

2. Statistische Entscheidungstheorie

- 2.1 Ziel und Aufgabe
- 2.2 Begriffe zur Wahrscheinlichkeit
- 2.3 Statistischer Klassifikator
- 2.4 A posteriori Wahrscheinlichkeit
- 2.5 Diskriminanzfunktionen und Merkmals-Teilräume
- 2.6 Fehlerwahrscheinlichkeiten
- 2.7 Datenpartitionierung und Kreuzvalidierung

2.1 Ziel und Aufgabe

Ziel in Anwendungsphase (Arbeitsphase):

- optimale Entscheidungen treffen (Klassifikation).

Aufgabe in der Trainingsphase (Lernphase):

- statistisch auswertbare Messungen/Beobachtungen (Trainingselemente) aus dem Problembereich erfassen, und
- statistisch beschreibbares Wissen über den Problembereich einbeziehen,

um einen Klassifikator zu erwerben.

In der Praxis sollten die Lern- und die Anwendungsphase in einem Zyklus ablaufen.

2.2 Begriffe zur Wahrscheinlichkeit

Überblick:

- Problembereich und Ereignisse
- Wahrscheinlichkeit und Verbundwahrscheinlichkeit
- Einschub: Gesetze der Wahrscheinlichkeit
- Bedingte Wahrscheinlichkeit
- Bayes-Formel
- Stochastische(r) Variable und Prozess
- Vektoren von Zufallsvariablen
- Erwartungswert versus Mittelwert

Problembereich und Ereignisse

- Ω : Problembereich.
Umfasst alle möglichen Situationen, Muster, Ereignisse, etc.
- $\Psi \subset \Omega$: Stochastisches Ereignis.
Ergebnis eines zufälligen Versuchs.
- $v \in \Psi$: Elementarereignis.
Ereignis, das nicht weiter zerlegbar ist.

$$\Omega := \bigcup_m v^m$$

Sprechweise:

Ψ ist eingetreten, wenn $v \in \Psi$ eingetreten ist.

Wahrscheinlichkeit und Verbundwahrscheinlichkeit

$P(\Psi)$: Wahrscheinlichkeit für Ereignis Ψ .

$P(\Psi^m, \Psi^n)$: Wahrscheinlichkeit, dass zwei Ereignisse Ψ^m, Ψ^n gleichzeitig eingetreten sind, d.h. $P(v \in \Psi^m \cap \Psi^n)$ ist Verbundwahrscheinlichkeit (joint probability).

Äquivalente Schreibweise:

$$P(\Psi^m, \Psi^n) = P(\Psi^m \cap \Psi^n)$$

Einschub: Gesetze der Wahrscheinlichkeit

$$0 \leq P(A) \leq 1$$

$$P(A \cup \bar{A}) = 1, \quad P(A \cap \bar{A}) = 0$$

$$P(A_1 \cup A_2) = P(A_1) + P(A_2) - P(A_1 \cap A_2)$$

Bedingte Wahrscheinlichkeit

$$P(\Psi^m | \Psi^n) := \frac{P(\Psi^m, \Psi^n)}{P(\Psi^n)}$$

$$\begin{aligned} P(\Psi^m | \Psi^n) \cdot P(\Psi^n) &= \\ P(\Psi^m \cap \Psi^n) &= P(\Psi^n \cap \Psi^m) \\ &= P(\Psi^n | \Psi^m) \cdot P(\Psi^m) \end{aligned}$$

Bayes-Formel

$$P(\Psi^m|\Psi^n) := \frac{P(\Psi^n|\Psi^m) \cdot P(\Psi^m)}{P(\Psi^n)}$$

Beispiel:

$P(\text{Fieber} \text{Grippe})$	Schätzung ist einfach
$P(\text{Grippe} \text{Fieber})$	schwierig schätzbar
$P(\text{Grippe})$	Existiert eine 'Grippewelle' ?
$P(\text{Fieber})$	Temperaturmessung

$$P(\text{Grippe}|\text{Fieber}) := \frac{P(\text{Fieber}|\text{Grippe}) \cdot P(\text{Grippe})}{P(\text{Fieber})}$$

Stochastische(r) Variable und Prozess

$\mathbf{X} : \Omega \rightarrow \mathbb{R}$, \mathbf{X} stochastische Variable (Zufallsvariable).

$x := \mathbf{X}(v)$, Realisierung einer Zufallsvariable.

Realisierungen werden durch Messungen/Beobachtungen erhalten.

