

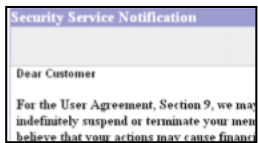
Introduction to Machine Learning

Linear Regression

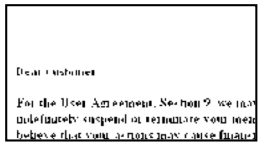
Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

Basic Supervised Learning Pipeline

Training Data



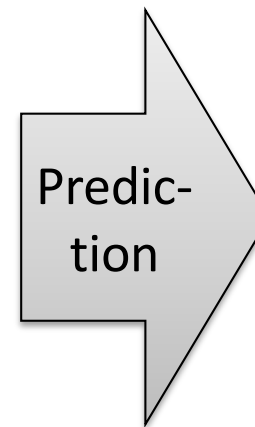
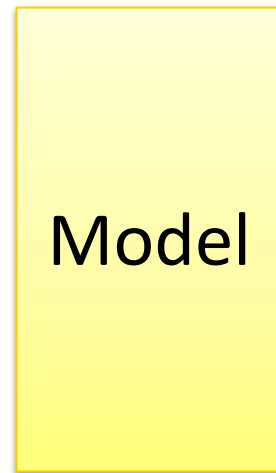
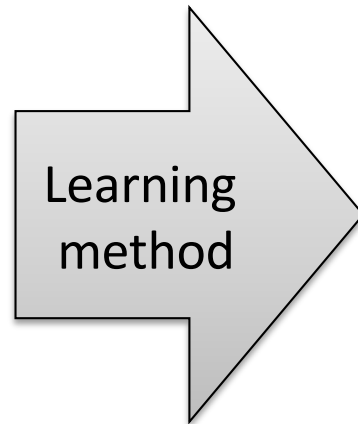
“spam”



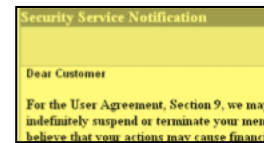
“ham”



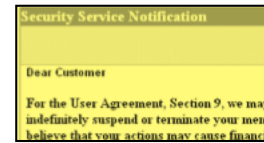
“spam”



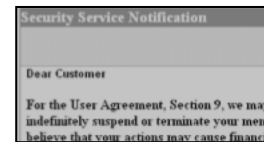
Test Data



?



?



?

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

Model fitting

Prediction/
Generalization

Regression

- Instance of supervised learning
- **Goal:** Predict **real valued** labels (possibly vectors)
- Examples:

X

Flight route

Real estate objects

Customer & ad features

Y

Delay (minutes)

