

Introduction to Machine Learning

Linear Regression

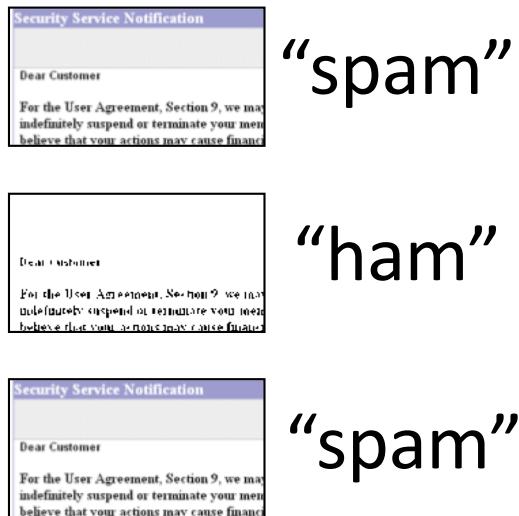
Prof. Andreas Krause
Learning and Adaptive Systems (las.ethz.ch)

Announcements

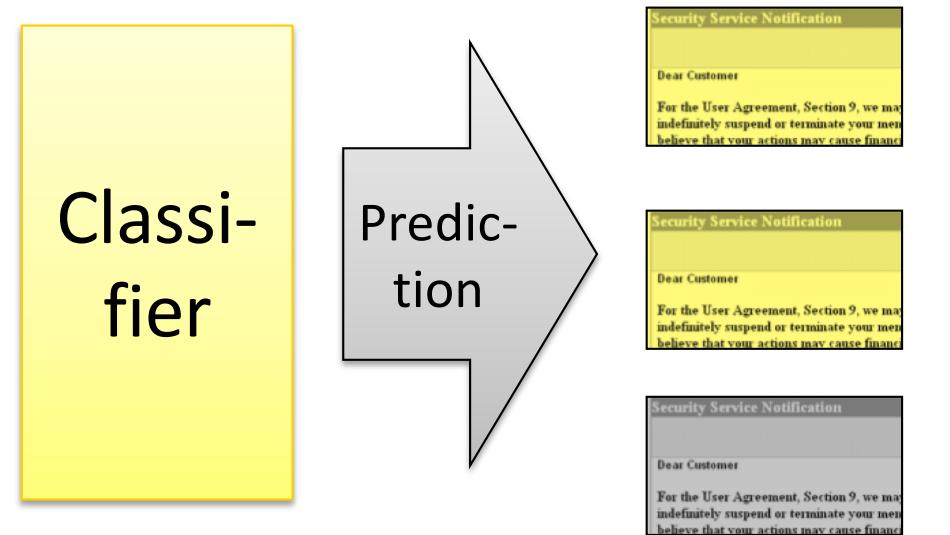
- First recordings online:
 - <https://www.video.ethz.ch/lectures/d-infk/2019/spring/252-0220-00L.html>
- Waitlist updates

Basic Supervised Learning Pipeline

Training Data



Test Data



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

Model fitting

Prediction/
Generalization

Regression

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

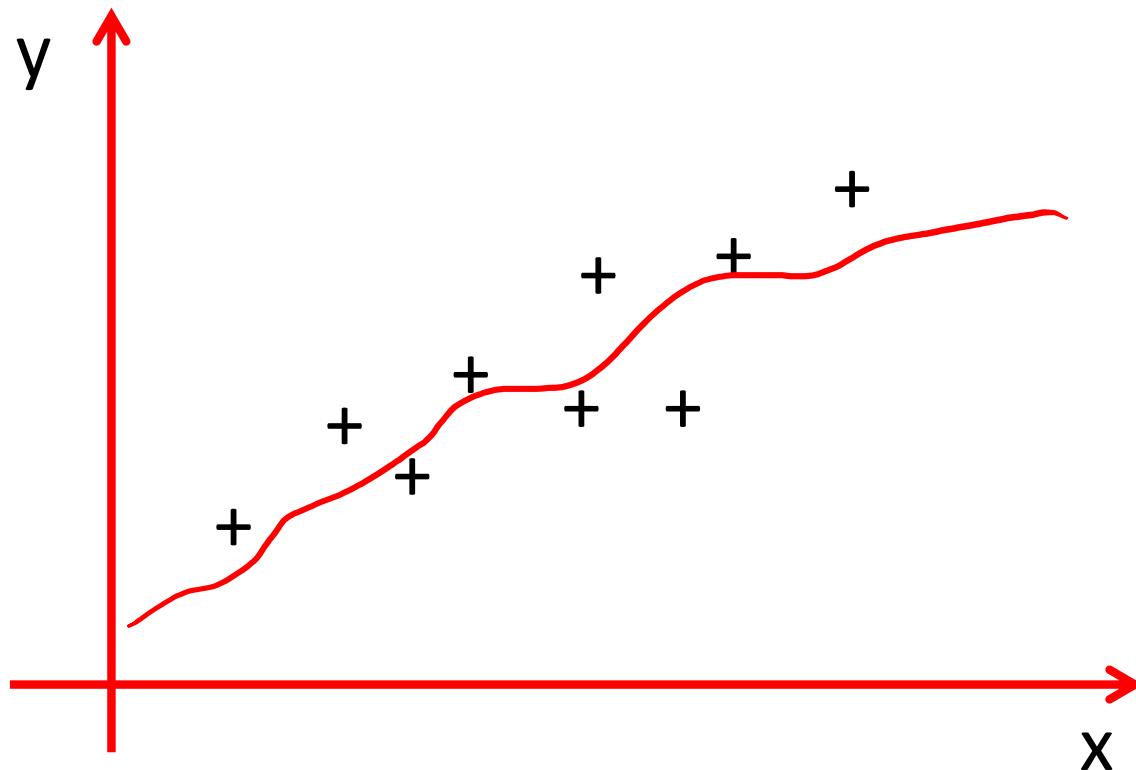
X	Y
Flight route	Delay (minutes)
Real estate objects	Price
Customer & ad features	Click-through probability