Zufallsvariable ist zusätzlich Funktion der Zeit, $t \in T$, weil die Messung selbst unstabil ist.

Neudefinition: $\mathbf{X}(v, t) : \Omega \times T \rightarrow \mathbb{R}$, stochastischer Prozess.

Vektoren von Zufallsvariablen

- Zum Problembereich erhält man somit via Messungen/ Beobachtungen im Computer eine Datenmenge.
- Realisierung x kann auch ein Vektor der Dimension I sein, dann wird durch $\mathbf{X} : \Omega \rightarrow \mathbb{R}^I$ ein Vektor von Zufallsvariablen definiert.
- Der Vektor von Zufallsvariablen kann auch zweigeteilt in einen sogenannten Eingabe- und einen Ausgabeteil sein.
Dies ist relevant für Aufgabenstellungen der Klassifikation oder Regression (siehe später).

Erwartungswert versus Mittelwert

Erwartungswert E einer Zufallsvariable X :

$$E(X) := \int_{-\infty}^{\infty} x \cdot P(x) dx$$

Spezialfall: Endlich viele Realisierungen x^1, \dots, x^M und $P(x^m) = \frac{1}{M}$, $\forall m$.

$$E(X) := \frac{1}{M} \sum_{m=1}^M x^m \quad \text{Mittelwert}$$

2.3 Statistischer Klassifikator

Überblick:

- Klassifizierte Trainingsmenge
- Repräsentation von Klassen
- Trainingselemente und Wahrscheinlichkeiten
- Klassenbedingte Verteilungsdichte
- Beispiele zu klassenbedingten Verteilungsdichten
- Ziel eines statistischen Klassifikators

Klassifizierte Trainingsmenge

Gegeben sei eine klassifizierte Trainingsmenge:

$$\{(x^m, y^{k_m}) | m = 1, \dots, M\}$$

I -dimensionaler Messvektor:

$$x^m := (x_1^m, \dots, x_I^m)^T$$

Klassen des Problembereichs:

$$y^{k_m} \in \{c^k | k = 1, \dots, K\}$$

Repräsentation von Klassen

1. Beispiel: $c^k := k$;

Probleme mit Metrik: Gewichte werden 'überladen'

2. Beispiel: $c^k := (\underbrace{0, \dots, 0, 1, 0, \dots, 0}_{K \text{ Komp., Wert 1 für } k. \text{ Komp.}})$

(1 aus K)-Codierung

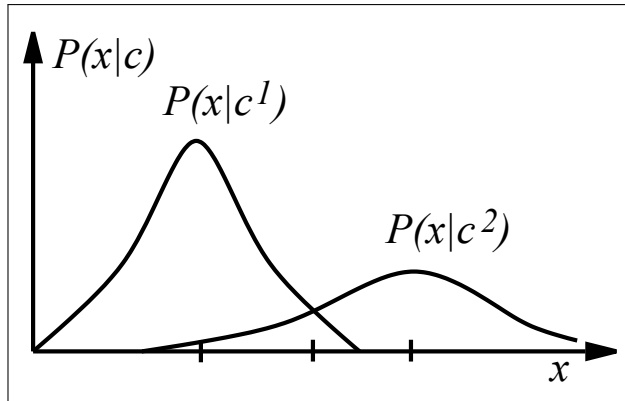
Trainingselemente und Wahrscheinlichkeiten

- Die Menge aller Trainingselemente ist das Ergebnis eines stochastischen Prozesses.
- Die Tupel (x, c^k) sind zufällig aber nicht regellos entstanden.
- Die Regelmäßigkeiten drücken sich durch die Verbundwahrscheinlichkeit $P(x, c^k)$ aus.
- Es gilt: $P(x, c^k) = P(x|c^k) \cdot P(c^k)$

Klassenbedingte Verteilungsdichte

- $P(x|c^k)$ drückt die Unsicherheit der Zugehörigkeit des Trainingselements x zu gegebener Klasse c^k aus.
- Jede Entscheidung beinhaltet ein Risiko, einen Fehler zu begehen.
- Form und Nachbarschaft der klassenbedingten Verteilungsdichten $P(x|c^k)$ wichtig.

Beispiele zu klassenbedingten Verteilungsdichten



Ziel eines statistischen Klassifikators

Bestimmung der Klassenzugehörigkeit von Messvektoren, sodass eine Minimierung der Wahrscheinlichkeit der Fehlklassifikation erreicht wird (minimising a loss function).

