist $f(x) = \lim_{n \rightarrow \infty} f_n(x)$ für alle $x \in X$. \llcorner

5-3.3 **Lemma** Sei k Kern auf X mit RKHS H und $k(x, \cdot)$ messbar für alle $x \in X$ und H separabel. Dann ist k messbar und die kanonische Featuremap $\Phi : X \rightarrow H$ ist messbar. \times

» Wir überspringen den Beweis. \llcorner

5-3.4 **Lemma** Sei X ein Hausdorff-Raum, k ein Kern auf X mit RKHS H . Dann sind äquivalent

(i) k ist beschränkt und separat stetig, d.h. $k(x, \cdot)$ ist stetig für alle $x \in X$.

(ii) Alle $f \in H$ sind stetig und beschränkt. In diesem Fall gilt

$$\|\text{id} : H \rightarrow C_b(X)\| = \|k\|_\infty < \infty. \quad \times$$

» „(i) \Rightarrow (ii)“: $H_{\text{pre}} \subset C_b(X)$ nach Voraussetzung. Sei $f \in H$, dann existiert eine Folge (f_n) in H_{pre} mit $\|f_n - f\|_H \rightarrow 0$. Da k beschränkt folgt mit Lemma 5-3.1, dass

$$\|f_n - f\|_\infty \rightarrow 0$$

und folglich ist $f \in C_b(X)$.

„(ii) \Rightarrow (i)“: $k(x, \cdot) \in H$ ist stetig nach Voraussetzung und beschränkt nach Lemma 5-3.1. Somit ist k beschränkt und $\|k\|_\infty = \|\text{id} : H \rightarrow C_b(X)\|$. «

5-3.5 **Definition** Sei k ein Kern auf X mit RKHS H und kanonischer Featuremap $\Phi : X \rightarrow H$. Dann heißt $d_k : X \times X \rightarrow [0, \infty)$

$$d_k(x, x') := \|\Phi(x) - \Phi(x')\|_H = \sqrt{k(x, x) + k(x', x') - 2k(x, x')}$$

die *Kernmetrik* von k . \times

Man sieht leicht ein, dass d_k genau dann eine Metrik ist, wenn Φ injektiv ist. Im Allgemeinen ist d_k lediglich eine Semi-Metrik.

5-3.6 **Lemma** Sei (X, τ) ein Hausdorff-Raum, k ein Kern auf X und H ein RKHS von k , $\Phi : X \rightarrow H$ die kanonische Featuremap. Dann sind äquivalent

(i) k ist stetig.

(ii) k ist separat stetig und $x \mapsto k(x, x)$ ist stetig.

(iii) $\Phi : X \rightarrow H$ ist stetig.

(iv) $\text{id} : (X, \tau) \rightarrow (X, d_k)$ ist stetig, d.h. jede d_k -offene Kugel ist τ -offen. \times

» „(i) \Rightarrow (ii)“: Trivial.

„(ii) \Rightarrow (iv)“: $d_k(\cdot, x) : (X, \tau) \rightarrow \mathbb{R}$ ist stetig und daher ist die Kugel

$$B_\varepsilon(x', d_k) = \{x' : d_k(x, x') < \varepsilon\}$$

τ -offen. Diese Kugeln erzeugen aber gerade die Topologie von d_k , d.h. die Topologie von τ ist größer oder gleich der Topologie von d_k .

„(iv) \Rightarrow (iii)“: $\Phi : (X, d_k) \rightarrow H$ ist stetig nach Konstruktion von d_k und somit ist $\Phi : (X, \tau) \rightarrow H$ stetig nach Voraussetzung.

„(iii) \Rightarrow (i)“: Seien (x_n) und (y_n) Folgen in X mit $x, y \in X$, so dass $x_n \rightarrow x$ und $y_n \rightarrow y$, so gilt

$$\begin{aligned} |k(x_n, y_n) - k(x, y)| &= |\langle \Phi(y_n), \Phi(x_n) - \Phi(x) \rangle + \langle \Phi(x), \Phi(y_n) - \Phi(y) \rangle| \\ &\leq \|\Phi(y_n)\| \|\Phi(x_n) - \Phi(x)\| + \|\Phi(x)\| \|\Phi(y_n) - \Phi(y)\| \\ &\leq \|\Phi(y_n)\| d_k(x_n, x) + \|\Phi(x)\| d_k(y_n, y). \end{aligned}$$

Nach Voraussetzung ist Φ stetig, d.h. $\|\Phi(y_n)\|$ ist beschränkt. «

5-3.7 **Satz** Sei X ein kompakter Hausdorffraum und k ein stetiger Kern mit RKHS H . Dann ist die Abbildung

$$\text{id} : H \rightarrow C(X)$$

wohldefiniert und kompakt. \times

» Da X kompakt und k stetig, ist k auch beschränkt. Folglich ist $\text{id} : H \rightarrow C(X)$ wohldefiniert und stetig (Lemma 5-3.4).

Sei $C(X, d_k)$ der Raum der d_k -stetigen Funktionen $f : X \rightarrow \mathbb{R}$. Offensichtlich ist $C_b(X, d_k) \subset l_\infty(X)$. Außerdem ist X kompakt und Φ stetig, d.h. $\Phi(X) \subset H$ ist kompakt. Da Φ Isometrie bezüglich d_k , ist (X, d_k) kompakt und folglich

$$C(X, d_k) = C_b(X, d_k).$$

Für $f \in B_H$ und $x, y \in X$ gilt

$$|f(x) - f(y)| = |\langle f, \Phi(x) - \Phi(y) \rangle_H| \leq \|f\| \|\Phi(x) - \Phi(y)\| \leq d_k(x, y).$$

Also sind alle Funktionen in B_H Lipschitz bezüglich d_k mit Lipschitz-Konstante ≤ 1 . Der Satz von Arzelà-Ascoli besagt nun, dass

$$\overline{B_H} \text{ kompakt ist in } C(X, d_k).$$

Somit ist $\text{id}_{H \rightarrow C(X, d_k)}(B_H)$ relativkompakt und daher $\text{id} : H \rightarrow C(X, d_k)$ kompakt. Da $C(X, d_k) \subset C(X)$ mit Normgleichheit und

$$H \xrightarrow{\text{id}} C(X, d_k) \xrightarrow{\text{stetig}} C(X)$$

ist $\text{id} : H \rightarrow C(X)$ kompakt. «

Bemerkung. Ist X nicht kompakt aber k stetig und beschränkt, so ist

$$\text{id} : H \rightarrow C_b(X)$$

wohldefiniert und stetig aber im Allgemeinen *nicht* kompakt.

Betrachte z.B. $X = \mathbb{R}$ und den Gauß-Kern $k_\sigma(x, y) = \exp(-\sigma^2 \|x - y\|_2^2)$. Seien $m, n \in \mathbb{N}$ mit $n \neq m$, so gilt

$$\begin{aligned} \|k_\sigma(m, \cdot) - k_\sigma(n, \cdot)\|_\infty &\geq |k_\sigma(m, m) - k_\sigma(n, m)| \\ &= 1 - \exp(-\sigma^2 |m - n|^2) \geq 1 - e^{-\sigma^2}, \end{aligned}$$

da $|m - n| \geq 1$. Es ist aber

$$\|k_\sigma(n, \cdot)\|_H = \sqrt{k_\sigma(n, n)} = 1$$

und folglich $\text{id} : H \rightarrow C_b(X)$ nicht kompakt. \rightarrow

5-3.8 **Lemma** Sei X ein separabler, metrischer Raum und k ein stetiger Kern mit RKHS H . Dann ist H separabel. \times

» $\Phi : X \rightarrow H$ ist stetig, da k stetig und somit ist $\Phi(X)$ separabel. Folglich ist auch $H_{\text{pre}} = \text{span } \Phi(X)$ separabel, denn man kann span durch Linearkombinationen mit rationalen Koeffizienten approximieren, und damit ist auch $H = \overline{H_{\text{pre}}}$ separabel. «

Wir wollen nun untersuchen, wie sich die Differenzierbarkeit des Kerns auf die Featuremap und den RKHS übertragen. Dazu betrachten wir einen Kern $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ und die Abbildung

$$\tilde{k} : \mathbb{R}^{2d} \rightarrow \mathbb{R}, \quad (x, y) \mapsto \tilde{k}((x, y)) = k(x, y).$$

Falls die partiellen Ableitungen von i und $i + d$ von \tilde{k} existieren, setzen wir

$$\partial_i \partial_{i+d} k := \partial_i \partial_{i+d} \tilde{k}.$$

Analog lässt sich dies definieren, wenn $X \subset \mathbb{R}^d$ offen.

5-3.9 **Lemma** Sei $X \subset \mathbb{R}^d$ offen, k ein Kern mit Featurespace H und Featuremap $\Phi : X \rightarrow H$, so dass $\partial_i \partial_{i+d} k$ für $1 \leq i \leq d$ existiert und stetig ist. Dann existiert auch

$$\partial_i \Phi : X \rightarrow H$$

und ist stetig und es gilt

$$\langle \partial_i \Phi(x), \partial_i \Phi(y) \rangle_H = \partial_i \partial_{i+d} k(x, y) = \partial_{i+d} \partial_i k(x, y). \quad \times$$

» Wir beweisen nur den Fall $X = \mathbb{R}^d$. Sei e_i der i -te Einheitsvektor, dann definiere für festes $h > 0$,

$$\Delta_h \Phi(x) := \Phi(x + h e_i) - \Phi(x).$$

Wir zeigen nun, dass $(h_n^{-1} \Delta_{h_n} \Phi(x))_n$ für alle Folgen $h_n \rightarrow 0$ mit $h_n \neq 0$ konvergiert. Dann ist $\lim_{n \rightarrow \infty} h_n^{-1} \Delta_{h_n} \Phi(x) := \partial_i \Phi(x)$ unabhängig von der Folge (h_n) . Da H vollständig ist, genügt es zu zeigen, dass $(h_n^{-1} \Delta_{h_n} \Phi(x))$ eine Cauchyfolge ist. Es ist

$$\begin{aligned} & \left\| h_n^{-1} \Delta_{h_n} \Phi(x) - h_m^{-1} \Delta_{h_m} \Phi(x) \right\|_H^2 \\ &= \left\| h_n^{-1} \Delta_{h_n} \Phi(x) \right\|_H^2 + \left\| h_m^{-1} \Delta_{h_m} \Phi(x) \right\|_H^2 \\ &\quad - 2 \langle h_n^{-1} \Delta_{h_n} \Phi(x), h_m^{-1} \Delta_{h_m} \Phi(x) \rangle \\ &\leq 2 \left(\varepsilon + \partial_i \partial_{i+d} k(x, x) - \langle h_n^{-1} \Delta_{h_n} \Phi(x), h_m^{-1} \Delta_{h_m} \Phi(x) \rangle \right) \end{aligned} \quad (*)$$

Wir definieren für festes n und x ,

$$K_{x,n}(\mathcal{Y}) = k(x + h_n e_i, \mathcal{Y}) - k(x, \mathcal{Y}),$$

dann ist

$$\begin{aligned} \langle \Delta_{h_n} \Phi(x), \Delta_{h_m} \Phi(\mathcal{Y}) \rangle_H &= \langle \Phi(x + h_n e_i) - \Phi(x), \Phi(\mathcal{Y} + h_m e_i) - \Phi(\mathcal{Y}) \rangle_H \\ &= k(x + h_n e_i, \mathcal{Y} + h_m e_i) - k(x, \mathcal{Y} + h_m e_i) \\ &\quad - k(x + h_n e_i, \mathcal{Y}) + k(x, \mathcal{Y}) \\ &= K_{x,n}(\mathcal{Y} + h_m e_i) - K_{x,n}(\mathcal{Y}). \end{aligned}$$

Wenden wir den Mittelwertsatz auf $K_{x,n}$ an, erhalten wir ein $\xi_{m,n} \in [-|h_m|, |h_m|]$ mit

$$\begin{aligned} h_m^{-1} \langle \Delta_{h_n} \Phi(x), \Delta_{h_m} \Phi(\mathcal{Y}) \rangle &= h_m^{-1} (K_{x,n}(\mathcal{Y} + h_m e_i) - K_{x,n}(\mathcal{Y})) \\ &= \partial_i K_{x,n}(\mathcal{Y} + \xi_{m,n} e_i) \\ &= \partial_{i+d} (k(x + h_n e_i, \mathcal{Y} + \xi_{m,n} e_i) - k(x, \mathcal{Y} + \xi_{m,n} e_i)) \end{aligned}$$

Wir wenden nun den Mittelwertsatz auf die erste Variable an und erhalten ein $\eta_{n,m} \in [-|h_n|, |h_n|]$, so dass

$$\langle h_n^{-1} \Delta_{h_n} \Phi(x), h_m^{-1} \Delta_{h_m} \Phi(\mathcal{Y}) \rangle_H = \partial_i \partial_{i+d} k(x + \eta_{n,m} e_i, \mathcal{Y} + \xi_{m,n} e_i)$$

Somit ist nach (*),

$$\begin{aligned} &\|h_n^{-1} \Delta_{h_n} \Phi(x) - h_m^{-1} \Delta_{h_m} \Phi(x)\|_H^2 \\ &\leq 2 (\varepsilon + \partial_i \partial_{i+d} k(x, x) - \partial_i \partial_{i+d} k(x + \eta_{n,m} e_i, x + \xi_{m,n} e_i)) \\ &< 4\varepsilon, \quad n, m \geq n_0, \end{aligned}$$

denn $\partial_i \partial_{i+d} k$ ist stetig. Also ist $(h_n^{-1} \Delta_{h_n} \Phi(x))$ Cauchyfolge und daher existiert $\partial_i \Phi(x)$ und ist stetig.