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

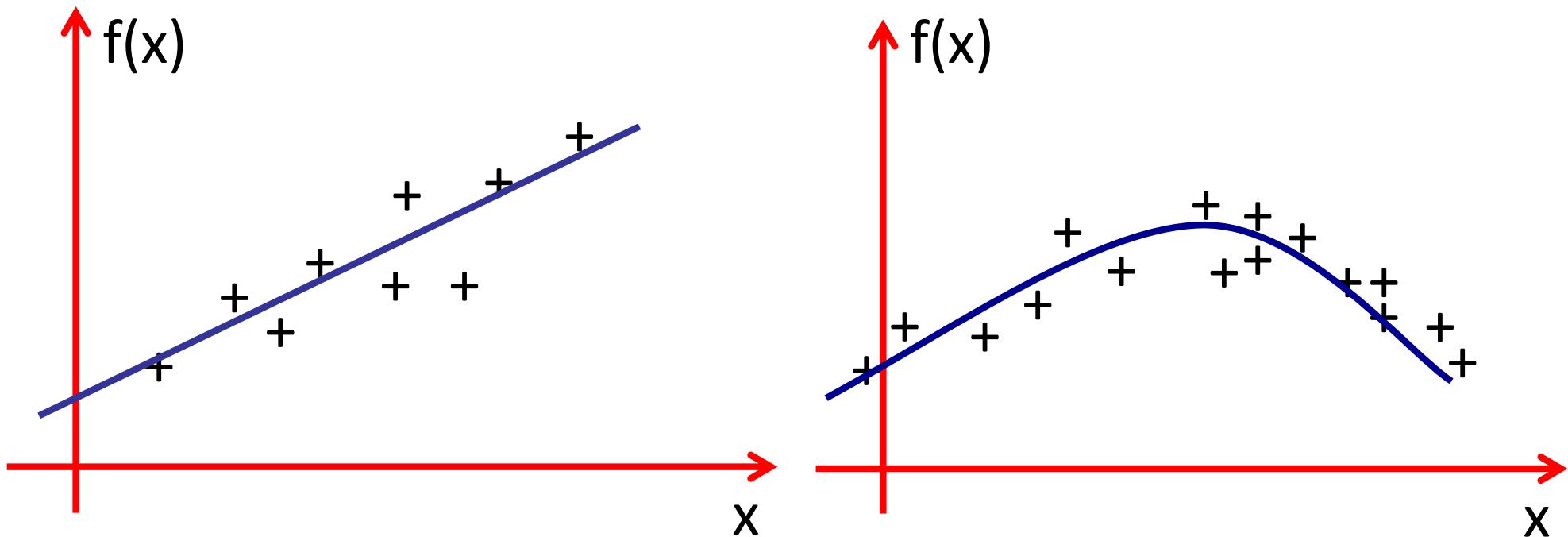
Regression



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

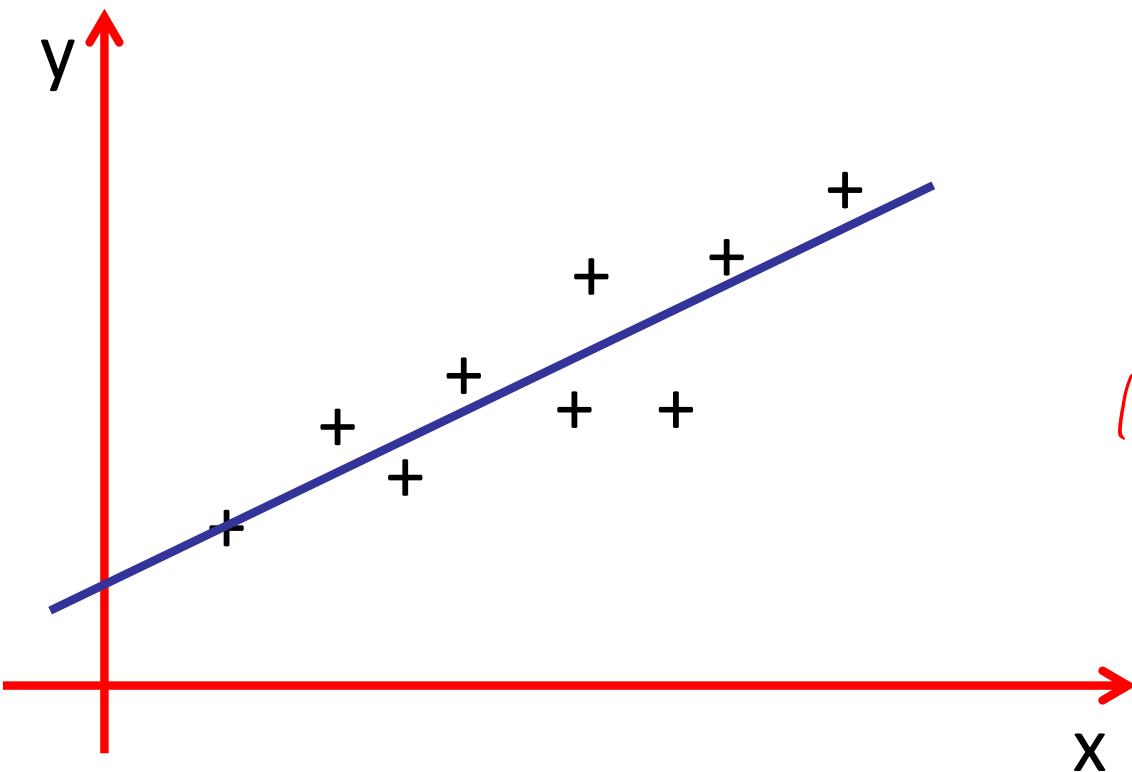
Important choices in regression

- What types of functions f should we consider? Examples



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

Example: linear regression



$$y \approx f(x)$$

f is linear
(affine) in x

$$\text{1-d: } f(x) = ax + b$$

$$2-d: \quad f(x) = a_1 x_1 + a_2 x_2 + c$$

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

$$= \sum_i w_i \cdot x_i + w_0$$

$$= w^T x + w_0$$

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