2.4 A posteriori Wahrscheinlichkeit

Überblick:

- Bayes-Formel zur Klassifikation
- Maximale a posteriori Wahrscheinlichkeit
- Bestandteile der Bayes-Formel
- Satz der totalen Wahrscheinlichkeit
- Normalisierte a posteriori Wahrscheinlichkeiten

Bayes-Formel zur Klassifikation

Bayes-Formel für die Klassifikation beschreibt, wie bei Beobachtung x , unter Einbezug von Erfahrung/Empirie, die a priori Wahrscheinlichkeit einer Klasse c^k in die a posteriori Wahrscheinlichkeit umgewandelt wird.

$$\underbrace{P(c^k|x)}_{\text{a posteriori}} := \frac{P(x|c^k)}{\underbrace{P(x)}_{\text{Erfahrung}}} \underbrace{P(c^k)}_{\text{a priori}}$$

Bestandteile der Bayes-Formel

$P(c^k)$: A priori Wahrscheinlichkeit der Klasse c^k .

$P(x|c^k)$: Klassenbedingte Wahrscheinlichkeit des Trainingselementes x für die Klasse c^k (Klassenbedingte Dichtefunktion).

$P(c^k|x)$: A posteriori Wahrscheinlichkeit dafür, dass x der Klasse c^k zuzuordnen ist.

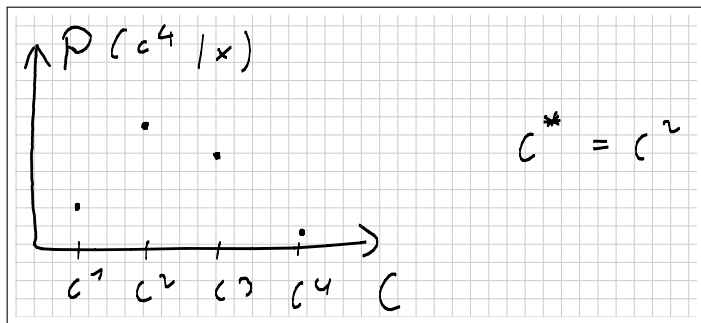
$P(x)$: Wahrscheinlichkeit für das Auftreten von x (Dichtefunktion).

Maximale a posteriori Wahrscheinlichkeit

Die Beobachtung x ist der Klasse c^k zuzuordnen, deren a posteriori Wahrscheinlichkeit maximal ist.

$$c^* := \arg \max_{k \in 1, \dots, K} \{P(c^k | x)\}$$

Beispiel:



Satz der totalen Wahrscheinlichkeit

Unter Verwendung von *Satz der totalen Wahrscheinlichkeit* gilt:

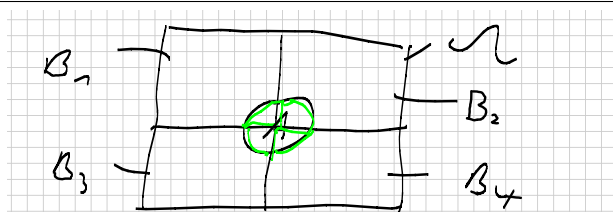
$$P(x) = \sum_{k=1}^K P(x|c^k) \cdot P(c^k)$$

Damit kann $P(x)$ als Normierungsfaktor in der Bayes-Formel dienen, so dass gilt:

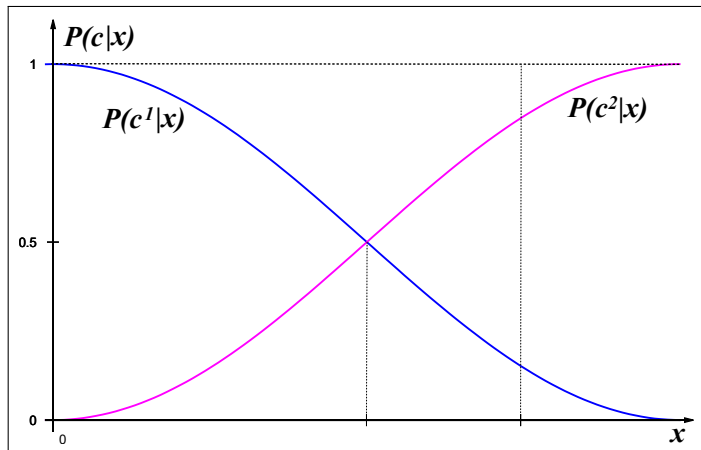
$$\sum_{k=1}^K P(c^k|x) = 1$$

Satz der totalen Wahrscheinlichkeit

Beispiel:


$$P(A) = \sum_{i=1}^4 P(A \cap B_i)$$
$$P(A | B_i) = \frac{P(A \cap B_i)}{P(B_i)}$$

