

# Wiederholung: Von Bayes'schem Schließen zu Maximum Likelihood

- Bayes'sche Datenmodellierung:

$$p(\mathbf{w} | D) = \frac{1}{p(D)} p(D | \mathbf{w}) p(\mathbf{w})$$

- MAP-Näherung: Verwende das Modell mit der größten a-posteriori-Wahrscheinlichkeit

$$\mathbf{w}^{MAP} = \arg \max_{\mathbf{w}} p(\mathbf{w} | D)$$

- I. d. Regel gute Näherung bei vielen Daten
- Übergang zu frequentistischer Sicht (relative Häufigkeiten, ein wahres Modell)

- Minimiere negativen log posterior:

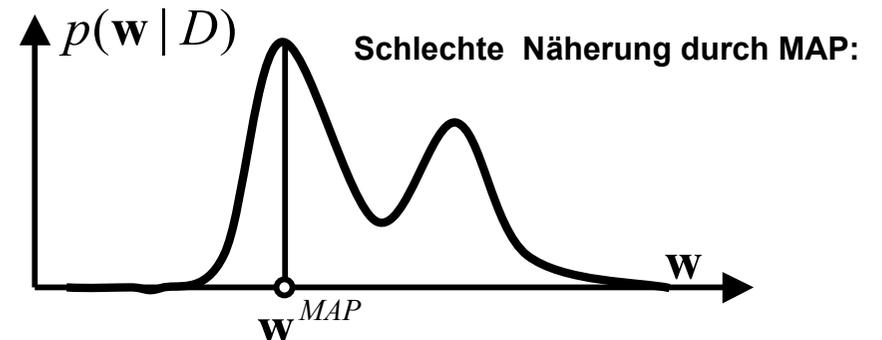
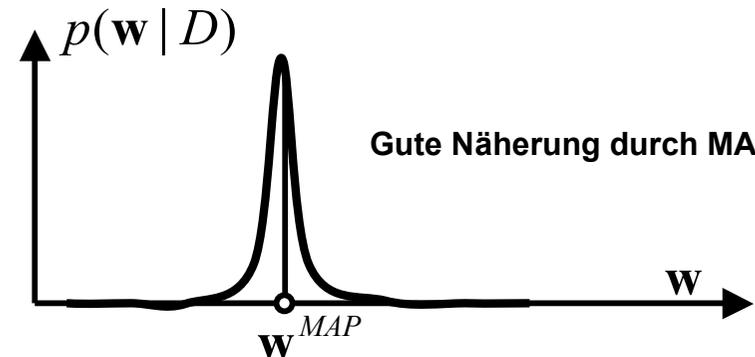
$$-\ln p(D | \mathbf{w}) - \ln p(\mathbf{w}) + const = \min$$



Maximum Likelihood (=min neg. log. likelihood)

Später:

Fehlerminimierung



# Maximum-Likelihood-Parameterschätzung

## Ziel:

- Maximiere Wahrscheinlichkeit, dass die Daten aus dem Modell erzeugt wurden  
=> Maximum Likelihood
- Bem: Maximum-Likelihood (ML) entspricht MAP ohne Vorwissen

## Prinzip:

- Likelihood:  $p(D | \mathbf{w}) = p(\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}\} | \mathbf{w})$
- Max. Likelihood:  $\mathbf{w}^{ML} = \arg \max_{\mathbf{w}} p(D | \mathbf{w})$
- Äquivalent:  $\mathbf{w}^{ML} = \arg \min_{\mathbf{w}} -\ln p(\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}\} | \mathbf{w}) =: \arg \min_{\mathbf{w}} L(\mathbf{w})$

$L(\mathbf{w})$  wird auch als negative log-likelihood oder „Likelihood-Funktion“ bezeichnet

- Für iid („independent, identically distributed“) gezogene Datenpunkte sind die Likelihood-Funktionen verschiedener Datenpunkte unabhängig:

$$L(\mathbf{w}) = -\ln p(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\} | \mathbf{w}) = -\ln \prod_{m=1}^M p(\mathbf{x}^{(m)} | \mathbf{w}) = -\sum_{m=1}^M \ln p(\mathbf{x}^{(m)} | \mathbf{w})$$

# Lernen von Datenmodellen

## Maximum-Likelihood Schätzung und Fehlerminimierung

- **Beispiele für Maximum-Likelihood Schätzung**
- **Maximum-Likelihood und Fehlerminimierung**
- **Kostenfunktionen**

## Maximum-Likelihood Schätzung für Regression:

Likelihood für Input-Output-Paare  $(\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$ ,  $m = 1, \dots, M$

$$p(D | \mathbf{w}) = \prod_{m=1}^M p(\mathbf{y}^{(m)}, \mathbf{x}^{(m)} | \mathbf{w}) = \prod_{m=1}^M p(\mathbf{y}^{(m)} | \mathbf{x}^{(m)}, \mathbf{w}) p(\mathbf{x}^{(m)}) \Rightarrow L(\mathbf{w}) = -\sum_{m=1}^M \ln p(\mathbf{y}^{(m)} | \mathbf{x}^{(m)}, \mathbf{w}) + c$$

Datenmodell Regression: Alles, was nicht durch  $f$  erklärt werden kann, muss Rauschen sein