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

Homogeneous representation

$$w^T x + w_0 = \tilde{w}^T \tilde{x}$$

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

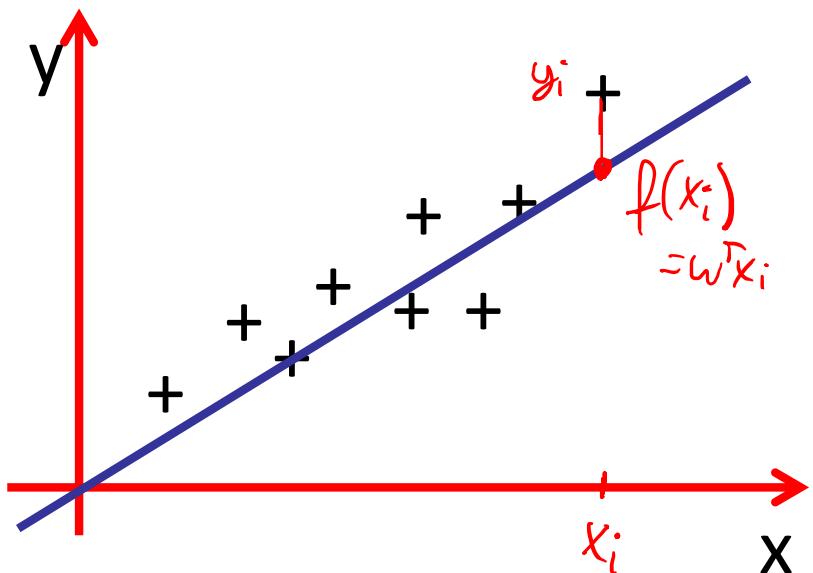
$$x = [x_1 \dots x_d]^T$$

$$\tilde{w} = [w_1 \dots w_d \ w_0]^T$$

$$\tilde{x} = [x_1 \dots x_d \ 1]^T$$

Quantifying goodness of fit

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



Residual $r_i = y_i - w^T x_i$

Cost $\hat{R}(w) = \sum_{i=1}^n r_i^2$
 $= \sum_{i=1}^n (y_i - w^T x_i)^2$

Note: We have made 2 decisions here:
1) quantify error via squared residuals
2) sum over these

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?