Price

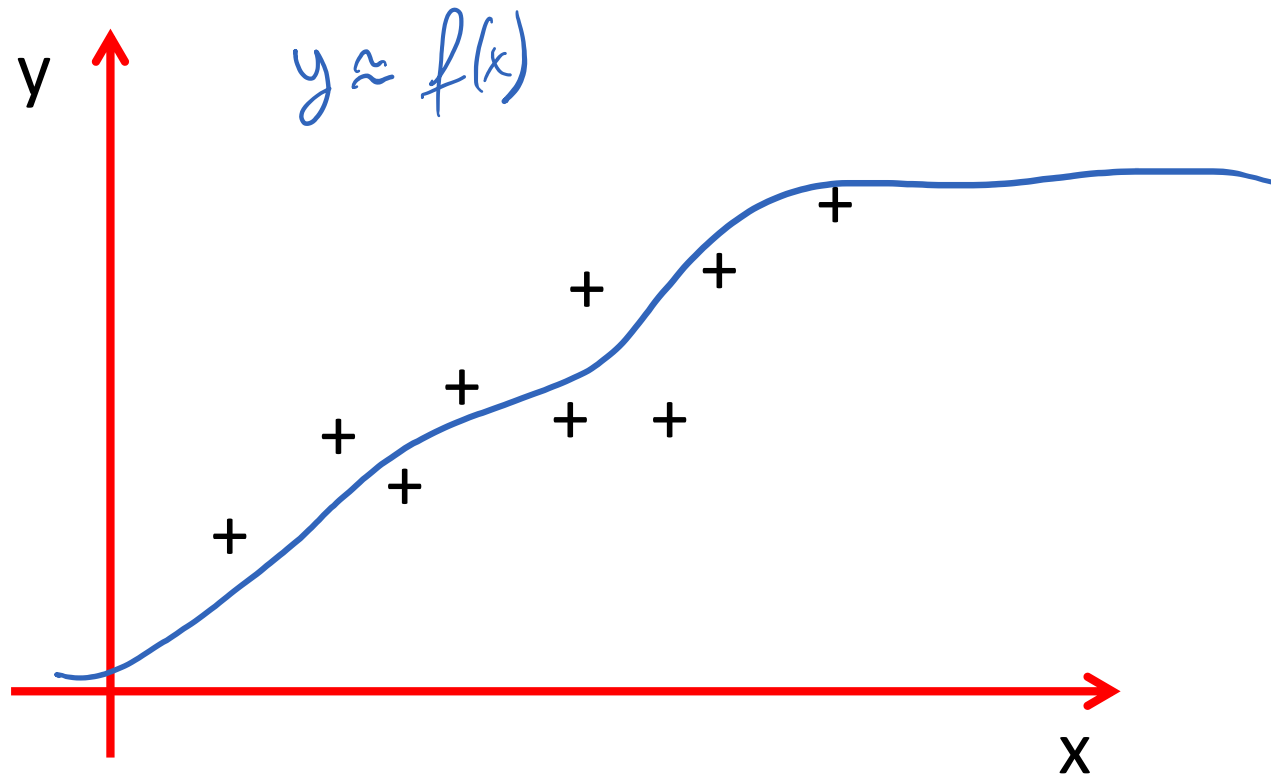
Click-through probability

Running example: Diabetes

[Efron et al '04]

- Features X :
 - Age
 - Sex
 - Body mass index
 - Average blood pressure
 - Six blood serum measurements (S1-S6)
- Label (target) Y
 - quantitative measure of disease progression

