

COMPUTATIONAL INTELLIGENCE

(INTRODUCTION TO MACHINE LEARNING) SS17

Lecture 2:

- Linear Regression
- Gradient Descent
- Non-linear basis functions

LINEAR REGRESSION

MOTIVATION

Why Linear Regression?

- Simplest machine learning algorithm for regression
 - Widely used **in biological, behavioral and social sciences** to describe and to extract relationships between variables from data
 - Prediction **of real-valued outputs**
 - Easy to implement, fast to execute
 - **Benchmark** algorithm for comparison with more complex algorithms
- Introduction to **notation and concepts** that we will need again later in the course
 - Data format, vector & matrix notation
 - Learning from data by minimizing a **cost function**
 - Gradient descent
 - Non-linear features and basis functions
 - Preparation for neural networks

Applications of (linear) regression

- Brain computer interfaces
 - <https://www.youtube.com/watch?v=Ae6En8-eaww>
- Neuroprosthetic control
 - https://www.youtube.com/watch?v=X_AI4MiY6L4

LINEAR REGRESSION WITH ONE INPUT

Linear regression with one input

$$\langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle$$

Training set

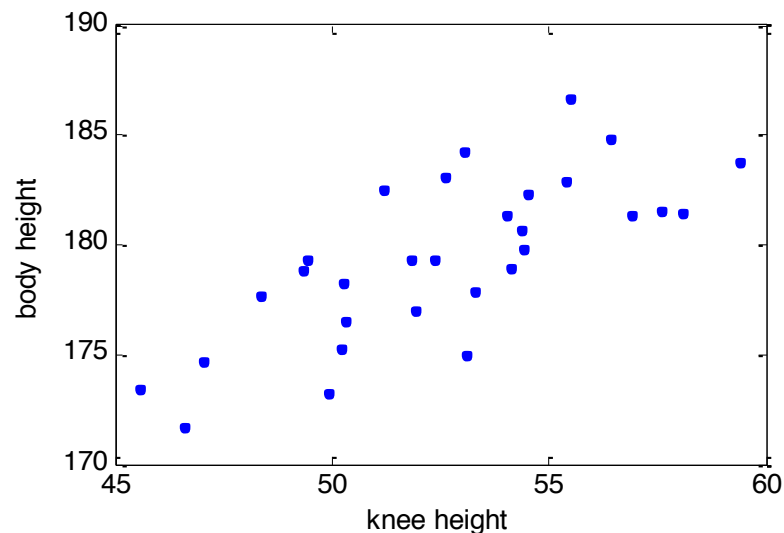
Learning algorithm?

Test input

x

„Hypothesis“
 h

Prediction



Hypothesis

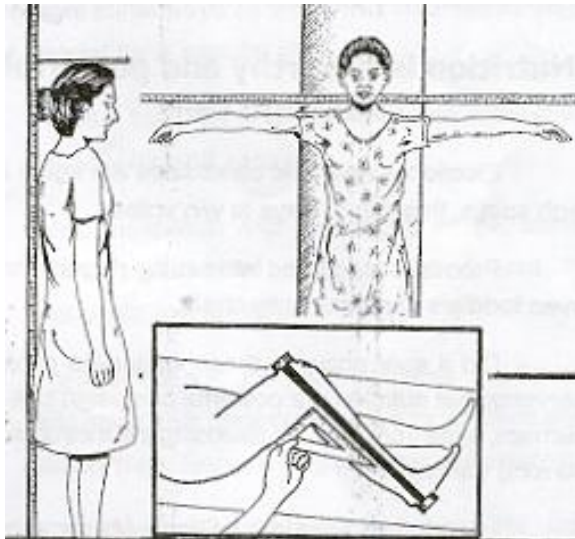
$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

Parameters?

$$\theta = (\theta_0, \theta_1)$$

A regression problem

- We want to learn to predict a **person's height** based on his/her **knee height** and/or **arm span**
- This is useful for patients who are **bed bound** or in a wheelchair and cannot stand to take an accurate measurement of their height

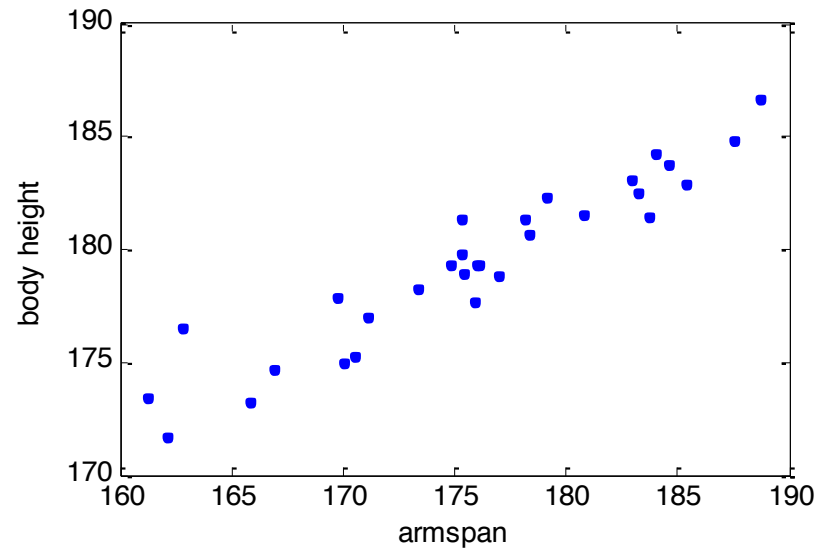
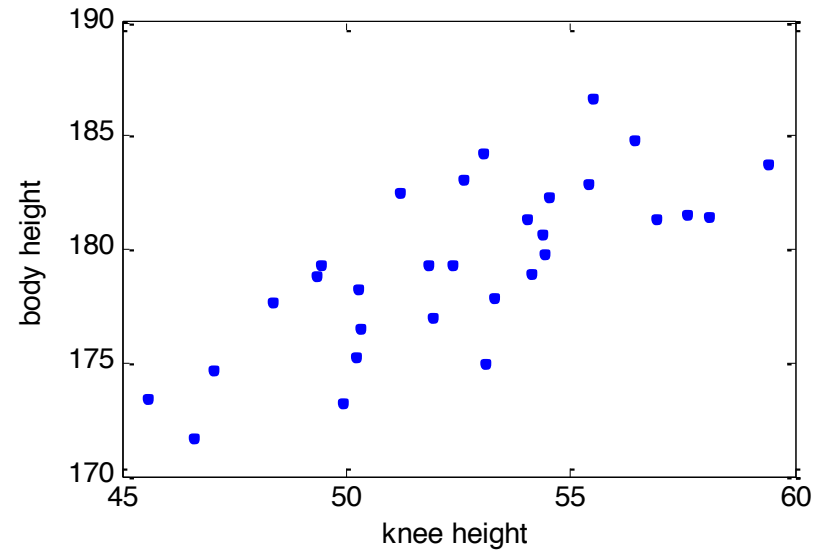


Knee Height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...

Example Data

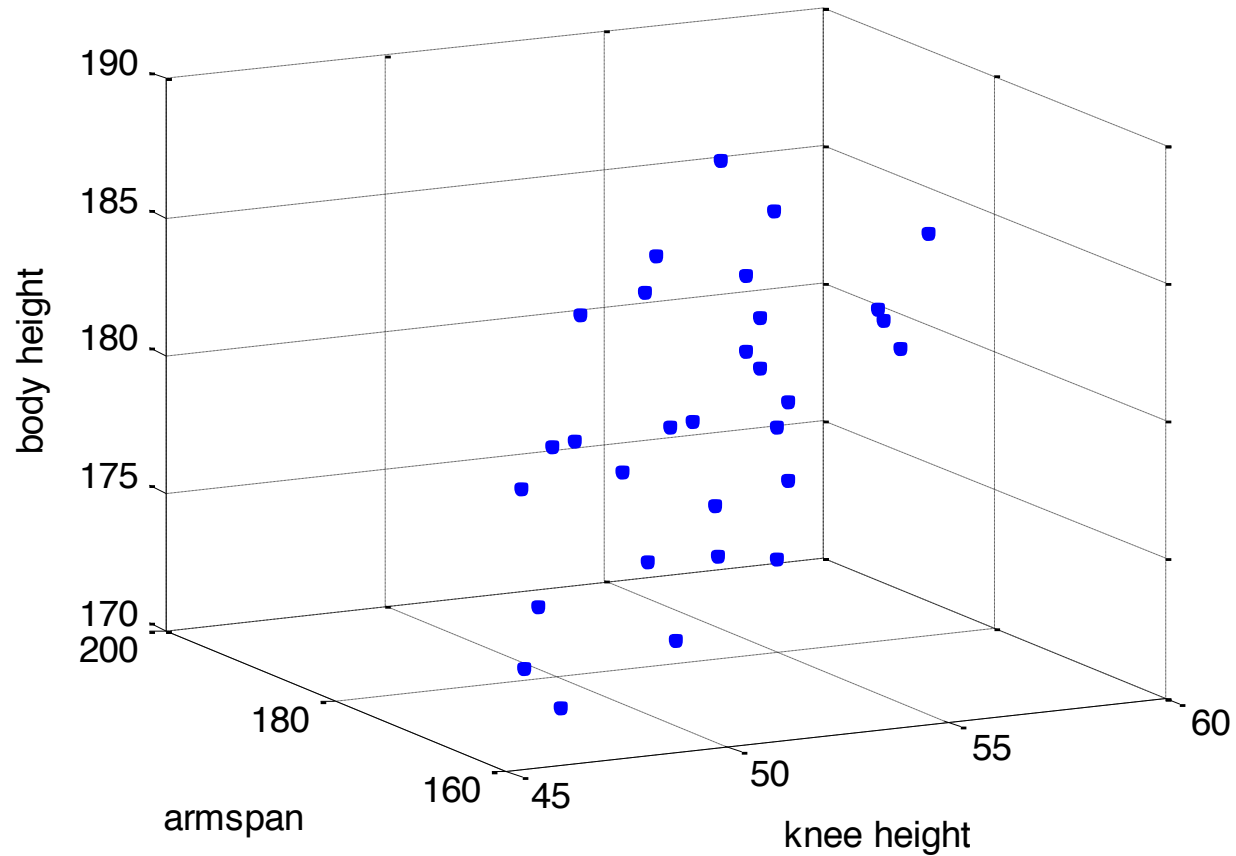
Knee height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...

m=30 data points



Example Data

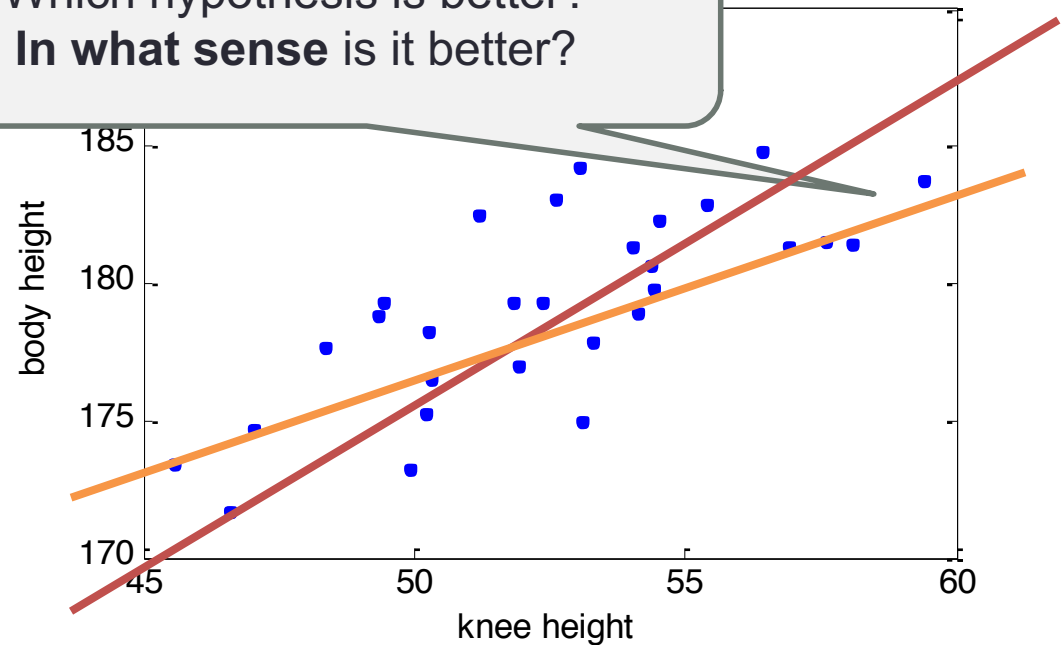
Knee Height [cm]	Arm span [cm]	Height [cm]
50	166	171
56	172	175
52	174	168
...