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

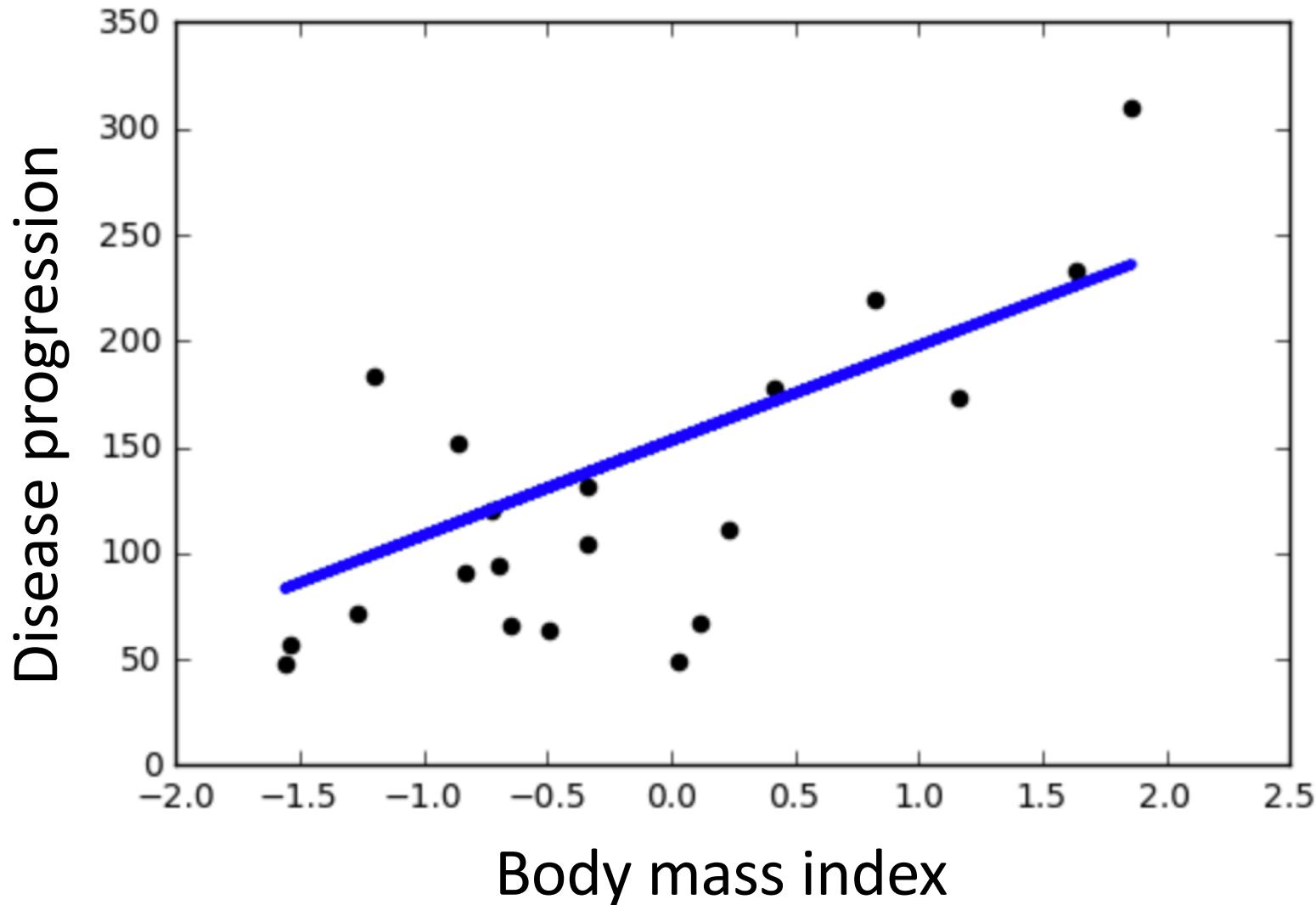
How to solve? Example: Scikit Learn

```
# 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 1: Closed form solution

- The problem

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

can be solved in **closed form**:

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

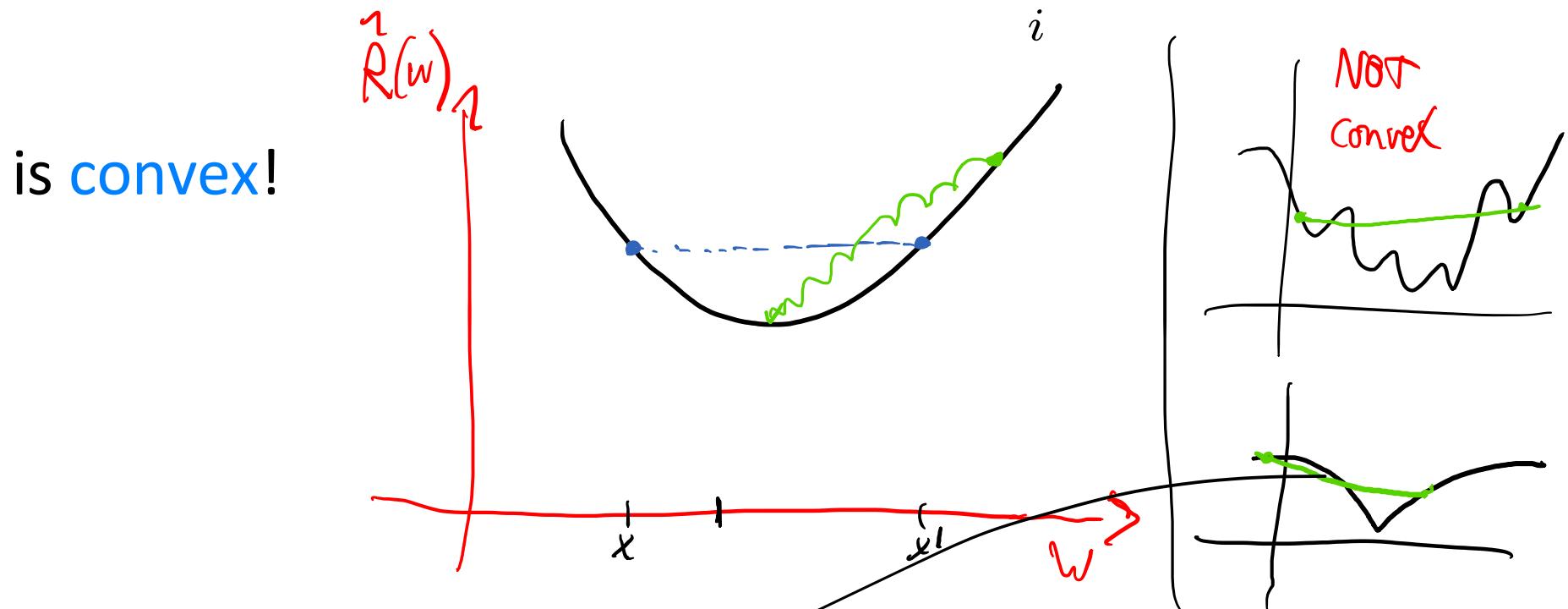
- Hereby:

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,d} \\ \vdots & & \vdots \\ x_{n,1} & \dots & x_{n,d} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Method 2: Optimization

- The objective function

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



