

Vorlesung 13 Linear Regression

Von Lehrer VL:

Gegeben sind Daten $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}^N$

Wir haben zwei Arten von Modellen

1. Deterministisches Modell:

$$f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^N, \quad \theta \in \mathbb{R}^P.$$

2. Statistisches Modell:

$$P_\theta(y|x)$$

mit $\theta \in \mathbb{R}^P$ und P_θ die Wkt.-Dichte von y gegeben x .

Strategien zur Parameterbestimmung waren:

1. ERM: (für det. Modell).

$$\theta \in \operatorname{argmin}_{\theta \in \mathbb{R}^P} R(\theta), \quad \text{wobei: } R(\theta) = \sum_{i=1}^n \ell(y_i, f_\theta(x_i))$$

Loss-Funktion.

2. MLE: (für stat. Modell)

$$\theta \in \operatorname{argmax}_{\theta \in \mathbb{R}^P} \prod_{i=1}^n P_\theta(y_i | x_i)$$

3. MAP (für stat. Modell)

$$\theta \in \operatorname{argmax} \alpha(\theta), \quad \alpha(\theta) = P(\theta | X, Y),$$

wobei θ selbst als zufällig angesehen.

$$X = [x_1, \dots, x_n]$$

$$Y = [y_1, \dots, y_n]$$

4. Bayes'sche Ansatz:

Die Wkt.-Dichte $P(\theta | X, Y)$ berechnen und θ

samplen oder $P(y|x, X, Y) = \int_{\Theta} P_{\theta}(y|x) P(\theta|X, Y) d\theta$
berechnen

Kurz zurück zur Modellauswahl: Den Bayes'schen Ansatz Parameter selbst als zufällig zu modellieren können wir auch hier benutzen.

Angenommen wir haben Modelle M_1, \dots, M_k zur Auswahl.

Wir setzen den Prior $P(M)$ auf die Menge der Modell (d.h. $P(M)$ bestimmt eine Wkt-Verteilung auf $\{M_1, \dots, M_k\}$). Z.B.: $P(M=M_i) = \frac{1}{k}$.

Dann haben wir die Posterior-Funktion $P(M|X, Y)$

Es gilt nach Bayes' Theorem:

$$\operatorname{argmax}_i P(M=M_i | X, Y) = \operatorname{argmax}_i P(M_i) P(X, Y | M_i)$$

$$\text{und } P(X, Y | M_i) = \int_{\Theta} P(X, Y | \theta) P(\theta | M_i)$$

└ Prior für Modell M_i

Lineare Regression

Definition 13.1

Das Lineare Modell $f_{\theta} : \mathbb{R}^D \rightarrow \mathbb{R}$ ist gegeben durch

$$\theta = (\theta_0, \theta_1, \dots, \theta_D)^T \in \mathbb{R}^{D+1}$$

mit

$$f_{\theta}(x) = x^T \theta' + \theta_0,$$

wobei $\theta' = (\theta_1, \dots, \theta_D)^T \in \mathbb{R}^D$.

Der quadratische Loss ist

$$\ell(y, \hat{y}) = (y - \hat{y})^2.$$

Das zugehörige ERM heißt das Problem der kleinsten Quadrate

bzw. least-squares.

Gegeben seien wieder Daten $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \mathbb{R}$.

Wir setzen wieder

$$X = [x_1, \dots, x_n]^T \in \mathbb{R}^{n \times D}, \quad Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$$

Wir setzen außerdem

$$\Omega = \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^{n \times (D+1)}$$

Ω heißt Feature-matrix.

Theorem 13.2

Sei $r(\Omega)$ der Rang von Ω . Dann gilt im obigen Setting:

1. Wenn $r(\Omega) < D+1$, hat $\arg\min_{\theta} R(\theta)$ unendlich viele Lösungen.
2. Wenn $r(\Omega) = D+1$, hat $\arg\min_{\theta} R(\theta)$ eine eindeutige Lösung:

$$\theta_{\text{ERM}} = \Omega^+ Y,$$

wobei Ω^+ die Pseudoinverse von Ω ist.