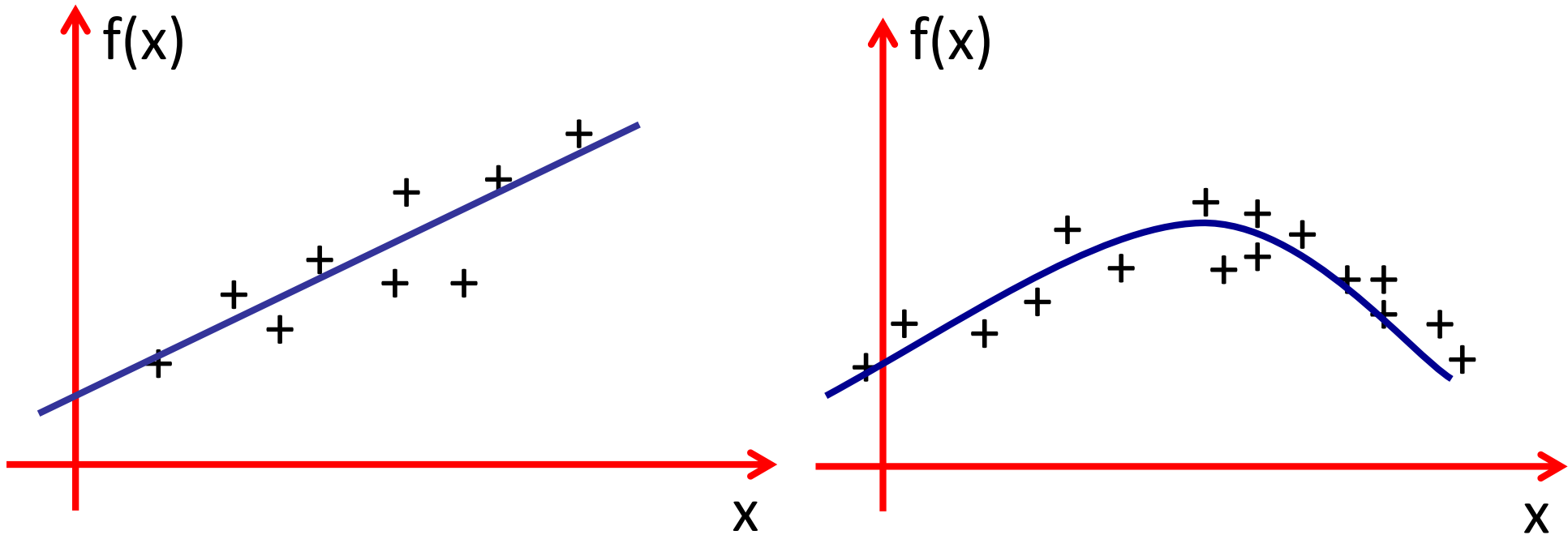
Regression



- **Goal:** learn real valued mapping $f : \mathbb{R}^d \rightarrow \mathbb{R}$

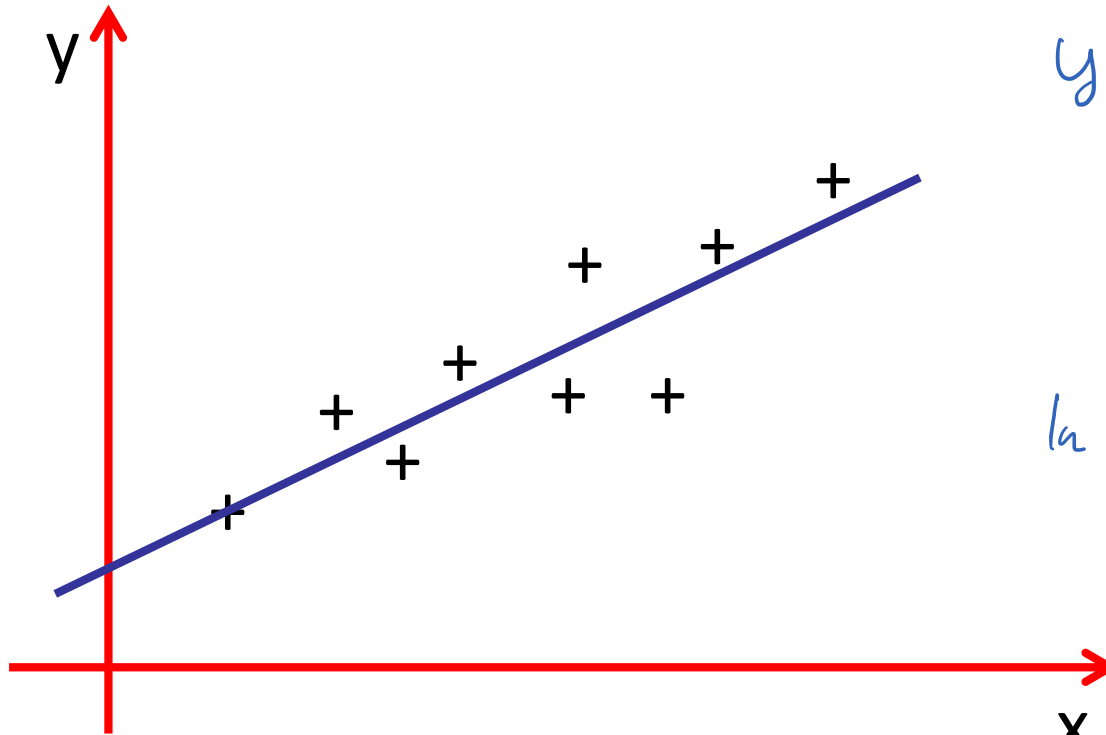
Important choices in regression

- What **types of functions f** should we consider? Examples



- How should we measure **goodness of fit**?

Example: linear regression



$$x = [x_1 \dots x_d]$$

$$w = [w_1 \dots w_d]$$

$$y \approx f(x)$$

f is linear (affine)