Linear regression with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...

Which hypothesis is better?
In what sense is it better?



Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

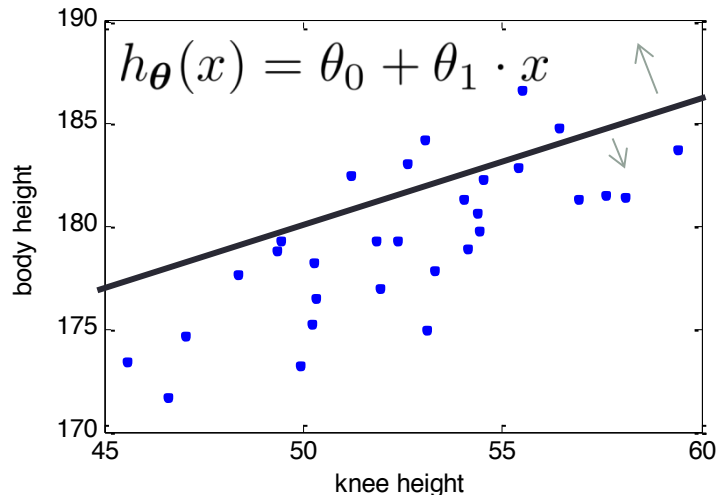
Parameters ?

$$\theta = (\theta_0, \theta_1)$$

Formalization of problem

$x^{(i)}$	Knee Height [cm]	Height [cm]	$y^{(i)}$
	50	171	
	56	175	
	52	168	
	

$m=30$ data points



- Given m training examples

$$\langle x^{(1)}, y^{(1)} \rangle \dots \langle x^{(m)}, y^{(m)} \rangle$$

- Goal: learn parameters

$$\theta = (\theta_0, \theta_1)$$

such that

$$h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 \cdot x^{(i)} \approx y^{(i)}$$

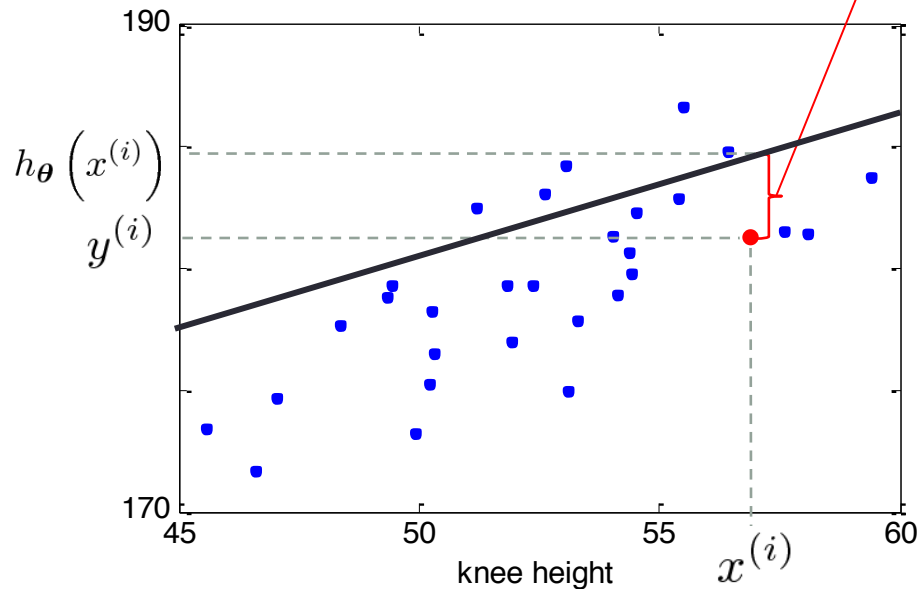
for all training examples $i=1 \dots 30$.

Least Squares Objective

- Minimize Error

$$J(\theta_0, \theta_1) = \left(\underbrace{h_{\theta} \left(x^{(i)} \right) - y^{(i)}}_{\text{error}} \right)^2$$

$$\begin{aligned}\theta_0 &= 150 \\ \theta_1 &= 0.6\end{aligned}$$



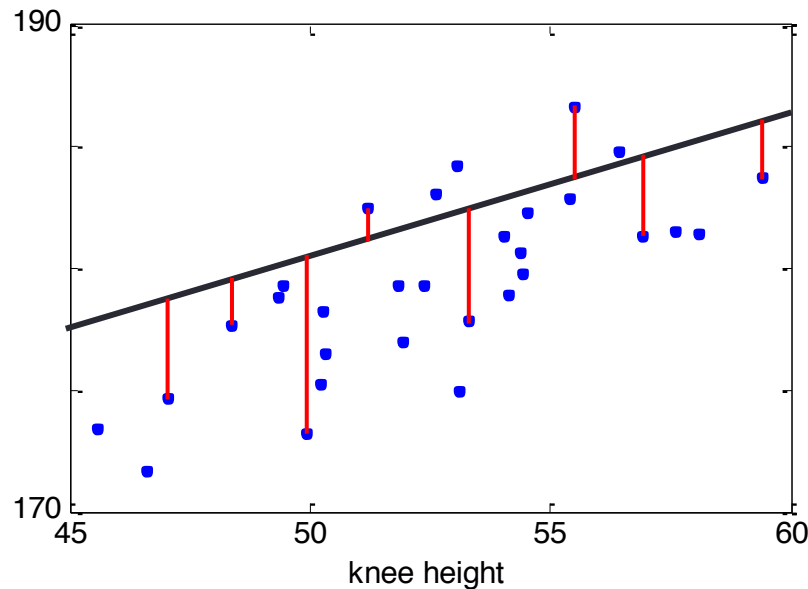
Least Squares Objective

- Minimize Error $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(\underbrace{h_{\theta} \left(x^{(i)} \right) - y^{(i)}}_{\text{mean squared error}} \right)^2$

cost function

mean squared error

$$\begin{aligned}\theta_0 &= 150 \\ \theta_1 &= 0.6 \\ J(\theta_0, \theta_1) &= 10.77\end{aligned}$$



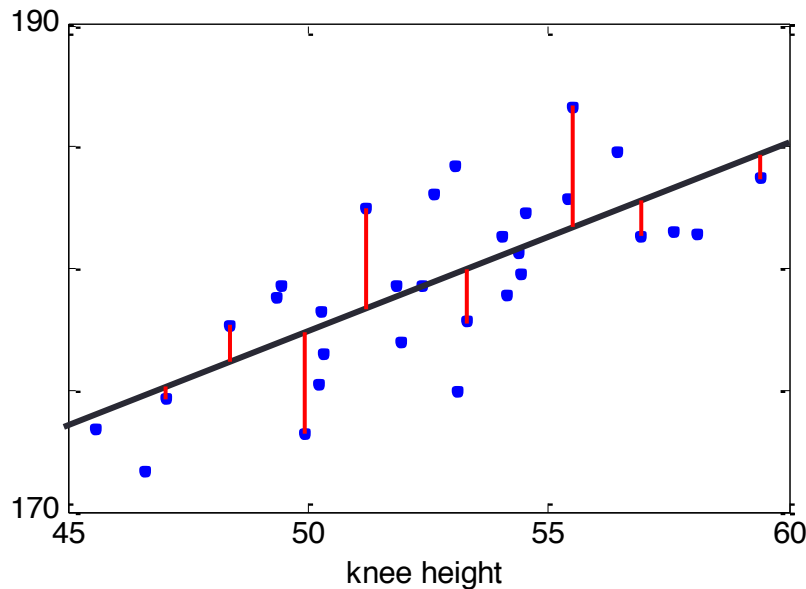
Least Squares Objective

- Minimize Error $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(\underbrace{h_{\theta} \left(x^{(i)} \right) - y^{(i)}}_{\text{mean squared error}} \right)^2$

cost function

mean squared error

$$\begin{aligned}\theta_0 &= 140 \\ \theta_1 &= 0.75 \\ J(\theta_0, \theta_1) &= 5.94\end{aligned}$$

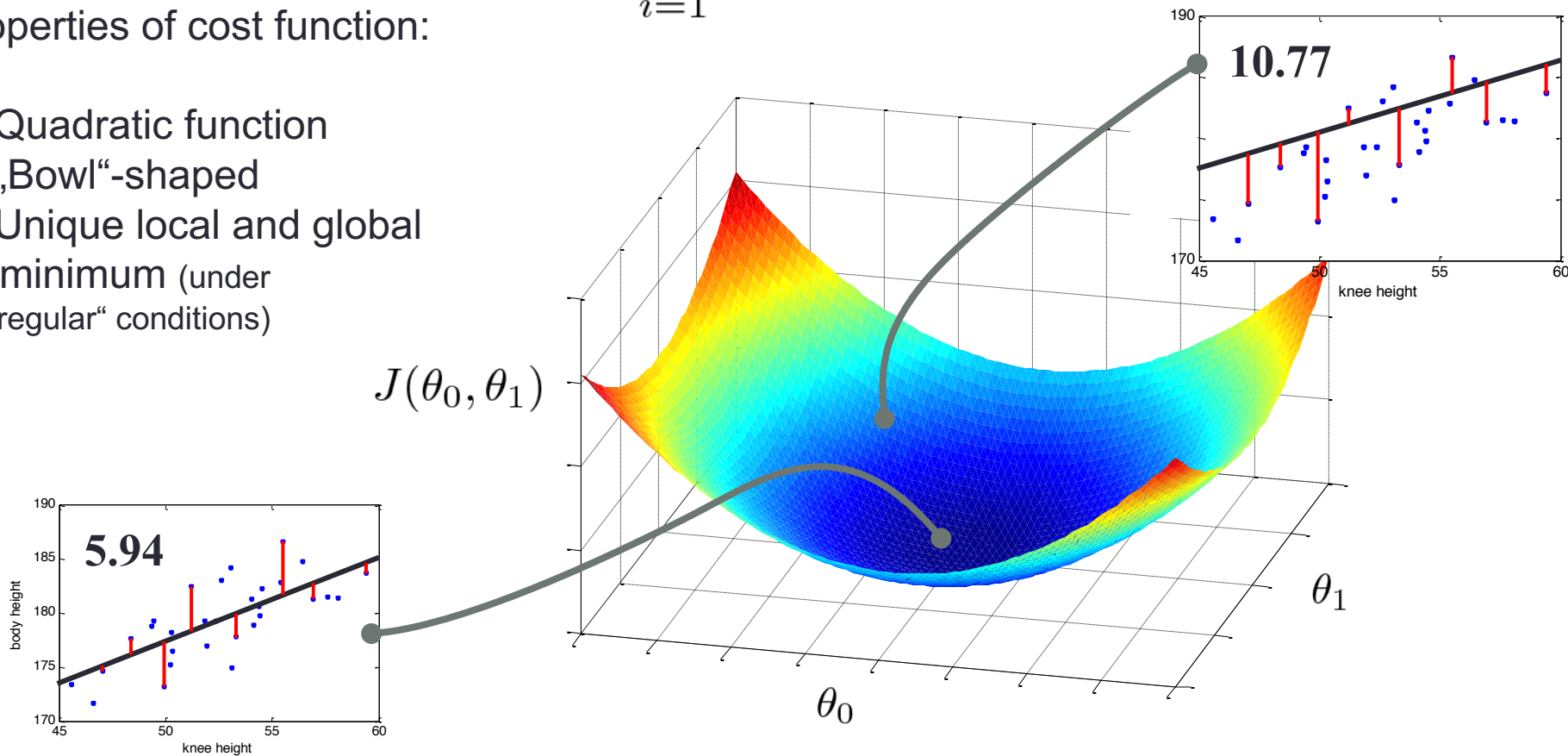


Cost function illustrated

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

Properties of cost function:

- Quadratic function
- „Bowl“-shaped
- Unique local and global minimum (under „regular“ conditions)

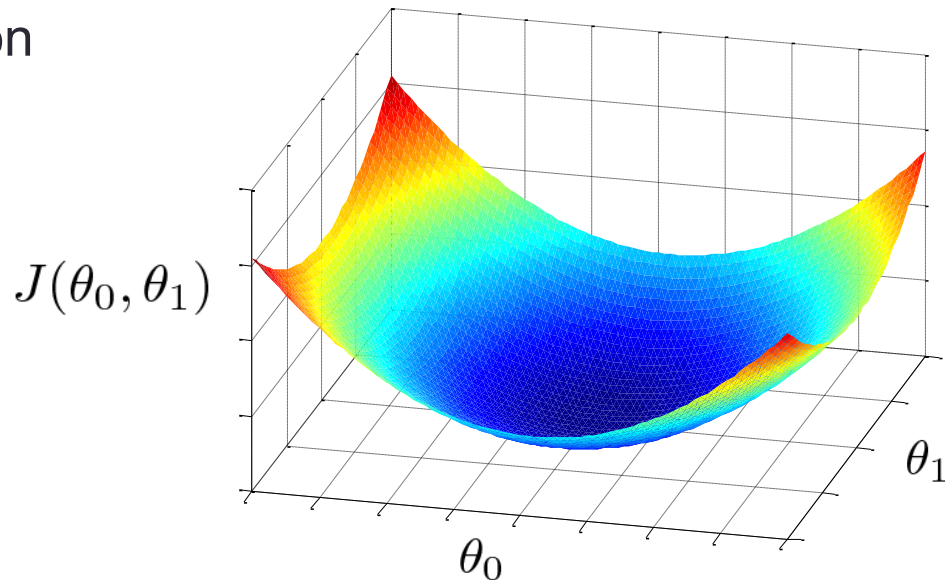


Minimizing the cost

- Two ways to find the parameters $\theta = (\theta_0, \theta_1)$ minimizing

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

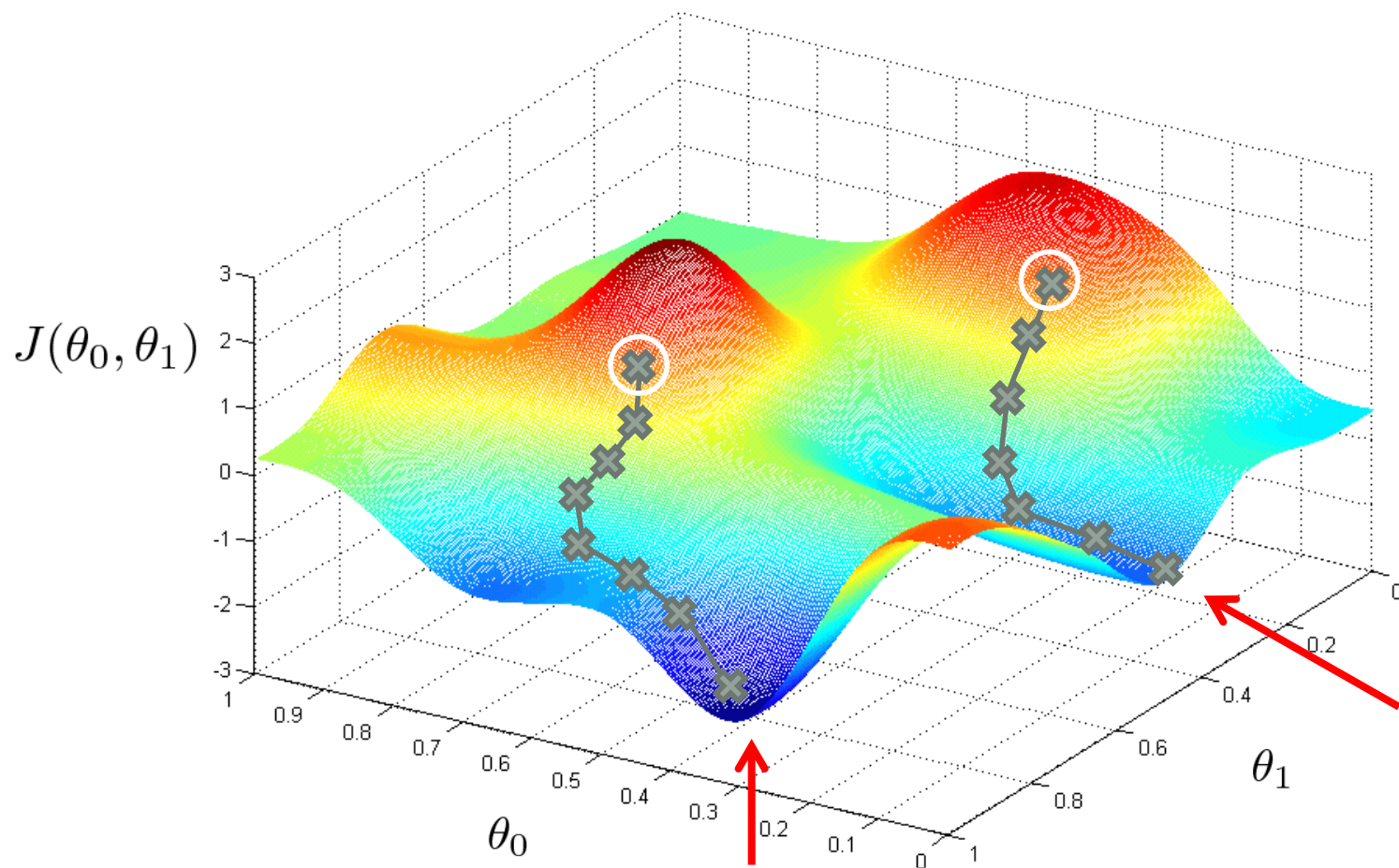
- Gradient descent
- Direct analytical solution (setting derivatives = 0)



GRADIENT DESCENT

Descending in the steepest direction

Gradient descent on some arbitrary cost function $J(\theta_0, \theta_1)$...



Gradient descent algorithm

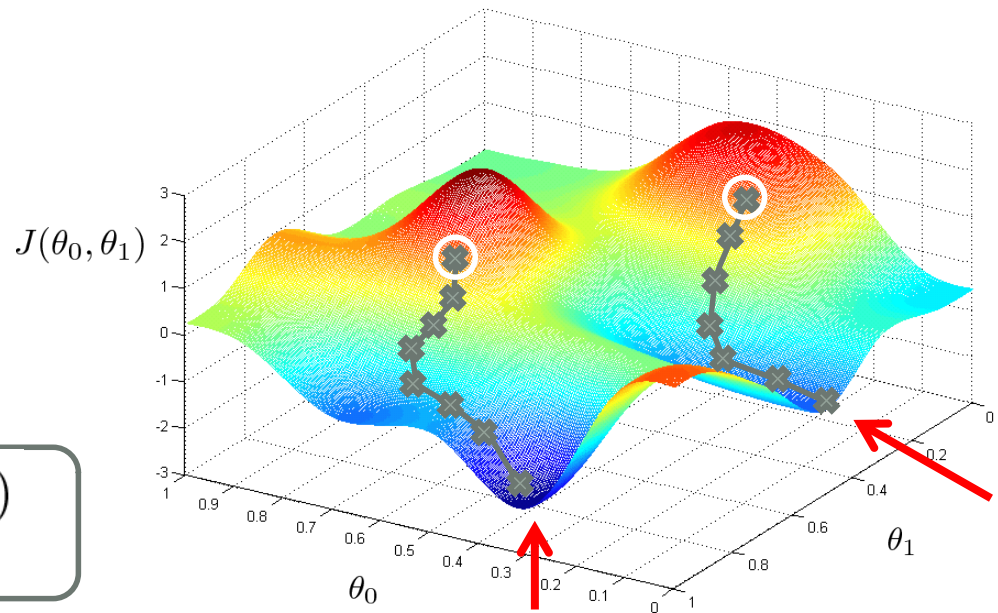
- Repeat until convergence

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{simultaneously updating } \theta_0 \text{ and } \theta_1)$$

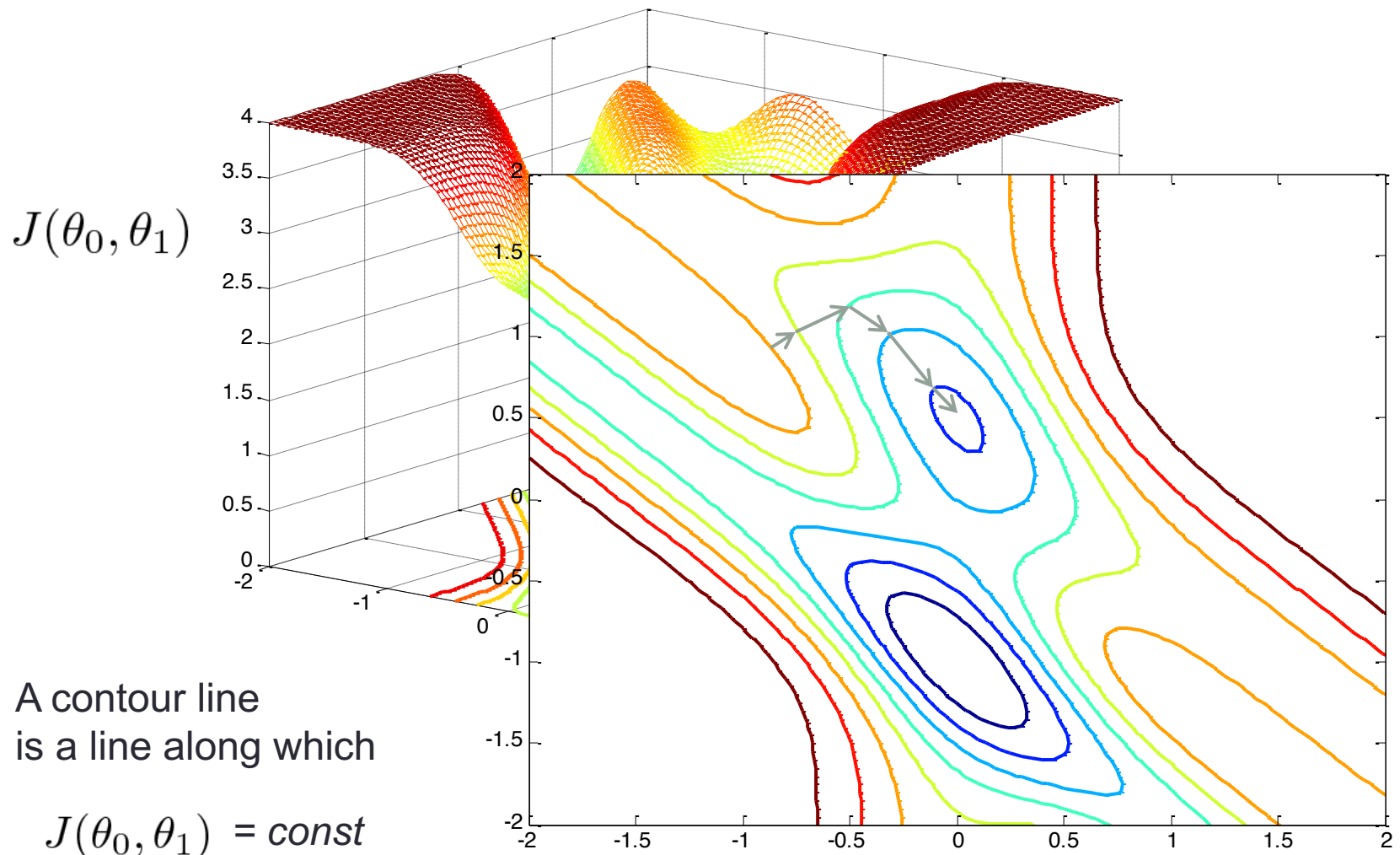
negative gradient =
descent

learning rate („eta“)

partial derivative of $J(\theta_0, \theta_1)$
with respect to θ_j

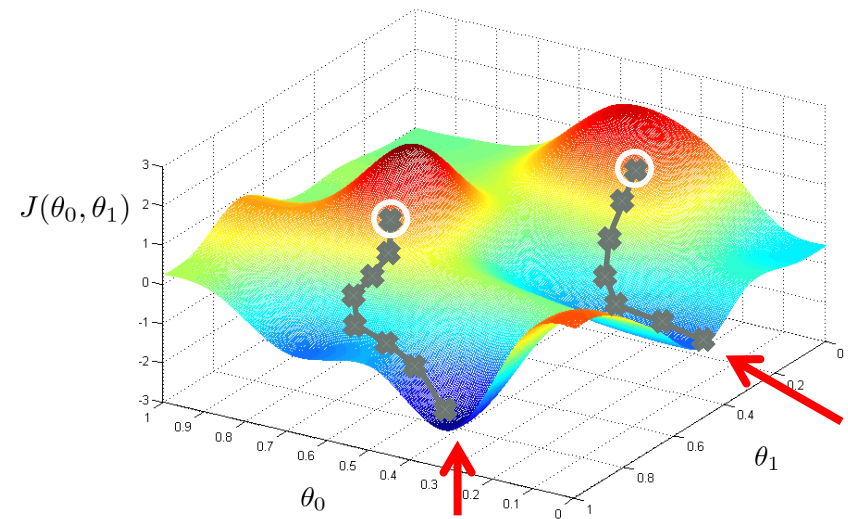


Gradient is orthogonal to contour lines

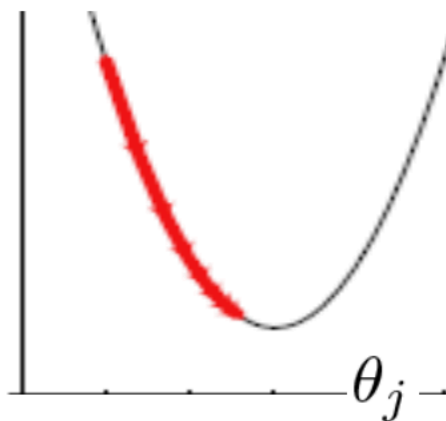


Potential issues with gradient descent

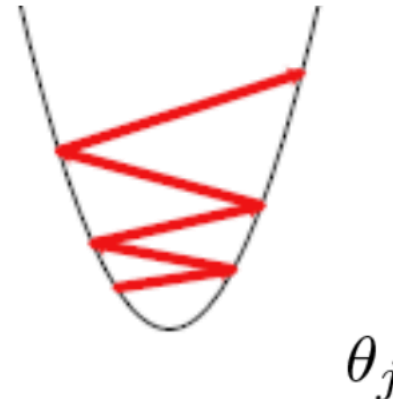
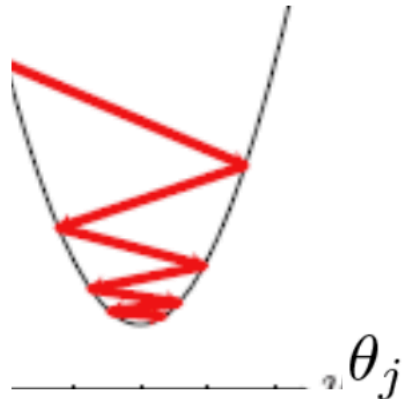
- May get stuck in local minima
- Learning rate too small: slow convergence
- Learning rate too large: oscillations, divergence



η too small



η too large



LINEAR REGRESSION WITH GRADIENT DESCENT

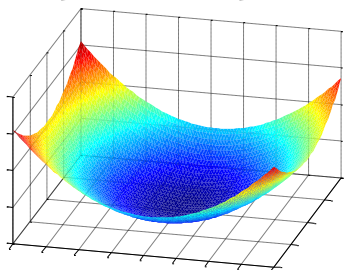
(ONE INPUT)

Application of gradient descent

- Linear regression cost

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$



- Gradient descent

$$\theta_j := \theta_j - \eta \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(simultaneous update)

$$\theta_0 := \theta_0 - \underline{2\eta} \cdot \frac{1}{m} \sum_{i=1}^m \left(\underline{h_{\theta} \left(x^{(i)} \right) - y^{(i)}} \right)$$

"learning rate"

$$\theta_1 := \theta_1 - \underline{2\eta} \cdot \frac{1}{m} \sum_{i=1}^m \left(\underline{h_{\theta} \left(x^{(i)} \right) - y^{(i)}} \right) \cdot \underline{x^{(i)}}$$

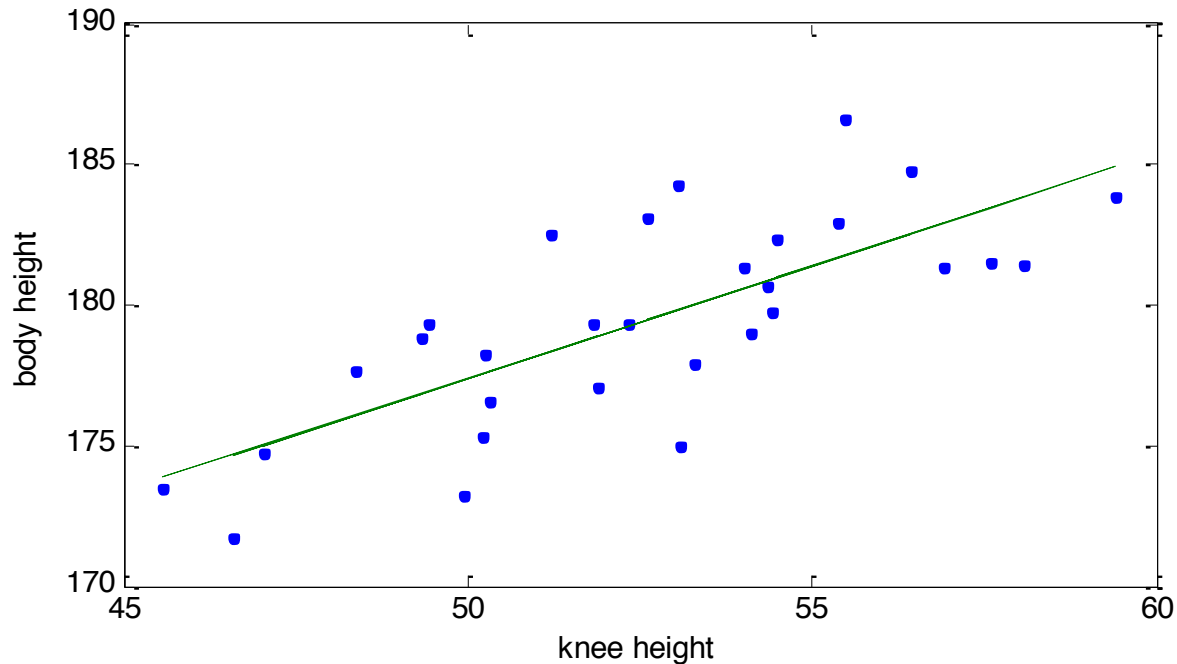
"error"

"input"

(simultaneous
update)

Predicting height from knee height

- Optimal fit to training data



$$\theta_0 = 137.4$$

$$\theta_1 = 0.8$$

LINEAR REGRESSION

MORE GENERAL FORMULATION: MULTIPLE FEATURES

Multiple inputs (features)

	Knee Height	Arm span	Age	Height
	x_1	x_2	x_3	y
m	50	166	32	171
	56	172	17	175
	52	174	62	168

	$n = 3$			

$$\mathbf{x}^{(2)} = \begin{pmatrix} 56 \\ 172 \\ 17 \end{pmatrix}$$

$$x_3^{(2)} = 17$$

- Notation:

m ... number of training examples

n ... number of features

$\mathbf{x}^{(i)}$... input features of i 'th training example (vector-valued)

$x_j^{(i)}$... value of feature j in i 'th training example

Linear hypothesis

- Hypothesis (one input):

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

- Hypothesis (multiple input features):

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n$$

Example: $h(x) = 50 + 0.5 \cdot \text{kneeheight} + 0.3 \cdot \text{armspan} + 0.1 \cdot \text{age}$

- More compact notation:

$$h_{\theta}(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}$$

$$\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix}$$

Introduce $x_0 = 1$
Why? Notation convenience!

Multiple inputs (features) revisited

	Knee Height	Arm span	Age	Height
x_0	x_1	x_2	x_3	y
1	50	166	32	171
1	56	172	17	175
1	52	174	62	168
1

m {

{ $n = 3$

- Notation:
 - m ... number of training examples
 - n ... number of features

$\mathbf{x}^{(i)}$... input features of i 'th training example (vector-valued)
 $x_j^{(i)}$... value of feature j in i 'th training example

$$\mathbf{x}^{(2)} = \begin{pmatrix} 1 \\ 56 \\ 172 \\ 17 \end{pmatrix}$$

$$x_0^{(2)} = 1$$

$$x_3^{(2)} = 17$$

Matrix and vector notation

	Knee Height	Arm span	Age	Height
x_0	x_1	x_2	x_3	y
1	50	166	32	171
1	56	172	17	175
1	52	174	62	168

$$\mathbf{X} = \begin{pmatrix} 1 & 50 & 166 & 32 \\ 1 & 56 & 172 & 17 \\ 1 & 52 & 174 & 62 \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} 171 \\ 175 \\ 168 \end{pmatrix}$$

$$\mathbf{x}^{(i)} = \begin{pmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix}$$

features of i 'th training example
 $(n+1) \times 1$

$$\mathbf{X} = \begin{pmatrix} \text{---} (\mathbf{x}^{(1)})^T \text{---} \\ \text{---} (\mathbf{x}^{(2)})^T \text{---} \\ \vdots \\ \text{---} (\mathbf{x}^{(m)})^T \text{---} \end{pmatrix}$$

design matrix
 $m \times (n+1)$

$$\mathbf{y} = \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

output/target vector
 $m \times 1$

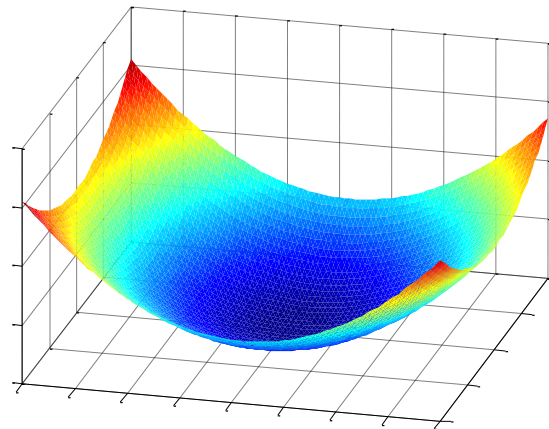
LINEAR REGRESSION WITH GRADIENT DESCENT

(GENERAL FORMULATION)

Linear regression problem statement

- Hypothesis: $h_{\theta}(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}$
- Cost function: $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$

high-dimensional quadratic
(„bowl“-shaped) function



Goal is to find parameters which minimize the cost

Gradient descent (multiple features)

with **one** input feature:

$$\begin{aligned}\theta_0 &:= \theta_0 - \underbrace{2\eta}_{\text{"learning rate"}} \cdot \frac{1}{m} \sum_{i=1}^m \underbrace{\left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)}_{\text{"error"}} \\ \theta_1 &:= \theta_1 - \underbrace{2\eta}_{\text{"learning rate"}} \cdot \frac{1}{m} \sum_{i=1}^m \underbrace{\left(h_{\theta} \left(x^{(i)} \right) - y^{(i)} \right)}_{\text{"error"}} \cdot \underbrace{x^{(i)}}_{\text{"input"}}\end{aligned}\quad \left. \vphantom{\sum_{i=1}^m} \right\} \begin{array}{l} \text{(simultaneous} \\ \text{update)}\end{array}$$

with **n** input features:

$$\theta_j := \theta_j - \underbrace{2\eta}_{\text{"learning rate"}} \cdot \frac{1}{m} \sum_{i=1}^m \underbrace{\left(h_{\theta} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)}_{\text{"error"}} \cdot \underbrace{x_j^{(i)}}_{\text{"input"}} \quad \left. \vphantom{\sum_{i=1}^m} \right\} \begin{array}{l} \text{(simultaneous} \\ \text{update for} \\ j=0\dots n) \end{array}$$

For $j = 0$: define for convenience $x_0^{(i)} = 1$

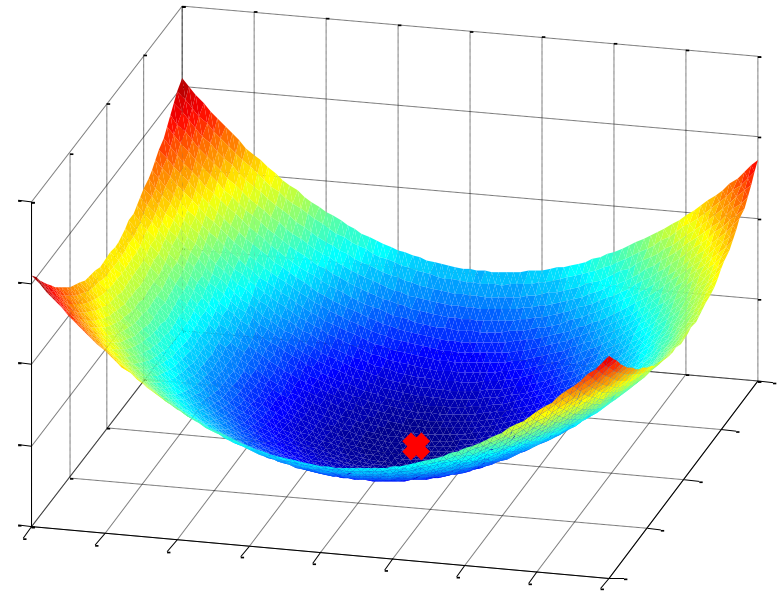
LINEAR REGRESSION ANALYTICAL SOLUTION

Analytical solution

- Set all partial derivatives of cost function $J(\boldsymbol{\theta}) = 0$
- Solving system of linear equations yields:

$$\boldsymbol{\theta}^* = \underbrace{\left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T}_{\text{Moore-Penrose Pseudoinverse of } \mathbf{X}} \mathbf{y}$$

Moore-Penrose Pseudoinverse of \mathbf{X}



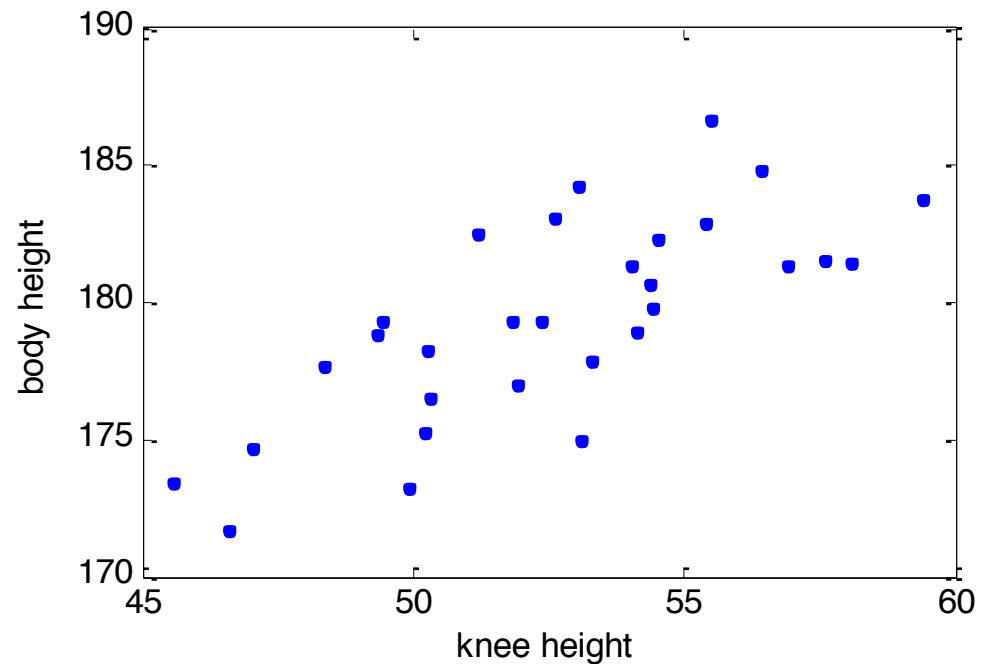
\mathbf{X} ... design matrix

\mathbf{y} ... output/target vector

- *Note: This analytical solution requires that columns of \mathbf{X} are linearly independent („regular“ conditions)*

Example: analytical solution applied to problem with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...



Example: analytical solution applied to problem with one input

Knee Height [cm]	Height [cm]
50	171
56	175
52	168
...	...

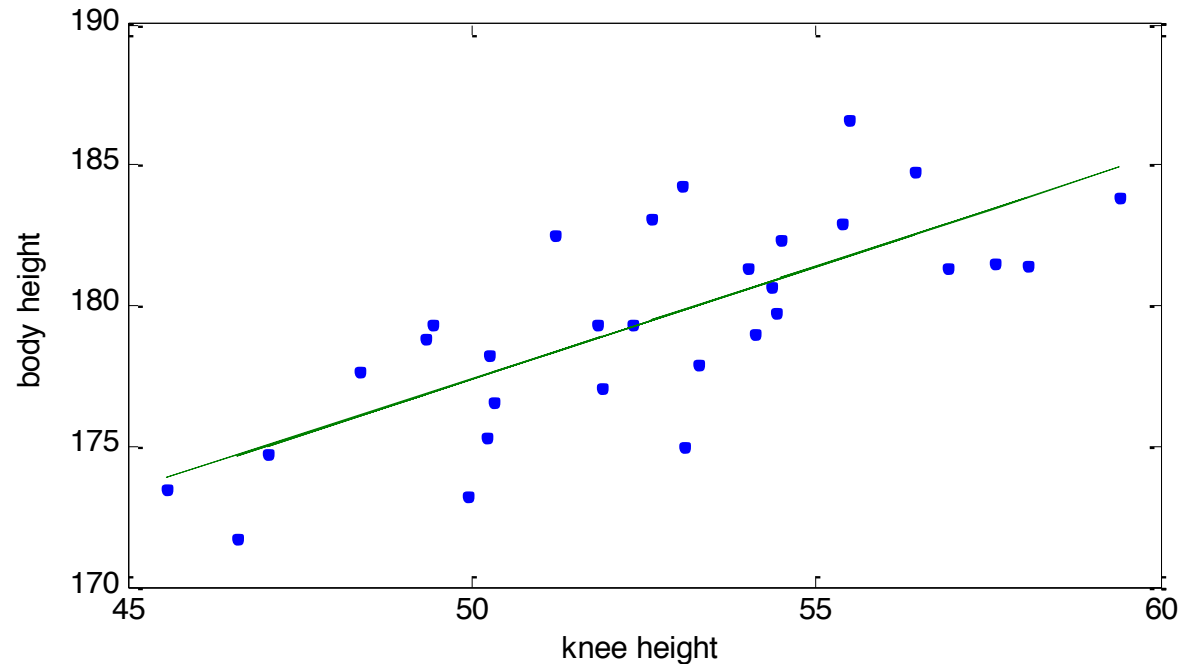
$$\mathbf{X} = \begin{pmatrix} 1 & 50 \\ 1 & 56 \\ 1 & 52 \\ \vdots & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} 171 \\ 175 \\ 168 \\ \vdots \end{pmatrix}$$

30×2 30×1

$$\begin{aligned} \theta^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \begin{pmatrix} 137.4 \\ 0.8 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \mathbf{X}^T \mathbf{X} &= \begin{pmatrix} 30 & 1577 \\ 1577 & 83222 \end{pmatrix} & 2 \times 2 \\ (\mathbf{X}^T \mathbf{X})^{-1} &= \begin{pmatrix} 7.994 & -0.152 \\ -0.152 & 0.003 \end{pmatrix} & 2 \times 2 \\ \mathbf{X}^T \mathbf{y} &= \begin{pmatrix} 5383 \\ 283210 \end{pmatrix} & 2 \times 1 \end{aligned}$$

Predicting height from knee height



$$\theta_0 = 137.4$$

$$\theta_1 = 0.8$$

$$\begin{aligned}\theta^* &= \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \\ &= \begin{pmatrix} 137.4 \\ 0.8 \end{pmatrix}\end{aligned}$$

Gradient descent

- Need to choose learning rate η
- Iterative algorithm (needs many iterations to converge)
- Works well even when number of input features is large n

Analytical solution

- No need to choose η
- Direct solution (no iteration)
- Slow if n is too large (inverting $n \times n$ matrix)

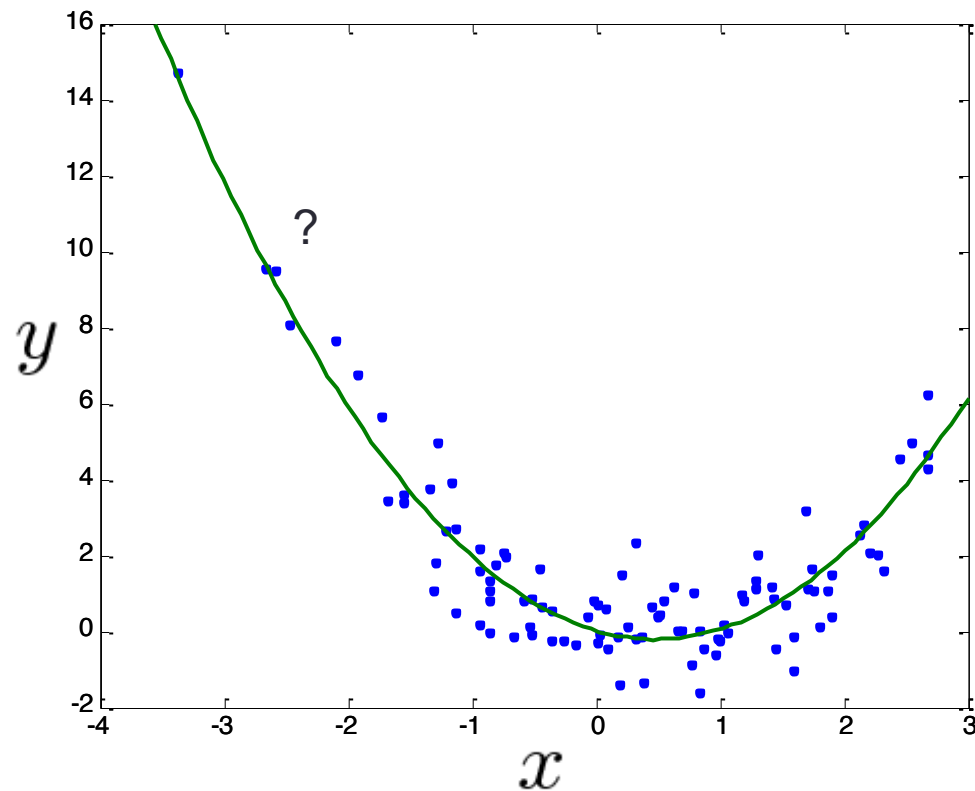
NON-LINEAR FEATURES

(NON-LINEAR BASIS FUNCTIONS)

Non-linear trends in data

- How can we learn non-linear hypotheses?

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...



$$h_{\theta}(x) = \overset{?}{\theta_0} + \overset{?}{\theta_1} \cdot x + \overset{?}{\theta_2} \cdot x^2$$

Linear fit to this “non-linear” data

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

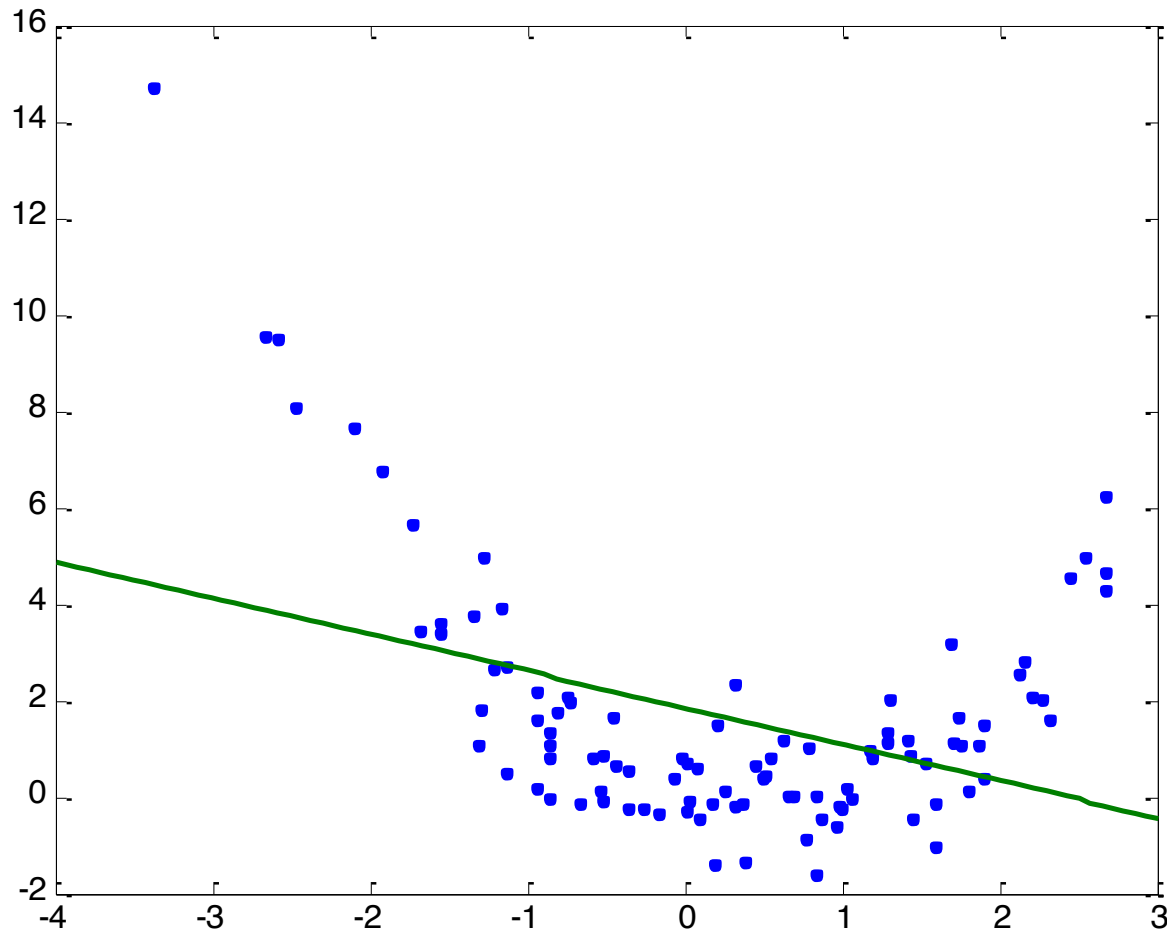
$$\mathbf{X} = \begin{pmatrix} 1 & 0.01 \\ 1 & -1.22 \\ 1 & 0.17 \\ & \vdots \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

standard design matrix

Hypothesis: $h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 \cdot x$

Optimal parameters: $\boldsymbol{\theta}^* = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$

Linear fit to this “non-linear” data



$$h_{\theta}(x) = 1.85 - 0.76 \cdot x$$

Non-linear (quadratic) fit

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

$$\phi_0 = 1 \quad \phi_1 = x \quad \phi_2 = x^2$$

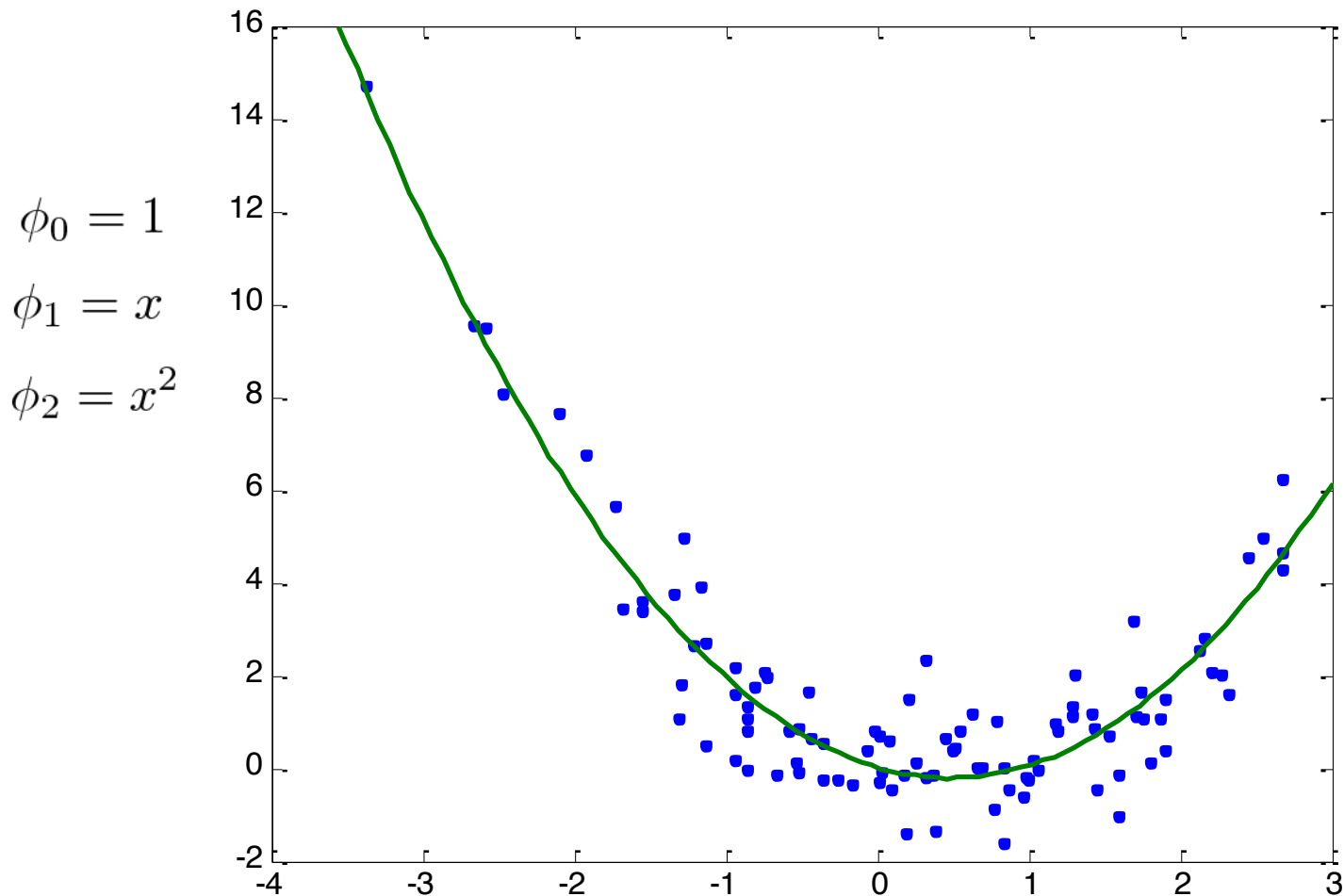
$$\Phi = \begin{pmatrix} 1 & 0.01 & 0.01^2 \\ 1 & -1.22 & (-1.22)^2 \\ 1 & 0.17 & (0.17)^2 \\ \vdots & & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

*design matrix with
non-linear features*

Hypothesis: $h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \theta_0 + \theta_1 \cdot \phi_1 + \theta_2 \cdot \phi_2$

Optimal parameters: $\boldsymbol{\theta}^* = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y}$

Non-linear (quadratic) fit



$$h_{\theta}(x) = 0.02 \cdot 1 - 0.95 \cdot x + 0.99 \cdot x^2$$

Non-linear (sinusoid) fit

x	y
0.01	-0.27
-1.22	2.63
0.17	-0.13
...	...

$$\phi_0 = 1 \quad \phi_1 = x \quad \phi_2 = \cos(x)$$

$$\Phi = \begin{pmatrix} 1 & 0.01 & \cos(0.01) \\ 1 & -1.22 & \cos(-1.22) \\ 1 & 0.17 & \cos(0.17) \\ \vdots & & \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} -0.27 \\ 2.63 \\ -0.13 \\ \vdots \end{pmatrix}$$

*design matrix with
non-linear features*

Hypothesis: $h_{\boldsymbol{\theta}}(\boldsymbol{\phi}) = \theta_0 + \theta_1 \cdot \phi_1 + \theta_2 \cdot \phi_2$

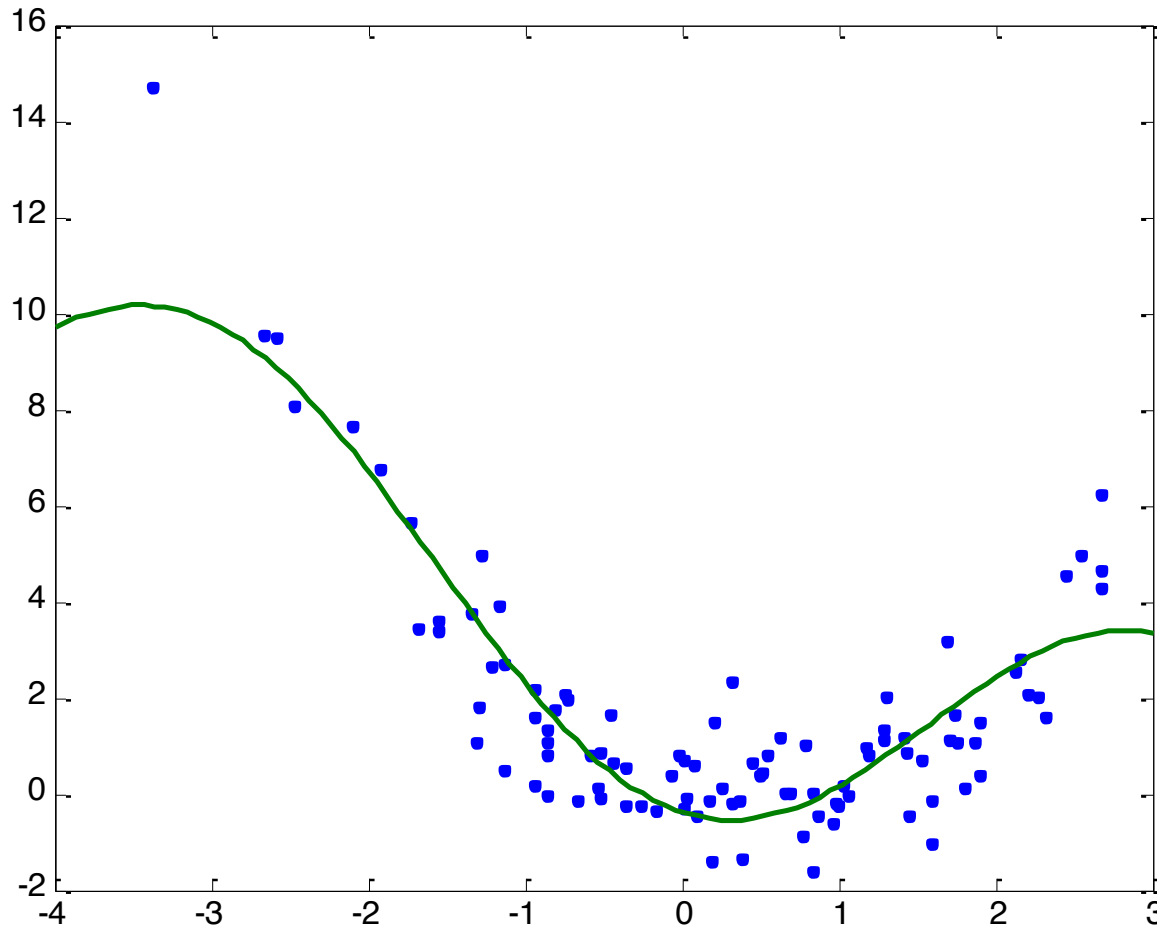
Optimal parameters: $\boldsymbol{\theta}^* = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y}$

Non-linear (sinusoidal) fit

$$\phi_0 = 1$$

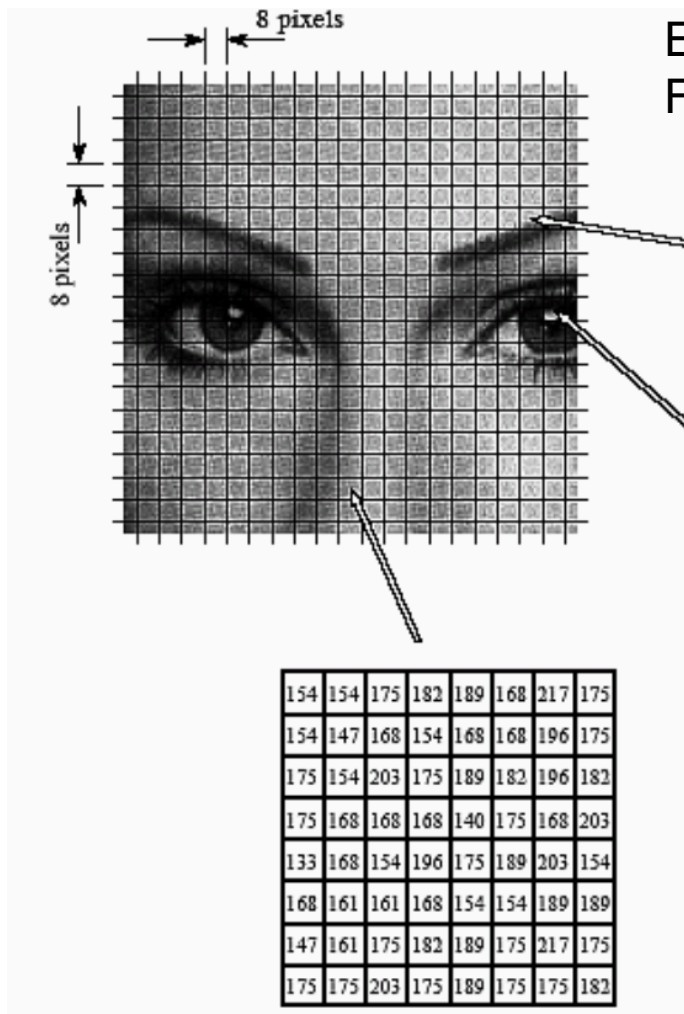
$$\phi_1 = x$$

$$\phi_2 = \cos(x)$$

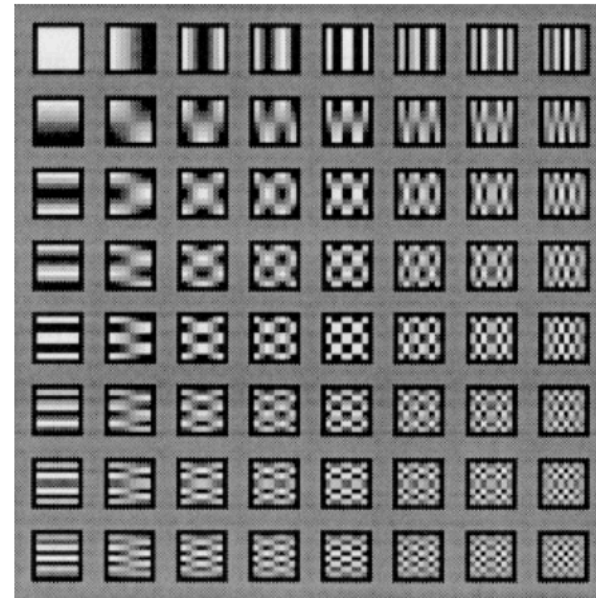


$$h_{\theta}(x) = 3.12 \cdot 1 - 1.07 \cdot x - 3.5 \cdot \cos(x)$$

Image: JPEG = cosine-basis



Each block of 8x8 pixels is represented in a Fourier basis of cosine filters



Better representation of edges and corners
Allows for compression

Non-linear input features (in general)

feature 2 of all training examples

$$\Phi = \begin{pmatrix} 1 & \phi_1^{(1)} & \phi_2^{(1)} & \dots & \phi_n^{(1)} \\ 1 & \phi_1^{(2)} & \phi_2^{(2)} & \dots & \phi_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1^{(m)} & \phi_2^{(m)} & \dots & \phi_n^{(m)} \end{pmatrix}$$

all features of 1st training example

- Feature 2 for each training example i is computed by applying a **non-linear basis function**:

$$\phi_2^{(i)} = \phi_2(\mathbf{x}^{(i)})$$

- Allows to learn a variety of non-linear functions with the same technique(s):
 - Analytical

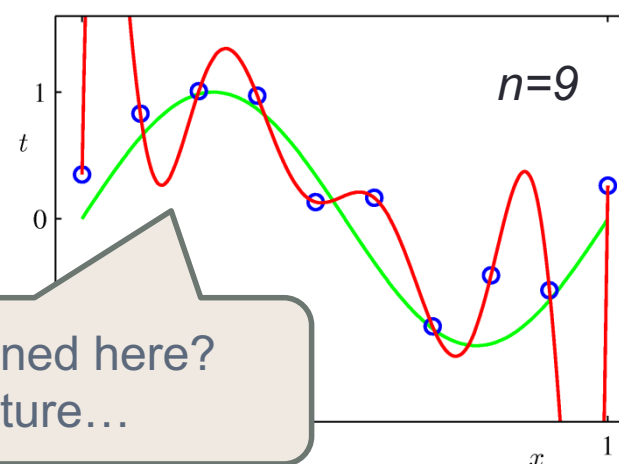
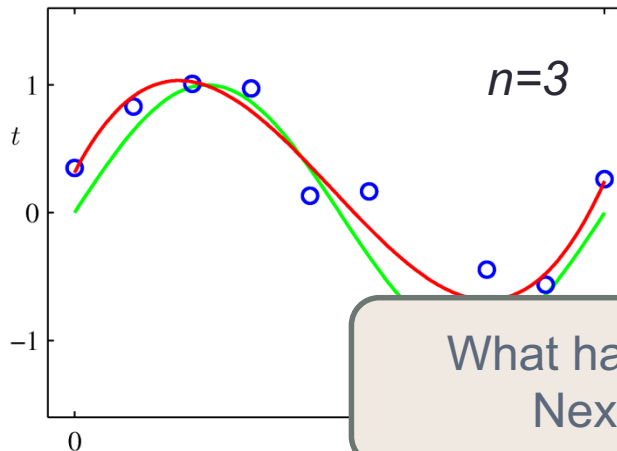
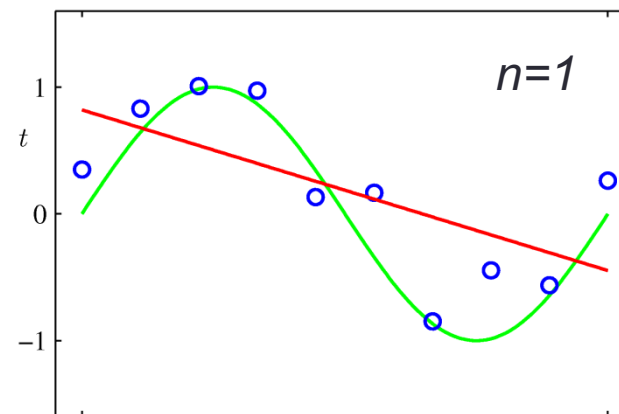
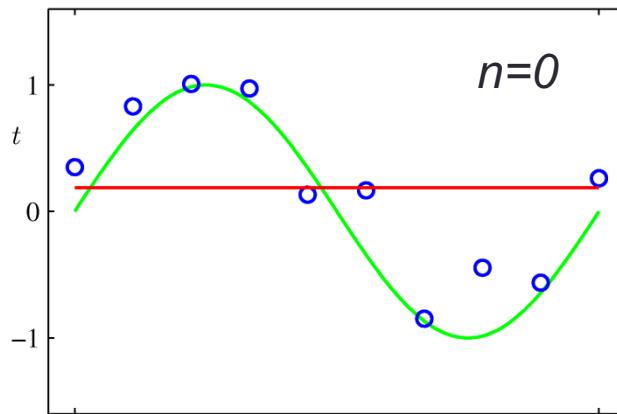
$$\theta^* = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} \quad \text{or gradient descent}$$

Polynomial regression

- Features are powers of x

n = degree of polynome
to be learned

$$\phi_0 = x^0, \phi_1 = x^1, \phi_2 = x^2, \dots, \phi_n = x^n$$

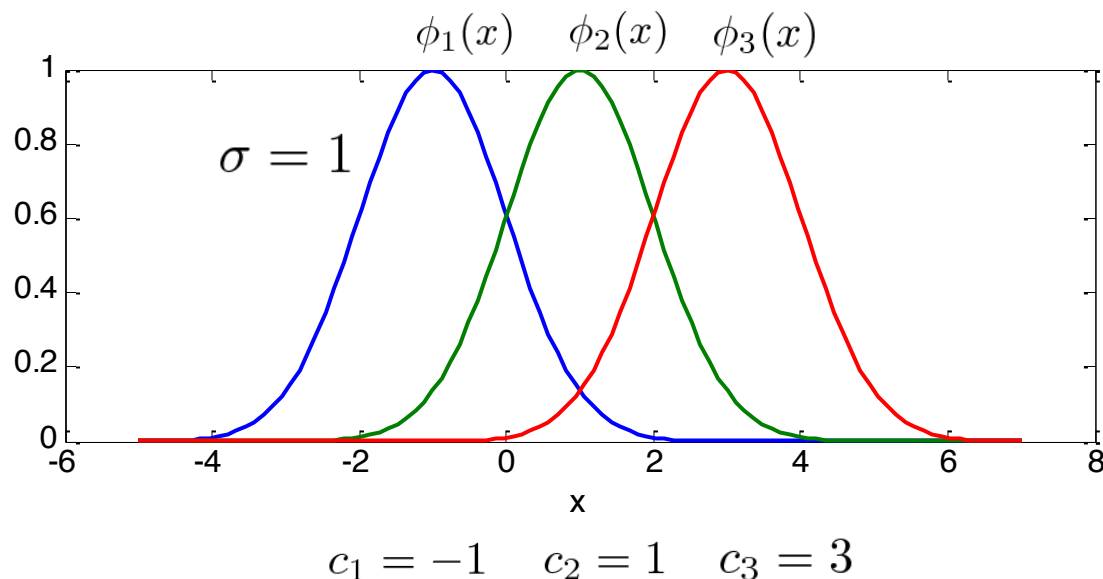


What happened here?
Next lecture...

Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ .

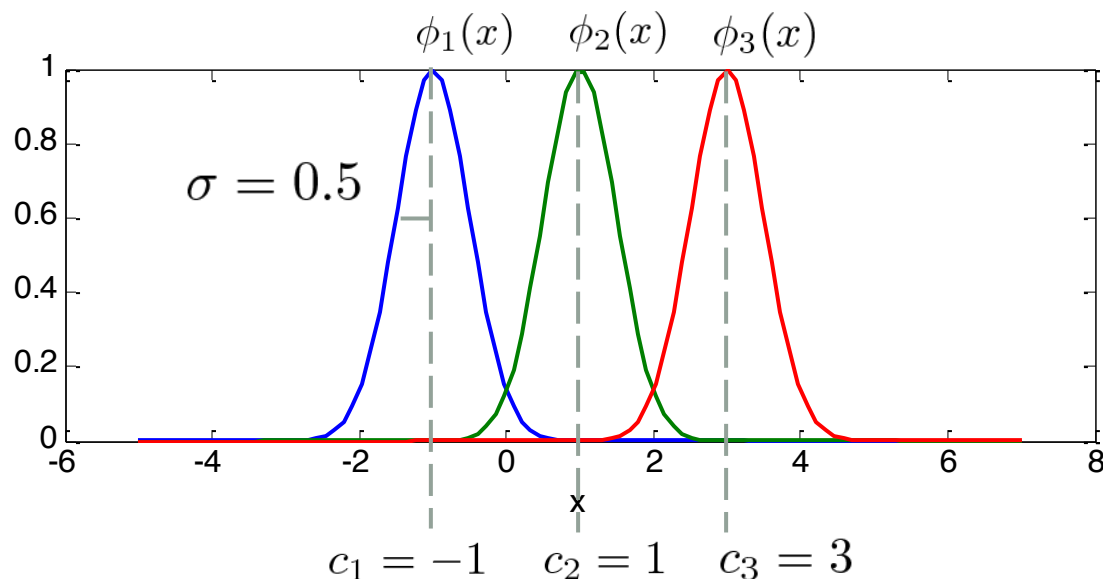
$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ .

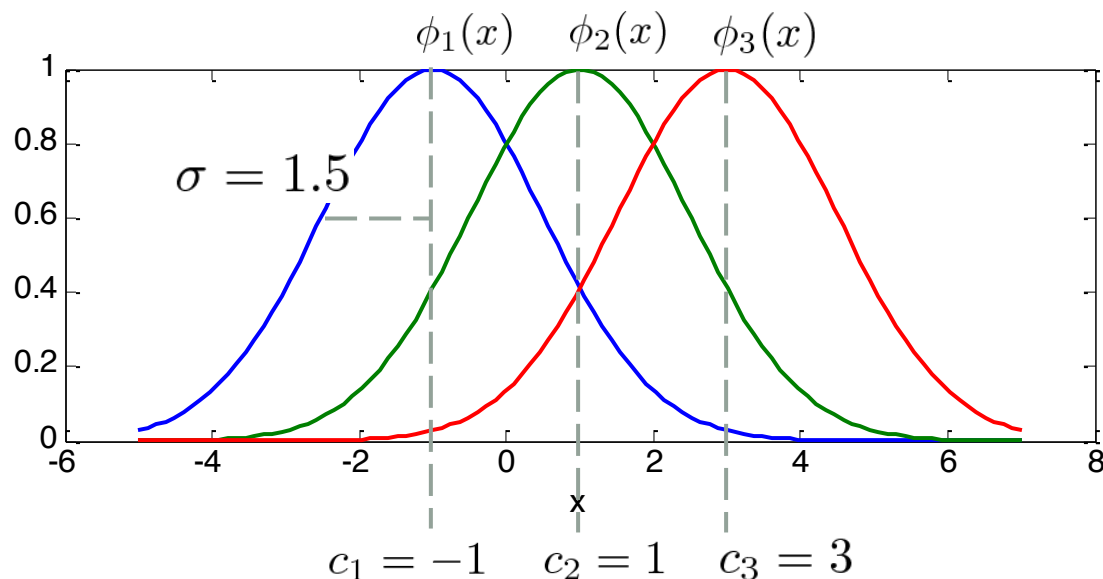
$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



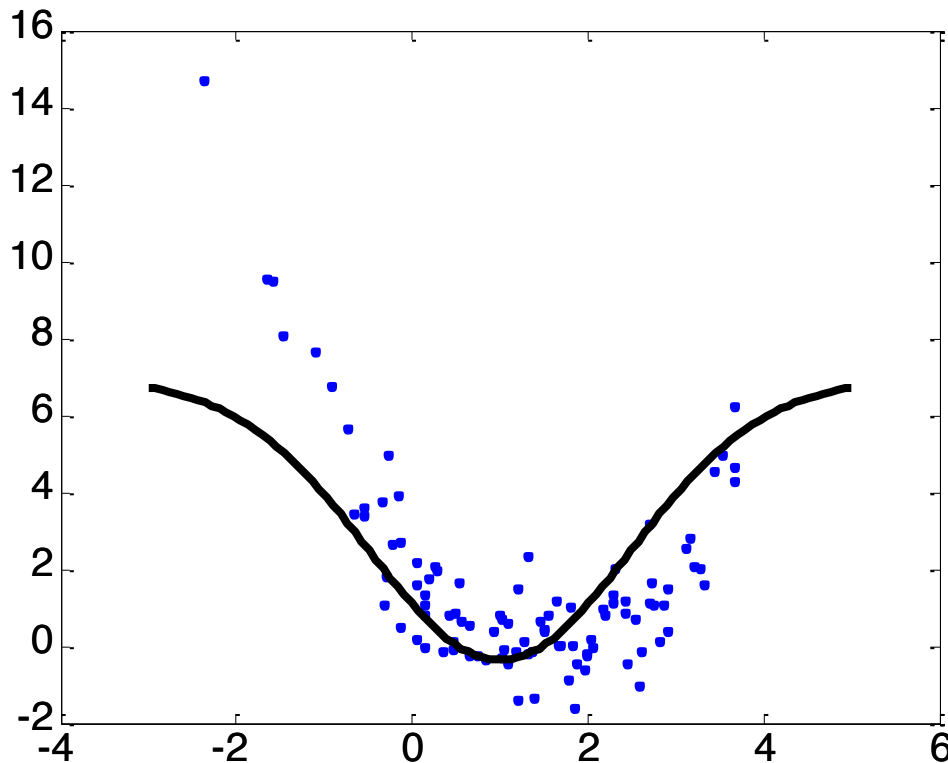
Radial basis functions

- „Gaussian“-shaped RBFs:
 - Each basis function j has a **center** \mathbf{c}_j in the input space
 - The **width** of the basis functions is determined by σ .

$$\phi_j(\mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{c}_j\|^2\right)$$



Fitting a single RBF to data

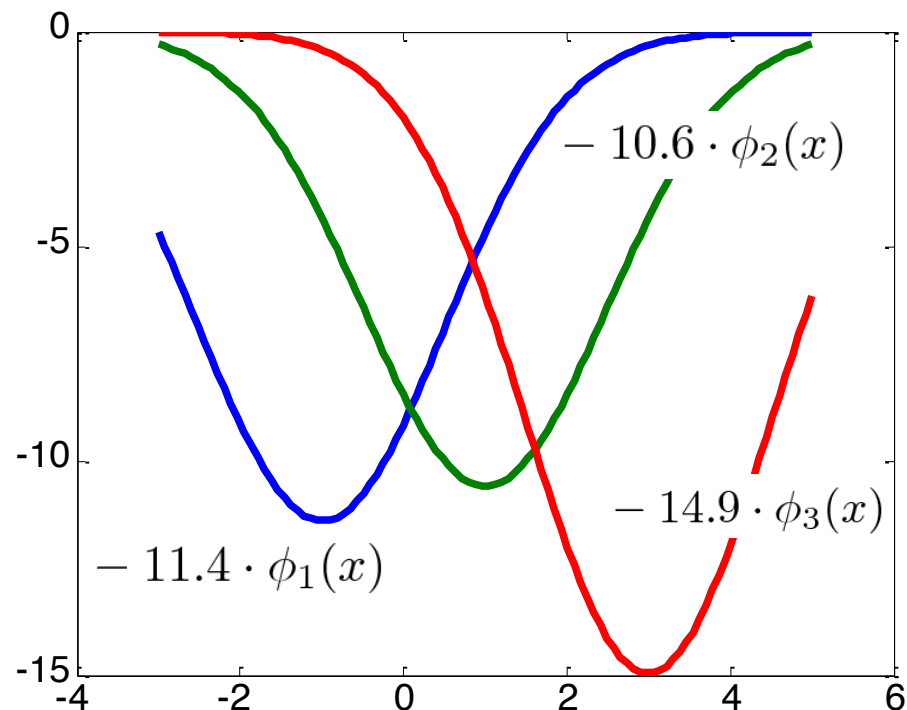
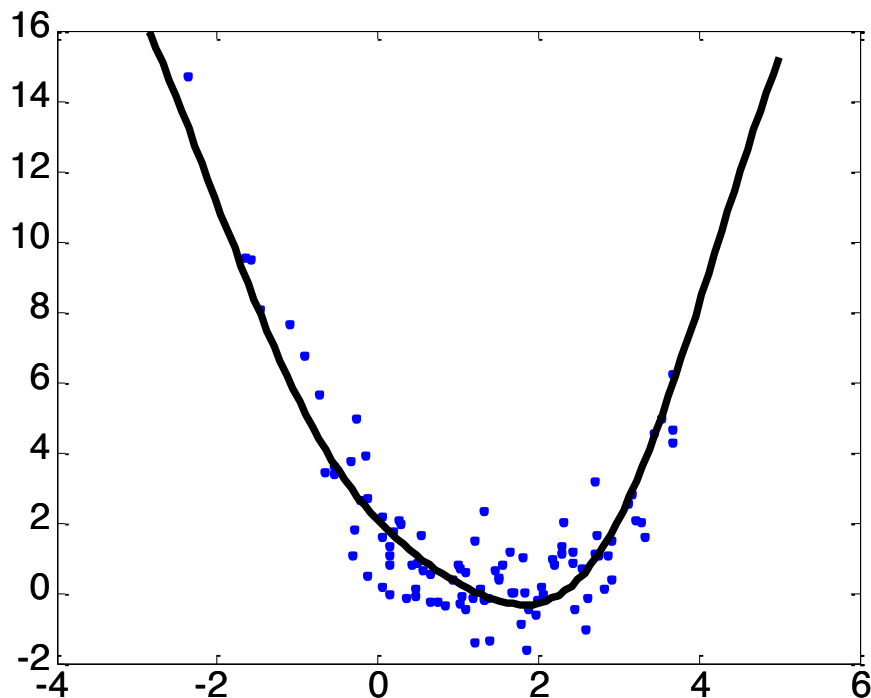


$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \phi_1(x)$$

$$h_{\theta}(x) = 6.9 - 7.3 \cdot \phi_1(x)$$

RBF with
 $c_1 = 1$
 $\sigma = 1.5$

Fitting RBFs to data



RBFs with $\sigma = 1.5$

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot \phi_1(x) + \theta_2 \cdot \phi_2(x) + \theta_3 \cdot \phi_3(x)$$

$$h_{\theta}(x) = 21.7 - 11.4 \cdot \phi_1(x) - 10.6 \cdot \phi_2(x) - 14.9 \cdot \phi_3(x)$$

SUMMARY (QUESTIONS)

Some questions...

- Hypothesis for linear regression = ?
- Cost function for linear regression = ?
- How many local minima may the cost function for lin. reg. have (under regular conditions)?
- Name two ways to minimize the cost function?
- General gradient descent formula?
- Linear regression with gradient descent formula?
- What issues can arise during gradient descent?
- What is the design matrix? What are its dimensions?
- Analytical solution for linear regression = ?
 - What are the components of the solution?
- Pros and Cons of gradient descent vs. analytical solution?
- How can one learn non-linear hypotheses with linear regression?
- What is polynomial regression?
- What are radial basis functions?

What is next?

- Classification with Logistic Regression
- Gradient descent tricks & more advanced optimization techniques
- Underfitting & Overfitting
- Model selection (Training, Validation and test set)