

## Vorlesung 14 Lineare & nicht-lineare Regression

Lineare Regression: wir haben das deterministische Modell

$$y = f_{\Theta}(x) = x^T \Theta' + \Theta_0, \quad x \in \mathbb{R}^D, y \in \mathbb{R}$$

wobei  $\Theta' = (\Theta_1, \dots, \Theta_D)^T \in \mathbb{R}^D$ ,  $\Theta_0 \in \mathbb{R}$ .

Wir haben das stat. Modell

$$y \sim N(f_{\Theta}(x), \sigma^2) \quad \text{mit } \sigma^2 \text{ fest} \quad (= \text{ein Hyperparameter})$$

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

Der MAP ist  $\Theta_{MAP}$  mit  $\Theta_{MAP} = \operatorname{argmax}_{\Theta} P(\Theta | X, Y)$ .

Wir wählen dazu den Prior  $\Theta \sim N(\mu, \Sigma)$ ,  $\mu \in \mathbb{R}^{D+1}$ ,  $\Sigma \in \mathbb{R}_{\text{pos. def.}}^{(D+1) \times (D+1)}$

### Theorem 14.1.

Sei  $\mu \in \mathbb{R}^{D+1}$ ,  $\Sigma \in \mathbb{R}^{(D+1) \times (D+1)}$  pos. definit. Für den Prior  $\Theta \sim N(\mu, \Sigma)$  haben wir:

$$\Theta_{MAP} = (\Sigma^T \Sigma + \sigma^2 \Sigma^{-1})^{-1} (\Sigma^T Y + \sigma^2 \Sigma^{-1} \mu).$$

mit

$$\Sigma = \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix}^T \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}. \quad \text{und unter der Annahme dass}$$

$\Sigma^T \Sigma + \sigma^2 \Sigma^{-1}$  invertierbar ist.

### Beweis

Sei  $\alpha(\Theta) = P(\Theta | X, Y)$  die Posteriori-Funktion.

$\Theta_{MAP} = \operatorname{argmax}_{\Theta} \alpha(\Theta)$ . Es gilt

$$\alpha(\Theta) = P(\Theta | X, Y) = P(\Theta) \frac{P(Y | X, \Theta)}{P(Y | X)}.$$

nach dem Satz von Bayes. Sei  $C := -\log P(Y | X)$ .

Dann gilt:

$$\log \alpha(\theta) = \log P(\theta) + \log P(Y|X, \theta) + c.$$

und  $P(\theta) = \frac{1}{\sqrt{(2\pi)^{D+1} \det(\Sigma)}} \exp(-\frac{1}{2} (\theta - \mu)^T \Sigma^{-1} (\theta - \mu))$

$$P(Y|X, \theta) = \prod_{i=1}^n P(y_i|x_i, \theta) = \prod_{i=1}^n P_\theta(y_i|x_i)$$

$$= \frac{1}{\sqrt{2\sigma^2 n}} \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y_i - f_\theta(x_i))^2\right)$$

$$= \frac{1}{\sqrt{2\sigma^2 n}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f_\theta(x_i))^2\right)$$

$$= \frac{1}{\sqrt{2\sigma^2 n}} \exp\left(-\frac{1}{2\sigma^2} \|Y - \Sigma\theta\|^2\right)$$

Wir erhalten:

$$\log \alpha(\theta) = -\frac{1}{2} (\theta - \mu)^T \Sigma^{-1} (\theta - \mu) - \frac{1}{2\sigma^2} \|Y - \Sigma\theta\|^2 + c'$$

mit  $c'$  unabh. von  $\theta$ .

Um  $\alpha$  zu maximieren sehen wir  $\nabla_\theta \log \alpha(\theta) = 0$ .

$$\nabla_\theta \log \alpha(\theta) = -\Sigma^{-1}(\theta - \mu) + \frac{1}{\sigma^2} \Sigma^T(Y - \Sigma\theta)$$

Mit  $\nabla_\theta \log \alpha(\theta) = 0$  gilt:

$$\Sigma^{-1}(\theta - \mu) = +\frac{1}{\sigma^2} \Sigma^T(Y - \Sigma\theta).$$

$$\Rightarrow \Sigma^{-1}\theta - \Sigma^{-1}\mu = +\frac{1}{\sigma^2} \Sigma^T Y - \frac{1}{\sigma^2} \Sigma^T \Sigma \theta$$

$$\Rightarrow \Sigma^{-1}\mu + \frac{1}{\sigma^2} \Sigma^T Y = \left(\frac{1}{\sigma^2} \Sigma^T \Sigma + \Sigma^{-1}\right)\theta$$

$$\Rightarrow \theta = (\Sigma^T \Sigma + \sigma^2 \Sigma^{-1})^{-1} (\Sigma^T Y + \sigma^2 \Sigma^{-1} \mu). \quad \square$$

In den Beweisen zu ERM, RR, MLE und MAP ist das Modell  $f_\theta$  in Form der Zeilen von  $\Sigma$  eingegangen.

D.h., wir können einfach das Linear Model zum nicht-linearen Modell erweitern, indem wir folgendes Modell definieren:

$$f_\theta(x) = \phi(x)^T \theta \quad \text{mit } \phi: \mathbb{R}^D \rightarrow \mathbb{R}^P \text{ nicht-linear.}$$

Die Feature-Matrix ist dann

$$\Sigma = [\phi(x_1) \dots \phi(x_n)]^T.$$

Z.B.: mit  $\phi(x) = \begin{pmatrix} 1 \\ x \end{pmatrix} \in \mathbb{R}^{D+1}$  erhalten wir das Lineare Modell

- für  $D=1$  und  $\phi(x) = (1, x, x^2, \dots, x^d)^T \in \mathbb{R}^{d+1}$  erhalten wir Polynomiale - Regression, da  $f_\theta(x) = \sum_{i=0}^d \theta_i x^i$ .