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

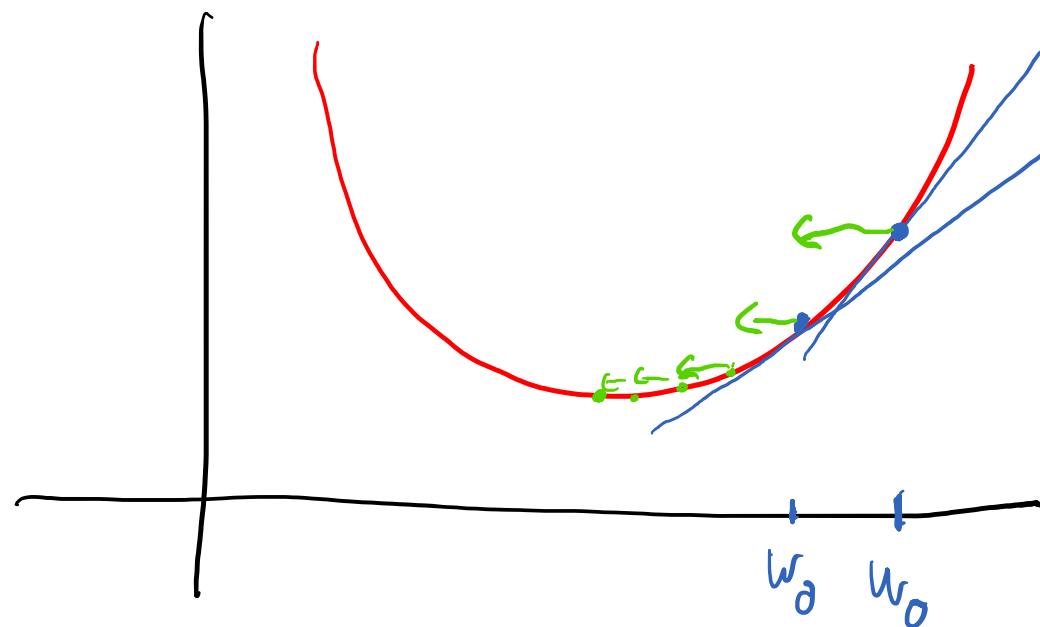
$$f(\lambda x + (1-\lambda)x') \leq \lambda f(x) + (1-\lambda)f(x')$$

Gradient Descent

$\frac{1}{2}$ for least squares

- 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 $\frac{1}{2}$ converges linearly

Can find x_t s.t. $f(x_t) - \min_x f(x) \leq \epsilon$ in $O(\log \frac{1}{\epsilon})$
iterations

Computing the gradient

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

(a) 1-d: $\nabla \hat{R}(w) = \frac{d}{dw} \hat{R}(w) = \frac{d}{dw} \sum_{i=1}^n (y_i - w \cdot x_i)^2$

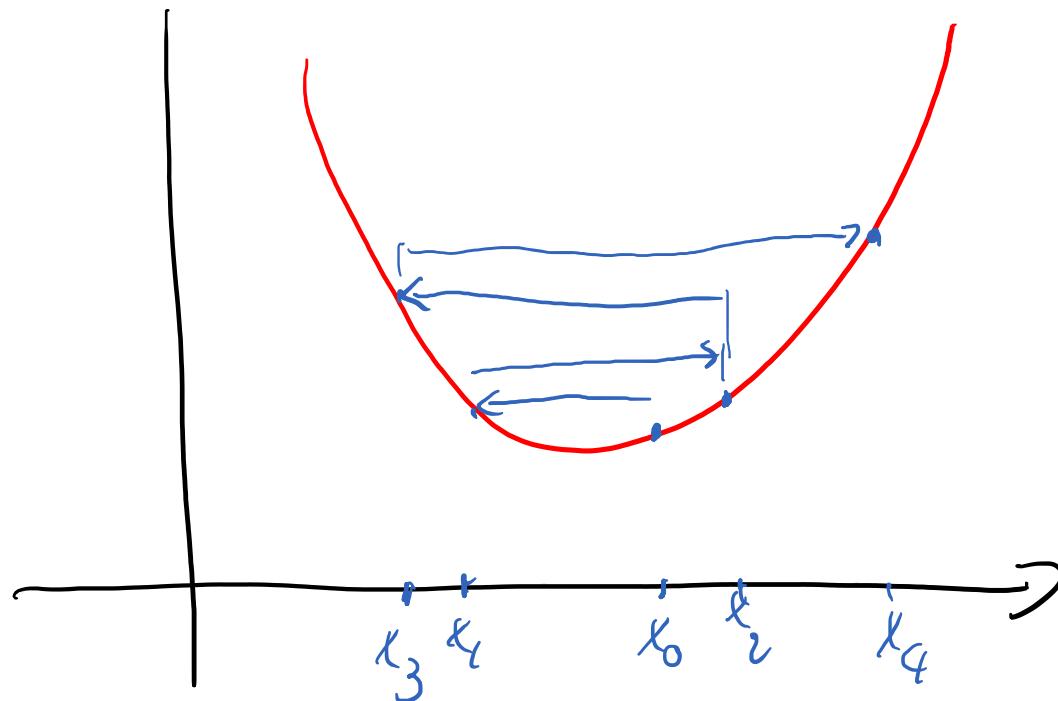
$$= \sum_{i=1}^n \frac{d}{dw} (y_i - w \cdot x_i)^2$$
$$= \sum_{i=1}^n 2 \underbrace{(y_i - w \cdot x_i)}_{r_i} (-x_i) = -2 \sum_i r_i \cdot x_i$$

(b) d-d: $\nabla \hat{R}(w) = -2 \sum_i r_i \cdot x_i^T$

Demo: Gradient descent

Choosing a stepsize

- What happens if we choose a poor stepsize?

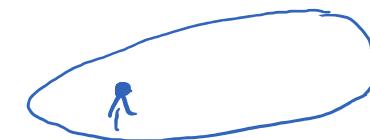


Adaptive step size

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

Sps. $\delta_t = \nabla \hat{R}(w_t)$

Pick $\gamma_t \in \arg\min_{\gamma \in (0, \infty)} \hat{R}(w_t - \gamma \cdot g_t)$



- 2) „Bold driver“ heuristic

- If function decreases, increase step size:

$$\text{If } \hat{R}(w_t - \gamma_t g_t) < \hat{R}(w_t) \Rightarrow \gamma_{t+1} = \gamma_t \cdot c_{\text{acc}}$$

e.g. $c_{\text{acc}} = 1.1$
 $c_{\text{dec}} = 0.5$

- If function increases, decrease step size:

$$\text{If } \hat{R}(\underbrace{w_t - \gamma_t g_t}_{w_{t+1}}) > \hat{R}(w_t) \Rightarrow \gamma_{t+1} = \gamma_t \cdot c_{\text{dec}}$$

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?

CF: $w^* = (X^T X)^{-1} (X^T y)$

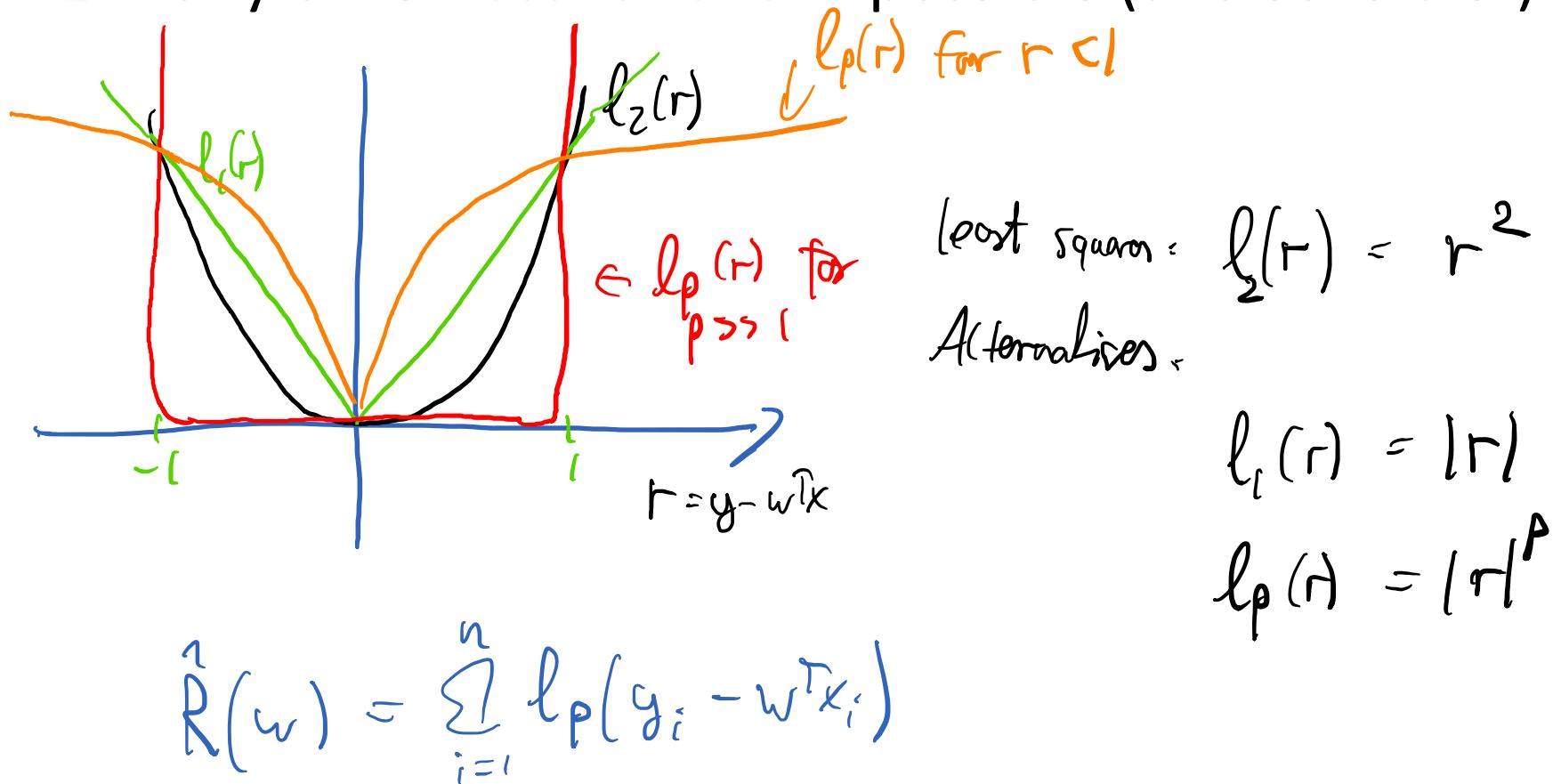
$\nwarrow O(n \cdot d^2)$, solve n -sys. $\rightarrow O(d^3)$

GD: Calc. $\nabla R(w) = \sum_i (y_i - \underline{w}^T x_i) x_i \Rightarrow O(n \cdot d)$, need $O(\log_{\frac{1}{\epsilon}})$ iters

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

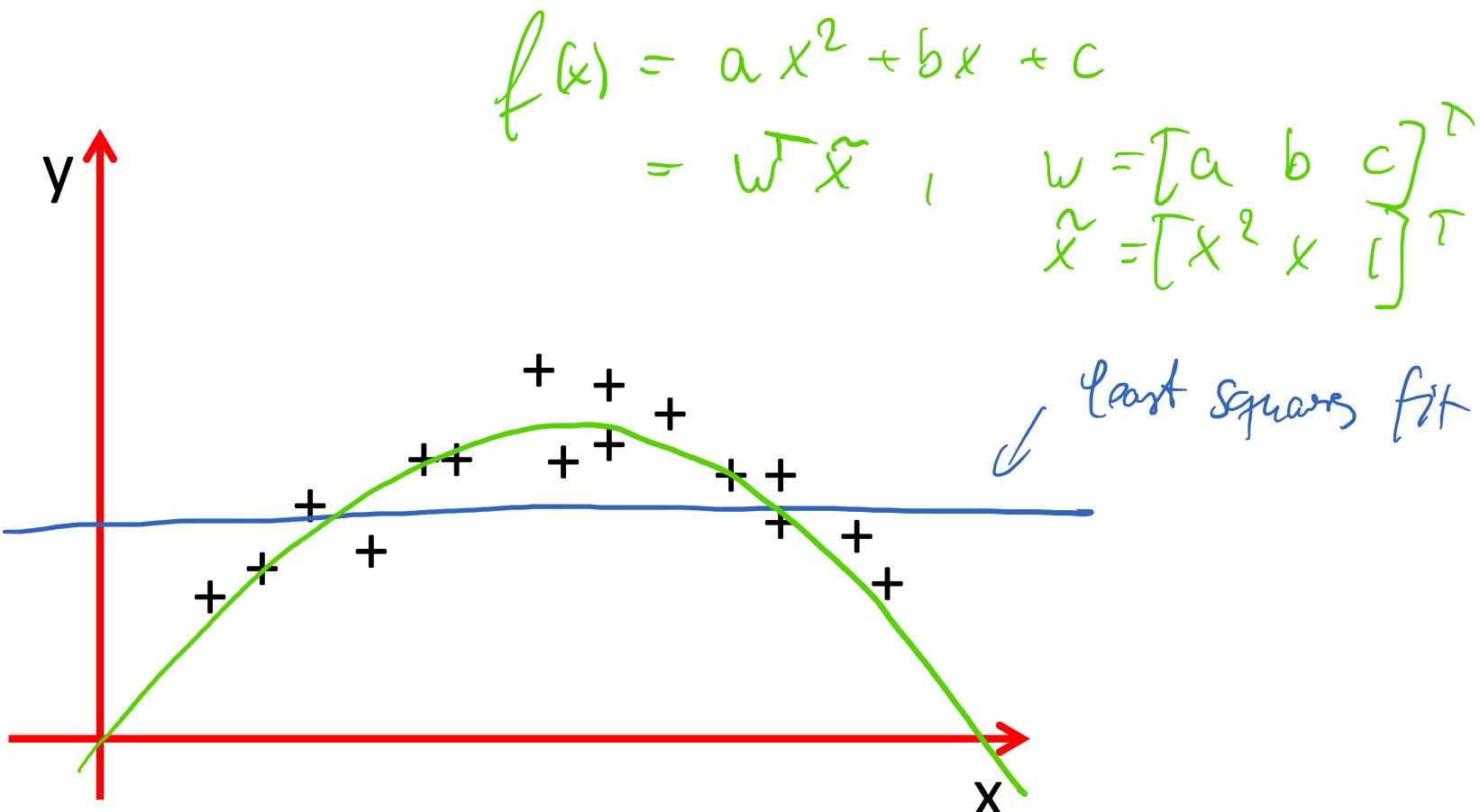
Other loss functions

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



Fitting nonlinear functions

- How about functions like this:



Linear regression for polynomials

We can fit non-linear functions via linear regression, using nonlinear features of our data (basis functions)

$$f(\mathbf{x}) = \sum_{i=1}^{\phi} w_i \phi_i(\mathbf{x})$$

$x \mapsto \tilde{x} = \phi(x)$

$$\Phi = \left| \phi(x) \right|$$

In 1-d: $\phi(x) = [1, x, x^2, \dots x^k]$

In 2-d: $\phi(x) = [1, x_1, x_2, x_1^2, x_2^2, x_1 \cdot x_2, \dots]$

In d-d, $\phi(x)$ = vector of all monomials of degree up to k in $x_1 \dots x_d$