Weiterhin gilt

$$\langle h_n^{-1} \Delta_{h_n} \Phi(x), h_m^{-1} \Delta_{h_m} \Phi(\mathcal{Y}) \rangle_H \rightarrow \partial_i \partial_{i+d} k(x, \mathcal{Y})$$

und da $h_n^{-1} \Delta_{h_n} \Phi(x) \rightarrow \partial_i \Phi(x)$, folgt die Formel. «

5-3.10 **Definition** Ein Kern k auf $X \subset \mathbb{R}^d$ offen heißt *m-fach stetig differenzierbar*, wenn

$$\partial^{\alpha, \alpha} k : X \times X \rightarrow \mathbb{R}$$

existiert und stetig ist für alle $\alpha \in \mathbb{N}_0^d$ mit $|\alpha| \leq m$. \times

5-3.11 **Korollar** Sei k ein *m-fach stetig differenzierbarer Kern* mit RKHS H . Dann sind alle $f \in H$ *m-fach stetig differenzierbar* und es gilt

$$|\partial^\alpha f(x)| \leq \|f\| \sqrt{\partial^{\alpha, \alpha} k(x, x)},$$

für alle $x \in X$ und $\alpha \in \mathbb{N}_0^d$ mit $|\alpha| \leq m$. \times

» Iteration von Lemma 5-3.9 impliziert, dass $\partial^\alpha \Phi : X \rightarrow H$ Featuremap von $\partial^{\alpha, \alpha} k$. Da $\langle f, \cdot \rangle$ stetig, folgt

$$\langle f, \partial^\alpha \Phi(x) \rangle_H = \partial^\alpha \langle f, \Phi(x) \rangle = \partial^\alpha f(x)$$

und nach Cauchy-Schwarz

$$|\partial^\alpha f(x)| = |\langle f, \partial^\alpha \Phi(x) \rangle| \leq \|f\|_H \|\partial^\alpha \Phi(x)\|_H = \|f\|_H \sqrt{\partial^{\alpha, \alpha} k(x, x)}. \quad \ll$$

5-4 Große RKHS

Wann kann ein RKHS „viele“ Funktionen approximieren?

5-4.1 **Definition** Sei (X, d) ein kompakter metrischer Raum, dann heißt ein stetiger Kern k auf X *universell*, wenn der RKHS H von k dicht in $C(X)$ ist, d.h.

$$\forall g \in C(X), \varepsilon > 0 \exists f \in H : \|f - g\|_\infty \leq \varepsilon. \quad \times$$

5-4.2 **Lemma** Sei (X, d) ein kompakter metrischer Raum und k ein universeller Kern auf X . Dann gelten

- 1) Jede Featuremap von k ist injektiv.
- 2) $k(x, x) > 0$ für alle $x \in X$.

3) Der Kern $k^* : X \times X \rightarrow \mathbb{R}$ gegeben durch

$$k^*(x, y) = \frac{k(x, y)}{\sqrt{k(x, x)k(y, y)}}$$

ist universell.

4) Ist $M \subset X$ abgeschlossen, dann ist $k|_{M \times M}$ universell. \times

» 1) Seien $x_1 \neq x_2 \in X$ und $g : X \rightarrow \mathbb{R}$ mit $x \in X$,

$$g(x) := \frac{d(x_1, x)}{d(x_1, x) + d(x_2, x)} - \frac{d(x_2, x)}{d(x_1, x) + d(x_2, x)},$$

$g(x_1) = -1$ und $g(x_2) = 1$ und g stetig auf X , also $g \in C(X)$. Sei nun H_0 ein Featurespace von k und $\Phi_0 : X \rightarrow H_0$ eine Featuremap und ferner H der RKHS von k . Da k universell, existiert ein $f \in H$ mit

$$\|f - g\|_\infty \leq \frac{1}{2} \Rightarrow f(x_1) \leq -\frac{1}{2} \text{ und } f(x_2) \geq \frac{1}{2}.$$

Nach Satz 5-2.4 existiert ein $\omega \in H_0$ mit $f = \langle \omega, \Phi_0(\cdot) \rangle$ und daher

$$\left. \begin{array}{l} \langle \omega, \Phi_0(x_1) \rangle_H \leq -\frac{1}{2}, \\ \langle \omega, \Phi_0(x_2) \rangle_H \geq \frac{1}{2}, \end{array} \right\} \Rightarrow \Phi_0(x_1) \neq \Phi_0(x_2).$$

2) Zeige $k(x, x) > 0$ für alle $x \in X$. In 1) haben wir gesehen, dass es für $x \in X$ ein $g \in C(X)$ mit $g(x) = 1$ gibt. Dazu gibt es ein $f \in H$ mit $\|f - g\|_\infty \leq \frac{1}{2}$ und folglich

$$\frac{1}{2} \leq f(x) = \langle f, \Phi(x) \rangle_H \Rightarrow \Phi(x) \neq 0.$$

Somit ist $k(x, x) = \langle \Phi(x), \Phi(x) \rangle_H > 0$.

3) Zeige, dass der normalisierte Kern universell ist. Schreibe dazu

$$\alpha(x) := (k(x, x))^{-1/2},$$

dann ist $\alpha\Phi : X \rightarrow H$ eine Featuremap von k^* . Sei nun $g \in C(X)$ und $\varepsilon > 0$. Sei weiterhin

$$c := \|\alpha\|_\infty < \infty.$$

Dann existiert ein $f \in H$ mit $\|f - \frac{g}{\alpha}\|_{\infty} \leq \frac{\varepsilon}{c}$ und damit

$$\| \langle f, \alpha \Phi(\cdot) \rangle_H - g \|_{\infty} \leq \underbrace{\|\alpha\|_{\infty}}_{\leq c} \left\| f - \frac{g}{\alpha} \right\|_{\infty} \leq \varepsilon.$$

4) Sei $M \subset X$ abgeschlossen, also kompakt, so ist $k|_{M \times M}$ universell. Tietzes Fortsetzungssatz besagt

$$\forall g \in C(M) \exists \hat{g} \in C(X) : \hat{g}|_M = g.$$

Zu jedem $g \in C(X)$ und $\varepsilon > 0$ existiert eine Fortsetzung $\hat{g} \in C(X)$ und ein $f \in H$ mit $\|\hat{g} - f\|_{\infty} \leq \varepsilon$. Dann ist auch

$$\|\hat{g}|_M - f|_M\|_{\infty} \leq \varepsilon,$$

wobei $f|_M$ im RKHS von $k|_{M \times M}$. «

5-4.3 **Satz** Sei (X, d) ein kompakter metrischer Raum, k stetiger Kern mit $k(x, x) > 0$ für alle $x \in X$. Ferner sei $\Phi : X \rightarrow l_2$ eine injektive Featuremap von k . Schreibe $\Phi_n : X \rightarrow \mathbb{R}$ für die n -te Komponente von Φ , d.h. $\Phi(x) = (\Phi_n(x))_{n \geq 1}$ für alle $x \in X$.

Schreibe $\mathcal{A} := \text{span} \{ \Phi_n : n \geq 1 \}$. Ist \mathcal{A} eine Algebra, dann ist k universell. ✕

Zum Beweis des Satzes benötigen wir den

Satz von Stone-Weierstraß Sei (X, d) ein kompakter metrischer Raum, $\mathcal{A} \subset C(X)$ Algebra mit

- 1) Zu jedem $x \in X$ existiert ein $f \in \mathcal{A}$ mit $f(x) \neq 0$. „ \mathcal{A} verschwindet nicht“
- 2) Zu $x \neq y \in X$ existiert ein $f \in \mathcal{A}$ mit $f(x) \neq f(y)$. „ \mathcal{A} separiert“.

Dann gilt $\mathcal{A} \subset C(X)$ ist dicht. ✕

- » 1) $\|(\Phi_n(x))_{n \geq 1}\|_{l_2}^2 = \langle \Phi(x), \Phi(x) \rangle_{l_2} = k(x, x) > 0$ und daher existiert ein $n \in \mathbb{N}$, so dass $\Phi_n(x) > 0$.
- 2) Sei nun $x \neq y$, dann ist $\Phi(x) \neq \Phi(y)$ nach Voraussetzung, d.h. es gibt ein $n \in \mathbb{N}$, so dass $\Phi_n(x) \neq \Phi_n(y)$.

- 3) Zu zeigen ist noch, dass $\mathcal{A} \subset C(X)$. Da k stetig, ist $\Phi : X \rightarrow l^2$ stetig und daher ist auch Φ_n stetig für alle $n \geq 1$, d.h. $\mathcal{A} \subset C(X)$. Damit zeigt Stone-Weierstraß, dass $\overline{\mathcal{A}}^{\|\cdot\|_\infty} = C(X)$. Für $g \in C(X)$ und $\varepsilon > 0$ existiert somit ein $f \in \mathcal{A}$ mit $\|f - g\|_\infty \leq \varepsilon$ und

$$f = \sum_{j=1}^m \alpha_j \Phi_{n_j}.$$

Setze

$$\omega_n := \begin{cases} \alpha_j, & \text{falls } n_j = n, \\ 0, & \text{sonst,} \end{cases}$$

und $\omega := (\omega_n)_{n \geq 1}$. Somit ist $\omega \in l^2$ und $f = \langle \omega, \Phi(\cdot) \rangle_{l^2}$ nach Konstruktion.
«

5-4.4 **Korollar** Für $r \in (0, \infty]$ sei $h : (-r, r) \rightarrow \mathbb{R}$ mit

$$h(t) = \sum_{n \geq 0} a_n t^n, \quad t \in (-r, r).$$

Falls $a_n > 0$ für alle $n \geq 0$, ist der Taylorkern universell auf allen kompakten Teilmengen $X \subset \sqrt{r} B_{\mathbb{R}^d}$. \times

» Wir haben schon gesehen, dass $\Phi : X \rightarrow l^2(\mathbb{N}_0^d)$,

$$\Phi(x) = \left(\underbrace{\sqrt{a_{j_1} \cdots a_{j_d}}}_{>0} \prod_{i=1}^d x_i^{j_i} \right)_{j_1, \dots, j_d \geq 0}$$

eine Featuremap von k ist.

- 1) k ist stetig, da h und $\langle \cdot, \cdot \rangle_{\mathbb{R}^d}$ stetig.
- 2) $a_0 > 0$ und daher ist $k(x, x) = \sum_{n \geq 0} a_n (\langle x, x \rangle)^n \geq a_0 > 0$.

3) Φ ist injektiv. Sei $x \neq y$ dann existiert ein $i \in \{1, \dots, d\}$ mit $x_i \neq y_i$. Nimm $j_i = 1$, alle anderen $j_l = 0$. Folglich ist $\Phi_{j_1, \dots, j_d}(x) \neq \Phi_{j_1, \dots, j_d}(y)$ und daher $\Phi(x) \neq \Phi(y)$.

Zu zeigen ist, dass $\mathcal{A} = \text{span} \{ \Phi_{j_1, \dots, j_d} : j_1, \dots, j_d \geq 0 \}$ ist Algebra. Da $a_n \geq 0$ ist \mathcal{A} Algebra, denn \mathcal{A} enthält alle Monome und Linearkombinationen davon. Also ist Satz 5-4.3 anwendbar. «

BSP Der *Exponentialkern*

$$h(x, x') := \exp(\langle x, x' \rangle_{\mathbb{R}^d})$$

ist universell auf allen Kompakta, da $a_n = \frac{1}{n!} > 0$. ■

BSP Für $\sigma > 0$ ist der *Gaußkern*

$$k_\sigma(x, x') = \exp(-\sigma^2 |x - x'|^2)$$

universell auf allen Kompakta. Da

$$h_\sigma(x, x') = \frac{\exp(2\sigma^2 \langle x, x' \rangle)}{\exp(\sigma^2 |x|^2) \exp(\sigma^2 |x'|^2)},$$

$\exp(2\sigma^2 \langle x, x' \rangle)$ universell und h_σ eine normalisierte Fassung davon. ■

Bemerkungen. A. Dahmen und Michelli zeigten 1987, dass ein Kern k genau dann universell ist, wenn

$$a_0 > 0 \text{ und } \sum_{a_{2n} > 0} \frac{1}{2n} = \sum_{a_{2n+1} > 0} \frac{1}{2n+1} = \infty.$$

B. Eine leichte Übung zeigt, dass jeder universelle Kern strikt positiv ist.

C. Pirkus zeigte 2004, dass Taylorkerne genau dann strikt positiv sind, wenn

$$a_0 > 0 \text{ und } |\{n : a_{2n} > 0\}| = |\{n : a_{2n+1} > 0\}| = \infty.$$

D. Man kann auch untersuchen, wann $H \subset L_p(\mu)$ dicht ist.

Für Gaußkerne ist $H_\sigma \subset L_p(\mu)$ dicht für alle $p \in [1, \infty)$ und alle W -Maße μ auf \mathbb{R}^d und $\mu = \lambda^d$ auf \mathbb{R}^d .

- E. Sei (X, d) kompakter metrischer Raum, dann existiert ein universeller Kern auf X .
- F. Sei (X, τ) kompakter Hausdorffraum, dann existiert genau dann ein universeller Kern auf X , wenn (X, τ) metrisierbar ist.
- G. In der Regel gilt für universelle RKHS H , dass $\dim H = \infty$, insbesondere gilt dies für $X \subset \mathbb{R}^d$ mit $\overset{\circ}{X} \neq \emptyset$. \rightarrow

6 Support Vector Machines (SVMs)

6-1 Definition, einfache Eigenschaften und Beispiele

6-1.1 **Definition** Sei H ein RKHS und $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ eine konvexe Verlustfunktion. Ferner sei $\lambda > 0$ „Regularisierungsparameter“. Dann heißt eine Lernmethode *Support Vector Machine (SVM)*, wenn

(a) Für jedes $D \in (X \times Y)^m$ die Entscheidungsfunktion $f_{D,\lambda} \in H$.

(b) $\lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,D}(f_{D,\lambda}) = \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f)$.