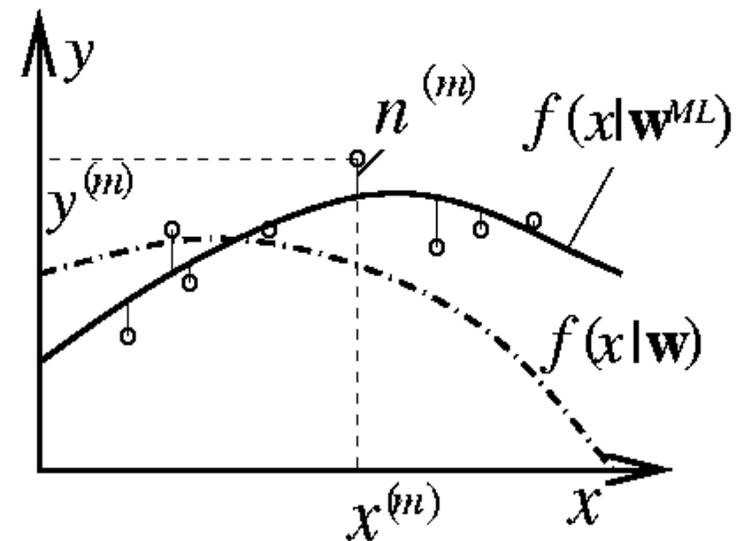
$$\mathbf{y}^{(m)} = \mathbf{f}(\mathbf{x}^{(m)} | \mathbf{w}) + \mathbf{n}^{(m)}, \Rightarrow \mathbf{n}^{(m)} = \mathbf{y}^{(m)} - \mathbf{f}(\mathbf{x}^{(m)} | \mathbf{w})$$

Bei Rauschverteilung  $p_n(\mathbf{n})$  (ohne Konstanten):

$$\begin{aligned} L(\mathbf{w}) &= -\sum_{m=1}^M \ln p(\mathbf{y}^{(m)} | \mathbf{x}^{(m)}, \mathbf{w}) \\ &= -\sum_{m=1}^M \ln p_n(\mathbf{y}^{(m)} - \mathbf{f}(\mathbf{x}^{(m)} | \mathbf{w})) \end{aligned}$$

Bem: Für Gauss'sches Rauschen gilt:

$$L(\mathbf{w}) = \frac{1}{2\sigma^2} \sum_{m=1}^M (\|\mathbf{y}^{(m)} - \mathbf{f}(\mathbf{x}^{(m)} | \mathbf{w})\|^2)$$



ML-Schätzung bei Gaussischem Rauschen entspricht dem Least-Squares-Fit

- **Bsp./Wiederholung: Maximum-Likelihood für das lineare Modell**

**Likelihood:**  $p(D | \mathbf{w}) \equiv p(\mathbf{y} | \mathbf{w}) = p_n(\mathbf{y} - \mathbf{H}\mathbf{w})$

**Annahme: Gaussisches weißes Rauschen**

$$p_n(\mathbf{n}) = \prod_m p_n(n^{(m)}) = \prod_m \frac{1}{(2\pi\sigma_n^2)^{1/2}} \exp\left(-\frac{n^{(m)2}}{2\sigma_n^2}\right)$$

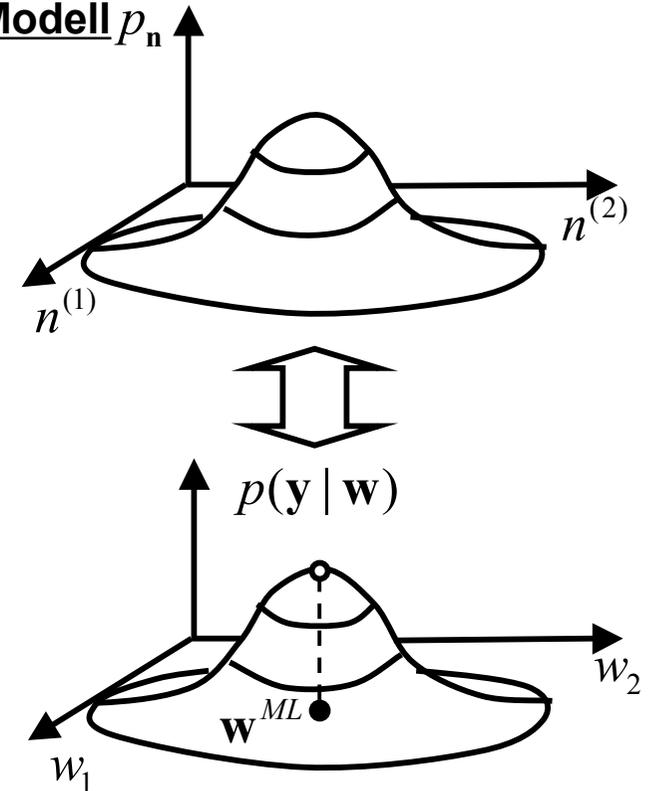
$$= \frac{1}{(2\pi\sigma_n^2)^{M/2}} \exp\left(-\frac{\mathbf{n}^T \mathbf{n}}{2\sigma_n^2}\right)$$

$$p(\mathbf{y} | \mathbf{w}) \propto \exp\left(-\frac{(\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})}{2\sigma_n^2}\right)$$