1-dim: $f(x) = ax + b$

2-dim: $f(x_1, x_2) = ax_1 + bx_2 + c$

x d -dim: $f(x) = w_1 x_1 + \dots + w_d x_d + w_0$

$$= \sum_{i=1}^d w_i x_i + w_0$$

$$= w^T x + w_0$$

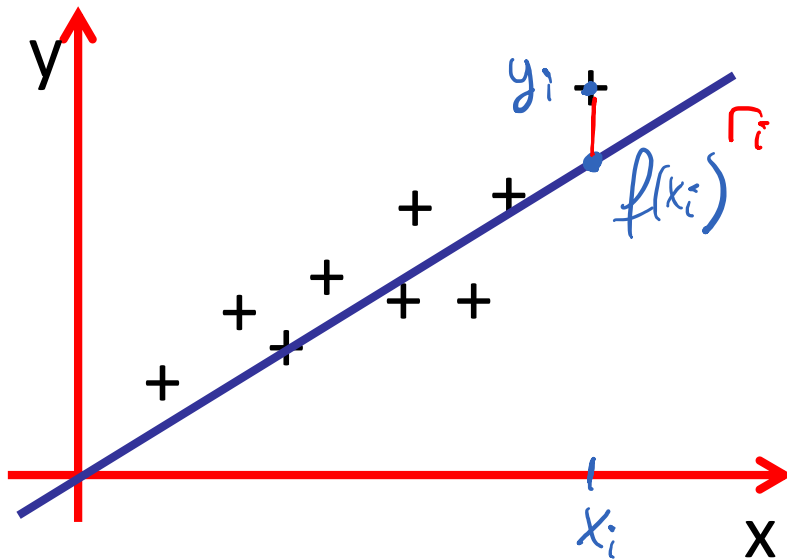
Homogeneous representation

$$\underbrace{w^T x + w_0}_{\in \mathbb{R}} = \underbrace{\tilde{w}^T \tilde{x}}_{\substack{\in \mathbb{R} \\ \in \mathbb{R}^{d+1}}}$$
$$\begin{array}{l} \overset{\in \mathbb{R}^d}{w} = [w_1 \dots w_d] \\ x = [x_1 \dots x_d] \end{array} \quad \Rightarrow \quad \begin{array}{l} \tilde{w} = [w_1 \dots w_d \ w_0] \\ \tilde{x} = [x_1 \dots x_d \ 1] \end{array}$$

$$\Rightarrow \text{w.log.} : f(x) = v^T x$$

Quantifying goodness of fit

$$D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \quad \mathbf{x}_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$



Residuals

$$\begin{aligned} r_i &= y_i - f(x_i) \\ &= y_i - w^T x_i \end{aligned}$$

$$\begin{aligned} \Rightarrow \text{Cost } \hat{R}(w) &= \sum_{i=1}^n r_i^2 \\ &= \sum_{i=1}^n (y_i - w^T x_i)^2 \end{aligned}$$

Note: We've made 2 decisions

- 1) quantify error for 1 point via squared residual
- 2) we sum over all points

Least-squares linear regression optimization

[Legendre 1805, Gauss 1809]

- Given data set $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- How do we find the **optimal weight vector**?

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

Method 1: Closed form solution

- The problem $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$

can be solved in **closed form**:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

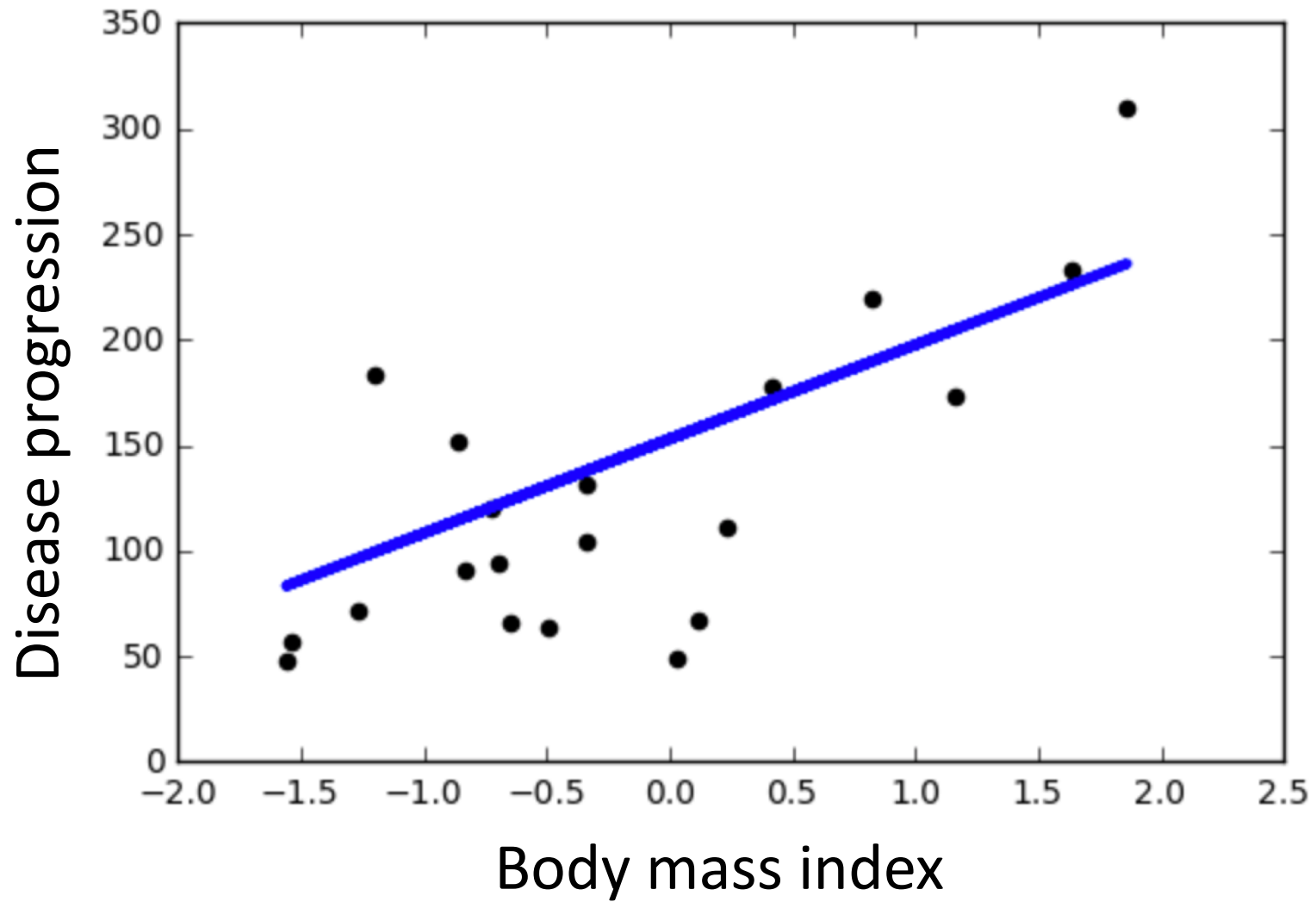
- Hereby:

$$\mathbf{X} = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

How to solve? Example: Scikit Learn

```
# Create linear regression object  
regr = linear_model.LinearRegression()  
  
# Train the model using the training set  
regr.fit(X_train, Y_train)  
  
# Make predictions on the testing set  
Y_pred = regr.predict(X_test)
```

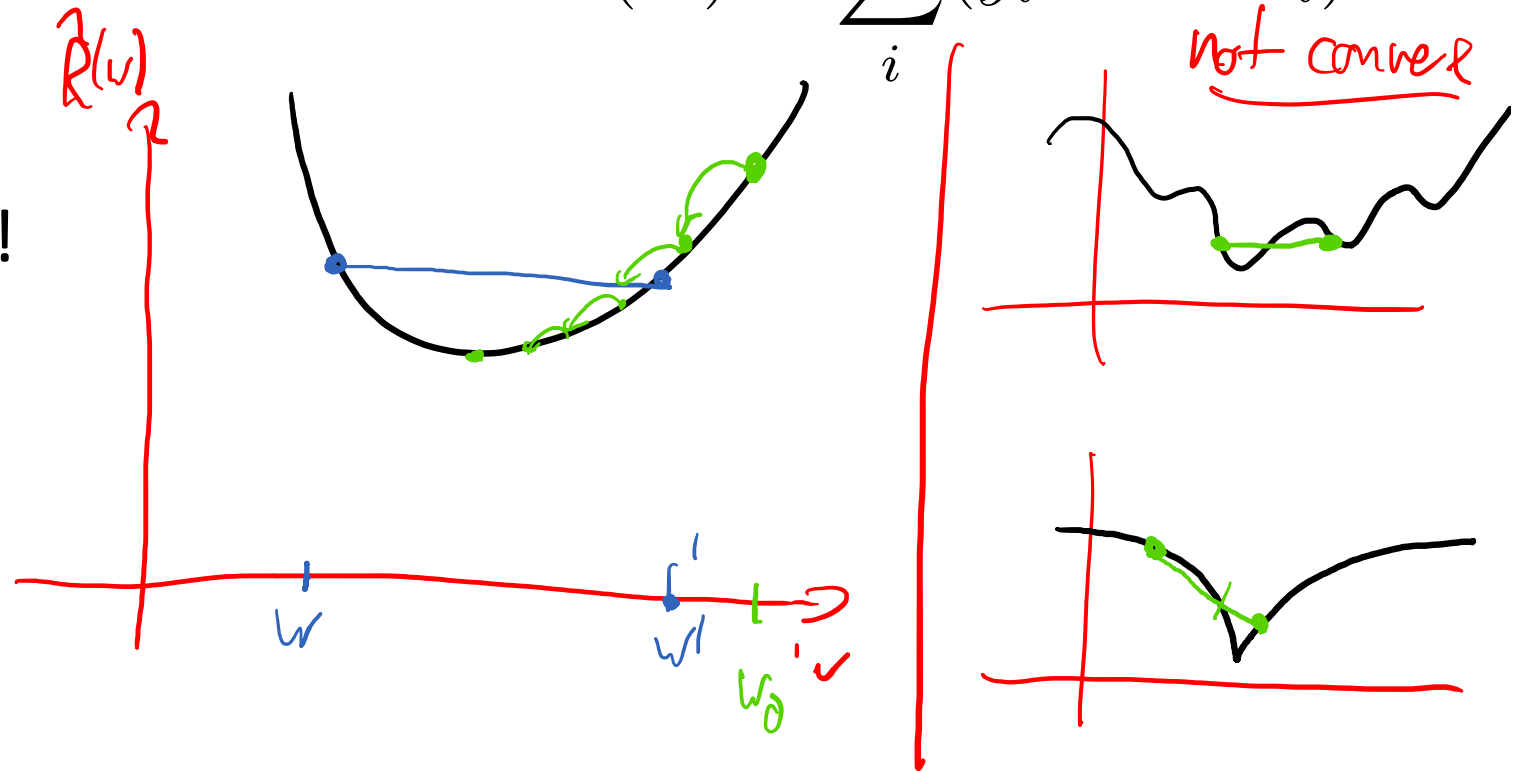
Demo



Method 2: Optimization

- The objective function $\hat{R}(\mathbf{w}) = \sum_i (y_i - \mathbf{w}^T \mathbf{x}_i)^2$

is convex!



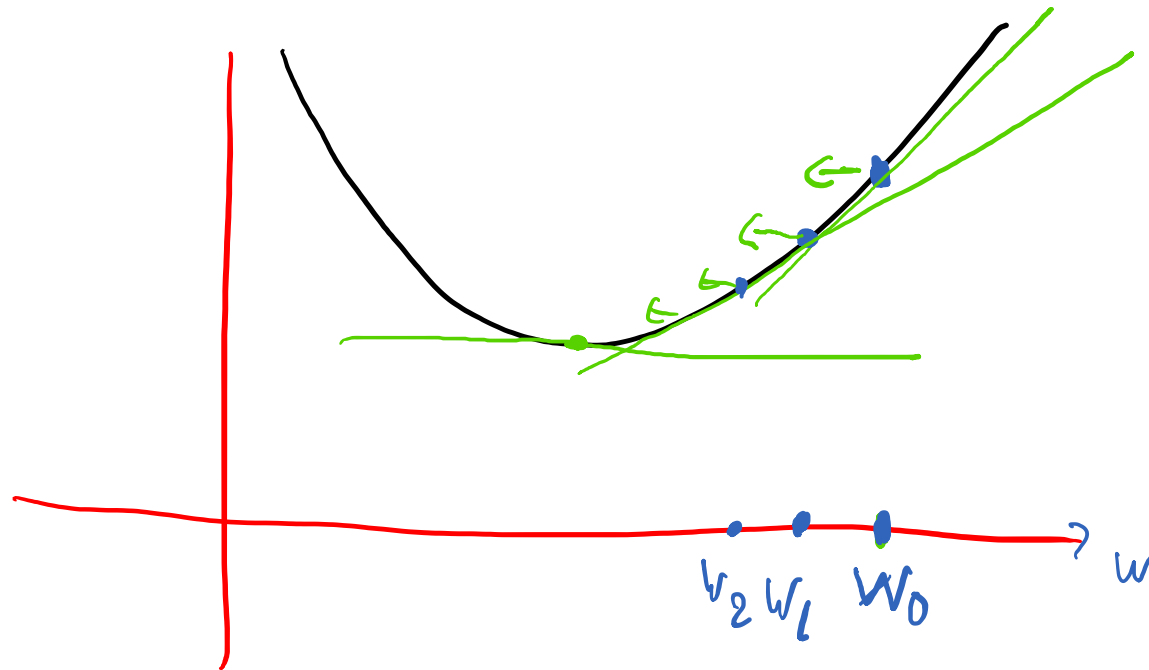
$f: \mathbb{R}^d \rightarrow \mathbb{R}$ convex iff $\forall x, x', \lambda \in [0, 1]$ it holds that

$$f(\lambda x + (1-\lambda)x') \leq \lambda f(x) + (1-\lambda)f(x')$$

Gradient Descent

- Start at an arbitrary $\mathbf{w}_0 \in \mathbb{R}^d$
- For $t=1, 2, \dots$ do $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla \hat{R}(\mathbf{w}_t)$

- Hereby, η_t is called learning rate



Convergence of gradient descent

- Under mild assumptions, if step size sufficiently small, gradient descent converges to a **stationary point** (gradient = 0)
- For convex objectives, it therefore finds the **optimal solution!**

- In the case of the squared loss, **constant stepsize** $\left(\frac{1}{2}\right)$ converges **linearly**

$$\Rightarrow \forall t \geq t_0 \exists \alpha < 1 \text{ s.t. } \underbrace{\left(\hat{R}(w_{t+\alpha}) - \hat{R}(w^*) \right)}_{\substack{\leq \epsilon \\ \text{"gap" at } t+\alpha}} \leq \alpha \underbrace{\left(\hat{R}(w_t) - \hat{R}(w^*) \right)}_{\substack{\text{"gap" at } t}} \leq \alpha \epsilon$$

\Rightarrow can find ϵ -optimal solution in $O\left(\ln \frac{1}{\epsilon}\right)$ iter.

