

# Linear Regression

Szymon Bobek

Institute of Applied Computer science  
AGH University of Science and Technology

Based on Carlos Guestrin and Emily Fox slides from  
Coursera Specialization on Machine Learning  
<http://geist.agh.edu.pl>



# Outline I

- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

# Presentation Outline

- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
- 5 Solution to linear regression problem

## 1 Regression

- Linear regression
- Ridge regression
- Bias, variance tradeoff

## 2 Classification

- Logistic regression
- Support Vector Machines
- Decision trees

# Presentation Outline

- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
- 5 Solution to linear regression problem

# Problem



# Problem



# Problem





# Presentation Outline

- 1 Roadmap
- 2 Use case
- 3 Model**
- 4 Simple linear regression
- 5 Solution to linear regression problem

# Model

## Data



$$(x_1 = 150 \text{ m}^2, y_1 = 100\,000\$)$$



$$(x_2 = 150 \text{ m}^2, y_2 = 100\,000\$)$$



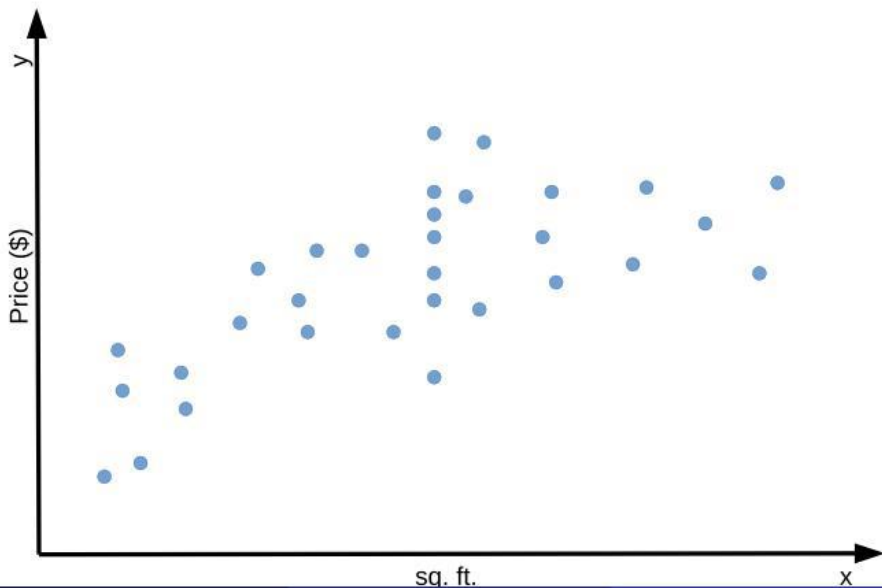
$$(x_3 = 150 \text{ m}^2, y_3 = 100\,000\$)$$

...