D.h. $f_{D,\lambda}$ minimiert das regularisierte empirische Risiko,

$$\lambda \|\cdot\|_{H^2} + \mathcal{R}_{L,D}(\cdot), \quad \text{über } H. \quad \times$$

6-1.2 **Eindeutigkeit** Für alle $D \in (X \times Y)^n$, $n \geq 1$, $\lambda > 0$ existiert höchstens ein $f_{D,\lambda} \in H$. \times

» Angenommen es gibt $f_1 \neq f_2 \in H$ und beide minimieren $\lambda \|\cdot\|_{H^2} + \mathcal{R}_{L,D}(\cdot)$ über H . Sei $f = \frac{1}{2}(f_1 + f_2)$. L konvex, dann ist $\mathcal{R}_{L,D}(\cdot)$ konvex

$$\mathcal{R}_{L,D}(f) \leq \frac{1}{2}\mathcal{R}_{L,D}(f_1) + \frac{1}{2}\mathcal{R}_{L,D}(f_2)$$

Parallelogrammgleichung in H :

$$\|f_1 + f_2\|_H^2 + \|f_1 - f_2\|_H^2 = 2\|f_1\|_H^2 + 2\|f_2\|_H^2.$$

Da $f_1 \neq f_2$ ist

$$\left\| \frac{1}{2}(f_1 + f_2) \right\|_H^2 < \frac{1}{2}\|f_1\|_H^2 + \frac{1}{2}\|f_2\|_H^2$$

und damit

$$\begin{aligned} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) &< \frac{1}{2} \left(\|f_1\|_H^2 + \mathcal{R}_{L,D}(f_1) \right) + \frac{1}{2} \left(\|f_2\|_H^2 + \mathcal{R}_{L,D}(f_2) \right) \\ &= \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f). \quad \neq \quad \ll \end{aligned}$$

6-1.3 **Satz** Ist $\dim H < \infty$, dann existiert ein $f_{D,\lambda} \in H$. \times

» L ist konvex, d.h. $L(x, y, \cdot) : \mathbb{R} \rightarrow [0, \infty)$ konvex und lokal lipschitz für alle x, y also ist L insbesondere stetig. Betrachte für $D \in (X \times Y)^n$ das Risikofunktional

$$\mathcal{R}_{L,D} : H \rightarrow [0, \infty), \quad f \mapsto \mathcal{R}_{L,D}(f).$$

Dieses ist konvex und stetig, da

$$f \mapsto \mathcal{R}_{L,D}(f) = \frac{1}{n} \sum_{i=1}^n L(x_i, y_i, f(x_i))$$

und die Konvergenz in H die punktweise Konvergenz impliziert. Ferner ist die Abbildung

$$H \rightarrow [0, \infty), \quad f \mapsto \lambda \|f\|_H^2$$

stetig und konvex (aufgrund der Parallelogrammungleichung sogar strikt konvex) und damit ist auch die Abbildung

$$H \rightarrow [0, \infty), \quad f \mapsto \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \tag{*}$$

stetig und konvex. Setze $m := \mathcal{R}_{L,D}(0) < \infty$ und

$$A = \left\{ f \in H : \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \leq m \right\}.$$

Dann gelten

- $0 \in A$.
- A ist abgeschlossen, da (*) stetig.

- A ist beschränkt, da für $f \in A$ gilt

$$\lambda \|f\|_H^2 \leq \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \leq m,$$

$$\text{ist } \|f\|_H \leq \sqrt{m/\lambda}.$$

Somit ist A nichtleer und kompakt, d.h. (*) besitzt ein Minimum $f^* \in A$ mit

$$\lambda \|f^*\|_H^2 + \mathcal{R}_{L,D}(f^*) \leq m$$

und für $f \notin A$ gilt

$$\lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) > m \geq \lambda \|f^*\|_H^2 + \mathcal{R}_{L,D}(f^*),$$

also ist f^* sogar ein globales Minimum. «

6-1.4 **Lemma** Falls $f_{D,\lambda} \in H$ existiert gelten

$$1) \|f_{D,\lambda}\|_H \leq (\lambda^{-1} \mathcal{R}_{L,D}(0))^{1/2}.$$

Inbesondere ist $L(x, y, 0) \leq 1$ für alle $(x, y) \in X \times Y$ und folglich

$$\|f_{D,\lambda}\|_H \leq \lambda^{-1/2}.$$

2) Falls ein $f^* \in H$ existiert mit $\mathcal{R}_{L,D}(f^*) \leq \mathcal{R}_{L,D}(0)$, ist $f_{D,\lambda} \neq 0$. ✕

» 1) Man rechnet direkt nach, dass

$$\begin{aligned} \lambda \|f_{D,\lambda}\|_H^2 &\leq \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,D}(f_{D,\lambda}) = \min_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \\ &\leq \lambda \|0\|_H^2 + \mathcal{R}_{L,D}(0) = \mathcal{R}_{L,D}(0). \end{aligned}$$

2) Für jedes $\alpha \in [0, 1]$ sei

$$h(\alpha) := 2\lambda\alpha \|f^*\|_H^2 + \mathcal{R}_{L,D}(f^*) - \mathcal{R}_{L,D}(0)$$

Da $\mathcal{R}_{L,D}(f^*) \leq \mathcal{R}_{L,D}(0)$ und $\alpha \mapsto h(\alpha)$ quadratisch, besitzt

$$h : [0, 1] \rightarrow [0, \infty)$$

ein globales Minimum $\alpha^* \in (0, 1]$, denn

$$\begin{aligned} 2\lambda\alpha \left\| f^* \right\|_H^2 + \mathcal{R}_{L,D}(f^*) - \mathcal{R}_{L,D}(0) &= 0 \\ \Leftrightarrow \alpha^* &= \frac{\mathcal{R}_{L,D}(0) - \mathcal{R}_{L,D}(f^*)}{2\lambda \left\| f^* \right\|_H^2}. \end{aligned}$$

Somit ist

$$\begin{aligned} \lambda \left\| \alpha^* f^* \right\|_H^2 + \mathcal{R}_{L,D}(\alpha^* f^*) &\leq h(\alpha^*) < h(0) = \mathcal{R}_{L,D}(0) \\ &= \lambda \left\| 0 \right\|_H^2 + \mathcal{R}_{L,D}(0). \end{aligned}$$

Also ist $f_{D,\lambda} \neq 0$. «

6-1.5 **Representer Theorem** Sei H ein RKHS mit Kern k , L eine konvexe Verlustfunktion und $\lambda > 0$, $D \in (X \times Y)^n$. Schreibe $D = ((x_1, y_1), \dots, (x_n, y_n))$, dann existieren ein $f_{D,\lambda}$ und $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ mit

$$f_{D,\lambda} = \sum_{i=1}^n \alpha_i k(x_i, \cdot). \quad \times$$

» Sei $X' := \{x_1, \dots, x_n\}$ und

$$H|_{X'} := \left\{ \sum_{i=1}^n \alpha_i k(x_i, \cdot) : \alpha_i \in \mathbb{R} \right\} = \text{span} \{k(x_i, \cdot) : i = 1, \dots, n\}.$$

Wir zeigen, dass $H|_{X'}$ der RKHS von $k|_{X' \times X'}$ ist.

Satz 5-2.4 sagte nämlich $\text{span} \{k(x_i, \cdot) : i = 1, \dots, n\}$ ist dicht im RKHS von $k|_{X' \times X'}$. Da $\dim \text{span} \{k(x_i, \cdot) : i = 1, \dots, n\} < \infty$, ist $\text{span} \{k(x_i, \cdot) : i = 1, \dots, n\}$ abgeschlossen im RKHS von $k|_{X' \times X'}$, d.h.

$$\text{span} \{k(x_i, \cdot)\} = \overline{\text{span} \{k(x_i, \cdot)\}} = \text{RKHS von } k|_{X' \times X'}.$$

Betrachte das SVM Problem bezüglich $k|_{X' \times X'}$,

$$\inf_{f \in H|_{X'}} \lambda \left\| f \right\|_H^2 + \mathcal{R}_{L,D}(f). \quad (*)$$

Nach Lemma 6-1.3 existiert ein $f_{D,\lambda,H}|_{X'} \in H|_{X'}$. Dieses ist Lösung von (*), denn $\dim H|_{X'} < \infty$. Da $H|_{X'} \subset H$ abgeschlossen, existiert das orthogonale Komplement $H|_{X'}^\perp$. Für $f \in H|_{X'}^\perp$ gilt dann

$$f(x, \cdot) = \langle f, k(x_i, \cdot) \rangle_H = 0.$$

Schreibe $P_{X'} : H \rightarrow H$ und $P_{X'}^\perp : H \rightarrow H$ für die orthogonale Projektion auf $H|_{X'}$ bzw. $H|_{X'}^\perp$. Für $f \in H$ gilt somit $f = P_{X'} f + P_{X'}^\perp f$ und nach (*) ist

$$P_{X'}^\perp f(x_i) = 0, \quad i = 1, \dots, n.$$

Somit ist auch

$$\mathcal{R}_{L,D}(f) = \mathcal{R}_{L,D}(P_{X'} f + P_{X'}^\perp f) = \mathcal{R}_{L,D}(P_{X'} f)$$

und da ferner $\|P_{X'} f\|_H^2 \leq \|f\|_H^2$, gilt

$$\begin{aligned} \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) &\leq \inf_{f \in H|_{X'}} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \\ &= \inf_{f \in H} \lambda \|P_{X'} f\|_H^2 + \mathcal{R}_{L,D}(P_{X'} f) \\ &\leq \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f). \end{aligned}$$

Insbesondere ist daher

$$\begin{aligned} \inf_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) &= \inf_{f \in H|_{X'}} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \\ &= \lambda \left\| f_{D,\lambda,H}|_{X'} \right\|_{H|_{X'}}^2 + \mathcal{R}_{L,D}(f_{D,\lambda,H}|_{X'}). \end{aligned}$$

Da $f_{D,\lambda,H}|_{X'} \in H$ löst $f_{D,\lambda,H}|_{X'}$ auch das original SVM-Problem und da $f_{D,\lambda}$ eindeutig und $f_{D,\lambda} \in H|_{X'}$ folgt die Darstellung. ◀

Als Konsequenz ist für $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$,

$$\begin{aligned} \min_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) &= \min_{\alpha_1, \dots, \alpha_n \in \mathbb{R}} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) \\ &= \min_{\alpha_1, \dots, \alpha_n \in \mathbb{R}} \lambda \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) + \mathcal{R}_{L,D} \left(\sum_{i=1}^n \alpha_i k(x_i, \cdot) \right) \end{aligned}$$

Dies ist ein endlichdimensionales und *konvexes* Optimierungsproblem.

Bsp Verlustfunktion der kleinsten Quadrate $L(x, y, t) = (y - t)^2$. Es ist

$$\begin{aligned} & \frac{\delta}{\delta \alpha_{i_0}} \left(\lambda \sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) + \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j) \right)^2 \right) \\ &= 2\lambda \sum_{i=1}^n \alpha_i k(x_i, x_{i_0}) + \frac{\delta}{\delta \alpha_{i_0}} \left(\frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j) \right)^2 \right), \quad (*) \end{aligned}$$

wobei

$$\begin{aligned} & \frac{\delta}{\delta \alpha_{i_0}} \left(\left(y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j) \right)^2 \right) \\ &= -2 \left(y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j) \right) \underbrace{\frac{\delta}{\delta \alpha_{i_0}} \left(\sum_{j=1}^n \alpha_j k(x_i, x_j) \right)}_{k(x_i, x_{i_0})}. \end{aligned}$$

Somit gilt

$$\begin{aligned} (*) &= 2\lambda \sum_{i=1}^n \alpha_i k(x_i, x_{i_0}) - \frac{2}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j) \right) k(x_i, x_{i_0}) \\ &\stackrel{!}{=} 0. \end{aligned}$$

Schreibe k als $n \times n$ -Matrix $k = (k(x_i, x_j))_{i,j=1}^n$. Dann ist $k = k^\top$ und die Optimierung ist gleichbedeutend mit

$$\begin{aligned} & 2\lambda k\alpha - \frac{2}{n} k\mathbf{y} + \frac{2}{n} k k \alpha \stackrel{!}{=} 0 \Leftrightarrow \lambda n k \alpha + k k \alpha = k\mathbf{y} \\ & \Leftrightarrow k(\lambda n E_n \alpha + k \alpha) = k\mathbf{y} \quad (**) \end{aligned}$$

wobei E_n die n -te Einheitsmatrix bezeichne. Falls $\lambda n E_n \alpha + k \alpha = \mathbf{y}$, dann ist $(**)$ erfüllt.

Da k ein Kern ist die Matrix k positiv semi-definit, d.h. k hat n nichtnegative Eigenwerte,

$$\lambda n \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix} + \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}.$$

$\lambda n E_n + k$ hat folglich strikt positive Eigenwerte und ist daher invertierbar.

$$\alpha = (\lambda n E_n + k)^{-1} y$$

ist die eindeutige Lösung und

$$f_{D,\lambda} = \sum_{i=1}^n \alpha_i k(x_i, \cdot). \quad \blacksquare$$

6-1-1 Ausflug in das Reich der konvexen Analysis

Sei E ein endlichdimensionaler Vektorraum, $P : E \rightarrow \mathbb{R}$ konvex und stetig differenzierbar sowie $h_i : E \rightarrow \mathbb{R}$ affin linear für $i = 1, \dots, m$.

Primales Optimierungsproblem *Wir suchen*

$$P^* = \inf_{\omega \in E} P(\omega)$$

unter den Nebenbedingungen $h_i(\omega) \leq 0$ für alle $i = 1, \dots, m$. \times

Wir nehmen für alles Weitere an, dass ein ω^* existiert mit $h_i(\omega^*) \leq 0$ für $i = 1, \dots, m$ und $P(\omega^*) = P^*$. Unser Ziel ist es das primale Optimierungsproblem auf ein duales Optimierungsproblem zurückzuführen.

Definition *Die Funktion*

$$L(\omega, \beta) := P(\omega) + \sum_{i=1}^m \beta_i h_i(\omega), \quad \omega \in E, \beta \in \mathbb{R}^m$$

heißt *Lagrangefunktion*.

$$D(\beta) := \inf_{\omega \in E} L(\omega, \beta)$$

heißt die zur Lagrangefunktion gehörige *duale Funktion*. \times

Da P konvex und die h_i affin linear, ist auch L konvex.

Duales Optimierungsproblem

$$D^* = \sup_{\beta \geq 0} D(\beta). \quad \times$$

Wie hängen das duale Problem mit seinen Lösungen mit dem primalen Problem zusammen?