Die Formeln für  $\theta_{\text{ERM}}$ ,  $\theta_{\text{RR}}$ ,  $\theta_{\text{MLE}}$  und  $\theta_{\text{MAP}}$  gelten genauso im nicht-linearen Modell.

Jetzt wollen wir den Bayes'schen Ansatz wählen und die Verteilung  $P(\theta | X, Y)$  direkt berechnen.

### Theorem 14.2

Mit obigen Annahmen gilt:

$$(\theta | X, Y) \sim N(\mu, S),$$

mit

$$S = (\sigma^{-2} \Sigma^T \Sigma + \Sigma^{-1})^{-1} \quad \text{und} \quad \mu = S \left( \frac{1}{\sigma^2} \Sigma^T Y + \Sigma^{-1} \mu \right)$$

Beweis Wir berechnen die Dichte  $P(\Theta | X, Y)$ .

Dazu berechnen wir  $\log P(\Theta | X, Y)$ .

$$\log P(\Theta | X, Y) = -\frac{1}{2} (\Theta - \mu)^T \Sigma^{-1} (\Theta - \mu) - \frac{1}{2\sigma^2} \|Y - \Sigma \Theta\|^2 + c$$

mit  $c$  unabh. von  $\Theta$ .

Wir setzen

$$Q = -(\Theta - \mu)^T \Sigma^{-1} (\Theta - \mu) - \frac{1}{\sigma^2} (Y - \Sigma \Theta)^T (Y - \Sigma \Theta)$$

so dass  $P(\Theta | X, Y) = \exp(\frac{1}{2} Q + c)$

Es gilt:

$$\begin{aligned} Q &= -\Theta^T \Sigma^{-1} \Theta + 2\mu^T \Sigma^{-1} \Theta - \mu^T \Sigma^{-1} \mu \\ &\quad - \frac{1}{\sigma^2} (Y^T Y - 2Y^T \Sigma \Theta + \Theta^T \Sigma^T \Sigma \Theta) \end{aligned}$$

$$= -\Theta^T \Sigma^{-1} \Theta + 2\mu^T \Sigma^{-1} \Theta + \frac{1}{\sigma^2} 2Y^T \Sigma \Theta - \frac{1}{\sigma^2} \Theta^T \Sigma^T \Sigma \Theta + c'$$

mit  $c'$  unabh. von  $\Theta$ . Wir setzen

$$A_1 = \frac{1}{\sigma^2} \Sigma^T \Sigma + \Sigma^{-1} = S.$$

$$\alpha = \frac{1}{\sigma^2} \Sigma^T Y + \Sigma^{-1} \mu$$

Damit:

$$\begin{aligned} Q &= \Theta^T A \Theta - 2\alpha^T \Theta + c' \\ &= (\Theta - b)^T A (\Theta - b) + c'' \end{aligned}$$

$$m = b = A^{-1} \alpha = S^{-1} \alpha = S^{-1} \left( \frac{1}{\sigma^2} \Sigma^T Y + \Sigma^{-1} \mu \right).$$

□

Wir können den Bayes'schen Ansatz verwenden um

- Den Parameter  $\Theta$  zu samplen, oder
- Den Parameter "auszuinkravieren".

Proposition 14.3

↓ Trainingsdaten

$$X = [x_1, \dots, x_n]^T \quad Y = [y_1, \dots, y_n]^T$$

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

Sei  $(x, y) \in \mathbb{R}^D \times \mathbb{R}$  ein weiterer Punkt. Die Verteilung von  $(y|x, X, Y)$  ist.

$$(y|x, X, Y) \sim N(\phi(x)^T \theta_m, \phi(x)^T \Sigma \phi(x) + \sigma^2)$$

(für das nicht-lineare Modell bzgl.  $\phi(x)$ ).

Beweisidee Verwende die Formel für Marginal-Density:

$$P(y|x, X, Y) = \int_{\mathbb{R}^P} P(y|x, X, Y, \theta) \cdot P(\theta|X, Y) d\theta$$

Die bisherige Diskussion handelt von (nicht)linearen Funktionen  $f_\theta: \mathbb{R}^D \rightarrow \mathbb{R}$ .

Wir können die Theorie zu Funktion

$$f_\theta: \mathbb{R}^D \rightarrow \mathbb{R}^N$$

erweitern, indem wir

$$(x) \quad f_\theta(x) = \begin{bmatrix} \phi(x)^T \theta_1 \\ \vdots \\ \phi(x)^T \theta_N \end{bmatrix} \quad \text{mit } \theta = [\theta_1, \dots, \theta_N] \in \mathbb{R}^{P \times N}$$

$$\phi: \mathbb{R}^D \rightarrow \mathbb{R}^P$$

setzen und die Loss Funktion

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

bzw das stat. Modell

$$(y|x) \sim N(f_\theta(x), \sigma^2 I_N).$$

Weil  $\|y - \hat{y}\|^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2$ , können wir jedes  $\Theta_i, i=1, \dots, N$  separat optimieren. D.h. für jeden Eintrag von  $f_\theta$  in  $\Theta$  können wir die 1-dim. Schicht von oben verwenden.