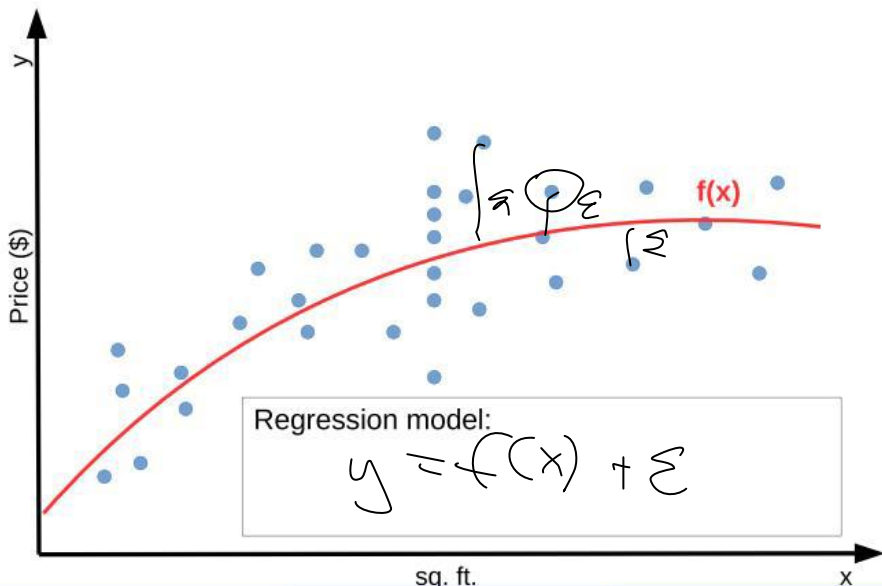


$$(x_m = 150 \text{ m}^2, y_m = 100\,000\$)$$

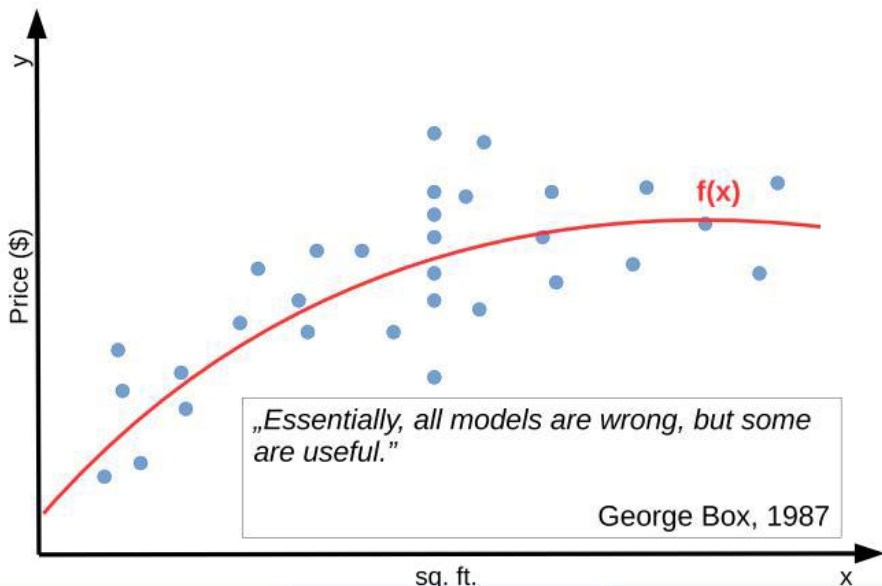
# Model



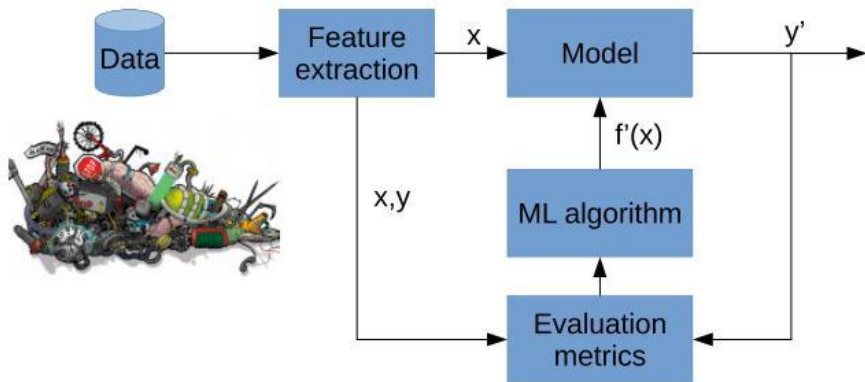
# Model



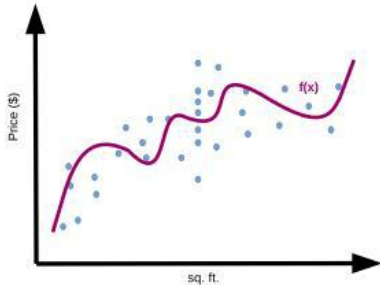
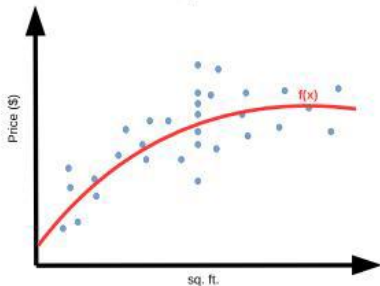
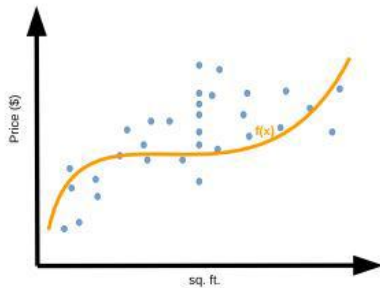
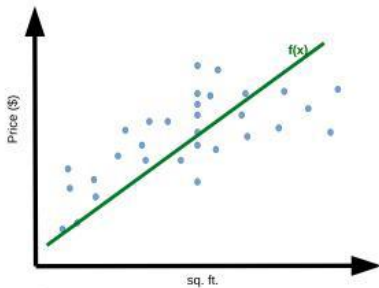
# Model