Normalisierte a posteriori Wahrscheinlichkeiten



2.5 Diskriminanzfunktionen und Merkmals-Teilräume

Überblick:

- Klassenbezogene Diskriminanzfunktionen
- Beispiele für Diskriminanzfunktionen
- Einschub: Logarithmus, Minimum Squared Error, Nearest Neighbor
- Perzeptron Diskriminanzfunktion
- Entscheidungsregionen

Klassenbezogene Diskriminanzfunktionen

Klassifikator wendet K Diskriminanzfunktionen d^k an und wählt diejenige Klasse c^l mit $d^l(x) > d^k(x)$ für alle $k \neq l$.

Diskriminanzfunktionen sind das Ergebnis eines Lernvorgangs.

Es werden Ungenauigkeiten bei Messungen zugelassen (weil unvermeidlich), und deren Auswirkungen werden bei der Klassifizierung minimiert.

Beispiele für Diskriminanzfunktionen

$d^k := P(c^k|x) \Rightarrow$ Maximum a posteriori
Wahrscheinlichkeit

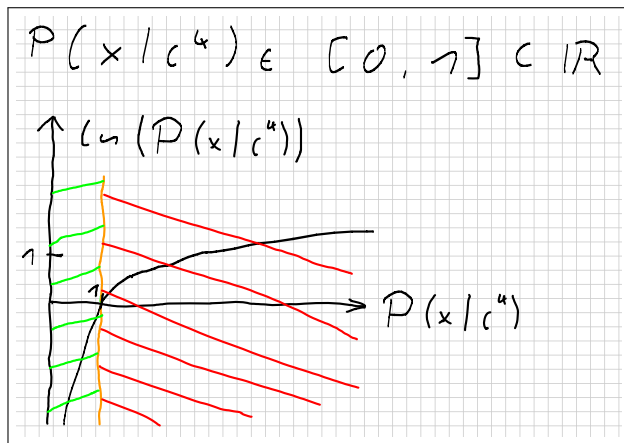
$d^k := P(x|c^k) \Rightarrow$ Maximum Likelihood

$d^k := \ln P(x|c^k) \Rightarrow$ Minimum Squared Error
(falls P eine Gauss-Verteilung)

$d^k := -R(c^k|x) \Rightarrow$ Risikominimierung

Einschub: Logarithmus-Funktion

Beispiel:



Einschub: Minimum Squared Error

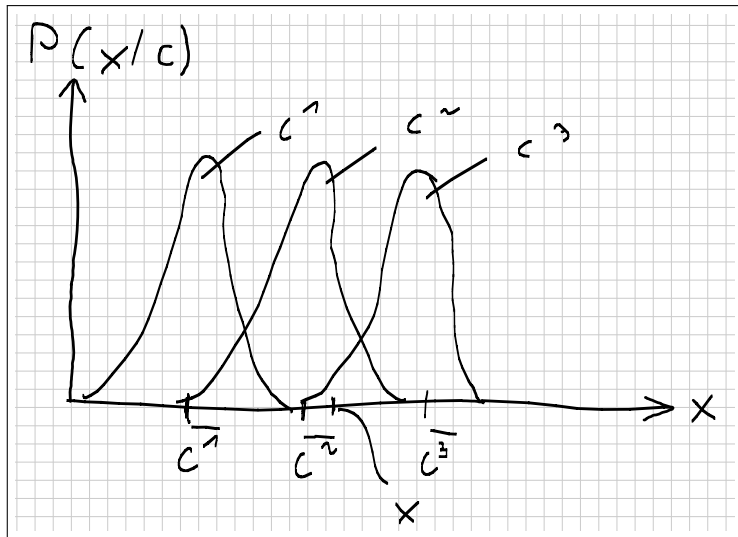
$$\text{Specialfall } \mathcal{D}(x|c^4) := e^{-(x - \bar{c}^4)^2}$$

$$\max_k \left\{ \ln(\mathcal{D}(x|c^k)) \right\} = \max_k \left\{ \ln(e^{-(x - \bar{c}^k)^2}) \right\}$$

$$= \max_k \left\{ -(x - \bar{c}^k)^2 \right\} = \min_k \left\{ (x - \bar{c}^k)^2 \right\}$$