Computing the gradient

$$\nabla \hat{R}(w) = \left[\frac{\partial}{\partial w_1} \hat{R}(w) \quad \dots \quad \frac{\partial}{\partial w_d} \hat{R}(w) \right]$$

$$\begin{aligned} \text{In 1-dim: } \nabla \hat{R}(w) &= \frac{d}{dw} \hat{R}(w) = \frac{d}{dw} \sum_{i=1}^n (y_i - w \cdot x_i)^2 \\ &= \sum_{i=1}^n \frac{d}{dw} (y_i - w \cdot x_i)^2 \\ &= \sum_{i=1}^n 2 \underbrace{(y_i - w \cdot x_i)}_{r_i} \cdot (-x_i) = -2 \sum_{i=1}^n r_i x_i \end{aligned}$$

$$\text{In d-dim: } \nabla \hat{R}(w) =$$

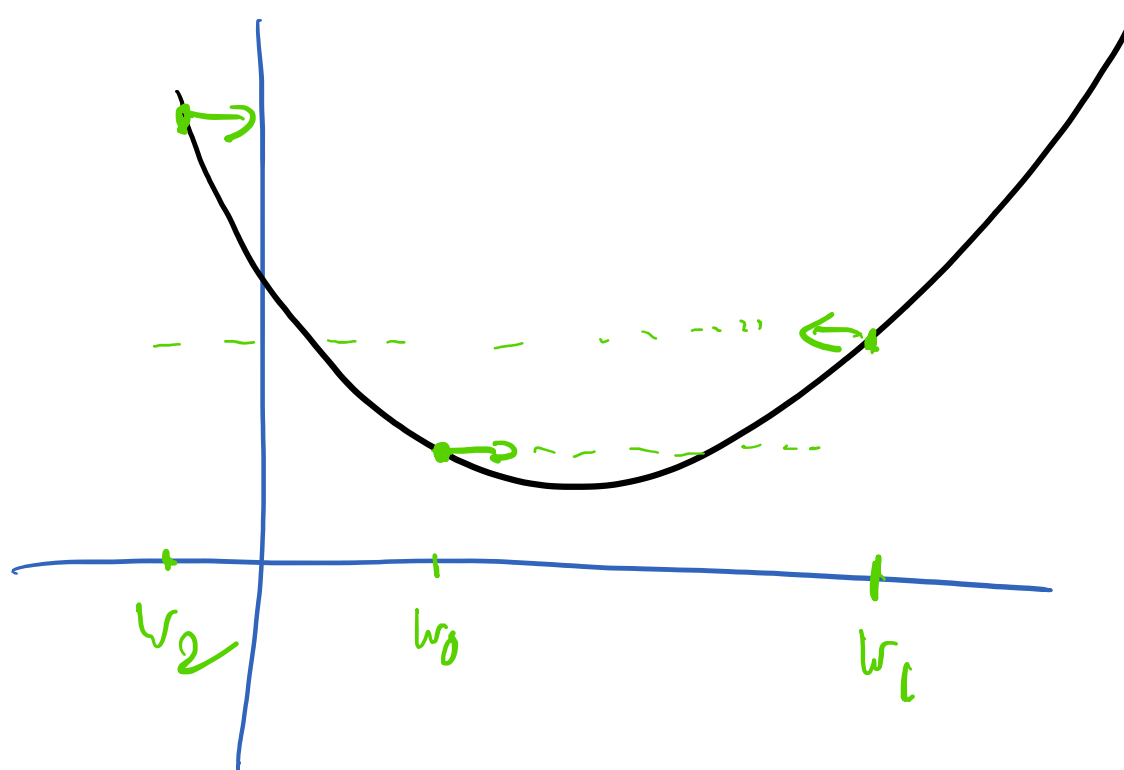
$$= -2 \sum_{i=1}^n r_i x_i$$

ER
↓
ER
↑
ER^d

Demo: Gradient descent

Choosing a stepsize

- What happens if we choose a poor stepsize?