Beweis

In unserem Fall ist das empirische Risiko:

$$\begin{aligned} R(\theta) &= \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) \\ &= \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 \\ &= \sum_{i=1}^n (y_i - x_i^T \theta + \theta_0)^2 \end{aligned}$$

$$= \left\| \begin{pmatrix} y_1 - x_1^T \theta' + \theta_0 \\ \vdots \\ y_n - x_n^T \theta' + \theta_0 \end{pmatrix} \right\|^2$$

$$= \| Y - \Omega \cdot \theta \|^2.$$

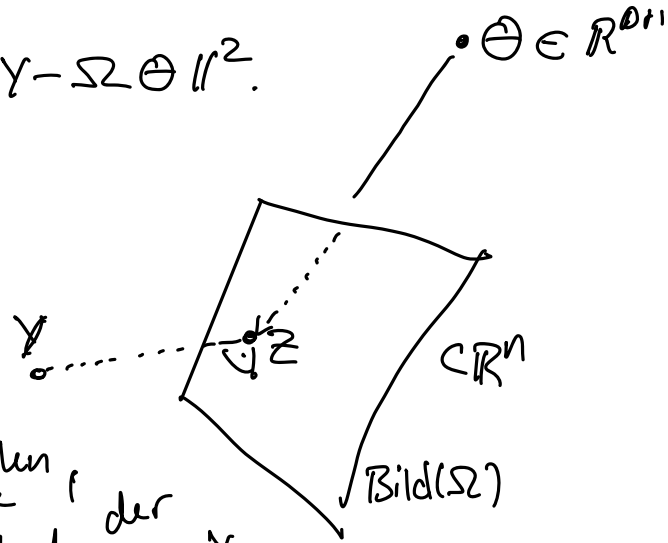
D.h. $\theta_{\text{ERM}} = \operatorname{argmin}_{\theta \in \mathbb{R}^{D+1}} \| Y - \Omega \theta \|^2$.

Sei $z := \Omega \Omega^+ Y$

Dann wissen wir

$$z = \operatorname{argmin}_{w \in \operatorname{Bild}(\Omega)} \| Y - w \|^2.$$

θ muss jetzt die Gleichung $z = \Omega \theta$ erfüllen,
weil $\| Y - z \|^2$ der minimale Abstand von Y nach $\operatorname{Bild}(\Omega)$ ist.



Falls $r(\Omega) < D+1$, hat $z = \Omega \theta$ unendl. viele Lösungen,
weil $\ker(\Omega) \neq \emptyset$.

Falls $r(\Omega) = D+1$, haben wir gezeigt, dass $\Omega^+ \Omega = I_{D+1}$.

$\Rightarrow \theta = \Omega^+ Y$ eindeutig definiert ist.
 $= \Omega^+ z$

Der kleinste-Quadrat-Schätzer führt oft zu Overfitting.

Um das zu verhindern wird oft ein Regularisierungs-Parameter λ eingeführt. Dieser definiert die regularisierte Loss-Funktion

$$\ell(y, \hat{y}) = (y - \hat{y})^2 + \lambda \cdot \| A \theta \|^2 \quad \text{mit } A \in \mathbb{R}^{(D+1) \times (D+1)}.$$

Der Spezialfall $A = I_{D+1}$ heißt Tikhonov-Regularisierung und das zugehörige Regressionsproblem heißt Ridge Regression.

Theorem 13.3

Sei $\Sigma \in \mathbb{R}^{n \times (D+1)}$ die Feature Matrix und $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$.

Für fast alle (= alle außer endlich viele) Werte von λ gibt es eine eindeutige Lösung

$$\Theta_{RR} = \arg \min_{\Theta} R(\Theta).$$

mit

$$\Theta_{RR} = (\Sigma^T \Sigma + n\lambda A^T A)^{-1} \Sigma^T y.$$

Bemerkung Thm 13.3. kann auch verwendet werden, wenn $r(\Sigma) < D+1$, d.h. z.B. falls $n < D+1$ (weniger Datenpunkte als Parameter).

Beweis von Thm 13.3.