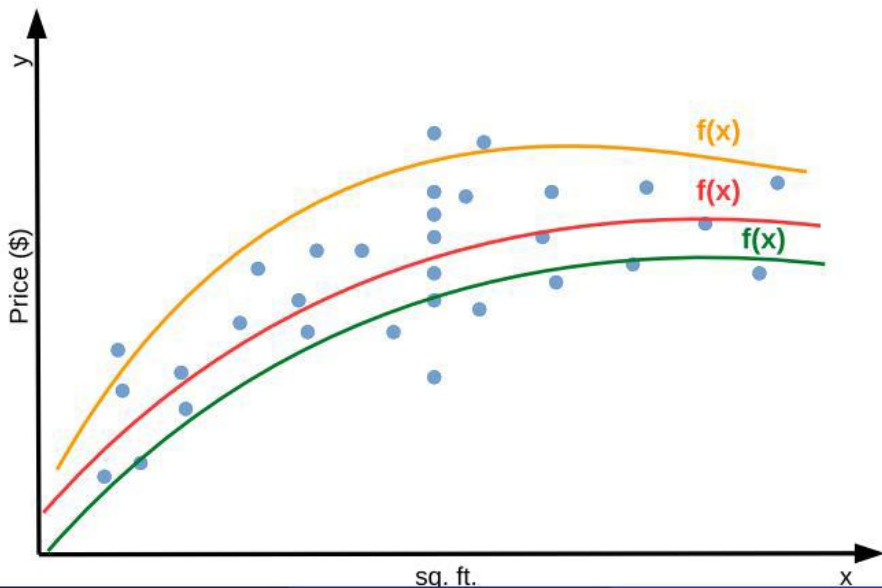
# How to get the model



# Which model?



# Which model?





# Presentation Outline

1 Roadmap

2 Use case

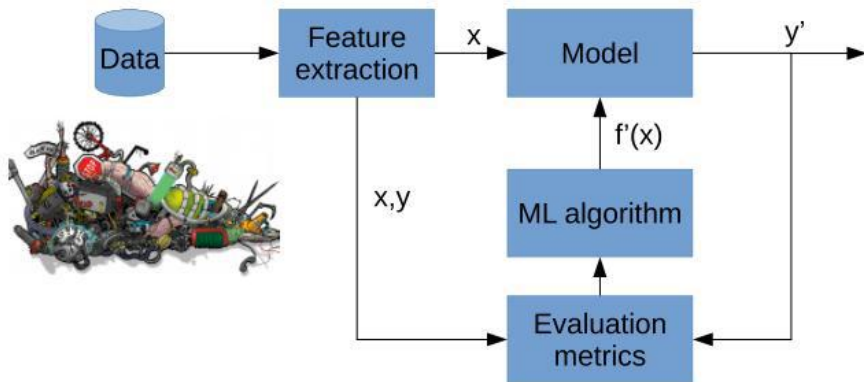
3 Model

4 Simple linear regression

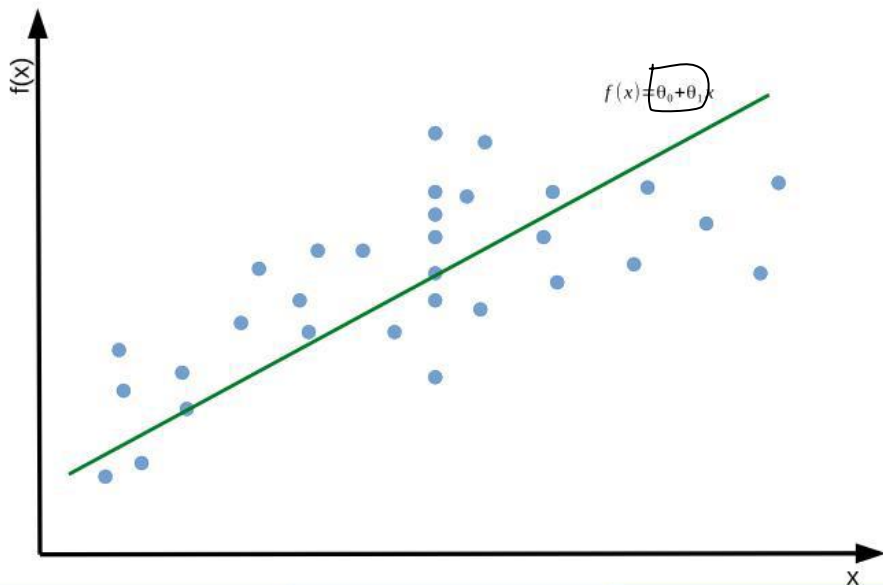
- General idea
- Why RSS?
- Other cost functions
- Linearity in linear regression

