

# Maschinelles Lernen

## Multivariate Lineare Regression mit Gradientenabstieg

Prof. Dr. Rainer Stollhoff

# Übersicht

- Motivation
- Multivariater Gradientenabstieg

# Supervised Learning

## Supervised Learning

### 1. Aufgabe A

Vorhersage  $\hat{Y} = A(X)$

### 2. Qualität Q

Verlustfunktion  $L(\hat{Y}, Y)$

### 3. Erfahrung E

Datensatz

$(x_i, y_i)$  für  $i = 1, \dots, n$

Eine Maschine **lernt** aus Erfahrung E eine Aufgabe A mit der Qualität Q, wenn die Qualität Q beim erfüllen der Aufgabe A mit Erfahrung E steigt (T. Mitchell, MIT, 1988)

# Einfache univariate Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(x_i, y_i)_{i=1}^n$

**Qualität:** Quadratische Verlustfunktion

$$L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = L(\theta)$$

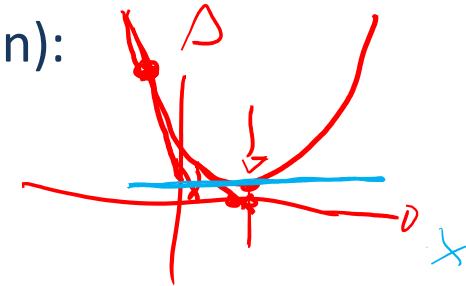
**Maschine:** vereinfachte lineare Regression mit

$$f(x; \theta) = \theta \cdot x$$

**Lernen:** Finde einen Wert für  $\theta$ , der die quadratische Verlustfunktion minimiert

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = 0$
2. Berechne Ableitung  $\frac{d}{d\theta} L(\theta) = \frac{d}{d\theta} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = \frac{d}{d\theta} \sum_{i=1}^n (y_i - x_i \cdot \theta)^2$ 
$$= \sum_{i=1}^n 2 \cdot (y_i - x_i \cdot \theta) \cdot (-x_i)$$
3. Update  $\theta^{t+1} = \theta^t + \alpha \cdot \frac{d}{d\theta} L(\theta^t) = \theta^t - \alpha \cdot 2 \sum_{i=1}^n (y_i - x_i \cdot \theta^t) \cdot x_i = \theta^t + 2 \sum_{i=1}^n (y_i - x_i \cdot \theta^t) \cdot x_i$



mit  $\alpha = -1$  (steilster Abstieg)

# Multivariate Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(x_i, y_i)_{i=1}^n$

**Qualität:** Quadratische Verlustfunktion

$$L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \underbrace{f(x_i; \theta)})^2 = L(\theta)$$

**Maschine:** vereinfachte lineare Regression mit

$$f(x_1, x_2; \theta_0, \theta_1, \theta_2) = \overline{\theta_0} + \underline{\theta_1} \cdot \underline{x_1} + \overline{\theta_2} \cdot \overline{x_2}$$

**Lernen:** Finde Werte für  $\theta = (\theta_0, \theta_1, \theta_2)$ , die die quadratische Verlustfunktion minimieren

# Multivariate Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(x_i, y_i)_{i=1}^n$

**Qualität:** Quadratische Verlustfunktion

$$L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = \underline{L(\theta)}$$

**Maschine:** vereinfachte lineare Regression mit

$$f(x_1, x_2; \theta_0, \theta_1, \theta_2) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2$$

**Lernen:** Finde Werte für  $\theta = (\theta_0, \theta_1, \theta_2)$ , die die quadratische Verlustfunktion minimieren

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = 0$

2. Berechne Ableitung  $\frac{d}{d\theta} L(\theta) = \frac{d}{d\theta} \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = \frac{d}{d\theta} \sum_{i=1}^n (y_i - x_i \cdot \theta)^2$   
 $= \sum_{i=1}^n 2 \cdot (y_i - x_i \cdot \theta) \cdot (-x_i)$

3. Update  $\theta^{t+1} = \theta^t + \alpha \cdot \frac{d}{d\theta} L(\theta^t) = \theta^k - \alpha \cdot 2 \sum_{i=1}^n (y_i - x_i \cdot \theta^t) \cdot x_i = \theta^k + 2 \sum_{i=1}^n (y_i - x_i \cdot \theta^t) \cdot x_i$

mit  $\alpha = -1$  (steilster Abstieg)