Lemma Für $\beta \geq 0$ und $\omega \in E$ mit $h_i(\omega) \leq 0$ für $i = 1, \dots, m$ gilt

$$D(\beta) \leq P(\omega)$$

und damit insbesondere $D^* \leq P^*$. \times

» Man rechnet direkt nach, dass

$$D(\beta) = \inf_{\omega' \in E} L(\omega', \beta) \leq L(\omega, \beta) = P(\omega) + \sum_{i=1}^m \beta_i h_i(\omega) \leq P(\omega). \quad \ll$$

Satz $D^* = P^*$. \times

» Wir überspringen den Beweis, da dieser zu weit in die konvexe Analysis hineinführt. \ll

Korollar Sei $\beta^* \geq 0$ mit $D(\beta^*) = D^*$, dann gilt

$$D^* = \max_{\beta \geq 0} \inf_{\omega} L(\omega, \beta) = \min_{\omega} \sup_{\beta \geq 0} L(\omega, \beta) = P^*. \quad \times$$

» „ \leq “: gilt immer.

„ \geq “: Sei $\omega^* \in E$ eine primale Lösung, d.h.

$$h_i(\omega^*) \leq 0, \quad i = 1, \dots, m, \quad P(\omega^*) = P^*.$$

Dann ist

$$D^* = D(\beta^*) = \max_{\beta \geq 0} \inf_{\omega \in E} L(\omega, \beta),$$

$$P^* = P(\omega^*) = \min_{\omega, h_i(\omega) \leq 0} P(\omega) = \min_{\omega \in E} \sup_{\beta \geq 0} P(\omega) + \sum_{i=1}^m \beta_i h_i(\omega) \quad (*)$$

denn

$$\sup_{\beta \geq 0} P(\omega) + \sum_{i=1}^m \beta_i h_i(\omega) = \begin{cases} P(\omega), & h_i(\omega) \leq 0, \quad i = 1, \dots, m, \\ \infty, & \text{sonst.} \end{cases}$$

und folglich ist

$$(*) = \min_{\omega \in E} \sup_{\beta \geq 0} L(\omega, \beta).$$

Die Behauptung folgt, dass $P^* = D^*$. «

Satz Sei β^* die duale Lösung, d.h. $\beta^* \geq 0$ und $D(\beta^*) = D^*$ und ω^* die primale Lösung, d.h. $h_i(\omega^*) \leq 0$ für $i = 1, \dots, m$ mit $P(\omega^*) = P^*$. Dann folgt

$$\max_{\beta \geq 0} L(\omega^*, \beta) = L(\omega^*, \beta^*) = \min_{\omega \in E} L(\omega, \beta^*)$$

und $L(\omega^*, \beta^*) = D^* = P^*$. ✕

(ω^*, β^*) ist Sattelpunkt von L .

» Es gilt

$$D^* = \max_{\beta \geq 0} \inf_{\omega \in E} L(\omega, \beta) = \min_{\omega \in E} \sup_{\beta \geq 0} L(\omega, \beta) = P(\omega^*) = \sup_{\beta \geq 0} L(\omega^*, \beta).$$

und

$$P^* = \min_{\omega \in E} \sup_{\beta \geq 0} L(\omega, \beta) = \max_{\beta \geq 0} \inf_{\omega \in E} L(\omega, \beta) = D^* = D(\beta^*) = \inf_{\omega \in E} L(\omega, \beta^*)$$

sowie

$$L(\omega^*, \beta^*) \geq \inf_{\omega \in E} L(\omega, \beta^*) = P^*,$$

$$L(\omega^*, \beta^*) \leq \sup_{\beta \geq 0} L(\omega^*, \beta) = D^*,$$

und folglich ist

$$L(\omega^*, \beta^*) = \sup_{\beta \geq 0} L(\omega^*, \beta) = \inf_{\omega \in E} L(\omega, \beta^*) = P^* = D^*$$

und sup ist max und inf ist min. «

Korollar Seien β^* und ω^* wie oben. Dann ist

$$\beta_i^* h_i(\omega^*) = 0 \text{ für alle } i = 1, \dots, m. \quad \times$$

$$\gg P(\omega^*) = P^* = L(\omega^*, \beta^*) = P(\omega^*) + \sum_{i=1}^m \beta_i h_i(\omega^*). \quad \ll$$

Korollar Sei $E = E_1 \times E_2$, $\omega = (\omega_1, \omega_2) \in E_1 \times E_2$ und L derart, dass:

Für alle $\beta^* \geq 0$ mit $D(\beta^*) = D^*$ existiert genau ein $\omega_1 \in E_1$ und ein (oder mehrere) $\omega_2 \in E_2$ mit $h_i(\omega_1, \omega_2) \leq 0$ für $i = 1, \dots, m$ und

$$L((\omega_1, \omega_2), \beta^*) = D(\beta^*) = \inf_{\omega \in E} L(\omega, \beta^*).$$

Dann gilt für jede primale Lösung $\omega^* = (\omega_1^*, \omega_2^*)$,

$$\omega_1 = \omega_1^*.$$

Insbesondere: Falls wir wissen, dass es eine primale Lösung ω^* gibt, so finden wir die erste Komponente ω_1^* indem wir

$$\inf_{\omega \in E} L(\omega, \beta^*)$$

lösen und für die Lösung $\omega = (\omega_1, \omega_2)$ gilt $\omega_1^* = \omega_1$. \times

\gg Sei ω^* eine Lösung, dann ist

$$h_i(\omega^*) = h_i(\omega_1^*, \omega_2^*) \leq 0, \quad i = 1, \dots, m$$

und der vorige Satz zeigte $L(\omega^*, \beta^*) = D^* = D(\beta^*)$. Mit der vorausgesetzten Eindeutigkeit in der ersten Komponente folgt $\omega_1^* = \omega_1$. \ll

Als Anwendung betrachten wir SVMs mit Hinge-loss.

Bsp SVMs mit Hinge-loss. Das Optimierungsproblem

$$\min_{f \in H} \lambda \|f\|_H^2 + \frac{1}{n} \sum_{i=1}^n \max\{0, 1 - y_i f(x)\}$$

ist ein SVM Optimierungsproblem, da $L_{\text{hinge}}(y, t) = \max\{0, 1 - yt\}$. Dieses ist äquivalent zu dem Optimierungsproblem

$$\min_{\substack{f \in H \\ \xi \in \mathbb{R}^n}} \frac{1}{2} \langle f, f \rangle + c \sum_{i=1}^n \xi_i, \quad c = \frac{1}{2\lambda n}$$

mit den Nebenbedingungen

$$\xi_i \geq 0, \quad \xi_i \geq 1 - y_i f(x_i) \Leftrightarrow \xi_i \geq \max\{0, 1 - y_i f(x_i)\}.$$

Die Lagrangefunktion des Problems ist

$$L(f, \xi, \beta, \gamma) = \frac{1}{2} \langle f, f \rangle + c \sum_{i=1}^n \xi_i - \sum_{i=1}^n \beta_i \xi_i + \sum_{i=1}^n \gamma_i (1 - y_i f(x_i) - \xi_i).$$

L ist differenzierbar und

$$\frac{\delta}{\delta f} L(f, \xi, \beta, \gamma) = f - \sum_{i=1}^n \gamma_i y_i k(x_i, \cdot) \stackrel{!}{=} 0$$

führt auf die eindeutige Lösung

$$f = \sum_{i=1}^n \gamma_i y_i k(x_i, \cdot).$$

Weiterhin führt

$$\frac{\delta}{\delta \xi_i} L(f, \xi, \beta, \gamma) = c - \beta_{i_0} - \gamma_{i_0} \stackrel{!}{=} 0$$

auf $\beta_{i_0} + \gamma_{i_0} = c$.

Ist $\beta + \gamma = c$, so minimiert jedes ξ . Ist andererseits $\beta + \gamma \neq c$, so gibt es keine Lösung.

Setzen wir dies in unser Optimierungsproblem ein, so erhalten wir

$$\begin{aligned} D(\beta, \gamma) &= \frac{1}{2} \sum_{i=1}^n \gamma_i \gamma_j y_i y_j k(x_i, x_j) + \sum_{i=1}^n \gamma_i - \sum_{i=1}^n \gamma_i \gamma_j y_i y_j k(x_i, x_j) \\ &= \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i=1}^n \gamma_i \gamma_j y_i y_j k(x_i, x_j) \end{aligned}$$

mit $\beta \geq 0$, $\gamma \geq 0$ und $\beta + \gamma = c$, d.h. $\gamma \in [0, c]^n$ und $\beta = c - \gamma$.

Damit ist das duale Problem

$$\max_{\gamma \in [0, c]^n} \sum_{i=1}^n \gamma_i - \frac{1}{2} \sum_{i,j=1}^n \gamma_i \gamma_j \gamma_i \gamma_j k(x_i, x_j). \quad (*)$$

- mit der Definition $\mathbf{k} = (\gamma_i \gamma_j k(x_i, x_j))_{i,j=1}^n$ -

$$\max_{\gamma \in [0, c]^n} \langle \gamma, \mathbf{1} \rangle - \frac{1}{2} \langle \gamma, \mathbf{k} \gamma \rangle.$$

Ist γ^* Lösung des dualen Problems, so ist

$$f_{D,x} = \sum_{i=1}^n \gamma_i \gamma_i^* k(x_i, \cdot)$$

die Lösung der ersten Komponente des primalen Problems.

$$\xi_i^* = \max \{0, 1 - \gamma_i f_{D,x}(x_i)\}$$

ist die zweite Komponente.

Bemerkung. Sei $\gamma \in [0, c]^n$ und $f_\gamma = \sum_{i=1}^n \gamma_i \gamma_i k(x_i, \cdot)$. Gilt nun

$$P(f_\gamma, \xi_\gamma) - D(\gamma) \leq \varepsilon,$$

dann ist auch

$$P(f_\gamma) - P^* \leq \varepsilon,$$

da $D(\gamma) \leq P^*$.

$$P(f_\gamma, \xi_\gamma) = \sum_{i,j=1}^n \gamma_i \gamma_j \gamma_i \gamma_j k(x_i, x_j) + c \sum_{i=1}^n \max \{0, 1 - \gamma_i f_\gamma(x_i)\} - \sum_{i=1}^n \gamma_i. \quad \rightarrow$$

Dann ist

$$\lambda \|f_\gamma\|_H^2 + \mathcal{R}_{L,D}(f_\gamma) - \min_{f \in H} (\lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f)) = 2\lambda (P(f_\gamma, \xi_\gamma) - P^*).$$

Gilt nun

$$P(f_Y, \xi_Y) - P^* \leq \frac{\varepsilon}{2\lambda} \quad (+)$$

so ist

$$\lambda \|f_Y\|_H^2 + \mathcal{R}_{L,D}(f_Y) \leq \inf_{f \in H} (\lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f)) + \varepsilon. \quad \blacksquare$$

Wie löst man das duale Problem?

- Standardsoftware.
- **Gradient ascent.** Zu y berechne den Gradienten von (*) und mache einen Step in Richtung des steilsten Anstiegs. Wiederhole den Vorgang, bis (+) erreicht ist.

Ein Problem bei diesem Verfahren ist, dass die Kernmatrix ein enormer Platzfresser ist. Es gibt jedoch zahlreiche Verfahren, die nicht die gesamte Kernmatrix in den RAM laden muss.

- Man betrachtet die Abbildung

$$y_i \mapsto D(y + y_i e_i).$$

Diese ist eine eindimensionale konkave und quadratische Funktion und ihr Optimierungsproblem ist explizit lösbar.

- (a) Suche Richtung i^* mit maximalem Gewinn in D bei Optimierung in Richtung i^* .
- (b) Optimierte in Richtung i^* .
- (c) Gehe zu (a), falls (*) nicht erfüllt ist.

Die meisten verfügbaren Programme machen dies so.

6-2 Orakelungleichungen für SVMs

6-2.1 **Satz** Sei (X, d) ein kompakter, metrischer Raum, L eine konvexe und lokal lipschitz-stetige Verlustfunktion mit $L(x, y, 0) \leq 1$ für alle $(x, y) \in X \times Y$. Ferner sei H

ein RKHS über X mit stetigem Kern k und $\|k\|_\infty \leq 1$ und P sei ein Wahrscheinlichkeitsmaß auf $X \times Y$.