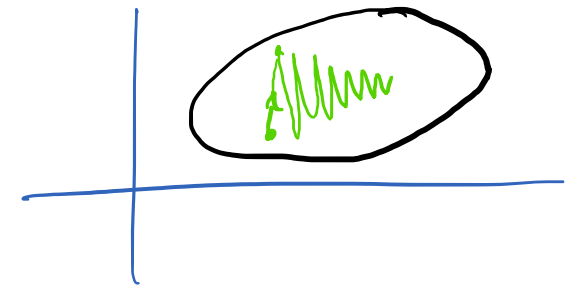


Adaptive step size

- Can update the step size adaptively. For example:
- 1) Via **line search** (optimizing step size every step)

Spes at iter t , have $w_t, g_t = \nabla \hat{R}(w_t)$

Define $\eta_t^* = \operatorname{argmin}_{\eta \in [0, \infty)} \hat{R}(w_t - \eta g_t)$

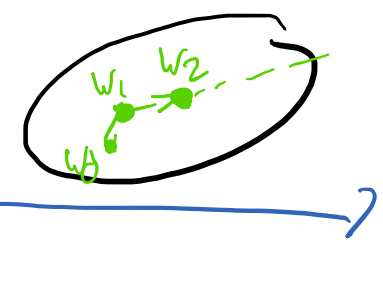


- 2) „Bold driver“ heuristic

- If function decreases, increase step size:

$$\text{If } \hat{R}(w_{t+1}) < \hat{R}(w_t) : \eta_{t+1} \leftarrow \eta_t \cdot c_{acc}$$

e.g. 1.1



- If function increases, decrease step size:

$$\text{If } \hat{R}(w_{t+1}) > \hat{R}(w_t) : \eta_{t+1} \leftarrow \eta_t \cdot c_{dec}$$

e.g. 0.5

Demo: Gradient Descent for Linear Regression

Gradient descent vs closed form

- Why would one ever consider performing gradient descent, when it is possible to find closed form solution?

Closed form: $\hat{w} = (X^T X)^{-1} (X^T y)$

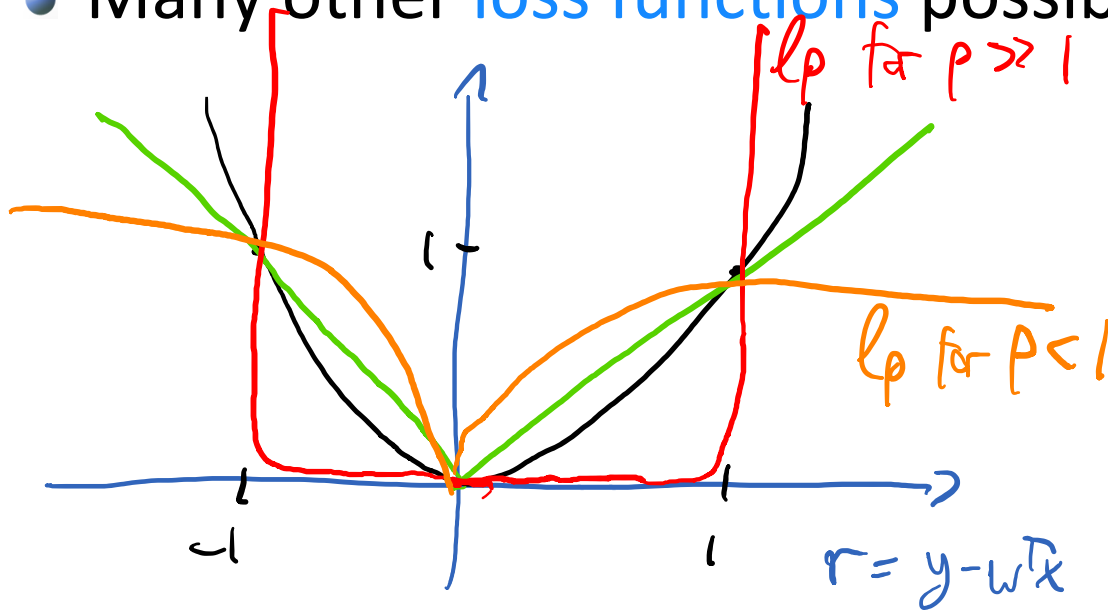
\uparrow
 $O(nd^2)$ solve lin sys. $O(d^3)$

Gradient descent: Calc. $\nabla R(w) = \sum_i (y_i - wx_i) x_i \Rightarrow O(nd) \cdot \underbrace{\ln(\frac{1}{\epsilon})}_{\text{iterations}}$

- **Computational complexity**
- **May not need an optimal solution**
- **Many problems don't admit closed form solution**

Other loss functions

- So far: Measure goodness of fit via squared error
- Many other **loss functions** possible (and sensible!)



least-squares

$$l_2(r) = r^2$$

Alternatives:

$$l_1(r) = |r|$$

$$l_p(r) = |r|^p$$

still convex for $p \geq 1$