- **Maximum-Likelihood Parameter**

$$0 = \nabla_{\mathbf{w}} \ln p(\mathbf{y} | \mathbf{w}) = -\frac{1}{2\sigma_n^2} \nabla_{\mathbf{w}} (\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{H}^T \mathbf{y} + \mathbf{w}^T \mathbf{H}^T \mathbf{H} \mathbf{w})$$

$$\Rightarrow \hat{\mathbf{w}} = \mathbf{w}^{ML} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y} \quad \text{analytisch berechenbar!}$$



$$= \frac{1}{\sigma_n^2} (\mathbf{H}^T \mathbf{y} - \mathbf{H}^T \mathbf{H} \mathbf{w})$$

$$= \frac{(\mathbf{H}^T \mathbf{H})}{\sigma_n^2} ((\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y} - \mathbf{w})$$

- Beispiel: Hyperparameter im linearen Modell :

Schätzer Rauschvarianz:  $\hat{\sigma}_n^2 = \frac{\hat{\mathbf{n}}^T \hat{\mathbf{n}}}{\text{tr}(\mathbf{R})}$        $\hat{\mathbf{n}} = \mathbf{y} - \mathbf{H}\hat{\mathbf{w}} = \mathbf{y} - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H} \mathbf{y} =: \mathbf{R} \mathbf{y}$

- Bem: Rausch-Schätzung entspricht Schätzung eines Likelihood-Hyperparameters

Likelihood-Hyperparam.:  $p(D | \mathbf{w}) = p(D | \mathbf{w}, \beta)$      $p(D | \beta) = \int p(D | \mathbf{w}, \beta) p(\mathbf{w} | \alpha) d\mathbf{w}$

$$p(\beta | D) \equiv \frac{p(D | \beta) p(\beta)}{\int p(D | \beta) p(\beta) d\beta}$$

Hier:  $p(D | \mathbf{w}, \beta) \equiv p(\mathbf{y} | \mathbf{w}, \sigma_n^2) = p_n(\mathbf{y} - \mathbf{H}\mathbf{w} | \sigma_n^2) \propto \exp\left(-\frac{(\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})}{2\sigma_n^2}\right)$

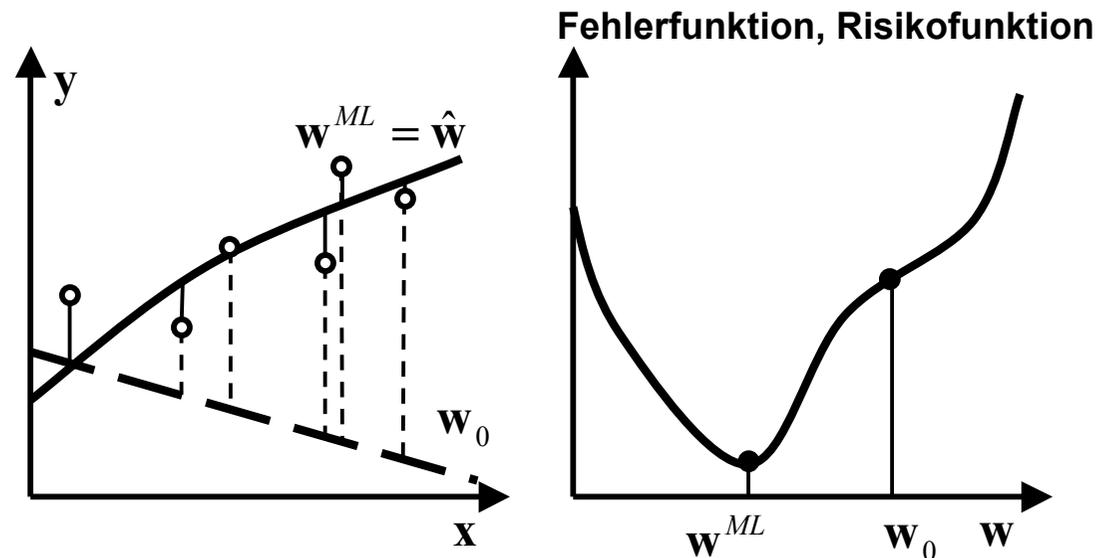
MAP/ML:  $p(D | \beta) \equiv p(\mathbf{y} | \sigma_n^2) = \int p(\mathbf{y} | \mathbf{w}, \sigma_n^2) p(\mathbf{w}) d\mathbf{w} \approx p(\mathbf{y} | \mathbf{w}^{ML}, \sigma_n^2)$

Maximiere bez:  $\sigma_n^2$      $\frac{1}{\sqrt{2\pi\sigma_n^2}} \exp\left(-\frac{\hat{\mathbf{n}}^T \hat{\mathbf{n}}}{2\sigma_n^2}\right)$       ( Siehe ML-Bsp):  $\hat{\sigma}_n^2 = \frac{\hat{\mathbf{n}}^T \hat{\mathbf{n}}}{M}$

# Maximum-Likelihood und Fehlerminimierung

Betrachte überwachtetes Lernproblem (Regression, Klassifikation)

- **Likelihood:** Gegeben ein Rauschmodell  $p_n(\mathbf{n})$ , wie wahrscheinlich ist eine Abweichung des tatsächlichen vom vorhergesagten Output?
- **Negative log-Likelihood:** Wie unwahrscheinlich ist die Abweichung  
Also: Wie groß ist der Fehler, den das Modell macht?
- Teilweise synonym (beobachtet):
  - ++ negative log-Likelihood
  - ++ neg. log. Rauschdichte
  - „Risikofunktion“
  - „Fehlerfunktion“
  - „Verlustfunktion“
  - „Energiefunktion“
- Anpassen („Fitten“) eines Modells durch Risiko/Fehlerminimierung



## Fehler und empirischer Fehler

- **Def. Fehlerfunktion:**  $l(\mathbf{x}, \mathbf{y}, f(\mathbf{x} | \mathbf{w})) \geq 0$  mit  $l(\mathbf{x}, \mathbf{y}, \mathbf{y}) = 0 \quad \forall \mathbf{x}, \mathbf{y}$

**Bem:** -- Es genügt, dass  $l$  am wahren Output sein Minimum annimmt

- **Mittlerer (erwarteter) Fehler:**

$$L_e(\mathbf{w}) = \int l(\mathbf{x}, \mathbf{y}, f(\mathbf{x} | \mathbf{w})) p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

- **Empirischer Fehler:**

**Dichte ist in der Regel unbekannt.**

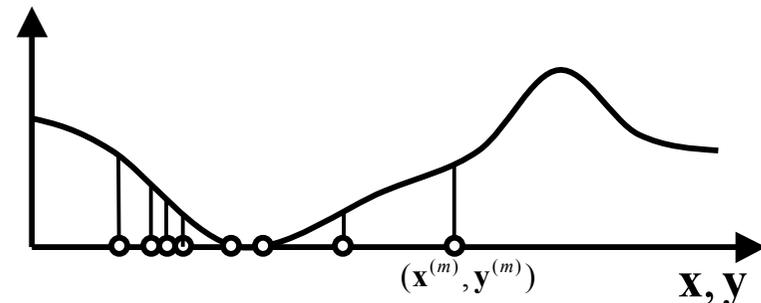
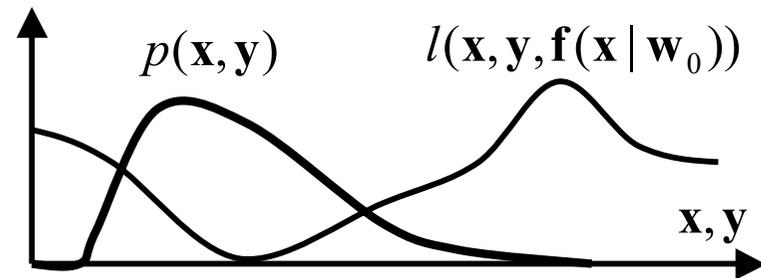
**Man hat Daten**  $(\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$ ,  $m = 1, \dots, M$

$$L(\mathbf{w}) = \frac{1}{M} \sum_{m=1}^M l(\mathbf{x}^{(m)}, \mathbf{y}^{(m)}, f(\mathbf{x}^{(m)} | \mathbf{w}))$$

- **Bsp: Quadratischer Fehler**

$$L(\mathbf{w}) = \frac{1}{M} \sum_{m=1}^M l(\mathbf{x}^{(m)}, \mathbf{y}^{(m)}, f(\mathbf{x}^{(m)} | \mathbf{w})) = \frac{1}{M} \sum_{m=1}^M (\mathbf{y}^{(m)} - f(\mathbf{x}^{(m)} | \mathbf{w}))^2$$

**Minimierung des quadratischen Fehlers= „Least Squares Fit“**



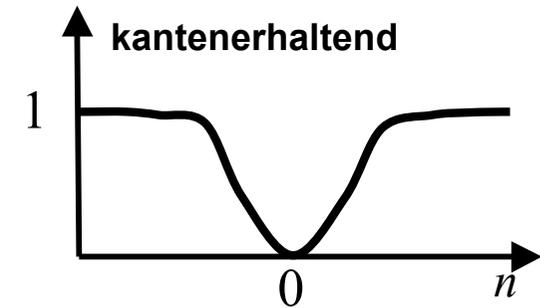
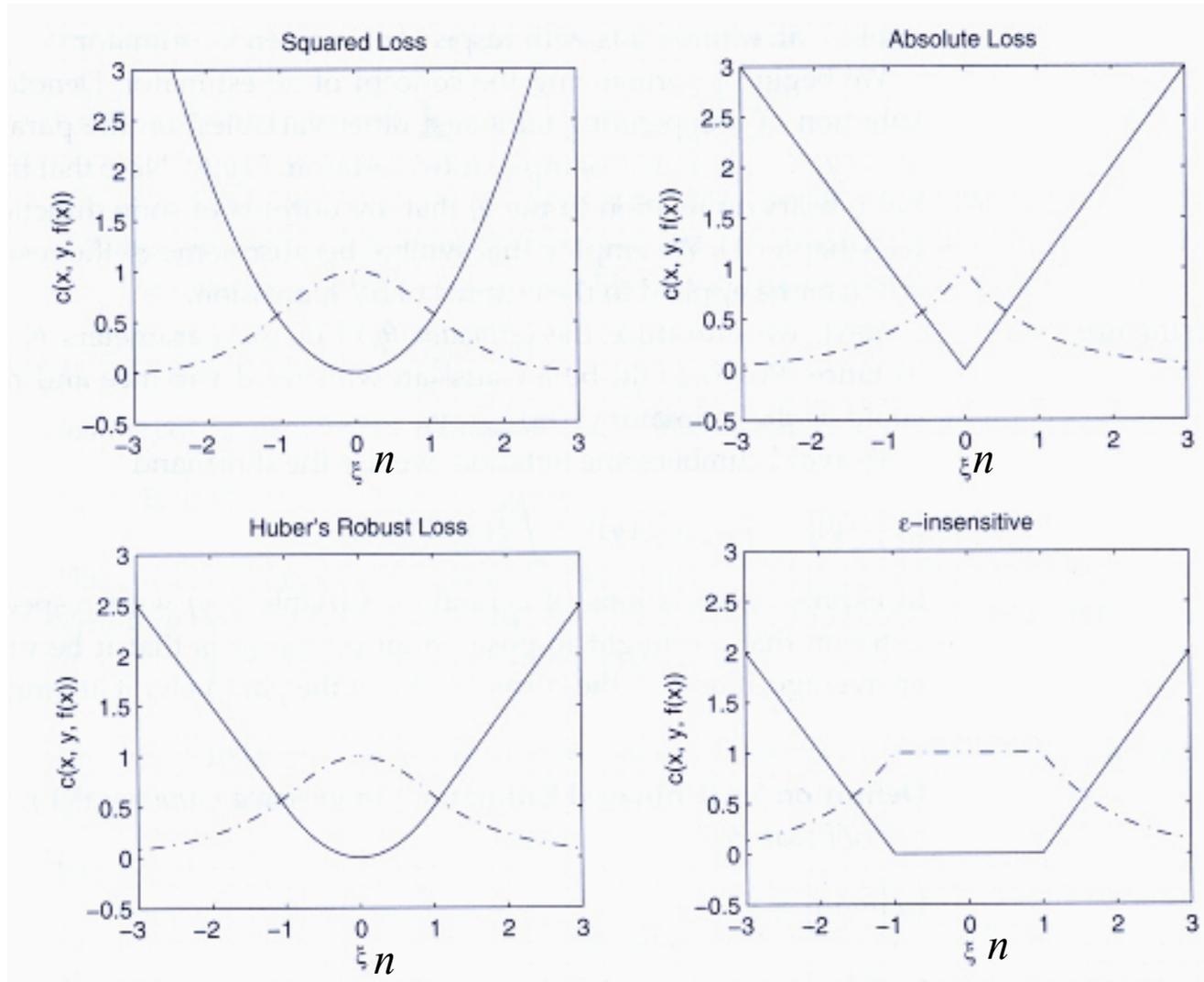
## Fehlerfunktionen für Regression

$$l(\mathbf{x}, \mathbf{y}, \mathbf{f}(\mathbf{x} | \mathbf{w})) = -\ln p(\mathbf{y} | \mathbf{x}, \mathbf{w}) + c = -\ln p_n(\mathbf{y} - \mathbf{f}(\mathbf{x} | \mathbf{w})) + c$$

| Rauschmodell              | Verteilung $p_n(n)$  | Fehler                    | Fehlerfunktion  |
|---------------------------|--|---------------------------|---|
| Gauss                     | $\frac{1}{\sqrt{2\pi}} \exp(-\frac{n^2}{2})$   | quadratisch               | $n^2$   |
| Laplace                   | $\frac{1}{2} \exp(- n )$   | absolut<br>(oszilliert)   | $ n $   |
| Huber's<br>robust         | $\infty \begin{cases} \exp(-\frac{n^2}{2\sigma}), &  n  \leq \sigma \\ \exp(\frac{\sigma}{2} -  n ) & \text{sonst} \end{cases}$                | Outlier-robust            | $\infty \begin{cases} \frac{1}{2\sigma} n^2 &  n  \leq \sigma \\  n  - \frac{\sigma}{2} & \text{sonst} \end{cases}$ |
| $\varepsilon$ -insensitiv | $\infty \begin{cases} 1/2(1 + \varepsilon), &  n  \leq \varepsilon \\ \exp( n  - \varepsilon) / 2(1 + \varepsilon) & \text{sonst} \end{cases}$ | $\varepsilon$ -insensitiv | $\infty \begin{cases} 0, &  n  \leq \varepsilon \\  n  - \varepsilon & \text{sonst} \end{cases}$                    |
|                           | —  | Kantenerhaltend           | $1 - \exp(-\frac{n^2}{2\sigma^2})$  |

**Bem: Es gibt Fehlerfunktionen, für die kein äquivalentes Rauschmodell existiert**

## Fehlerfunktionen: Beispiele für Regression



## Fehlerfunktionen für binäre Klassifikation

- **Likelihood:** Wie groß ist die Wahrscheinlichkeit einer korrekten Klassenzuteilung?

$$\Pr(y = 1 | f(\mathbf{x} | \mathbf{w})) \equiv \Pr(y | \mathbf{x}, \mathbf{w})$$

- **Fehlerfunktion:**

$$l(\mathbf{x}, y, f(\mathbf{x} | \mathbf{w})) = -\ln \Pr(y | \mathbf{x}, \mathbf{w})$$

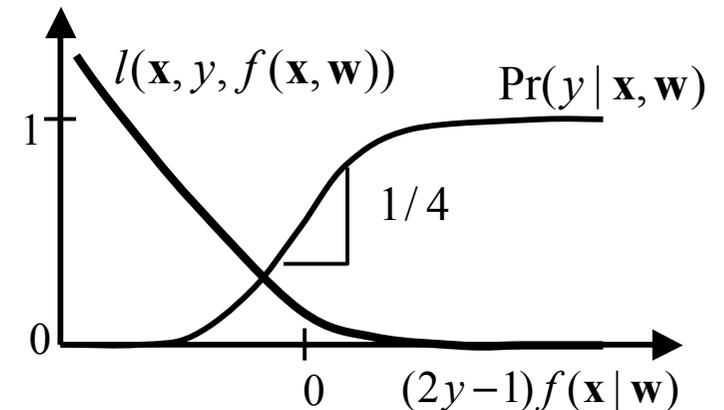
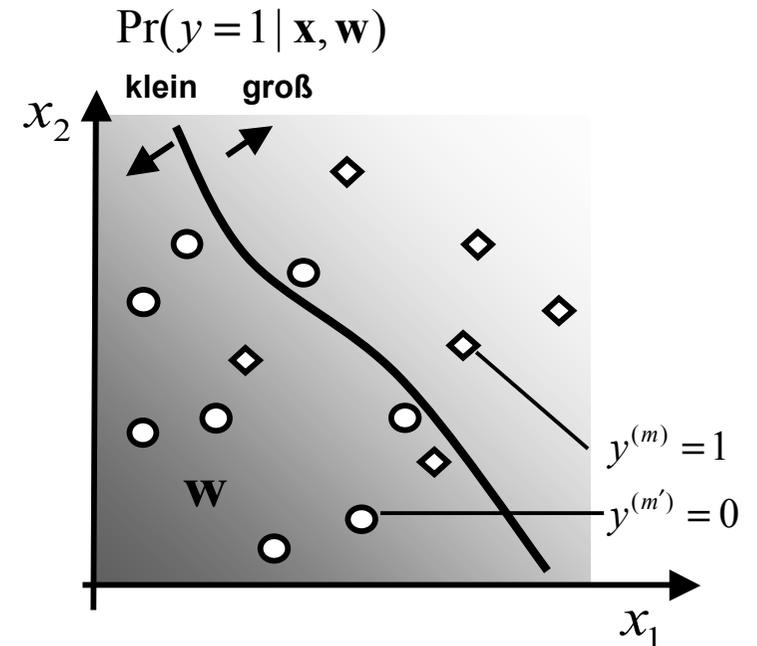
- **Logistische Fehlerfunktion:**

$$\Pr(y = 1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-f(\mathbf{x} | \mathbf{w}))}$$

$$l(\mathbf{x}, y = 1, f(\mathbf{x} | \mathbf{w})) = \ln(1 + \exp(-f(\mathbf{x} | \mathbf{w})))$$

- **Logistischer empirischer Fehler:**

$$L(\mathbf{w}) = \sum_{y^{(m)}=1} \ln(1 + \exp(-f(\mathbf{x}^{(m)} | \mathbf{w}))) + \sum_{y^{(m)}=-1} \ln(1 + \exp(+f(\mathbf{x}^{(m)} | \mathbf{w}))) = \sum_{m=1}^M \ln(1 + \exp(-(2y^{(m)} - 1)f(\mathbf{x}^{(m)} | \mathbf{w})))$$



## Weitere Fehlerfunktionen

- **Falsch-Klassifikations-Zähler**

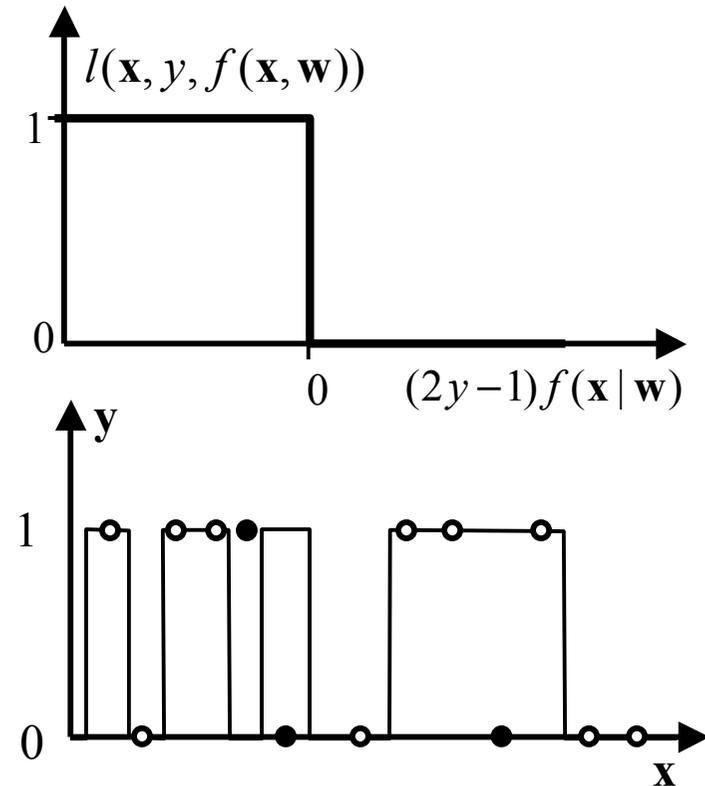
$$l(\mathbf{x}, y, f(\mathbf{x} | \mathbf{w})) = \begin{cases} 0, & (2y-1)f(\mathbf{x} | \mathbf{w}) \geq 0 \\ 1, & \text{sonst} \end{cases}$$

