

# Mustererkennung: Diskriminatives Lernen

D. Schlesinger – TUD/INF/KI/IS

Gegeben sie eine (parametrisierte) Klasse der Wahrscheinlichkeitsverteilungen.

Jeder konkreten Wahrscheinlichkeitsverteilung entspricht ein Klassifikator.

Das eigentliche Ziel ist die Erkennung – die Anwendung des Klassifikators.

---

**Generatives** Lernen – der Ausgangspunkt ist die Klasse der WV:

1. Die Parameter der WV werden gesucht (z.B. nach dem ML).
  2. Der Klassifikator wird daraus abgeleitet (z.B. nach dem Bayesschen Prinzip).
  3. Der Klassifikator wird angewendet – Erkennung.
- 

**Diskriminatives** Lernen – der Ausgangspunkt ist die Klasse der Klassifikatoren:

1. Die Parameter des Klassifikators werden „direkt“ gesucht.
2. Der Klassifikator wird angewendet – Erkennung.

Ist die (parametrisierte) Klasse der Klassifikatoren explizit angegeben, so wird die Klasse der WV nicht mehr benötigt.

„Die Wahrheit (Parameter des darunterliegenden Modells) ist egal. Wichtig ist nur, dass alles am Ende gut funktioniert (die vorgegebene Lernstichprobe wird richtig klassifiziert)“.

Beispiel:

zwei Klassen, Gausssiane gleicher Varianz als bedingte Wahrscheinlichkeiten für  $x \in \mathbb{R}^n$ :

$$p(x, k) = p(k) \cdot \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left[-\frac{\|x - \mu^k\|^2}{2\sigma^2}\right]$$

Klassifikator ist eine Hyperebene

→ suche gleich nach einer „guten“ Hyperebene  $\langle w, x \rangle < b$ .

Vergleiche:  $2 \cdot n + 2$  freie Parameter der WV gegen  $n$  freie Parameter des Klassifikators.

Konsequenz: einem Klassifikator entsprechen viele WV.

Beispiel: Die Lage der Hyperebene im obigen Beispiel hängt von  $\sigma$  gar nicht ab. Zentren  $\mu^1$  und  $\mu^2$  sind nicht relevant, sondern die Differenz  $\mu^1 - \mu^2$ .

Wie ist eine gute Hyperebene zu finden?

Das Bayessche Risiko:

$$R_b(e) = \sum_x \sum_k p(x, k) C(e(x), k) \rightarrow \min_e$$

Jetzt hat man aber keine  $p(x, k)$ . Man hat nur eine Lernstichprobe  $L = ((x^l, k^l) \dots)$   
Das Bayessche Risiko wird durch das **Empirische** Risiko ersetzt:

$$R_e(e) = \sum_l C(e(x^l), k^l) \rightarrow \min_{e \in \mathcal{E}}$$

wobei  $\mathcal{E}$  die Klasse der Klassifikatoren ist.

Beispiel:

- lineare Klassifikatoren, d.h.  $\langle w, x \rangle < b$  mit unbekanntem  $w \in \mathbb{R}^n$  und  $b \in \mathbb{R}$
- $\mathbb{I}(k \neq k')$  als Kostenfunktion,
- Annahme: es existiert ein Klassifikator, der die vorgegebene Klassifikation wiedergibt

→ Perzeptron Algorithmus.

Die Frage: wie genau ist das Lernen?

Satz (Vapnik, Chervonenkis):

notwendige und hinreichende Bedingungen für Konsistenz des empirischen Risikos:

$$\lim_{|L| \rightarrow \infty} \left\{ \sup_e |R_e - R_b| > \varepsilon \right\} = 0 \quad \forall \varepsilon > 0$$

Obere Schranke für die Differenz der Fehler:

$$P \left\{ |R_b - R_e| < \sqrt{\frac{h(\log(2N/h) + 1) - \log(\delta/4)}{N}} \right\} > 1 - \delta$$

mit VC-Dimension  $h$ .

(notw. und hinr. Bedingungen für Konsistenz des empirischen Risikos:  $h < \infty$ ).

**VC-Dimension:**

Die kleinste Zahl  $n$  so, dass kein  $(n + 1)$ -Tupel von Punkten beliebig Klassifizierbar ist.

Beispiel: lineare Klassifikatoren im  $\mathbb{R}^n$ ,  $VC = n + 1$ .

Eine etwas andere Vorgehensweise:

$p(x)$  ist bei der Erkennung irrelevant, wichtig ist nur  $p(k|x)$

→ definiere die Klasse der **bedingten** WV explizit ( $p(x)$  ist „beliebig“).

Maximum-Likelihood:

$$p(L; \Theta) = \prod_l [p(x^l) \cdot p(k^l|x^l; \Theta)]$$

$$\ln p(L; \Theta) = \sum_l \ln p(x^l) + \sum_l \ln p(k^l|x^l; \Theta)$$

Beispiel: wieder die Gaussiane gleicher Streuung:

$$\begin{aligned} p(k=1|x) &= \frac{p(1)p(x|1)}{p(1)p(x|1) + p(2)p(x|2)} = \frac{1}{1 + \frac{p(2)p(x|2)}{p(1)p(x|1)}} \\ &= \frac{1}{1 + \exp\left[-\frac{\|x-\mu^2\|^2}{2\sigma^2} + \frac{\|x-\mu^1\|^2}{2\sigma^2} + \ln p(2) - \ln p(1)\right]} \\ &= \frac{1}{1 + \exp(\langle w, x \rangle + b)} \end{aligned}$$

Fragen/Nachteile: kein unüberwachtes Lernen möglich, VC =  $\infty$  (?).