# Maschinelles Lernen

## Multivariate Analysis - Einführung

**Prof. Dr. Rainer Stollhoff**

# Univariate Lineare Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(x_i, y_i)_{i=1}^n$

**Qualität:** Verlustfunktion:  $L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = L(\theta)$

**Maschine:** Regression mit  $f(x; \theta) = \underline{\theta} \cdot \underline{x}$

**Lernen:** Finde Werte für  $\theta = (\underline{\theta})$ , die die quadratische Verlustfunktion minimieren

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = (\underline{1})$

2. Berechne Gradienten

$$\underline{\nabla L(\theta^0)} = \left( \frac{\partial}{\partial \underline{\theta}} L(\theta^0) \right) = \underline{\left( \sum_{i=1}^n (y_i - (\underline{\theta} \cdot \underline{x}_i)) \cdot (-2 \cdot \underline{x}_i) \right)}$$

3. Update  $\theta^{t+1} = \theta^t + \underline{\alpha} \cdot \nabla L(\theta^t)$



# Bivariate Lineare Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(\underline{x}_i, y_i)_{i=1}^n$  mit  $\underline{x}_i = (x_{i,1}, x_{i,2})$

**Qualität:** Verlustfunktion:  $L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = L(\theta)$

**Maschine:** Regression mit  $f(x_1, x_2; \theta_0, \theta_1, \theta_2) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2$

**Lernen:** Finde Werte für  $\theta = (\theta_0, \theta_1, \theta_2)$ , die die quadratische Verlustfunktion minimieren

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$

2. Berechne Gradienten

$$\underline{\nabla L(\theta^0)} = \begin{pmatrix} \frac{\partial}{\partial \theta_0} L(\theta^0) \\ \frac{\partial}{\partial \theta_1} L(\theta^0) \\ \frac{\partial}{\partial \theta_2} L(\theta^0) \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n \left[ (y_i - \overbrace{(\theta_0^0 + \theta_1^0 \cdot x_{i,1} + \theta_2^0 \cdot x_{i,2})}^{f(x_i, \theta)}) \cdot (-2 \cdot 1) \right] \\ \sum_{i=1}^n \left[ (y_i - (\theta_0^0 + \theta_1^0 \cdot x_{i,1} + \theta_2^0 \cdot x_{i,2})) \cdot (-2 \cdot x_{i,1}) \right] \\ \sum_{i=1}^n \left[ (y_i - (\theta_0^0 + \theta_1^0 \cdot x_{i,1} + \theta_2^0 \cdot x_{i,2})) \cdot (-2 \cdot x_{i,2}) \right] \end{pmatrix}$$

3. Update  $\theta^{t+1} = \theta^t + \alpha \cdot \nabla L(\theta^t)$

gleiche Schrittweite  
 $\alpha \leadsto A$  spezifische Schrittweite

# Multivariate Lineare Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(\mathbf{x}_i, y_i)_{i=1}^n$  mit  $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, \mathbf{x}_{i,n})$

**Qualität:** Verlustfunktion:  $L(y, \hat{y}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 = L(\theta)$

**Maschine:** Regression mit  $f(\mathbf{x}; \boldsymbol{\theta}) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \boldsymbol{\theta}_n \cdot \mathbf{x}_n$

**Lernen:** Finde Werte für  $\theta = (\theta_0, \theta_1, \theta_2, \dots, \boldsymbol{\theta}_n)$ , die die quadratische Verlustfunktion minimieren

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ \mathbf{1} \end{pmatrix}$

2. Berechne Gradienten

$$\nabla L(\theta^0) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} L(\theta^0) \\ \frac{\partial}{\partial \theta_1} L(\theta^0) \\ \vdots \\ \frac{\partial}{\partial \theta_n} L(\theta^0) \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \boldsymbol{\theta}_n \cdot \mathbf{x}_n)) \cdot (-2) \cdot 1 \\ \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \boldsymbol{\theta}_n \cdot \mathbf{x}_n)) \cdot (-2 \cdot x_{i,1}) \\ \vdots \\ \sum_{i=1}^n (y_i - (\theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \boldsymbol{\theta}_n \cdot \mathbf{x}_n)) \cdot (-2 \cdot x_{i,n}) \end{pmatrix}$$