Für alle $\lambda > 0$, $n \geq 1$, $\varepsilon > 0$ und $\tau > 0$ gilt mit einer Wahrscheinlichkeit P^n nicht kleiner als $1 - \varepsilon^{-\tau}$:

$$\begin{aligned} \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}_{L,P}^* &< \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \right) \\ &+ 4\varepsilon |L|_{2\lambda^{-1/2},1} + \left(2 |L|_{2\lambda^{-1/2},1} \lambda^{-1/2} + 1 \right) \\ &\cdot \sqrt{\frac{2\tau + 2 \log 2\mathcal{N}(B_H, \|\cdot\|_\infty, \frac{1}{2}\lambda^{1/2}\varepsilon)}{n}}. \quad \times \end{aligned}$$

» Satz 5-3.7 zeigte, dass die Abbildung

$$\text{id} : H \rightarrow C(X)$$

kompakt ist. Daraus folgt

$$\log \mathcal{N}(B_H, \|\cdot\|_\infty, \varepsilon) < \varepsilon, \quad \varepsilon > 0.$$

Lemma 6-1.4 und 5-3.1 zeigten

$$\begin{aligned} \|f_{D,\lambda}\|_H &\leq \lambda^{-1/2}, \quad \lambda > 0, \\ \|f_{D,\lambda}\|_\infty &\leq \lambda^{-1/2}, \quad \text{da } \|\text{id} : H \rightarrow C(X)\| = \|k\|_\infty \leq 1. \end{aligned}$$

Sei nun $\delta > 0$, dann existiert ein $f_\delta \in H$ mit

$$\lambda \|f_\delta\|_H^2 + \mathcal{R}_{L,P}(f_\delta) \leq \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) + \delta,$$

und damit

$$\begin{aligned} \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) \\ \leq \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \lambda \|f_\delta\|_H^2 - \mathcal{R}_{L,P}(f_\delta) + \delta \\ - \mathcal{R}_{L,D}(f_{D,\lambda}) + \mathcal{R}_{L,D}(f_{D,\lambda}) \end{aligned} \quad (*)$$

Weiterhin ist

$$\lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,D}(f_D) \leq \lambda \|f_\delta\|_H^2 + \mathcal{R}_{L,D}(f_\delta),$$

da $f_{D,\lambda}$ diese Funktion minimiert. Es folgt

$$\begin{aligned}
 (*) &\leq \lambda \|f_\delta\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \lambda \|f_\delta\|_H^2 - \mathcal{R}_{L,P}(f_\delta) + \delta \\
 &\quad - \mathcal{R}_{L,D}(f_{D,\lambda}) + \mathcal{R}_{L,D}(f_\delta) \\
 &= \mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}_{L,D}(f_{D,\lambda}) - (\mathcal{R}_{L,P}(f_\delta) - \mathcal{R}_{L,D}(f_\delta)) + \delta \\
 &\leq |\mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}_{L,D}(f_{D,\lambda})| + |(\mathcal{R}_{L,P}(f_\delta) - \mathcal{R}_{L,D}(f_\delta))| + \delta. \quad (**)
 \end{aligned}$$

Dann gilt ferner

$$\begin{aligned}
 \lambda \|f_\delta\|_H^2 &\leq \lambda \|f_\delta\|_H^2 + \mathcal{R}_{L,P}(f_\delta) \leq \inf_{f \in H} (\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f)) + \varepsilon \\
 &\leq \mathcal{R}_{L,P}(0) + \delta \leq 1 + \delta \leq 4,
 \end{aligned}$$

falls $\delta \leq 3$. Damit ist

$$\begin{aligned}
 f_\delta &\in 2\lambda^{-1/2}B_H, \\
 f_{D,\lambda} &\in \lambda^{-1/2}B_H \subset 2\lambda^{-1/2}B_H.
 \end{aligned}$$

Folglich ist

$$(***) \leq 2 \sup_{f \in 2\lambda^{-1/2}B_H} |\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f)| + \delta$$

für jedes $\delta > 0$ also gilt auch

$$(***) \leq 2 \sup_{f \in 2\lambda^{-1/2}B_H} |\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,D}(f)|.$$

Definiere nun

$$B := 2\lambda^{-1/2} |L|_{2\lambda^{-1/2}} + 1.$$

Für $f \in 2\lambda^{-1/2}B_H$ gilt dann

$$\begin{aligned}
 L(x, \gamma, f(x)) &\leq |L(x, \gamma, f(x)) - L(x, \gamma, 0)| + \underbrace{L(x, \gamma, 0)}_{\leq 1} \\
 &\leq |L|_{2\lambda^{-1/2}, 1} |f(x) - 0| + 1 \leq |L|_{2\lambda^{-1/2}, 1} 2\lambda^{-1/2} + 1 \leq B.
 \end{aligned}$$

Sein nun \mathcal{F}_ε ein ε -Netz von $2\lambda^{-1/2}B_H$ und

$$|\mathcal{F}_\varepsilon| = \mathcal{N}(2\lambda^{-1/2}B_H, \|\cdot\|_\infty, \varepsilon) = \mathcal{N}(B_H, \|\cdot\|_\infty, \frac{1}{2}\lambda^{1/2}\varepsilon).$$

Für $f \in 2\lambda^{-1/2}B_H$ gibt es $g \in \mathcal{F}_\varepsilon$ mit $\|f - g\|_\infty < \varepsilon$ und

$$\begin{aligned} |\mathcal{R}_{L,D}(f) - \mathcal{R}_{L,P}(f)| &\leq |\mathcal{R}_{L,D}(f) - \mathcal{R}_{L,D}(g)| + |\mathcal{R}_{L,D}(g) - \mathcal{R}_{L,P}(g)| \\ &\quad + |\mathcal{R}_{L,P}(g) - \mathcal{R}_{L,P}(f)| \\ &\leq 2\varepsilon |L|_{2\lambda^{-1/2},1} + |\mathcal{R}_{L,D}(g) - \mathcal{R}_{L,P}(g)|. \end{aligned}$$

Somit ist

$$\begin{aligned} \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \inf_{f \in H} (\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f)) \\ \leq 4\varepsilon |L|_{2\lambda^{-1/2},1} + 2 \sup_{g \in \mathcal{F}_\varepsilon} |\mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g)|. \end{aligned}$$

Wir schätzen weiterhin ab,

$$\begin{aligned} P^n \left(\left\{ D : \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \inf_{f \in H} (\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f)) \right. \right. \\ \left. \left. \geq B \sqrt{\frac{2\tau}{m}} + 4\varepsilon |L|_{2\lambda^{-1/2}} \right\} \right) \\ \leq P^n \left(\left\{ D : 2 \sup_{g \in \mathcal{F}_\varepsilon} |\mathcal{R}_{L,P}(g) - \mathcal{R}_{L,D}(g)| \geq B \sqrt{\frac{2\tau}{m}} \right\} \right) \\ \stackrel{!}{\leq} 2 |\mathcal{F}_\varepsilon| e^{-\tau} = \exp(-\tau + \log 2\mathcal{N}(B_H, \|\cdot\|_\infty, \frac{1}{2}\lambda^{1/2}\varepsilon)), \end{aligned}$$

wobei der union bound und die Hoeffdings Ungleichung verwendet wurden. Mit einer Variablentransformation für τ folgt schließlich die Behauptung. «

Bemerkung. Mit etwas mehr Theorie über Hilberträume kann man zeigen, dass das inf tatsächlich angenommen wird. →

6-2.2 **Korollar** *Es gelten die Voraussetzungen des Satzes 6-2.1 und außerdem sei L lipschitz stetig mit $|L|_1 \leq 1$. Ferner existieren Konstanten $a \geq 1$, $p > 0$ mit*

$$\log \mathcal{N}(B_H, \|\cdot\|_\infty, \varepsilon) \leq a\varepsilon^{-2p}, \quad \varepsilon > 0,$$

dann folgt, dass für alle $\lambda \in (0, 1]$, $\tau > 0$ und $n \geq 1$ mit einer Wahrscheinlichkeit P^n nicht kleiner als $1 - e^{-\tau}$ gilt,

$$\lambda \left\| f_{D,\lambda} \right\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}_{L,P}^* \leq \inf_{f \in H} \lambda \left\| f \right\|_H^2 + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* + 2\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 4\lambda^{-1/2} \sqrt{\frac{2\tau + 2}{n}}. \quad \times$$

» Eine einfache Rechnung zeigt ...

$$\sqrt{2\tau + 2 \log(2\mathcal{N}(B_H, \|\cdot\|_\infty, \frac{1}{2}\lambda^{1/2}\varepsilon))} \leq \sqrt{2\tau + 2 \log 2 + a \left(\frac{1}{2}\lambda^{1/2}\varepsilon \right)^{-2p}}$$

und da $2 \log 2 \leq 2$, $\sqrt{s+r} \leq \sqrt{s} + \sqrt{r}$ für $s, r > 0$,

$$\leq \sqrt{2(\tau + 1)} + a^{1/2} 2^p \lambda^{-p/2} \varepsilon^{-p}.$$

Da $4\varepsilon |L|_{2\lambda^{-1/2},1} \leq 4\varepsilon$ und analog

$$2 |L|_{2\lambda^{-1/2},1} \lambda^{-1/2} + 1 \leq 2\lambda^{-1/2} + 1 \leq 4\lambda^{-1/2},$$

erhalten wir

$$4\varepsilon |L|_{2\lambda^{-1/2},1} + \left(2 |L|_{2\lambda^{-1/2},1} \lambda^{-1/2} + 1 \right) \cdot \sqrt{\frac{2\tau + 2 \log 2\mathcal{N}(B_H, \|\cdot\|_\infty, \frac{1}{2}\lambda^{1/2}\varepsilon)}{n}} \leq 4\varepsilon + 4\lambda^{-1/2} \left(\sqrt{\frac{2(\tau + 1)}{n}} + \frac{a^{1/2} 2^p \lambda^{-p/2} \varepsilon^{-p}}{\sqrt{n}} \right).$$

Betrachten wir nun

$$h(\varepsilon) := \varepsilon + \lambda^{-1/2} 2^p \lambda^{-p/2} \varepsilon^{-p} \left(\frac{a}{n} \right)^{1/2},$$

so führt das Lösen von $h'(\varepsilon) = 0$ auf die Minimalstelle

$$\varepsilon^* = p^{(1+p)^{-1}} 2^{(1+p)^{-1}} \lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}}.$$

Setzen wir diese ein, ergibt sich,

$$h(\varepsilon^*) = \underbrace{2^{\frac{p}{1+p}}}_{\leq 2} p^{\frac{1-p}{1+p}} \lambda^{-1/2} \left(\frac{a}{n} \right)^{\frac{1}{2+2p}}$$

Weiterhin besitzt

$$q(p) = p^{\frac{1-p}{1+p}},$$

$$q'(p) = -\frac{p^{-\frac{2p}{1+p}} (2 \ln p^p - 1 + p^2)}{(1+p)^2}$$

ein Maximum in $p = 1$. Folglich ist

$$h(\varepsilon^*) \leq 2\lambda^{-1/2} \left(\frac{a}{n}\right)^{\frac{1}{2+2p}}.$$

Setzen wir dies in die Orakelungleichung 6-2.1 ein, folgt die Behauptung. «

BSP Sei $X \subset \mathbb{R}^n$ kompakt, $O = X^\circ$ nichtleer.

$$C_b^m(\overline{O}) := \{f \in C(X) \cap C^m(O) : f^{(k)} \text{ stetig fortsetzbar auf } X, 1 \leq k \leq m\}.$$

Weiterhin sei $k \in C_b^m(\overline{O})$ ein Kern mit RKHS H . Mit Korollar 5-3.11 folgt, dass

$$\text{id} : H \rightarrow C_b^m(\overline{O})$$

eine stetige Abbildung ist. Aus den Bemerkungen in Kapitel 4-2 folgt für die Entropiezahl e_n ,

$$e_n(\text{id} : C_b^m(\overline{O}) \rightarrow l_\infty(\overline{O})) \leq cn^{-\frac{m}{d}}, \quad n \geq 1$$

mit einer Konstanten c .

Aus der Multiplikativität von Entropiezahlen und $e_i(\cdot) = \|\cdot\|$ folgt, weiter

$$e_n(\text{id} : H \rightarrow l_\infty(\overline{O})) \leq cn^{-\frac{m}{d}} \|\text{id} : H \rightarrow C_b^m(\overline{O})\|.$$

Wir können nun Lemma 4-2.3 anwenden und erhalten

$$\log \mathcal{N}(B_H, \|\cdot\|_\infty, \varepsilon) \leq \log 4 \left(\frac{c}{\varepsilon}\right)^{d/m}.$$

Mit $a = \ln 4c^{d/m}$ und $2p = \frac{d}{m}$ ist nun Korollar 6-2.2 anwendbar.

Betrachten wir nun den Gauß-Kern

$$\Omega_\sigma(x, x') = \exp(-\sigma^2 \|x - x'\|^2)$$

mit RKHS H_Ω , dann können wir m beliebig groß wählen. Es gilt dann

$$e_n(\text{id} : H_\sigma \rightarrow l_\infty(X)) \leq c \sigma^m n^{-m/d}, \quad \sigma \geq 1, \quad m \geq 1, \quad n \geq 1$$

und c einer von n, m, σ unabhängigen Konstanten. Folglich ist

$$\log \mathcal{N}(B_\sigma, \|\cdot\|_\infty, \varepsilon) \leq \log 4 \left(\frac{c \sigma^m}{\varepsilon} \right)^{\frac{d}{m}} = \underbrace{c^{d/m} \log(4) \sigma^d}_{=: a} \varepsilon^{-\frac{d}{m}}.$$

Die rechte Seite der Orakelungleichung ist somit bis auf Konstanten, die von n, τ, m, λ unabhängig sind

$$\lambda^{-1/2} \sigma^{d \cdot \frac{1}{2+2\frac{d}{m}} n^{-\frac{1}{2+2\frac{d}{m}}} + 4\lambda^{-1/2} \sqrt{\frac{2\tau+2}{n}} = \lambda^{-1/2} \sigma^{\frac{dm}{2m+d}} n^{-\frac{m}{2m+d}} + 4\lambda^{-1/2} \sqrt{\frac{2\tau+2}{n}}$$

Für $m \gg 1$ entspricht dies ungefähr

$$\lambda^{-1/2} \sigma^{d/2} n^{-1/2} + 4\lambda^{-1/2} \sqrt{\frac{2\tau+2}{n}}. \quad \blacksquare$$

6-3 Die Funktion $A(\lambda)$

Wir definieren die Funktion A durch

$$A(\lambda) := \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) - \mathcal{R}_{L,P}^*, \quad \lambda > 0.$$

Im Folgenden untersuchen wir, wann $A(\lambda) \approx 0$.

6-3.1 **Satz** Sei P ein W -Maß auf $X \times Y$ und L eine P -integrierbare NVF. Dann gilt für jedes $p \in (0, \infty]$

$$\mathcal{R}_{L,P}^* = \inf_{f \in \mathcal{L}_p(P_X)} \mathcal{R}_{L,P}(f) = \mathcal{R}_{L,P,\mathcal{L}_p(P_X)}^*. \quad \times$$

Es genügt also beschränkte, messbare Funktionen zur Berechnung des Bayes-Risikos zu betrachten.

» „ $p = \infty$ “: Sei $f : Y \rightarrow \mathbb{R}$ messbar mit $\mathcal{R}_{L,p}(f) < \infty$. (Falls es keine solche Funktion gibt, sind wir fertig.) Setze nun

$$f_n := f \cdot \chi_{\{|f| \leq n\}},$$

dann ist $f_n \in \mathcal{L}_\infty$ für alle $n \geq 1$. Weiterhin gilt

$$\begin{aligned} |\mathcal{R}_{L,p}(f_n) - \mathcal{R}_{L,p}(f)| &\leq \int_{X \times Y} |L(x, y, f_n(x)) - L(x, y, f(x))| \, dP(x, y) \\ &= \int_{\{|f| > n\}} |L(x, y, 0) - L(x, y, f(x))| \, dP(x, y) \\ &\leq \int_{\{|f| > n\}} L(x, y, 0) + L(x, y, f(x)) \, dP(x, y) \end{aligned}$$

und da $L(x, y, t) \leq b(x, y) + h(|t|)$ folgt

$$\leq \int_{\{|f| > n\}} \underbrace{b(x, y) + h(0)}_{(*)} + \underbrace{L(x, y, f(x))}_{(**)} \, dP(x, y).$$

(*) $\in L^1(P)$, denn L eine P -integrierbare NVE, (**) $\in L^1(P)$, da $\mathcal{R}_{L,p}(f) < \infty$ und folglich

$$\leq \int_{X \times Y} \chi_{\{|f| > n\}} g \, dP \rightarrow 0, \quad n \rightarrow \infty,$$

denn $g_n \leq g$ und $g_n \rightarrow 0$. Für $\varepsilon > 0$ gibt es daher eine beschränkte Funktion f_ε mit

$$|\mathcal{R}_{L,p}(f_\varepsilon) - \mathcal{R}_{L,p}(f)| \leq \varepsilon.$$

Dann ist

$$\mathcal{R}_{L,p}(f) - \varepsilon \leq \mathcal{R}_{L,p}(f_\varepsilon) \leq \mathcal{R}_{L,p}(f) + \varepsilon$$

und die Behauptung folgt.

» $p < \infty$ “: $\mathcal{L}_\infty(P_Y) \subset \mathcal{L}_p(P_X)$. Damit sind wir fertig. «

6-3.2 **Satz** Sei P ein W -Maß auf $X \times Y$ und L eine P -integrierbare NVF der Ordnung $p \in [1, \infty)$ und ist H ein universeller RKHS (also insbesondere ist X kompakt). Dann gilt

$$\mathcal{R}_{L,P}^* = \mathcal{R}_{L,P,H}^* := \inf_{f \in H} \mathcal{R}_{L,P}(f). \quad \times$$

» Satz 2-1.6 besagt, dass die Abbildung

$$\mathcal{R}_{L,P} : \mathcal{L}_p(P_X) \rightarrow \mathbb{R}$$

stetig ist; weiterhin sind die Abbildungen

$$\text{id} : H \rightarrow C(X), \quad \text{id} : C(X) \rightarrow \mathcal{L}_p(P_X)$$

stetig und haben dichtes Bild. Folglich existiert zu $g \in \mathcal{L}_p(P_X)$ eine Folge (f_n) in H mit $\|f_n - g\|_{\mathcal{L}_p(P_X)} \rightarrow 0$. Aufgrund der Stetigkeit von $\mathcal{R}_{L,P}$ ist

$$\lim_{n \rightarrow \infty} \mathcal{R}_{L,P}(f_n) = \mathcal{R}_{L,P}(g)$$

und mit Satz 6-3.1 folgt die Behauptung

$$\mathcal{R}_{L,P,H}^* = \mathcal{R}_{L,P,\mathcal{L}_p(X)}^* = \mathcal{R}_{L,P}^*. \quad \ll$$

Korollar Falls die Voraussetzungen des Satzes 6-3.2 erfüllt sind, ist

$$A(\lambda) = \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) - \mathcal{R}_{L,P,H}^*, \quad \lambda \geq 0. \quad \times$$

Bemerkung. Satz 6-3.2 gilt auch für nicht universelle Kerne, falls der RKHS H dicht in $\mathcal{L}_p(P_X)$ liegt. \rightarrow

6-3.3 **Lemma** Sei L eine Verlustfunktion, H ein RKHS über X und P ein W -Maß auf $X \times Y$, so dass $\mathcal{R}_{L,P,H}^* = \mathcal{R}_{L,P}^*$. Dann erfüllt die Funktion $A : [0, \infty) \rightarrow [0, \infty]$ die folgenden Eigenschaften:

(i) A ist monoton steigend, stetig und konkav.

(ii) $A(0) = 0$.

(iii) $\mu^{-1}A(\mu) \leq \lambda^{-1}A(\lambda)$ für $0 < \lambda < \mu$.

(iv) $A(\lambda) \leq \mathcal{R}_{L,P}(0) - \mathcal{R}_{L,P,H}^*$ für alle $\lambda \geq 0$.

(v) A ist subadditiv, d.h. $A(\lambda + \mu) \leq A(\lambda) + A(\mu)$ für $\lambda, \mu \geq 0$.

(vi) Falls eine Funktion $h : [0, 1] \rightarrow [0, \infty)$ mit $\lim_{\lambda \rightarrow 0} h(\lambda) = 0$ und $A(\lambda) \leq \lambda h(\lambda)$ für $\lambda \in [0, 1]$ existiert, dann ist

$$A(\lambda) = 0, \quad \lambda \geq 0. \quad \times$$

» „(ii)“: Man rechnet direkt nach, dass

$$A(0) = \inf_{f \in H} \left(0 \cdot \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) - \mathcal{R}_{L,P,H}^* = \inf_{f \in H} \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P,H}^* = 0.$$

„(i)“: A ist konkav. Zu $\lambda \in [0, \infty)$ und $f \in H$ definiere

$$h_f(\lambda) := \lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P,H}^*.$$

Dann ist h_f affin linear und

$$A(\lambda) = \inf_{f \in H} h_f(\lambda).$$

Seien $\varepsilon > 0$, $\lambda, \mu \geq 0$ und $\alpha \in [0, 1]$. Wir zeigen nun, dass

$$\alpha A(\lambda) + (1 - \alpha)A(\mu) \leq A(\alpha\lambda + (1 - \alpha)\mu).$$

Da A das Infimum über h_f ist, existieren $f_1, f_2, f_3 \in H$, so dass

$$A(\lambda) \leq h_{f_1}(\lambda) \leq A(\lambda) + \varepsilon,$$

$$A(\mu) \leq h_{f_2}(\mu) \leq A(\mu) + \varepsilon,$$

$$h_{f_3}(\alpha\lambda + (1 - \alpha)\mu) \leq A(\alpha\lambda + (1 - \alpha)\mu) + \varepsilon$$

Somit gilt

$$\begin{aligned} \alpha A(\lambda) + (1 - \alpha)A(\mu) &\leq \alpha h_{f_1}(\lambda) + (1 - \alpha)h_{f_2}(\mu) \\ &\leq \alpha A(\lambda) + (1 - \alpha)A(\mu) + \varepsilon \\ &\leq \alpha h_{f_3}(\lambda) + (1 - \alpha)h_{f_3}(\mu) + \varepsilon \end{aligned}$$

und da h_f affin linear also konkav folgt

$$\dots \leq h_{f_\varepsilon}(\alpha\lambda + (1 - \alpha)\mu) + \varepsilon \leq A(\alpha\lambda + (1 - \alpha)\mu) + 2\varepsilon.$$

Grenzübergang $\varepsilon \rightarrow 0$ liefert, dass A konkav ist.

A ist stetig in 0. Sei $\varepsilon > 0$, dann existiert ein $f_\varepsilon \in H$ mit

$$\mathcal{R}_{L,P}(f_\varepsilon) - \mathcal{R}_{L,P,H}^* \leq \varepsilon.$$

Ohne Einschränkung ist $f_\varepsilon \neq 0$, denn sonst ist $A \equiv 0$. Für $\lambda \leq \|f_\varepsilon\|_H^{-2} \varepsilon$ folgt daher,

$$0 \leq A(\lambda) = \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) - \mathcal{R}_{L,P,H}^* \leq \underbrace{\lambda \|f_\varepsilon\|_H^2}_{\leq \varepsilon} + \underbrace{\mathcal{R}_{L,P}(f_\varepsilon) - \mathcal{R}_{L,P,H}^*}_{\leq \varepsilon} \leq 2\varepsilon.$$

Somit folgt $A(\lambda) \rightarrow 0 = A(0)$ für $\lambda \rightarrow 0$.

A ist monoton steigend. Seien $0 \leq \lambda \leq \mu$, dann ist

$$h_f(\lambda) \leq h_f(\mu), \quad f \in H$$

und folglich

$$A(\lambda) = \inf_{f' \in H} h_{f'}(\lambda) \leq h_f(\mu)$$

Dies gilt für jedes f , also auch für das Minimum

$$A(\lambda) \leq A(\mu).$$

„(iii)“: Sei $\lambda \leq \mu$, dann ist

$$\begin{aligned} A(\lambda) &= A\left(\frac{\lambda}{\mu}\mu + \left(1 - \frac{\lambda}{\mu}\right) \cdot 0\right) \geq \frac{\lambda}{\mu}A(\mu) + \left(1 - \frac{\lambda}{\mu}\right)A(0) = \frac{\lambda}{\mu}A(\mu) \\ &\Rightarrow \frac{1}{\mu}A(\mu) \leq \frac{1}{\lambda}A(\lambda). \end{aligned}$$

„(v)“: Ohne Einschränkung sei $\lambda \leq \mu$ und folglich,

$$A(\lambda + \mu) \leq \frac{\lambda + \mu}{\mu}A(\mu) = \frac{\lambda}{\mu}A(\mu) + A(\mu) \leq A(\lambda) + A(\mu)$$

A ist stetig. Sei $\lambda > 0$, dann gilt für $\mu \geq 0$,

$$A(\mu) \leq A(\lambda + \mu) \leq A(\lambda) + A(\mu) \rightarrow A(\mu), \quad \lambda \rightarrow 0.$$

„(iv)“:

$$\begin{aligned} A(\lambda) &= \inf_{f \in H} h_f(\lambda) = \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) + \mathcal{R}_{L,P,H}^* \\ &\leq \mathcal{R}_{L,P}(0) - \mathcal{R}_{L,P,H}^*. \end{aligned}$$

„(vi)“: Für $\lambda \in (0, 1]$ gilt nach (iii)

$$A(1) \leq \lambda^{-1} A(\lambda) \leq h(\lambda).$$

Konvergiert $h(\lambda) \rightarrow 0$ für $\lambda \rightarrow 0$, so ist $A(1) = 0$. Da A konkav und nichtnegativ folgt $A \equiv 0$. «

6-3.4 **Satz** Existiert ein $f^* \in H$ mit

$$\mathcal{R}_{L,P}(f^*) = \mathcal{R}_{L,P}^* = \mathcal{R}_{L,P,H}^*.$$

Dann ist

$$A(\lambda) \leq \lambda \|f^*\|_H^2, \quad \lambda \geq 0. \quad \times$$

» Sei $\lambda \geq 0$, dann ist

$$A(\lambda) = \inf_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right) - \mathcal{R}_{L,P}^* \leq \lambda \|f^*\|_H^2. \quad \llcorner$$

Bemerkungen. A. Sei L konvex und k beschränkt, dann gilt

$$\exists c \geq 0 \forall \lambda \geq 0 : A(\lambda) \leq c\lambda \Leftrightarrow \exists f^* \in H : \mathcal{R}_{L,P}(f^*) = \mathcal{R}_{L,P}^*. \quad (*)$$

B. Ein schnelleres Konvergenzverhalten als in (*) ist nur in trivialen Situationen möglich. Lemma 6-3.3 zeigte, dass dann $A(\lambda) = 0$ für $\lambda \geq 0$, d.h.

$$\mathcal{R}_{L,P}(0) = \mathcal{R}_{L,P}^*.$$

c. In vielen Situationen ist $A(\lambda) \leq c\lambda^\beta$ für $\lambda \geq 0$, mit $c \geq 0$ und $\beta \in (0, 1]$ (z.B. gilt dies für die Verlustfunktion der kleinsten Quadrate), wobei

$$\exists c \geq 0 \forall \lambda \geq 0 : A(\lambda) \leq c\lambda^\beta, \quad \beta \in (0, 1] \Leftrightarrow \exists f_{L,P}^* \in [L_2(P_X), H]_{\beta, \infty} . \quad \rightarrow$$

6-3.5 **Lemma** Sei $\mathcal{R}_{L,P}^* = \mathcal{R}_{L,P,H}^*$ und $I \subset (0, \infty)$ ein beschränktes Intervall und Λ sei ein endliches ε -Netz von I .

Für Konstanten $\alpha, c \in (0, \infty)$ gilt dann

$$\min_{\lambda \in \Lambda} (A(\lambda) + c\lambda^{-\alpha}) \leq A(2\varepsilon) + \inf_{\lambda \in I} (A(\lambda) + c\lambda^{-\alpha}) . \quad \times$$

» Ohne Einschränkung ist $\Lambda = \{\lambda_1, \dots, \lambda_m\}$ mit $\lambda_{i-1} \leq \lambda_i$ und $\lambda_0 := \inf I$.

Es gilt $0 < \lambda_i - \lambda_{i-1} \leq 2\varepsilon$ für alle $i = 1, \dots, m$, denn $B_\varepsilon(\lambda_i) \cap B_\varepsilon(\lambda_{i+1}) \neq \emptyset$, sonst wäre Λ kein ε -Netz.

Sei nun $\delta > 0$, dann existiert ein $\lambda^* \in I$ mit

$$A(\lambda^*) + c(\lambda^*)^{-\alpha} \leq \inf_{\lambda \in I} (A(\lambda) + c\lambda^{-\alpha}) + \delta$$

und folglich existiert ein i mit $\lambda_{i-1} \leq \lambda^* \leq \lambda_i$, also $\lambda_i \leq \lambda^* + 2\varepsilon$. Dann gilt

$$\begin{aligned} \min_{\lambda \in \Lambda} (A(\lambda) + c\lambda^{-\alpha}) &\leq A(\lambda_i) + c(\lambda_i)^{-\alpha} \leq A(\lambda^* + 2\varepsilon) + c(\lambda^*)^{-\alpha} \\ &\leq A(2\varepsilon) + \inf_{\lambda \in I} (A(\lambda) + c\lambda^{-\alpha}) + \delta. \end{aligned}$$

Im Limes für $\delta \rightarrow 0$ folgt die Behauptung. «

6-4 Konsistenz und Lernraten für SVMs

6-4.1 **Satz** Sei (X, d) ein kompakter, metrischer Raum, L eine konvexe, lipschitz-stetige Verlustfunktion mit $|L|_1 \leq 1$ und $L(x, y, 0) \leq 1$ für alle $(x, y) \in X \times Y$. Ferner sei H ein universeller RKHS mit Kern k , $\|k\|_\infty \leq 1$ und

$$\log \mathcal{N}(B_H, \|\cdot\|_\infty, \varepsilon) \leq a\varepsilon^{-2p}, \quad \varepsilon > 0,$$

wobei $a \geq 1$ und $p > 0$ Konstanten. Dann ist für alle Folgen (λ_n) in $[0, 1]$ mit

$$\lambda_n \rightarrow 0 \quad \text{und} \quad n\lambda_n^{1+p} \rightarrow \infty$$

die Lernmethode

$$(X, Y)^n \ni D \mapsto f_{D, \lambda_n}$$

universell konsistent. \times

SVMs mit geeigneter Regularisierung sind also universell konsistent.

» L ist konvex, lipschitz-stetig und erfüllt $L(x, y, 0) \leq 1$. Somit ist L eine P -integrierbare NVF der Ordnung 1 und es folgt

$$\mathcal{R}_{L, P, H}^* = \mathcal{R}_{L, P}^*.$$

Korollar 6-2.2 sichert weiterhin, dass

$$\mathcal{R}_{L, P}(f_{D, \lambda_n}) - \mathcal{R}_{L, P}^* \leq A(\lambda_n) + 2\lambda_n^{-1/2} \left(\frac{a}{n}\right)^{\frac{1}{2+2p}} + 4\lambda_n^{-1/2} \sqrt{\frac{2\tau + 2}{n}}$$

mit einer Wahrscheinlichkeit nicht kleiner als $1 - e^{-\tau}$.

Da $A(\lambda_n) \rightarrow 0$ und $n\lambda^{1+p} \rightarrow \infty$ folgt, dass

$$\lambda_n^{-1/2} \left(\frac{1}{n}\right)^{\frac{1}{2+2p}} = \left(\frac{1}{\lambda_n^{1+p} n}\right)^{\frac{1}{2+2p}} \rightarrow 0,$$

sowie

$$4\lambda_n^{-1/2} \sqrt{\frac{2\tau + 2}{n}} = 4\sqrt{2\tau + 2} \frac{1}{\sqrt{\lambda_n n}} \leq 4\sqrt{2\tau + 2} \frac{1}{\sqrt{\lambda_n^{1+p} n}} \rightarrow 0. \quad \ll$$

BSP Betrachte die Verlustfunktion $L = L_{\text{hinge}}$, dann besagt Zhang's Ungleichung 2-2.8

$$\mathcal{R}_{L_{\text{class}}, P}(f) - \mathcal{R}_{L_{\text{class}}, P}^* \leq \mathcal{R}_{L, P}(f) - \mathcal{R}_{L, P}^*.$$

Damit ist die SVM, die L benutzt, unter den Voraussetzungen des Satzes 6-4.1 nicht nur universell konsistent bezüglich L , sondern auch universell klassifikationskonsistent. ∞

6-4.2 **Satz** Es gelten die Voraussetzungen von Satz 6-4.1. Zudem seien $c > 0$ und $\beta \in (0, 1]$ Konstanten mit

$$A(\lambda) \leq c\lambda^\beta.$$

Für $n \geq 1$ definiere dann

$$\lambda_n := n^{-\frac{1}{(1+p)(1+2\beta)}}.$$

Dann lernt die Lernmethode

$$(X \times Y)^n \ni D \mapsto f_{D, \lambda_n}$$

in Verteilung mit der Rate $n^{-\frac{1}{(1+p)(1+2\beta)}}$. \times

» Mit Korollar 6-2.2 folgt,

$$\lambda \left\| f_{D, \lambda_n} \right\|_H^2 + \mathcal{R}_{L, P}(f_{D, \lambda_n}) - \mathcal{R}_{L, P}^* \leq c\lambda_n^\beta + 2\lambda_n^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 4\lambda_n^{-1/2} \sqrt{\frac{2\tau + 2}{n}}$$

mit einer Wahrscheinlichkeit nicht kleiner als $1 - \varepsilon^{-\tau}$. Einsetzen ergibt die Behauptung. \llcorner

Bemerkungen. A. Die Definition von λ_n ist asymptotisch (bezgl. n) für die Orakelungleichung aus Korollar 6-2.2 die beste Wahl.

B. Um λ_n so definieren zu können, benötigt man Wissen über β . In der Regel verfügt man darüber jedoch nicht! \rightarrow

6-4.3 **Definition** Sei L eine Verlustfunktion, die bei 1 abgeschnitten werden kann, H ein RKHS über X und $\Lambda = (\Lambda_n)$ eine Familie von endlichen Teilmengen von $(0, 1]$. Für $n \geq 3$ und $D = ((x_1, y_1), \dots, (x_n, y_n)) \in (X \times Y)^n$ definiere

$$m := \left\lfloor \frac{n}{2} \right\rfloor,$$

$$D_1 = ((x_1, y_1), \dots, (x_m, y_m)),$$

$$D_2 = ((x_{m+1}, y_{m+1}), \dots, (x_n, y_n)),$$

und betrachte

$$f_{D_1, \lambda} = \operatorname{argmin}_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L, D_1}(f) \right),$$

$$\lambda_{D_2} \in \operatorname{argmin}_{\lambda \in \Lambda_n} \mathcal{R}_{L, D_2}(\widehat{f}_{D_1, \lambda}).$$

Dann heißt $D \mapsto f_{D, \lambda_{D_2}}$ *TV-SVM (Training/Validation-Support vector machine)*. \times

6-4.4 **Satz** Es gelten die Voraussetzungen von Satz 6-4.1 und L sei bei 1 abschneidbar. Ferner sei Λ_n ein ε_n -Netz von $(0, 1]$, wobei $\varepsilon_n > 0$. Für $\tau \geq 1$ sei ferner

$$\tau_n = 2\tau + 4 \log |\Lambda_n| + 1.$$

Dann gilt mit einer Wahrscheinlichkeit P^n nicht kleiner als $1 - e^{-\tau}$,

$$\begin{aligned} \mathcal{R}_{L, P}(\widehat{f}_{D_1, \lambda_2}) - \mathcal{R}_{L, P}^* &\leq \inf_{\lambda \in (0, 1]} \left(A(\lambda) + 3\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} \right. \\ &\quad \left. + 12\lambda^{-1/2} \sqrt{\frac{\tau_n}{n}} \right) + A(2\varepsilon_n). \end{aligned}$$

Insbesondere ist die TV-SVM universell konsistent, falls $\varepsilon_n \rightarrow 0$ und $n^{-1} \log |\Lambda_n| \rightarrow 0$. Ist außerdem $A(\lambda) \leq c\lambda^\beta$ für $\lambda \geq 0$ und geeignetes $c \geq 0$ und $\beta \in (0, 1]$, so lernt die TV-SVM in Verteilung mit Rate $n^{-\frac{\beta}{(1+p)(1+2\beta)}}$, sofern $\varepsilon \leq \min \{n^{-1}, |\Lambda_n|\}$, verhält sich die Lernrate polynomial in n . \times

» Da $m = \lfloor \frac{n}{2} \rfloor$ ist $m > \frac{n}{2}$. Mit Korollar 6-2.2 folgt somit

$$\begin{aligned} \mathcal{R}_{L, P}(f_{D_1, \lambda}) - \mathcal{R}_{L, P}^* &\leq A(\lambda) + 2\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 4\lambda^{-1/2} \sqrt{\frac{2\tau + 2}{m}} \\ &\leq A(\lambda) + 3\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 8\lambda^{-1/2} \sqrt{\frac{\tau + 1}{n}} \end{aligned}$$

mit einer Wahrscheinlichkeit nicht kleiner als $1 - e^{-\tau}$. Mit dem union bound folgt, dass diese Abschätzung für alle $\lambda \in \Lambda_n$ simultan mit einer Wahrscheinlichkeit P^n nicht kleiner als $1 - |\Lambda_n| e^{-\tau}$ gilt. Ferner gilt

$$L(x, y, \widehat{t}) \leq \underbrace{|L|_1}_{\leq 1} + \underbrace{L(x, y, 0)}_{\leq 1} \leq 2 =: B.$$

Satz 4-1.2 impliziert mit $n - m \geq n/2 - 1 \geq n/4$, dass

$$\mathcal{R}_{L,P}(f_{D_1}, \lambda_{D_2}) \leq \inf_{\lambda \in \Lambda_n} \mathcal{R}_{L,P}(\widehat{f}_{D,\lambda}) + 4\sqrt{\frac{2\tau + 2\log 2 |\Lambda_n|}{n}}$$

mit Wahrscheinlichkeit P^{n-m} nicht kleiner als $1 - e^{-\tau}$. Man sieht mit Lemma 6-3.5 leicht ein, dass somit auch mit einer Wahrscheinlichkeit P^n nicht kleiner als $1 - (1 + |\Lambda_n|)e^{-\tau}$ gilt

$$\begin{aligned} \mathcal{R}_{L,P}(\widehat{f}_{D_1, \lambda_{D_2}}) - \mathcal{R}_{L,P}^* &\leq \inf_{\lambda \in \Lambda_n} \left(\mathcal{R}_{L,P}(\widehat{f}_{D,\lambda}) - \mathcal{R}_{L,P}^* \right) + 4\sqrt{\frac{2\tau + 2\log 2 |\Lambda_n|}{n}} \\ &\leq \inf_{\lambda \in \Lambda_n} \left(A(\lambda) + 3\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 8\lambda^{-1/2} \sqrt{\frac{\tau + 1}{n}} \right) \\ &\quad + 4\sqrt{\frac{2\tau + 2\log 2 |\Lambda_n|}{n}} \\ &\leq \inf_{\lambda \in (0,1]} \left(A(\lambda) + 3\lambda^{-1/2} \left(\frac{a}{n} \right)^{(2+2p)^{-1}} + 8\lambda^{-1/2} \sqrt{\frac{\tau + 1}{n}} \right) \\ &\quad + A(2\varepsilon) + 4\sqrt{\frac{2\tau + 2\log 2 |\Lambda_n|}{n}}. \end{aligned}$$

Weiterhin ist

$$2\tau + 2\log 2 |\Lambda_n| \leq 2\tau + 2 + 2\log |\Lambda_n| \leq \tau_n$$

sowie $\tau + 1 \leq \tau_n$ und $\lambda^{-1/2} \geq 1$ für $\lambda \in (0, 1]$. Somit folgt die Behauptung. «

Bsp Wähle Λ_n als $\frac{1}{n}$ -Netz mit Kardinalität $\leq 2 - n$. Dann ist die Konstruktion von β unabhängig von D und die TV-SVM lernt mit Rate $n^{-\frac{\beta}{(1+p)(1+2\beta)}}$. ■

Bemerkungen. A. Die TV-SVM kann bei geeigneter Modifikation auch weitere Parameter adaptiv bestimmen, wie z.B. den Kernparameter σ des Gaußkerns.

B. In der Praxis werden kleinere Netze Λ_n benutzt, z.B. $10 \leq |\Lambda_n| \leq 20$ mit geometrischer Verteilung.

c. Die benutzte Technik mit $\|\cdot\|_\infty$ -Überdeckungszahlen war Stand der Forschung bis ca. 2002. Dies kann noch deutlich verbessert werden, denn in Satz 6-2.1 haben wir ausgenutzt, dass $f_{D,\lambda} \in \lambda^{-1/2}B_H$. Dies führte zu $b \approx \lambda^{-1/2}$ in Hoeffdings Ungleichung. Man kann dies aus mehreren Gründen noch deutlich verbessern

1.) Sei $f_{D,\lambda} \in \operatorname{argmin}_{f \in H} \left(\lambda \|f\|_H^2 + \mathcal{R}_{L,P}(f) \right)$, dann gilt

$$\lambda \|f_{D,\lambda}\|_H^2 \leq \lambda \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda}) - \mathcal{R}_{L,P}^* = A(\lambda).$$

Für $A(\lambda) \leq c\lambda^\beta$ folgt somit

$$\|f_{D,\lambda}\|_H \leq \sqrt{c} \lambda^{\frac{\beta-1}{2}}, \quad \lambda \geq 0.$$

Wir sind bisher also stets von $\beta = 0$ und $c = 1$ ausgegangen...

2.) Sei λ_n eine Nullfolge mit $\lambda_n^{-1/2} n^{-(2+2p)^{-1}} \rightarrow 0$ polynomial. Dann besagt Korollar 6-2.2, dass mit hoher Wahrscheinlichkeit gilt

$$\begin{aligned} \lambda_n \|f_{D,\lambda_n}\|_H^2 &\leq \lambda_n \|f_{D,\lambda_n}\|_H^2 + \mathcal{R}_{L,P}(f_{D,\lambda_n}) - \mathcal{R}_{L,P}^* \\ &\leq A(\lambda_n) + \lambda_n^{-1/2} n^{-(2+2p)}, \end{aligned}$$

was polynomial gegen Null konvergiert, falls $A(\lambda) \leq c\lambda^\beta$. Sei α der entsprechende Exponent, dann folgt

$$\|f_{D,\lambda_n}\|_H^2 \leq \kappa \lambda^{-1/2} n^{-\alpha}, \quad (*)$$

mit hoher Wahrscheinlichkeit. Betrachte nun im Beweis von Satz 6-2.1 nur noch die Datensätze für die (*) gilt. Dies ergibt eine bessere Orakelungleichung, welche wiederum auf eine Verbesserung von (*) führt. Iteration dieses Arguments liefert eine deutlich verbesserte Orakelungleichung.

3.) Kann L abgeschnitten werden, z.B. bei $M = 1$, so folgt

$$\|\widehat{f}_{D,\lambda}\|_\infty \leq 1,$$

wobei wir bisher verwendet haben, dass

$$\|f_{D,\lambda}\|_\infty \leq \lambda^{-1/2}.$$

Man sollte daher eigentlich

$$\mathcal{R}_{L,P}(\widehat{f}_{D,\lambda}) - \mathcal{R}_{L,P}^*$$

abschätzen.

4.) In vielen Fällen liegt ein sogenanntes **variance bound** vor,

$$\mathbf{E}(L \circ f - L \circ f_{L,P}^*)^2 \leq V(\mathbf{E}(L \circ f - L \circ f_{L,P}^*))^\vartheta,$$

wobei $L \circ f(x, \mathcal{Y}) = L(x, \mathcal{Y}, f(x))$ und $V > 0, \vartheta \in [0, 1]$ Konstanten.