↗
nearest neighbor

Einschub: Nearest Neighbor



Perzeptron Diskriminanzfunktion

Spezialfall: 2-Klassen-Aufgabe

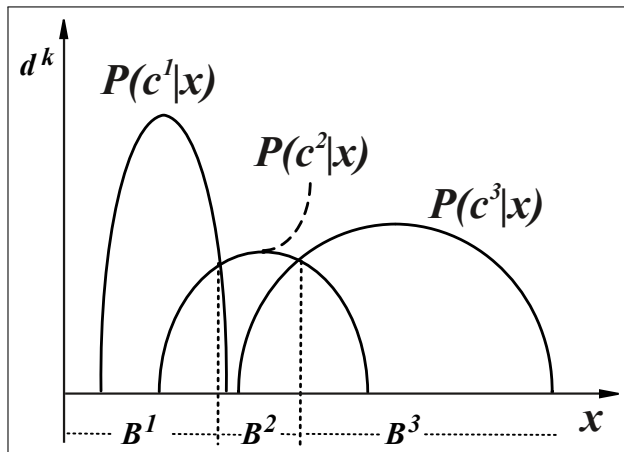
Entscheidung durch Perzeptron (Vorgriff auf Kap.3):

$$\begin{aligned}x &\rightarrow c^1, & \text{wenn } d(x) > 0 \\x &\rightarrow c^2, & \text{sonst}\end{aligned}$$

$$\text{mit } d^1(x) - d^2(x) =: d(x) := w^T x - \Theta$$

Entscheidungsregionen

Klassifikator unterteilt Merkmalsraum in Teilräume B^1, \dots, B^K .



2.6 Fehlerwahrscheinlichkeiten

Überblick:

- Wahrscheinlichkeit von Fehlklassifikationen
- Erwartungswert der Fehlerwahrscheinlichkeit
- Illustration der Fehlerwahrscheinlichkeit
- Fehlerwahrscheinlichkeit bei einer Zwei-Klassen-Aufgabe
- Praktisch relevante Evaluationskriterien

Wahrscheinlichkeit von Fehlklassifikationen

Wahrscheinlichkeit der Fehlklassifikation f (Ereignis), wenn ein neuronales Netz den Messvektor x der Klasse c^k zuordnet, sich also für diese Klasse entscheidet ?

Wahrscheinlichkeit von Fehlklassifikationen

1. Beispiel: Zwei-Klassen-Aufgabe

$$P(f|x) := P(c^l|x), l \neq k$$

2. Beispiel: Mehr-Klassen-Aufgabe

$$\begin{aligned} P(f|x) &:= \sum_{\substack{l=1 \\ l \neq k}}^K P(c^l|x) \\ &:= 1 - P(c^k|x) \end{aligned}$$

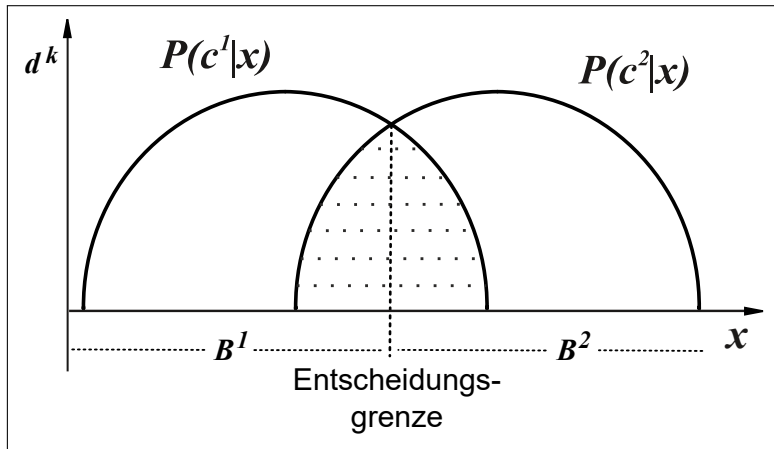
Erwartungswert der Fehlerwahrscheinlichkeit

Erwartungswert der Wahrscheinlichkeit, einen Fehler bei der Klassifizierung zu begehen:

Integration über alle möglichen Messvektoren.

$$P(f) := E(P(f|x)) = \int_{-\infty}^{\infty} P(f|x) \cdot P(x) dx$$

Illustration der Fehlerwahrscheinlichkeit



Fehlerwahrscheinlichkeit bei einer Zwei-Klassen-Aufgabe

$$\begin{aligned} P(f) &:= P(x \in B^1, c^2) + P(x \in B^2, c^1) \\ &:= \int_{B^1} P(x|c^2) \cdot P(c^2) + \\ &\quad \int_{B^2} P(x|c^1) \cdot P(c^1) \end{aligned}$$