Das empirische Risiko ist:

$$\begin{aligned} R(\Theta) &= \sum_{i=1}^n \ell(y_i, f_{\Theta}(x_i)) \\ &= \sum_{i=1}^n (y_i - f_{\Theta}(x_i))^2 + \lambda \|A\Theta\|^2 \\ &= \|y - \Sigma\Theta\|^2 + n\lambda \|A\Theta\|^2. \end{aligned}$$

Θ_{RR} minimiert $R(\Theta)$. Um Θ_{RR} zu finden sehen wir

$$\frac{d}{d\Theta} R(\Theta) = 0.$$

$$\text{D.h. } 0 = \frac{d}{d\Theta} \left((y - \Sigma\Theta)^T (y - \Sigma\Theta) + n\lambda \Theta^T A^T A \Theta \right)$$

$$= \sum \Omega^T (Y - \Omega \Theta) + 2n\lambda \bar{A}^T \bar{A} \Theta$$

$$\rightarrow \underbrace{(\Omega^T \Omega - n\lambda A^T A)}_{\in \mathbb{R}^{(D+1) \times (D+1)}} \Theta = \Omega^T Y.$$

Es gilt: $\Omega^T \Omega - n\lambda A^T A$ ist invertierbar, falls

$$f(\lambda) = \det(\Omega^T \Omega - n\lambda A^T A) \neq 0.$$

Aber $f(\lambda)$ hat höchstens $D+1$ Nullstellen, weil es ein Polynom vom Grad $D+1$ ist.

\leadsto Bis auf endlich viele Werte von λ ist $\Omega^T \Omega - n\lambda A^T A$ invertierbar und es gilt:

$$\Theta_{RR} = (\Omega^T \Omega - n\lambda A^T A)^{-1} \Omega^T Y.$$

Jetzt: Statistisches Modell. Wir wählen folgende Verteilung für y gegeben x :

$$P_{\Theta}(y|x) = \mathcal{I}(y | f_{\Theta}(x), \sigma^2).$$

$$\text{D.h. } (y|x) = f_{\Theta}(x) + \varepsilon, \text{ mit } \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

(wieder ist σ^2 ein Parameter, den wir nicht optimieren, sondern vorher festlegen, d.h. σ^2 ist ein hyperparameter).

Wir wählen wieder die ^{lineare} Funktion $f_{\Theta}(x) = x^T \Theta' + \Theta_0$

Theorem 13.4

Sei wieder $\Omega \in \mathbb{R}^{n \times (D+1)}$ die Feature matrix und $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$.
Der Maximum-Likelihood Schätzer ist

1. nicht eindeutig, falls $r(\Omega) < D+1$.
2. eindeutig bestimmt als

$$\Theta_{MLE} = \Omega^+ Y$$

falls $r(\Omega) = D+1$.

Beweis

Es ist $\Theta_{MLE} \in \operatorname{argmax} L(\Theta)$ mit $L(\Theta) = \prod_{i=1}^n P_{\Theta}(y_i | x_i)$.

Wir betrachten

$$\ell(\Theta) = \log L(\Theta) = \sum_{i=1}^n \log P_{\Theta}(y_i | x_i).$$

Anstatt $L(\Theta)$ können wir $\ell(\Theta)$ maximieren.

$$\begin{aligned} \text{Es ist } \ell(\Theta) &= \sum_{i=1}^n \log \Phi(y_i | f_{\Theta}(x_i), \sigma^2) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - f_{\Theta}(x_i))^2\right) \\ &= n \log \frac{1}{\sqrt{2\pi\sigma^2}} + \sum_{i=1}^n \left(-\frac{1}{2\sigma^2}(y_i - f_{\Theta}(x_i))^2\right) \\ &= n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f_{\Theta}(x_i))^2 \\ &= n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \|\gamma - \Omega\Theta\|^2. \end{aligned}$$

$\leadsto \Theta_{MLE} = \operatorname{argmin}_{\Theta} \|\gamma - \Omega\Theta\|^2$. Die Aussage folgt nun wie im Beweis von Thm. 13.2.

□