Vergleichen wir nun die Fehlerterme von Bernsteins-

$$\sqrt{\frac{\sigma^2 \tau}{n}} + \frac{B}{n}$$

und Hoeffdings-Ungleichung

$$B\sqrt{\frac{\tau}{n}},$$

so können wir σ^2 durch $\mathbf{E}(L \circ f - L \circ f_{L,P}^*)^2$ „ersetzen“. Dies führt zu einer neuen Orakelungleichung, die zeigt, dass

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \quad (**)$$

mit hoher Wahrscheinlichkeit klein ist (z.B. $O(n^{-\alpha\vartheta})$). Das variance bound zeigt dann, dass $\sigma^2 \hat{=} \mathbf{E}(L \circ f - L \circ f_{L,P}^*)^2$ ebenfalls mit hoher Wahrscheinlichkeit klein ist und folglich ist (**) noch kleiner (z.B. $O(n^{-\alpha\vartheta-1/2})$). Iteration dieses Arguments ergibt eine deutliche Verbesserung.

- 5.) Dies zusammengenommen kann Lernraten bis zu $O(n^{-1})$ ergeben. Vergleiche dies mit unserem Ergebnis $O(n^{-\frac{\beta}{(1+p)(1+2\beta)}})$, welches nie kleiner ist als $O(n^{-1/3})$.
- 6.) Bisher haben wir nur Überdeckungszahlen bezüglich der Supremumsnorm verwendet. Dies erfordert, dass der Eingaberaum kompakt ist. Oftmals sind die Daten jedoch auf ganz \mathbb{R}^n verteilt (z.B. Gauß-verteilt). Man kann anstatt $\mathcal{N}(B_H, \|\cdot\|_\infty, \varepsilon)$ bessere Überdeckungszahlen wie $\mathcal{N}(B_H, \|\cdot\|_{L^2(P_X)}, \varepsilon)$ betrachten. Dies geht jedoch nicht mehr mit elementaren Methoden. Dafür erhält man z.B. für die Verlustfunktion der kleinsten Quadrate optimale Lernraten. Resultate in dieser Richtung sind relativ neu (2007-2009). \rightarrow

7 Verlustfunktionen II

Bisher haben wir die Verlustfunktionen studiert, die das Lernziel gut beschreiben (L_{class} , L_{hinge} , L_{LS} , ...). Wir haben aber auch schon gesehen, dass L_{hinge} in einem günstigen Verhältnis zu L_{class} steht (Zhang's Ungleichung)

$$\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^* \leq \mathcal{R}_{L_{\text{hinge}},P}(f) - \mathcal{R}_{L_{\text{hinge}},P}^*.$$

Während die Minimierung des Überschussrisikos von L_{class} ein kombinatorisches Optimierungsproblem darstellt und somit äußerst diffizil werden kann, beschreibt die Minimierung des Überschussrisikos von L_{hinge} ein konvexes Optimierungsproblem, zu dessen Lösung wir bereits Methoden kennengelernt haben (siehe Kapitel 6-1). Dies motiviert den Einsatz von L_{hinge} z.B. in SMVs. Weiterhin haben wir in den Übungen gesehen, dass

$$\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^* \leq \sqrt{\mathcal{R}_{L_{\text{LS}},P}(f) - \mathcal{R}_{L_{\text{LS}},P}^*}.$$

Ziel dieses Abschnitts ist es nun zu untersuchen, wie man eine Zielverlustfunktion L , welche das Lernziel beschreibt, durch eine Ersatzverlustfunktion ersetzen kann, die einer „solchen“ Ungleichung genügt.

7-1 Marginbasierte Verlustfunktionen

7-1.1 **Satz von Bartlett, Jordan, Mc Auliffe (2007)** Sei L eine konvexe, marginbasierte Verlustfunktion, die durch

$$\phi : \mathbb{R} \rightarrow [0, \infty)$$

dargestellt wird ($L(y, t) = \phi(y - t)$). Dann sind äquivalent

- (i) ϕ ist in 0 differenzierbar mit $\phi' \Big|_0 = 0$.

(ii) Es existiert eine Funktion $Y : [0, 1] \rightarrow [0, \infty)$ streng monoton und

$$Y(\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^*) \leq \mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* \quad (*)$$

für alle W -Maße P auf $X \times \{-1, 1\}$ und $f : X \rightarrow \mathbb{R}$.

In diesem Fall gilt

$$Y(t) = \phi(0) - \inf_{s \in \mathbb{R}} \left(\frac{t+1}{2} \phi(s) + \left(1 - \frac{t+1}{2}\right) \phi(-s) \right) \\ = \inf_{\eta = \frac{t+1}{2}} \phi(0) - \inf_{s \in \mathbb{R}} (\eta \phi(s) + (1 - \eta) \phi(-s)). \quad \times$$

BSP Die Äquivalenz gilt für L_{hinge} und L_{LS} mit

$$L_{\text{hinge}}, \quad Y(t) = t, \quad \text{für alle } t \in [0, 1], \\ L_{\text{LS}}, \quad Y(t) = t^2, \quad \text{für alle } t \in [0, 1]. \quad \blacksquare$$

Satz Existieren $q \in [0, \infty]$, $c_p > 0$ mit

$$P_X(\{x \in X : |2\eta(x) - 1| < t\}) \leq (c_p t)^q, \quad t \in [0, \infty), \quad (**)$$

und gilt für das in (*) definiert Y , $Y(t) \geq ct^p$ für ein $p > 1$. Dann gilt

$$\mathcal{R}_{L_{\text{class}},P}(f) - \mathcal{R}_{L_{\text{class}},P}^* \leq 2c^{-\frac{q+1}{q+p}} c_p^{\frac{q(p-1)}{q+p}} (\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^*)^{\frac{q+1}{q+p}}. \quad \times$$

(**) heißt **Tsybakov's Noise Assumption** (2004). Es ist

$$\eta = \frac{1}{2} \Leftrightarrow |2\eta - 1| = 0, \quad \eta \approx \frac{1}{2} \Leftrightarrow |2\eta - 1| \approx 0.$$

Die linke Seite in (**) misst die Masse der x , bei denen $\eta(x)$ nahe bei $\frac{1}{2}$ ist. Die rechte Seite beschränkt diese Masse. Für $t \rightarrow 0$ sagt die Noise Assumption

$$P_X(\{x \in X : |2\eta(x) - 1| < t\}) = O(t^q).$$

Wir definieren

$$L_{\alpha\text{-class}}(\mathcal{Y}, t) := \begin{cases} 1 - \alpha, & \mathcal{Y} = 1, \quad t < 0, \\ \alpha, & \mathcal{Y} = -1, \quad t \geq 0, \\ 0, & \text{sonst.} \end{cases}$$

7-1.2 **Satz** Sei L eine konvexe Verlustfunktion, die durch ϕ dargestellt werden kann, mit

$$L_{\alpha\text{-class}}(y, t) = \begin{cases} 1 - \alpha\phi(t), & y = 1, \\ \alpha\phi(t), & y = -1. \end{cases}$$

Dann sind äquivalent

(i) ϕ ist in 0 differenzierbar mit $\phi'(0) = 0$.

(ii) Es gibt eine streng monoton steigende Funktion $Y : [0, 1] \rightarrow [0, \infty)$ mit

$$Y\left(\mathcal{R}_{\alpha\text{-class}, P}(f) - \mathcal{R}_{\alpha\text{-class}, P}^*\right) \leq \mathcal{R}_{L, P}(f) - \mathcal{R}_{L, P}^*. \quad \times$$

BSP Für $0 < \alpha \leq \frac{1}{2}$ sind

$$\begin{aligned} L_{\text{hinge}}, & & Y(t) &= t, \\ L_{\text{LS}}, & & Y(t) &= \frac{t^2}{2\alpha(1-\alpha) + (1-2\alpha)}. \quad \blacksquare \end{aligned}$$

7-2 Distanzbasierte Verlustfunktionen

Für die Verlustfunktion der kleinsten Quadrate L_{LS} ist

$$\|f - f_{L, P}^*\|_{L^2(P_X)}^2 = \mathcal{R}_{L, P}(f) - \mathcal{R}_{L, P}^*, \quad f_{L, P}^* = \mathbf{E}(Y | X).$$

Wie wollen nun den bedingten Median schätzen. Dazu nehmen wir an, dass

$$\text{Median}(Y | x) = \{f^*(x)\}, \quad P_X - \text{f.s.}$$

7-2.1 **Satz** Sei $L(y, t) = |y - t|$ für $y \in \mathbb{R}$ und $t \in \mathbb{R}$. Dann sind für $f : X \rightarrow \mathbb{R}$ äquivalent

(i) $\mathcal{R}_{L, P}(f) = \mathcal{R}_{L, P}^*$,

(ii) $f = f^*$ P_X -f.s.

D.h. der Minimierer des Überschussrisikos entspricht dem Median. \times

Wenn wir bereits wissen, dass $\mathcal{R}_{L,p}(f) \approx \mathcal{R}_{L,p}(f^*)$, wie nahe ist dann f an f^* ?

Definition Eine Verteilung Q auf $[-1, 1]$ hat den **Mediantyp** $q \in (1, \infty)$, wenn $\alpha_Q \in [0, 2]$ und ein $b_Q > 0$ existieren, so dass für $t^* = \text{Median}(Q)$ und jedes $s \in [0, \alpha_Q]$ gelten

$$Q((t^* - s, t^*)) \geq b_Q s^{q-1}, \quad Q((t^*, t^* + s)) \geq b_Q s^{q-1}. \quad \times$$

Definition Sei $p \in (0, \infty]$, $q \in (1, \infty)$ und P ein W -Maß auf $X \times [-1, 1]$. Dann hat P einen **p -mittelbaren Mediantyp** q , wenn folgende Eigenschaften erfüllt sind

- (a) $P(Y | x)$ hat Mediantyp q P_X -f.s..
- (b) Die Abbildung $\gamma : X \rightarrow \mathbb{R}$, $x \mapsto b_{P(Y|x)} \alpha_{P(Y|x)}^{q-1}$ erfüllt $\gamma^{-1} \in L_p(P_X)$. \times

BSP Alle Dichten von $P(Y | x)$ erfüllen für $p = \infty$ und $q = 2$,

$$\min \{1 + f^*(x), 1 - f^*(x)\} \geq \alpha^*.$$

Unter diesen Voraussetzungen gilt

$$r := \frac{p_q}{p+1},$$

$$\|f - f^*\|_{L_r(P_X)} \leq c (\mathcal{R}_{L,p}(f) - \mathcal{R}_{L,p}^*)^{1/q}$$

für alle $f : X \rightarrow [-1, 1]$ und $r = q = 2$. ■

Index

- P*-konsistent, 24
- ε -Netz, 67
- Überdeckungszahl, 67
- Überschussrisiko, 19
- Überwachtes Lernen, 5

- Active Learning, 11
- Atom, 37
- Ausreißeridentifikation, 8

- Batch Learning, 10
- Bayes
 - Entscheidungsfunktion, 18
 - Risiko, 18
- bedingter Median, 22
- Bernsteins Ungleichung, 60

- Clusteranalyse, 9

- Dichteschätzung, 9
- Dimensionsreduzierung, 9

- empirisches Maß, 33
- Entropiezahl, 67
- ERM, 65
- Exponential Kern, 76

- Gauß Kern, 77

- Halbüberwachtes Lernen, 7
- Hilbertfunktionenraum, 80
- Histogrammregel, 11
 - universelle Konsistenz, 31
- Hoeffdings Ungleichung, 62

- Kern, 73
 - stetig differenzierbar, 93
 - universell, 93
- Kernmetrik, 88
- Klassifikation, 6

- Lagrangefunktion, 105
- Lemma
 - für Polynom-Kerne, 75
 - für Taylor-Kerne, 75
- Lernmethode, 17
- Lernrate
 - im Mittel, 36
 - in Verteilung, 36
- Lernraten, 36
- Lernverfahren, 17

- Moving window / kernel rules, 12
- Multi-armed Bandit Learning, 10

- Nearest Neighbor, 12

- No-free-lunch Theorem, 37
- Online Learning, 10
- Optimierungsproblem
 - Duales, 106
 - Primales, 105
- Orakelungleichung
 - ERM, 65
 - ERM (2), 69
 - SVM, 111
- positiv definit, 77
 - strikt, 77
- Regression, 7
- Reinforcement Learning, 10
- reproduzierende Eigenschaft, 80
- reproduzierender Kern, 80
- Risiko, 17, 18
- Satz
 - von Lyapunov, 37
 - von Stone-Weierstraß, 95
- smart, 45
- SVM, 99
 - Eindeutigkeit, 99
 - TV, 126
- symmetrisch, 77
- Theorem
 - Representer, 102
- Unüberwachtes Lernen, 8
- Ungleichung
 - Zhang's-, 55
- Union bound, 63
- universell konsistent, 24
- Verlustfunktion, 17, 18
 - P -integrierbar, 49
 - distanzbasiert, 56
 - konvexe, 48
 - Nemitski-, 49
 - stetige, 48
 - strikt überwachte, 22
 - symmetrisch, 56
 - unüberwachte, 22

Literaturverzeichnis

- [1] F. Cucker and D.X. Zhou, Learning Theory: An Approximation Theory Viewpoint, Cambridge University Press (2007)
- [2] L. Devroye, L. Györfi und G. Lugosi: A Probabilistic Theory of Pattern Recognition, Springer (1996)
- [3] L. Györfi, M. Kohler, A. Krzyzak und Harro Walk, A Distribution-Free Theory of Nonparametric Regression, Springer (2002)
- [4] T. Hastie, R. Tibshirani und J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition, Springer (2009)
- [5] I. Steinwart und A. Christmann, Support Vector Machines, Springer (2008)
- [6] Bauer, H. Wahrscheinlichkeitstheorie; de Gruyter, Berlin (1990), 4. Aufl.
- [7] Bauer, H. Maß- und Integrationstheorie; de Gruyter, Berlin (1998), 2. Aufl.
- [8] Werner, Dirk: Funktionalanalysis. 6. Auflage, Springer 2007