Praktisch relevante Evaluationskriterien

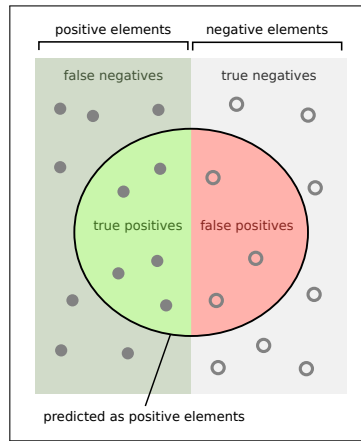
Beispiel für eine Klassifikation:

TP: True Positives

FP: False Positives

TN: True Negatives

FN: False Negatives



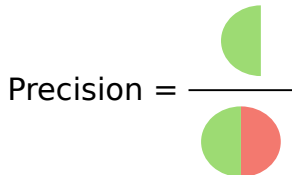
(adaptiert aus Wikipedia, Precision and Recall)

Praktisch relevante Evaluationskriterien

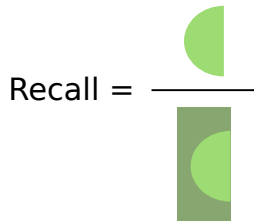
Definition von Precision und Recall:

$$\text{Precision} := \frac{TP}{TP+FP} , \quad \text{Recall} := \frac{TP}{TP+FN}$$

How many predicted positive elements are true?



How many true positive elements are predicted?



(adaptiert aus Wikipedia, Precision and Recall)

2.7 Datenpartitionierung und Kreuzvalidierung

Überblick:

- Lernmenge, Evaluationsmenge, Anwendungsmenge
- Trainingsmenge, Validierungsmenge
- Kreuzvalidierung

Lernmenge, Evaluationsmenge, Anwendungsmenge

Partitionierung des Problembereichs (drei disjunkte Datenmengen):

$$\Omega := \Omega_L \cup \Omega_E \cup \Omega_A$$

Ω_L : Lernmenge zum Lernen eines Klassifikators

$$\Omega_L := \{(x^m, c^{k_m}) | m \in \{1, \dots, M_L\}\}$$

Klassifizierte Elemente, d.h.

Messvektoren und Klassenzugehörigkeiten.

Grundlage für das Lernen des Klassifikators.

Lernmenge, Evaluationsmenge, Anwendungsmenge

Ω_E : Evaluationsmenge zum Evaluieren eines Klassifikators

$$\Omega_E := \{(x^m, c^{k_m}) | m \in \{1, \dots, M_E\}\}$$

Klassifizierte Elemente, d.h.

Messvektoren und Klassenzugehörigkeiten.

Evaluation des Klassifikators und Publikation.

Ω_A : Anwendungsmenge bei Anwendung eines Klassifikators

$$\Omega_A := \{x^m | m \in \{1, \dots, M_A\}\}$$

Unklassifizierte Elemente, d.h.

nur Messvektoren.

Prädiktion der Klassenzugehörigkeiten.

Trainingsmenge, Validierungsmenge

Weitere disjunkte Partitionierung der Lernmenge $\Omega_L := \Omega_T \cup \Omega_V$,
in Trainingsmenge Ω_T und Validierungsmenge Ω_V .

Ω_T : Trainiere Klassifikator so lange, bis die Entscheidung $e : x \rightarrow c^k$
mit einem akzeptablen Minimum an Fehlentscheidungen erfolgt.

Ω_V : Validiere, ob Fehlerrate des Klassifikators bei Ω_V höher ist als
die Fehlerrate bei Ω_T :

nein \rightarrow Beendigung des Trainings.

ja \rightarrow Fortsetzung des Trainings, evtl. Adaption von
Netztopologie und/oder Hyperparameter erforderlich.

Kreuzvalidierung

Engl.: Cross Validation

Praktikabler Lernansatz:

- Variables Aufspalten der Lernmenge $\Omega_L := \Omega_{T_j} \cup \Omega_{V_j}$.
- Iteratives Verändern der Zusammensetzung von Trainings- und Validierungsmenge.
- Jeweils trainieren mit aktueller Trainingsmenge und Evaluation bezüglich aktueller Validierungsmenge.
- Zusammenführen aller Ergebnisse.

Kreuzvalidierung

Beispiel:

