

Wissensentdeckung in Datenbanken

Modellklassen, Verlustfunktionen

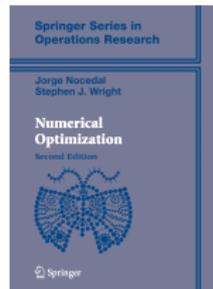
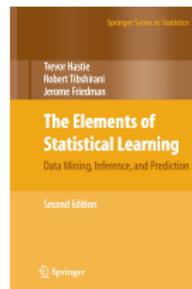
Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

02.05.2017

Literatur

- Trevor Hastie, Robert Tibshirani, Jerome Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd edition. Springer Series in Statistics. Springer. 2009
 - Im Moodle zum Download verfügbar.
- Jorge Nocedal. Stephen Wright. Numerical Optimization. 2nd edition. Springer Series in Operations Research and Financial Engineering. Springer-Verlag New York. 2006
 - In der Zentralbibliothek verfügbar.



Daten—und dann?



- Personendaten
- Medizinische Daten
- Konto- und Zahlungsdaten
- Verbindungsdaten
- Soziale Netzwerke



Martina Mustermann
 BARMER
 123456789
 Versicherung

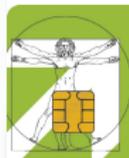
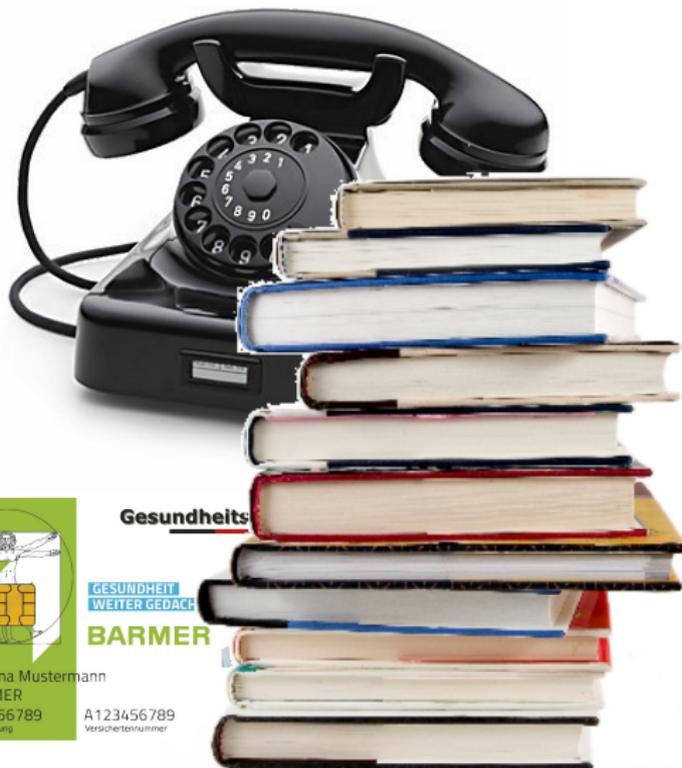
Gesundheits
 GESUNDHEIT
 WETTER GEDACHT
BARMER

A123456789
 Versicherungsnummer

Realisierung von Zufallsvektoren X



$$X = \begin{pmatrix} \text{Alter} \\ \text{Geschlecht} \\ \text{Krankheitstage} \\ \text{Medikation} \\ \text{Kontostand} \\ \text{Kredite} \\ \text{Webseiten} \\ \dots \end{pmatrix}$$



Martina Mustermann
 BARMER
 123456789
 Versicherung

Gesundheits
 GESUNDHEIT
 WEITER GEDACHT
 BARMER

A123456789
 Versichertennummer

Begriffe

- Daten $\mathcal{D} = \{(\mathbf{x}^1, \mathbf{y}^1), (\mathbf{x}^2, \mathbf{y}^2), \dots, (\mathbf{x}^N, \mathbf{y}^N)\}$
 - Realisierungen von Zufallsvariable X mit n -dimensionaler Domäne \mathcal{X}
 - Multimenge; $\#_{\mathcal{D}}(\mathbf{x}, \mathbf{y}) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{N}$
 - Messfehler, Rauschen $\epsilon \sim \mathbb{P}$
 - Fehlende Werte; $\mathbf{x}_i \in \mathcal{X}_i \cup \{?\}, 1 \leq i \leq n$
- Modell f aus Modellklasse \mathcal{M}
 - Funktionen $\mathcal{M} \subseteq \mathcal{F}$
 - Koeffizienten / Parameter $\mathcal{M} \subseteq \mathbb{R}^d$
 - Datenpunkte $\mathcal{M} \subseteq \mathcal{D}, \mathcal{M} \subseteq \mathcal{X} \times \mathcal{Y}$
- Verlustfunktion (Güte) $\ell : (f; \mathcal{D}) \mapsto \mathbb{R}$



Begriffe

- Daten $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)\}$
 - Realisierungen von Zufallsvariable X mit n -dimensionaler Domäne \mathcal{X}
 - Multimenge; $\#_{\mathcal{D}}(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{N}$
 - Messfehler, Rauschen $\epsilon \sim \mathbb{P}$
 - Fehlende Werte; $x_i \in \mathcal{X}_i \cup \{?\}, 1 \leq i \leq n$

- Modell f aus Modellklasse \mathcal{M}
 - Funktionen $\mathcal{M} \subseteq \mathcal{F}$
 - Koeffizienten / Parameter $\mathcal{M} \subseteq \mathbb{R}^d$
 - Datenpunkte $\mathcal{M} \subseteq \mathcal{D}, \mathcal{M} \subseteq \mathcal{X} \times \mathcal{Y}$

- Verlustfunktion (Güte) $\ell : (f; \mathcal{D}) \mapsto \mathbb{R}$



Begriffe

- Daten $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)\}$
 - Realisierungen von Zufallsvariable X mit n -dimensionaler Domäne \mathcal{X}
 - Multimenge; $\#_{\mathcal{D}}(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{N}$
 - Messfehler, Rauschen $\epsilon \sim \mathbb{P}$
 - Fehlende Werte; $x_i \in \mathcal{X}_i \cup \{?\}, 1 \leq i \leq n$

- Modell f aus Modellklasse \mathcal{M}
 - Funktionen $\mathcal{M} \subseteq \mathcal{F}$
 - Koeffizienten / Parameter $\mathcal{M} \subseteq \mathbb{R}^d$
 - Datenpunkte $\mathcal{M} \subseteq \mathcal{D}, \mathcal{M} \subseteq \mathcal{X} \times \mathcal{Y}$

- Verlustfunktion (Güte) $\ell : (f; \mathcal{D}) \mapsto \mathbb{R}$



Allgemeines Vorgehen

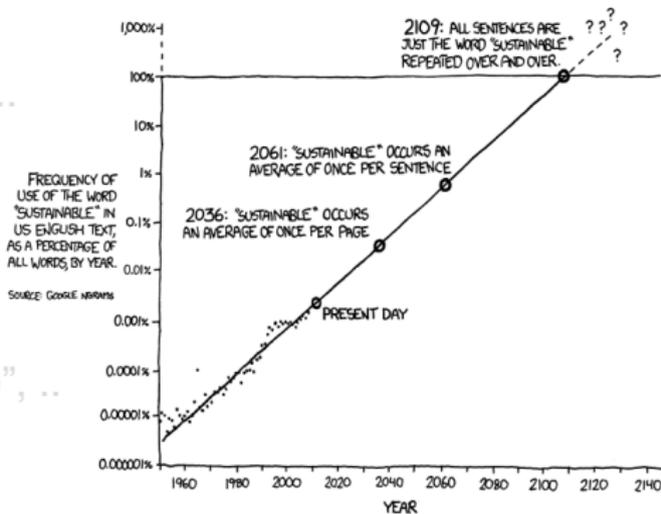
- Daten $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)\}$
- Datenpunkt $x \in \mathcal{X}$, Label (Klasse) $y \in \mathcal{Y}$
- Modelle \mathcal{M}

Modell "lernen", Datenanalyse, ..

$$f^* = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D})$$

Modellanwendung, "Vorhersage", ..

$$\hat{y} = f^*(x)$$



Allgemeines Vorgehen

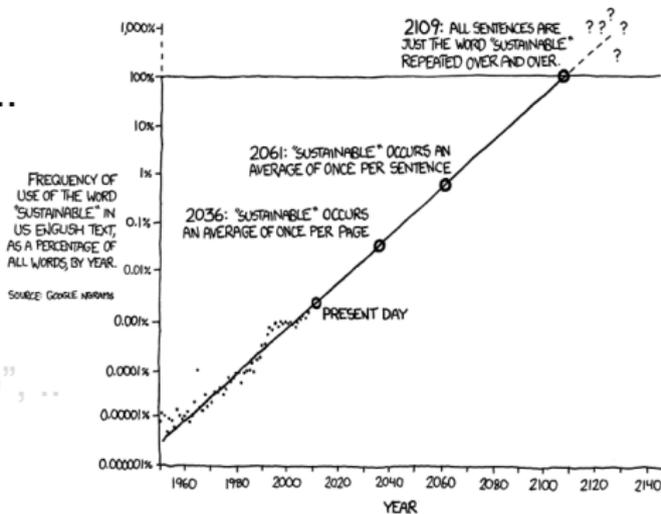
- Daten $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)\}$
- Datenpunkt $x \in \mathcal{X}$, Label (Klasse) $y \in \mathcal{Y}$
- Modelle \mathcal{M}

Modell "lernen", Datenanalyse, ..

$$f^* = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D})$$

Modellanwendung, "Vorhersage", ..

$$\hat{y} = f^*(x)$$



Allgemeines Vorgehen

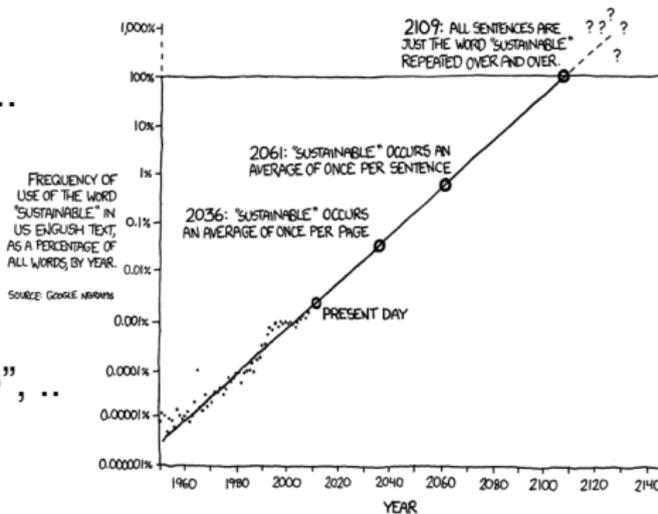
- Daten $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^N, y^N)\}$
- Datenpunkt $x \in \mathcal{X}$, Label (Klasse) $y \in \mathcal{Y}$
- Modelle \mathcal{M}

Modell "lernen", Datenanalyse, ..

$$f^* = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D})$$

Modellanwendung, "Vorhersage", ..

$$\hat{y} = f^*(x)$$



Modellklassen \mathcal{M}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- Linear

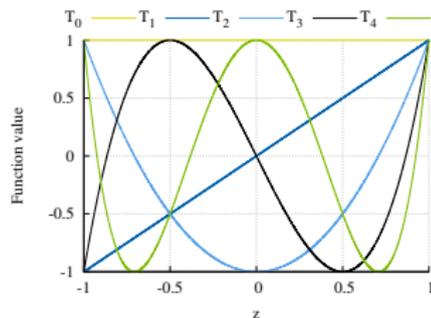
$$f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$$

- Polynom

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \beta_i x_i + \sum_{i=1}^n \sum_{j=1}^n \beta_{i,j} x_i x_j + \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \beta_{i,j,k} x_i x_j x_k + \dots$$

- Trigonometrisch/Periodisch

$$f(\mathbf{x}) = \sum_{i=0}^k \theta_i \cos(i \arccos(\langle \boldsymbol{\beta}^i, \mathbf{x} \rangle))$$



Modellklassen \mathcal{M}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- Linear

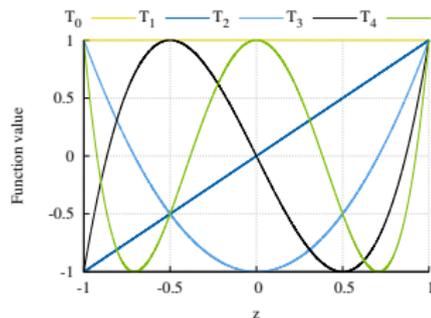
$$f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$$

- Polynom

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \beta_i \mathbf{x}_i + \sum_{i=1}^n \sum_{j=1}^n \beta_{i,j} \mathbf{x}_i \mathbf{x}_j + \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \beta_{i,j,k} \mathbf{x}_i \mathbf{x}_j \mathbf{x}_k + \dots$$

- Trigonometrisch/Periodisch

$$f(\mathbf{x}) = \sum_{i=0}^k \theta_i \cos(i \arccos(\langle \boldsymbol{\beta}^i, \mathbf{x} \rangle))$$



Modellklassen \mathcal{M}

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- Linear

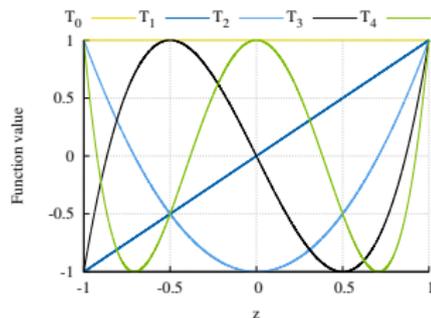
$$f(\mathbf{x}) = \beta_0 + \langle \boldsymbol{\beta}, \mathbf{x} \rangle$$

- Polynom

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \beta_i \mathbf{x}_i + \sum_{i=1}^n \sum_{j=1}^n \beta_{i,j} \mathbf{x}_i \mathbf{x}_j + \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \beta_{i,j,k} \mathbf{x}_i \mathbf{x}_j \mathbf{x}_k + \dots$$

- Trigonometrisch/Periodisch

$$f(\mathbf{x}) = \sum_{i=0}^k \theta_i \cos(i \arccos(\langle \boldsymbol{\beta}^i, \mathbf{x} \rangle))$$





Modellklassen \mathcal{M} (II)

- Probabilistisch/Exponentialfamilie/Bayesianisch

$$p_{\text{Gauss}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$p_{\text{Posterior}}(\mathbf{x}) = p_{\text{Likelihood}}(\mathbf{x}) p_{\text{Prior}_1}(\alpha) p_{\text{Prior}_2}(\gamma) \dots$$

- Unstetig/Piecewise/Thresholding

$$f_v(\mathbf{x}) = \begin{cases} f_{\text{Left}(v)}(\mathbf{x}) & g(\mathbf{x}) \geq \rho_v \\ f_{\text{Right}(v)}(\mathbf{x}) & g(\mathbf{x}) < \rho_v \end{cases}$$

Modellklassen \mathcal{M} (II)

- Probabilistisch/Exponentialfamilie/Bayesianisch

$$p_{\text{Gauss}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$p_{\text{Posterior}}(\mathbf{x}) = p_{\text{Likelihood}}(\mathbf{x}) p_{\text{Prior}_1}(\alpha) p_{\text{Prior}_2}(\gamma) \dots$$

- Unstetig/Piecewise/Thresholding

$$f_v(\mathbf{x}) = \begin{cases} f_{\text{Left}(v)}(\mathbf{x}) & g(\mathbf{x}) \geq \rho_v \\ f_{\text{Right}(v)}(\mathbf{x}) & g(\mathbf{x}) < \rho_v \end{cases}$$

Beispiel: Fehlerwahrscheinlichkeit

- $\mathcal{Y} = \{c_1, c_2, \dots, c_k\}$, Modell f
- Wahrscheinlichkeitsdichte falscher Vorhersagen:

$$\begin{aligned} p(f(\mathbf{X}) \neq \mathbf{Y}) &= \mathbb{E}_{\mathbf{X}, \mathbf{Y}}[\mathbb{1}_{\{f(\mathbf{X}) \neq \mathbf{Y}\}}] = \mathbb{E}_{\mathbf{X}} \mathbb{E}_{\mathbf{Y} | \mathbf{X}}[\mathbb{1}_{\{f(\mathbf{X}) \neq \mathbf{Y}\}}] \\ &= \mathbb{E}_{\mathbf{X}} \left[\sum_{i=1}^k \mathbb{1}_{\{f(\mathbf{X}) \neq c_i\}} \mathbb{P}(\mathbf{Y} = c_i | \mathbf{X}) \right] \end{aligned}$$

- Wahl von f , so dass innere Summe minimiert wird..

$$\begin{aligned} f(\mathbf{x}) &= \arg \min_{j=1}^k \sum_{i=1}^k \mathbb{1}_{\{c_j \neq c_i\}} \mathbb{P}(\mathbf{Y} = c_i | \mathbf{X} = \mathbf{x}) \\ &= \arg \min_{j=1}^k (1 - \mathbb{P}(\mathbf{Y} = c_j | \mathbf{X} = \mathbf{x})) \\ &= \arg \max_{j=1}^k \mathbb{P}(\mathbf{Y} = c_j | \mathbf{X} = \mathbf{x}) \end{aligned}$$



Beispiel: Fehlerwahrscheinlichkeit

- $\mathcal{Y} = \{c_1, c_2, \dots, c_k\}$, Modell f
- Wahrscheinlichkeitsdichte falscher Vorhersagen:

$$\begin{aligned} p(f(\mathbf{X}) \neq \mathbf{Y}) &= \mathbb{E}_{\mathbf{X}, \mathbf{Y}}[\mathbb{1}_{\{f(\mathbf{X}) \neq \mathbf{Y}\}}] = \mathbb{E}_{\mathbf{X}} \mathbb{E}_{\mathbf{Y} | \mathbf{X}}[\mathbb{1}_{\{f(\mathbf{X}) \neq \mathbf{Y}\}}] \\ &= \mathbb{E}_{\mathbf{X}} \left[\sum_{i=1}^k \mathbb{1}_{\{f(\mathbf{X}) \neq c_i\}} \mathbb{P}(\mathbf{Y} = c_i | \mathbf{X}) \right] \end{aligned}$$

- Wahl von f , so dass innere Summe minimiert wird..

$$\begin{aligned} f(\mathbf{x}) &= \arg \min_{j=1}^k \sum_{i=1}^k \mathbb{1}_{\{c_j \neq c_i\}} \mathbb{P}(\mathbf{Y} = c_i | \mathbf{X} = \mathbf{x}) \\ &= \arg \min_{j=1}^k (1 - \mathbb{P}(\mathbf{Y} = c_j | \mathbf{X} = \mathbf{x})) \\ &= \arg \max_{j=1}^k \mathbb{P}(\mathbf{Y} = c_j | \mathbf{X} = \mathbf{x}) \end{aligned}$$

Verlustfunktionen ℓ

- Absoluter Fehler/SSE/MSE/RMSE

$$\text{Err}(f; \mathcal{D}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}} |y - f(\mathbf{x})|$$

$$\text{RMSE}(f; \mathcal{D}) = \sqrt{\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} (y - f(\mathbf{x}))^2}$$

- Hinge Loss

$$\text{Hinge}(f; \mathcal{D}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}} \max\{0, 1 - yf(\mathbf{x})\}$$

Verlustfunktionen ℓ

- Absoluter Fehler/SSE/MSE/RMSE

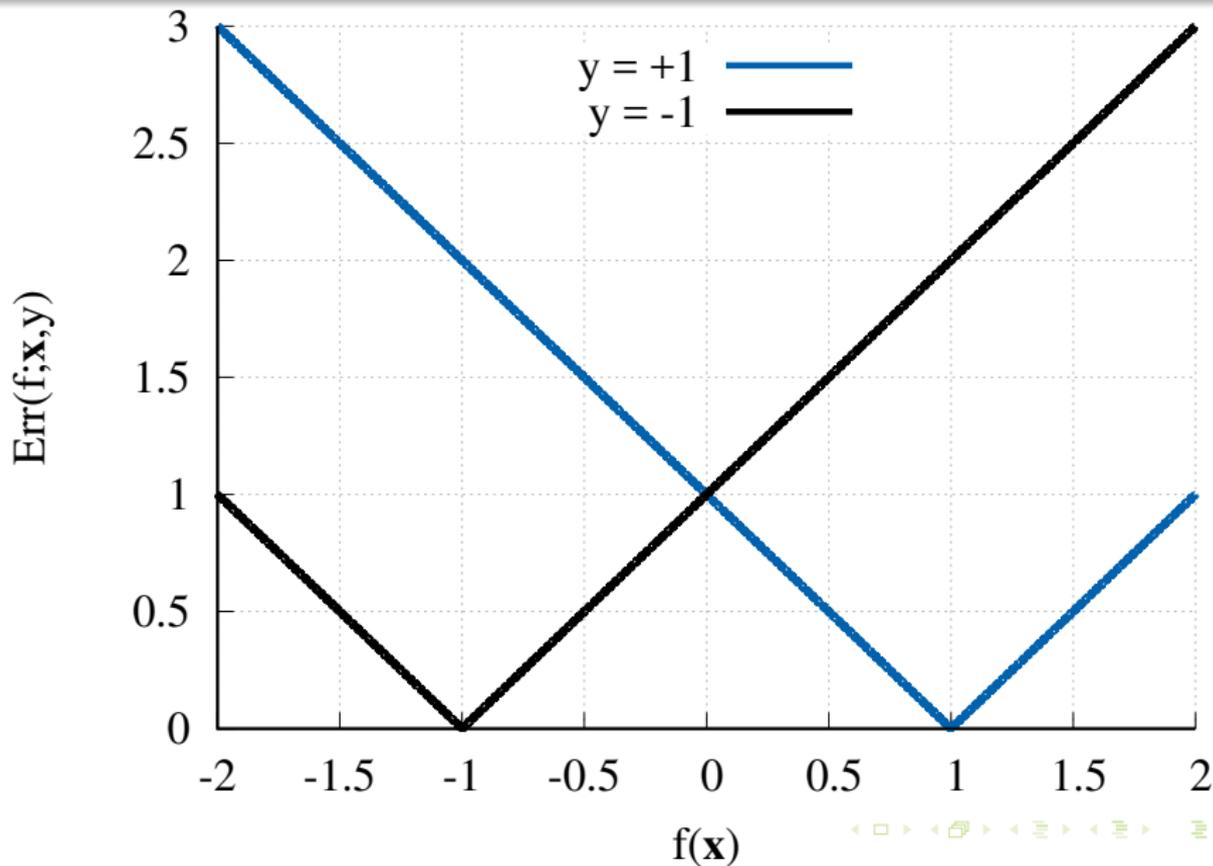
$$\text{Err}(f; \mathcal{D}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}} |y - f(\mathbf{x})|$$

$$\text{RMSE}(f; \mathcal{D}) = \sqrt{\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} (y - f(\mathbf{x}))^2}$$

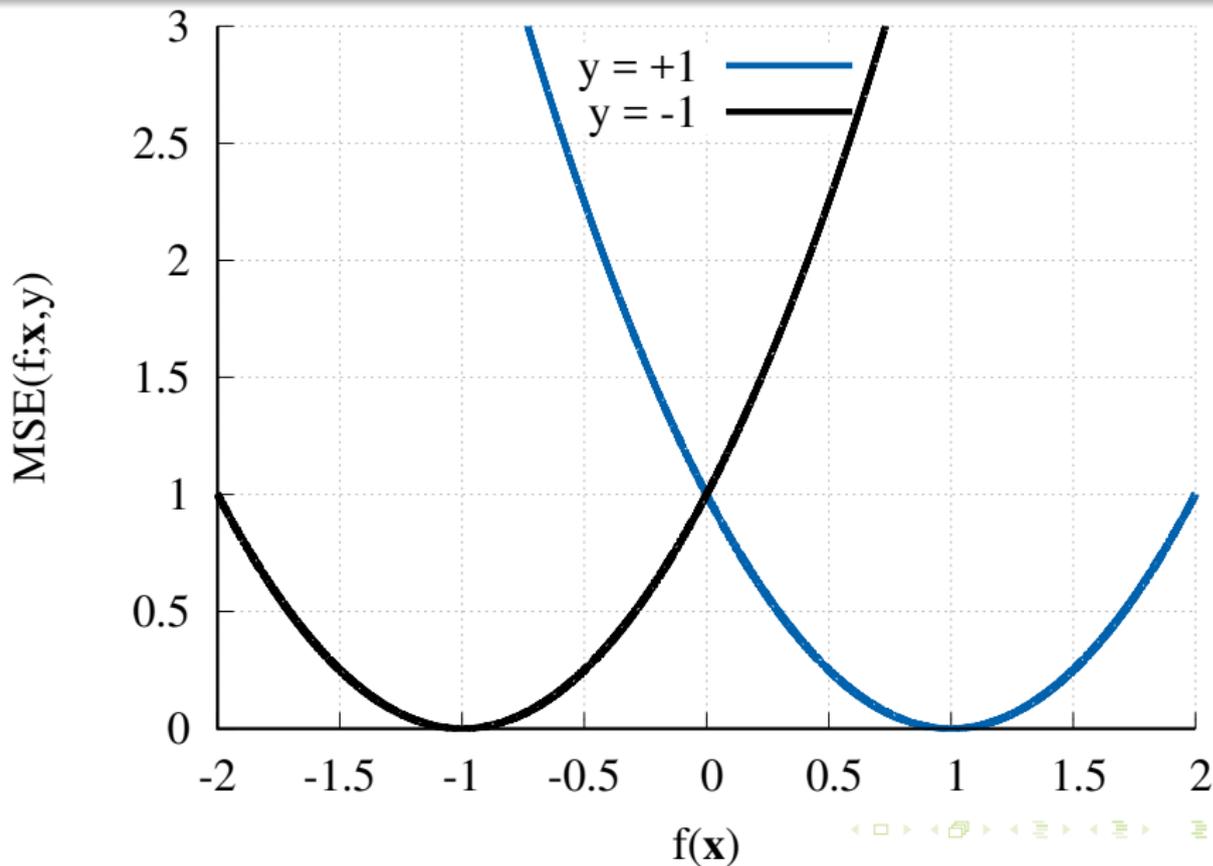
- Hinge Loss

$$\text{Hinge}(f; \mathcal{D}) = \sum_{(\mathbf{x}, y) \in \mathcal{D}} \max\{0, 1 - yf(\mathbf{x})\}$$

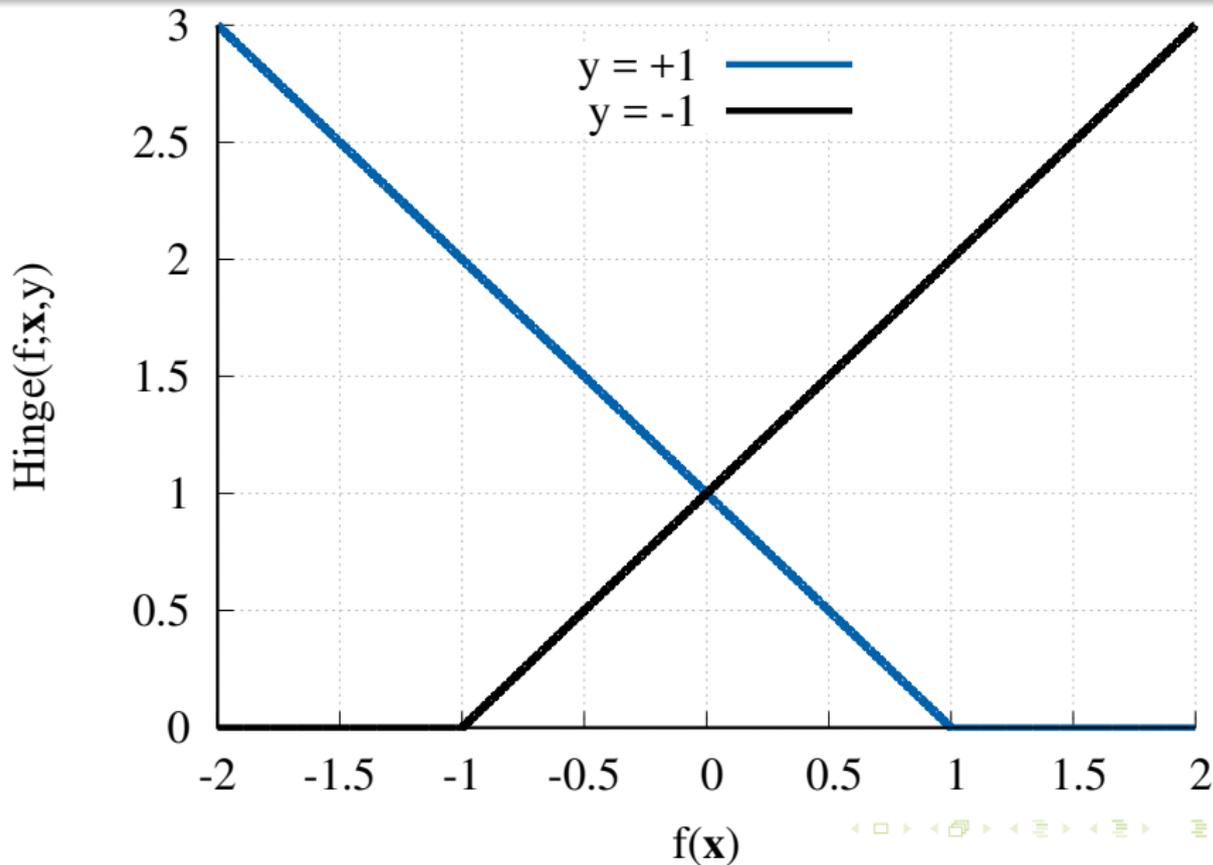
Absoluter Fehler, $|y - f(x)|$



Quadratischer Fehler, $(y - f(x))^2$



Hinge Fehler, $\max\{0, 1 - y - f(x)\}$



Verlustfunktionen ℓ (II) - Likelihood Varianten

- Likelihood

$$\mathcal{L}(p; \mathcal{D}) = \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{x}, \mathbf{y})$$

- Log-Likelihood

$$\log \mathcal{L}(p; \mathcal{D}) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$

- Neg. Avg. Log-Likelihood

$$\ell(p; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$



Verlustfunktionen ℓ (II) - Likelihood Varianten

- Likelihood

$$\mathcal{L}(p; \mathcal{D}) = \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{x}, \mathbf{y})$$

- Log-Likelihood

$$\log \mathcal{L}(p; \mathcal{D}) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$

- Neg. Avg. Log-Likelihood

$$\ell(p; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$



Verlustfunktionen ℓ (II) - Likelihood Varianten

- Likelihood

$$\mathcal{L}(p; \mathcal{D}) = \prod_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} p(\mathbf{x}, \mathbf{y})$$

- Log-Likelihood

$$\log \mathcal{L}(p; \mathcal{D}) = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$

- Neg. Avg. Log-Likelihood

$$\ell(p; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \log p(\mathbf{x}, \mathbf{y})$$

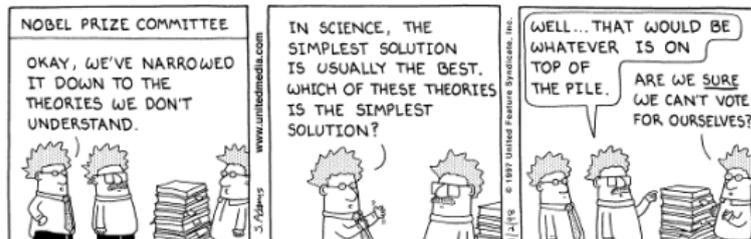


Verlustfunktionen ℓ (III) - MDL

- Minimum Description Length (MDL)
 - Formalisierung von Ockhams Rasiermesser

$$\min_{C \in \mathcal{M}} L(C) + L(\mathcal{D} | C)$$

- Intuition:
 - Nur wenige Objekte können kurze Codes haben
 - Nur wenige Objekte können hohe Wahrscheinlichkeit haben
- Formal: Alphabet $\mathcal{A} = \{1, 2, \dots, m\}$, Codierung C , Codelängen $L_C(1), L_C(2), \dots, L_C(m)$



$$\Leftrightarrow \sum_{a \in \mathcal{A}} 2^{-L_C(a)} \leq 1$$

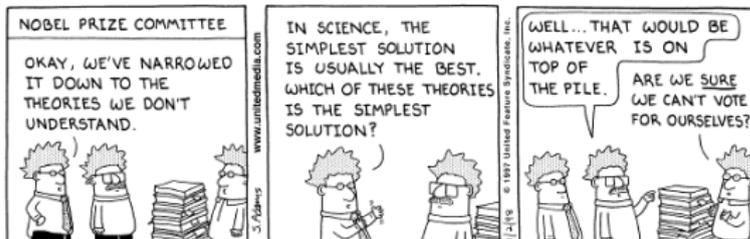
Kraft, 1949

Verlustfunktionen ℓ (III) - MDL

- Minimum Description Length (MDL)
 - Formalisierung von Ockhams Rasiermesser

$$\min_{C \in \mathcal{M}} L(C) + L(\mathcal{D} | C)$$

- Intuition:
 - Nur wenige Objekte können kurze Codes haben
 - Nur wenige Objekte können hohe Wahrscheinlichkeit haben
- Formal: Alphabet $\mathcal{A} = \{1, 2, \dots, m\}$, Codierung C , Codelängen $L_C(1), L_C(2), \dots, L_C(m)$



$$\Leftrightarrow \sum_{a \in \mathcal{A}} 2^{-L_C(a)} \leq 1$$

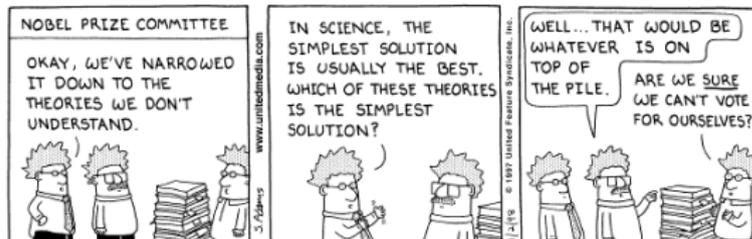
Kraft, 1949

Verlustfunktionen ℓ (III) - MDL

- Minimum Description Length (MDL)
 - Formalisierung von Ockhams Rasiermesser

$$\min_{C \in \mathcal{M}} L(C) + L(\mathcal{D} | C)$$

- Intuition:
 - Nur wenige Objekte können kurze Codes haben
 - Nur wenige Objekte können hohe Wahrscheinlichkeit haben
- Formal: Alphabet $\mathcal{A} = \{1, 2, \dots, m\}$, Codierung C , Codelängen $L_C(1), L_C(2), \dots, L_C(m)$



$$\Leftrightarrow \sum_{a \in \mathcal{A}} 2^{-L_C(a)} \leq 1$$

Kraft, 1949



Verlustfunktionen ℓ (IV) - Clustering

- Datensatz $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$ ohne Label
- Modell $\subset \mathbb{R}^n$
- k -Means (Intra-Cluster Varianz)

$$\text{ICV}(k; \mathcal{D}) = \min_{C \subset \mathbb{R}^n, |C|=k} \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|_2^2$$

- Mixture Modelle (Expectation Maximization)
 - Example: Gaussian

$$\min_{\{(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\}_{1 \leq i \leq k} \subset \mathbb{R}^n \times S_{++}^n} - \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_{i=1}^k p_{\text{Gauss}}(\mathbf{x} | i) p(i)$$



Verlustfunktionen ℓ (IV) - Clustering

- Datensatz $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^N\}$ ohne Label
- Modell $\subset \mathbb{R}^n$
- k -Means (Intra-Cluster Varianz)

$$\text{ICV}(k; \mathcal{D}) = \min_{C \subset \mathbb{R}^n, |C|=k} \sum_{\mathbf{x} \in \mathcal{D}} \min_{\mathbf{c} \in C} \|\mathbf{x} - \mathbf{c}\|_2^2$$

- Mixture Modelle (Expectation Maximization)
 - Example: Gaussian

$$\min_{\{(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\}_{1 \leq i \leq k} \subset \mathbb{R}^n \times S_{++}^n} - \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_{i=1}^k p_{\text{Gauss}}(\mathbf{x} \mid i) p(i)$$