- **Inputabhängiger Falsch-Klassifikations-Zähler**  
(Bsp: Klassifikation Steine und Diamanten)

$$l(\mathbf{x}, y, f(\mathbf{x} | \mathbf{w})) = \begin{cases} 0, & y = f(\mathbf{x} | \mathbf{w}) \\ l_0(\mathbf{x}), & \text{sonst} \end{cases}$$

- **Outputabhängiger Falsch-Klassifikations-Zähler**  
(Bsp: Klassifikation Krebs/gesund)

$$l(\mathbf{x}, y, f(\mathbf{x} | \mathbf{w})) = \begin{cases} 0, & y = f(\mathbf{x} | \mathbf{w}) \\ l_1, & y = 1, f(\mathbf{x} | \mathbf{w}) = 0 \\ l_2, & y = 0, f(\mathbf{x} | \mathbf{w}) = 1 \end{cases}$$



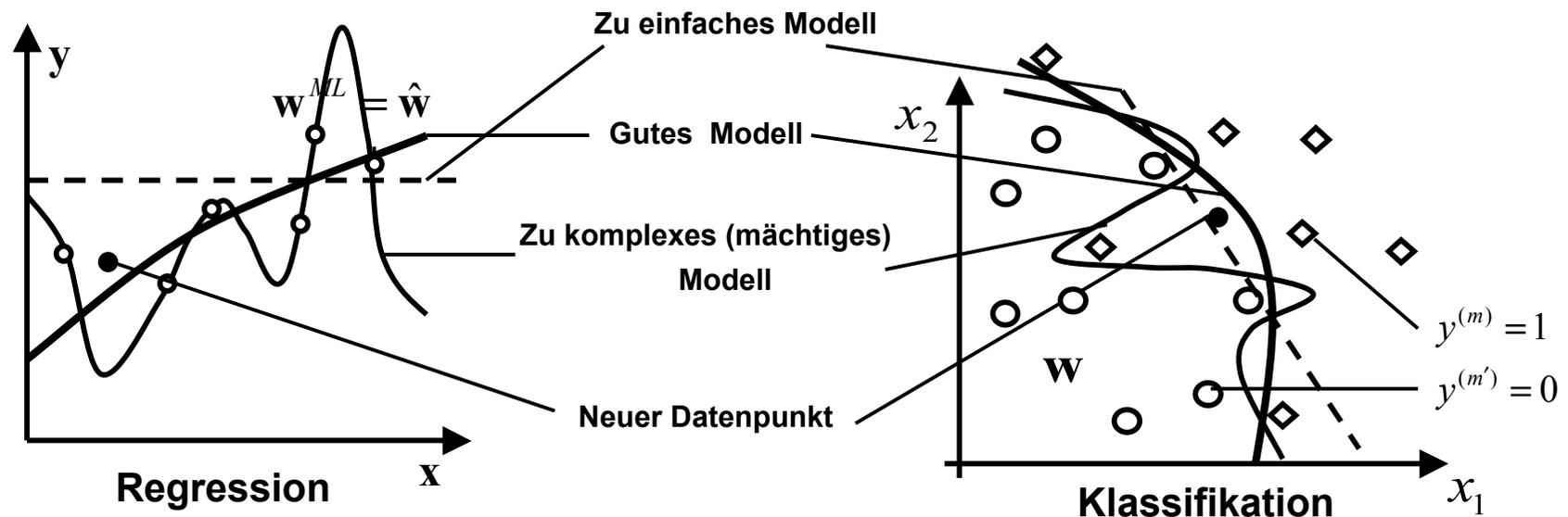
# Lernen von Datenmodellen

## Generalisierung und Regularisierung

- **Modell-Komplexität und Varianz-Bias Problem**
- **Generalisierung durch Regularisierung**
- **Kreuzvalidierung: Optimierung der Modellkomplexität**
- **Anwendungsbeispiel für Regularisierung: Funktionelle Kernspin-Bilder**

# Generalisierung und Regularisierung

- Betrachte überwachtetes Lernproblem mit endlichem Datensatz  $(\mathbf{x}^{(m)}, y^{(m)})$ ,  $m = 1, \dots, M$



- **Beobachtungen:**
  - Ein genügend komplexes Modell kann die Trainingsdaten beliebig genau erklären (kleiner Trainingsfehler)
  - Neue Datenpunkte werden dann aber schlechter erklärt (großer Testfehler, schlechte Generalisierungsfähigkeit, „Overfitting“)
- **Finden der statistischen Struktur: Erkläre Daten mit möglichst einfachem Modell**

## Etwas formaler: Das Bias-Varianz Problem

- **Ziel der Datenmodellierung:**
  - Modellierung der statistischen Struktur  $h$
  - nicht Modellierung eines speziellen Datensatzes
- **Betrachte mittleres Verhalten über viele Datensätze:**

$$\mathbf{D} = \{D_1, D_2, \dots\}$$
- **Bias:** Wie stark weicht das über alle Datensätze gemittelte Modell von dem wahren Modell ab?
- **Varianz:** Wie stark variiert das Modell (wie stark hängt es vom einzelnen Datensatz ab?)
- Sei  $E_{\mathbf{D}}[\cdot]$  der Erwartungswert über alle Datensätze  
Für Regression und quadratischen Fehler:

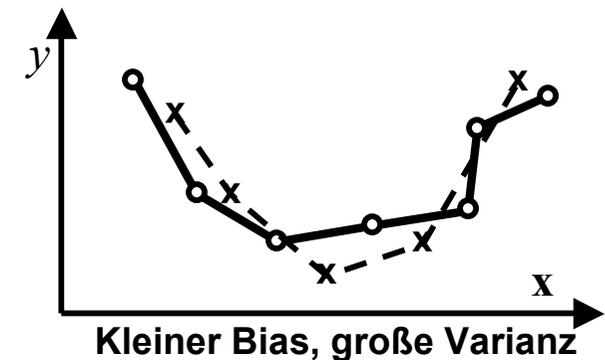
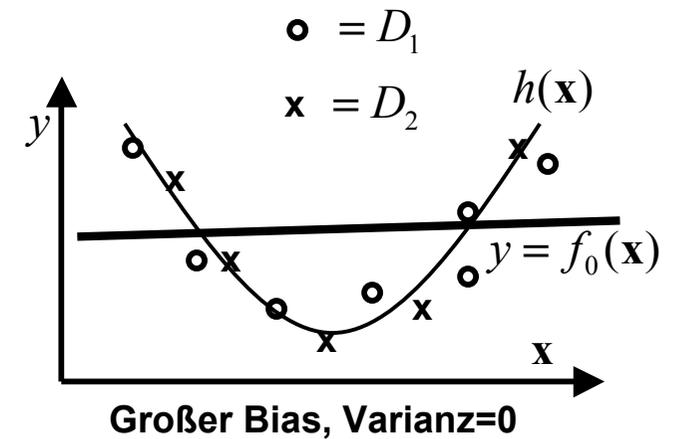
$$(f(\mathbf{x} | \mathbf{w}) - h(\mathbf{x}))^2 = (f(\mathbf{x} | \mathbf{w}) - E_{\mathbf{D}}[f(\mathbf{x} | \mathbf{w})] + E_{\mathbf{D}}[f(\mathbf{x} | \mathbf{w})] - h(\mathbf{x}))^2$$

$$E_{\mathbf{D}}[(f(\mathbf{x} | \mathbf{w}) - h(\mathbf{x}))^2] = (E_{\mathbf{D}}[f(\mathbf{x} | \mathbf{w})] - h(\mathbf{x}))^2 + E_{\mathbf{D}}[(f(\mathbf{x} | \mathbf{w}) - E_{\mathbf{D}}[f(\mathbf{x} | \mathbf{w})])^2]$$

**Minimiere gemeinsam:**

**Bias**

**Varianz**



## Bias-Varianz-Trade-off in Bayes-Modellen

- **Betrachte Bayes-Schätzer:**  $p(\mathbf{w} | D) = \frac{1}{p(D)} p(D | \mathbf{w}) p(\mathbf{w})$
- **MAP:**
- **Maximierung der Likelihood: Erzielung eines kleinen Bias (guter Fit der Datenpunkte)**
- **Der Prior kann zu komplexe Modelle bestrafen („Occam’s Razor“)**
  - Maximierung des Prior favorisiert einfache Modelle
  - Erzielung einer kleinen Varianz (kein overfitting)
- **Maximierung der Evidenz: Bestimmung der Hyperparameter**  
z. B. Optimierung der relativen Gewichtung von Bias und Varianz-Minimierung
- **Bsp: Weight-Decay Prior**  $p(\mathbf{w}) = \frac{1}{Z(\alpha)} \exp\left(-\frac{\alpha}{2} \|\mathbf{w}\|^2\right)$
- **Favorisiert kleine Gewichtswerte**  
=> favorisiert glattere Kurven ( $y$  ändert sich „langsamer“ mit  $x$ )
- **Negativer Log-Prior entspricht Regularisierungsterm**  
Hyperparameter  $\alpha$  gewichtet den Regularisator  
$$-\ln p(\mathbf{w}) \equiv \alpha R(\mathbf{w}) = \frac{\alpha}{2} \|\mathbf{w}\|^2 + c$$

## Regularisierung: Favorisierung von Modellen mit bestimmter (z.B. niedriger) Komplexität

### Bayes-Formalismus

- negative log-Likelihood
- negativer log-Prior
- Max. a Posteriori

### Schätztheorie

- Fehlerfunktion  $L$
- Regularisierungsfunktion  $R$
- $F(\mathbf{w}) = L(\mathbf{w}) + \alpha R(\mathbf{w}) \stackrel{!}{=} \min$

### Beobachtung

- Rauschen enthält alle (auch hohe) Frequenzen
- Deterministische Zusammenhänge sind oft glatt (niedrige Frequenzen)
- Sinnvolle Regularisierung: Bevorzuge „glatte“ Modelle
- Glattere Modelle  $\Leftrightarrow$  weniger Parameter, oder kleinere Gewichte

### Beispiele für Regularisierungsfunktionen

- **Weight-Decay**  $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$
- **Kurven-basierte Regularisierung**  $R(\mathbf{w}) = \sum_{m=1}^M \sum_{j=1}^c \sum_{i=1}^d \left( \frac{\partial^2 f(\mathbf{x} | \mathbf{w})}{\partial x_i \partial x_j} \Big|_{\mathbf{x}^{(m)}} \right)^2$