oder  $A$

3. Update  $\theta^{t+1} = \theta^t + \alpha \cdot \nabla L(\theta^t)$

# Exkurs: Multivariate Lineare Regression – Analytisch / Lineare Algebra

- Datensatz  $(x_i, y_i)_{i=1}^n$  in Matrixschreibweise  $(X, Y)$  ↖ n Zeilen
- Verlustfunktion und Ableitung in Matrixschreibweise

$$\begin{aligned} \nabla \sum_{i=1}^n (y_i - \hat{f}(x_i; \theta))^2 &= \nabla \sum_{i=1}^n (y_i - \underbrace{x_i}_{1 \times 1} \underbrace{\theta}_{1 \times 1})^2 \quad \left( \begin{matrix} \text{J} \left( \begin{matrix} y_1 - x_1 \theta \\ y_2 - x_2 \theta \\ \vdots \\ y_n - x_n \theta \end{matrix} \right) \end{matrix} \right) \\ &= \nabla ((Y - X\theta)^T (Y - X\theta)) = 2X^T(Y - X\theta) \end{aligned}$$

$$f(x; \theta) = x \cdot \theta = (x_{i,1}, \dots, x_{i,n}) \cdot \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix}$$

- Algebraische Lösung in Matrixschreibweise

$$\theta = (X^T X)^{-1} X^T Y$$

$$\begin{aligned} X X^T \cdot (Y - X\theta) &= 0 \\ X^T Y &= X^T X \cdot \theta \quad | \cdot (X^T X)^{-1} \end{aligned}$$

Inverse der Varianz  $X$  Kovarianz( $X, Y$ )

# Exkurs: Gradientenabstieg oder Lineare Algebra?

## Gradientenabstieg

- Iterative Berechnung

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \alpha \cdot \nabla L(\boldsymbol{\theta}^t)$$

- Wahl der Lernrate  $\alpha$  bzw. Matrix A
  - Fest
  - Adaptiv
- Benötigt u.U. viele Iterationen
- Funktioniert auch für großes  $n$ , d.h. viele Daten

## Lineare Algebra

- Algebraische Lösung

$$\boldsymbol{\theta} = (X'X)^{-1}X'Y$$

- keine Meta-Parameter
- Direkte Lösung, keine Iterationen
- Benötigt Berechnung von  $(X'X)^{-1}$ , d.h.
  - Inverse muss existieren – insbesondere keine linear abhängigen Variablen!
  - Rechenintensives Invertieren einer  $n \times n$  Matrix  $\sim O(n^3)$
- Langsam für große  $n$

# Ausblick: Multivariate Regression – Gradientenabstieg

**Aufgabe:** Regression, d.h. Vorhersage  $\hat{y} = \hat{y}(x) = f(x)$

**Erfahrung:** Datensatz  $(x_i, y_i)_{i=1}^n$  mit  $x_i = (x_{i,1}, x_{i,2}, \dots, x_{i,m})$

**Qualität:** Verlustfunktion:  $L(y, \hat{y}) = L(\theta) = \mathcal{L}(f(\theta))$

**Maschine:** Regression mit  $f(x; \theta)$

**Lernen:** Finde Werte für  $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_m)$ , die die quadratische Verlustfunktion minimieren

Durch geeignete Wahl von  $\theta$  in einem iterativen Prozess (Gradientenabstiegsverfahren):

1. Wähle Startwert z.B.  $\theta^0 = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$  }  $m+1$  äußere Abl.      innere Abl.
2. Berechne Gradienten  
$$\nabla L(\theta^0) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} L(\theta^0) \\ \frac{\partial}{\partial \theta_1} L(\theta^0) \\ \vdots \\ \frac{\partial}{\partial \theta_m} L(\theta^0) \end{pmatrix} = \begin{pmatrix} \frac{d}{df} L(f(\theta^0)) \cdot \frac{\partial}{\partial \theta_0} f(x; \theta^0) \\ \frac{d}{df} L(f(\theta^0)) \cdot \frac{\partial}{\partial \theta_1} f(x; \theta^0) \\ \vdots \\ \frac{d}{df} L(f(\theta^0)) \cdot \frac{\partial}{\partial \theta_m} f(x; \theta^0) \end{pmatrix}$$
 }  $m+1$  Parameter
3. Update  $\theta^{t+1} = \theta^t + \mathbf{A} \cdot \nabla L(\theta^t)$