5 Solution to linear regression problem

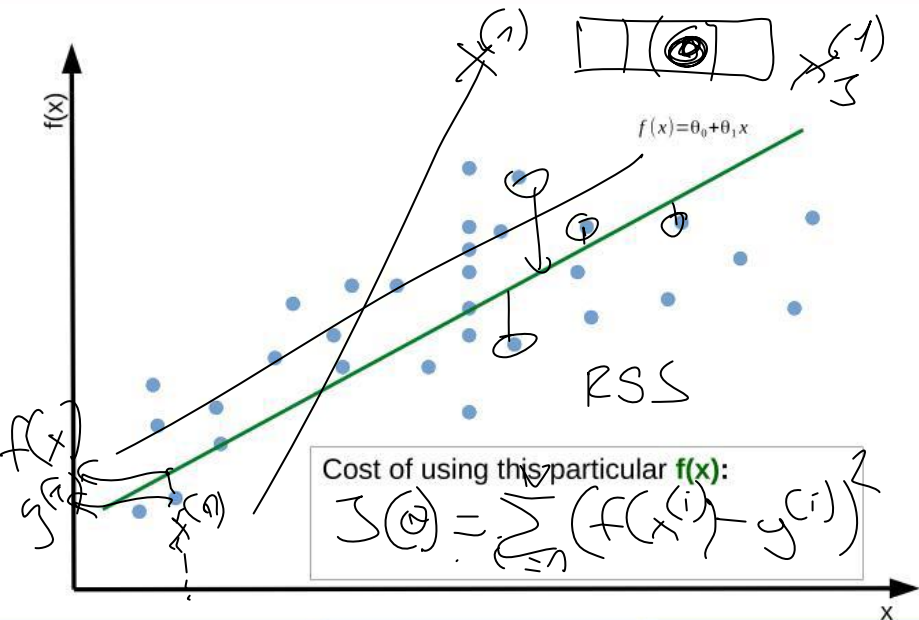
# Representing a model



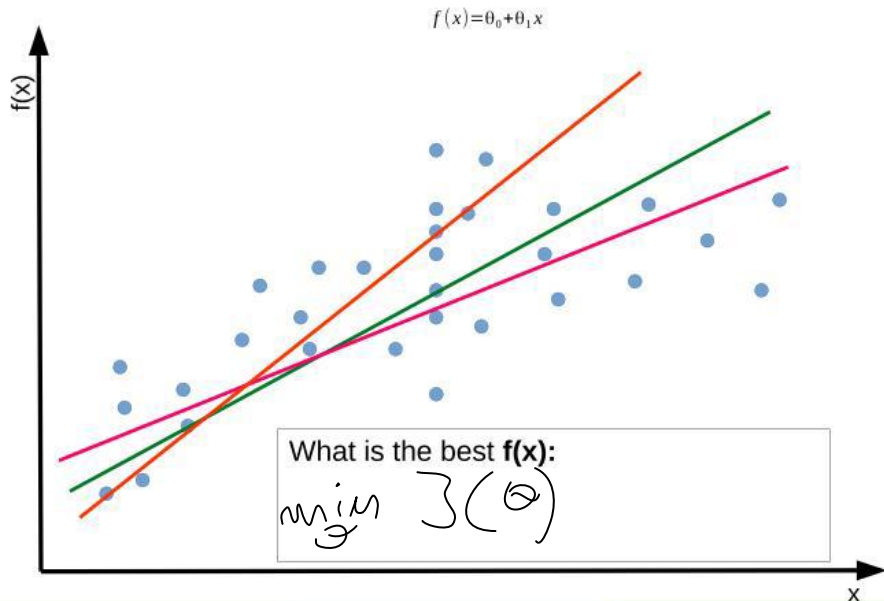
# Representing a model



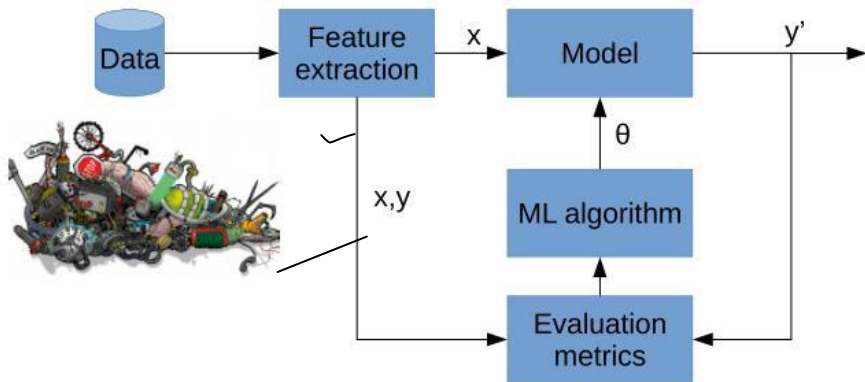
## Cost of using a model



# Cost of using a model



# Cost of using a model



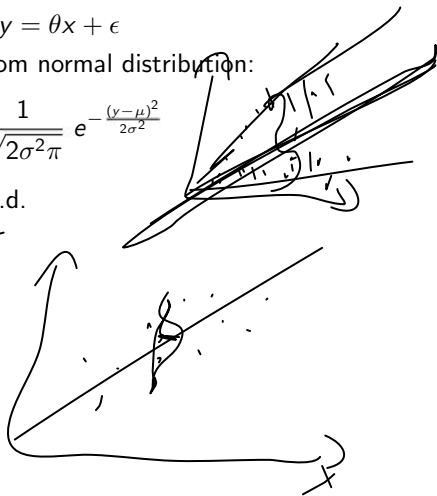
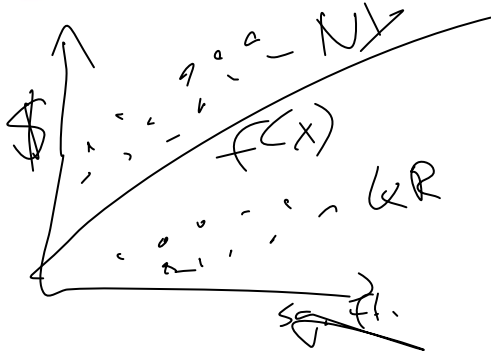
- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

# Gaussian interpretation

- Prediction is linear function and noise:  $y = \theta x + \epsilon$
- We assume that the noise  $\epsilon$  is drawn from normal distribution:

$$\mathcal{N}(y \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

- We assume that training samples are i.i.d.





# Gaussian interpretation

- We want to learn

$$P(y | \theta, x, \sigma^2) = \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{(y - \sum_i \theta_i x_i)^2}{2\sigma^2}}$$

$f(x) = \theta_0 + \theta_1 x_1$   
~~Wrong~~

- The best  $P$  is when probability is max for every training set:


maximize  $\ln P(D | \theta, \sigma^2) = \frac{1}{\sqrt{2\sigma^2\pi}} \prod_j e^{-\frac{(y^{(j)} - \sum_i \theta_i x_i^{(j)})^2}{2\sigma^2}}$

~~$= \ln \frac{1}{\sqrt{2\sigma^2\pi}}$~~

$\sum \ln e^{-\frac{(y^{(i)} - \dots)^2}{2\sigma^2}}$

$+\sum \frac{(y^{(i)} - \dots)^2}{2\sigma^2} = \text{RSS}$

~~$2\sigma^2$~~



# Other cost functions

- Residual sum of squares:

$$RSS(w) = \sum_i^N (y^{(i)} - \theta x^{(i)})^2$$

- Mean squared error:

$$\frac{1}{N} \sum_i^N (y^{(i)} - \theta x^{(i)})^2$$

- Log error:

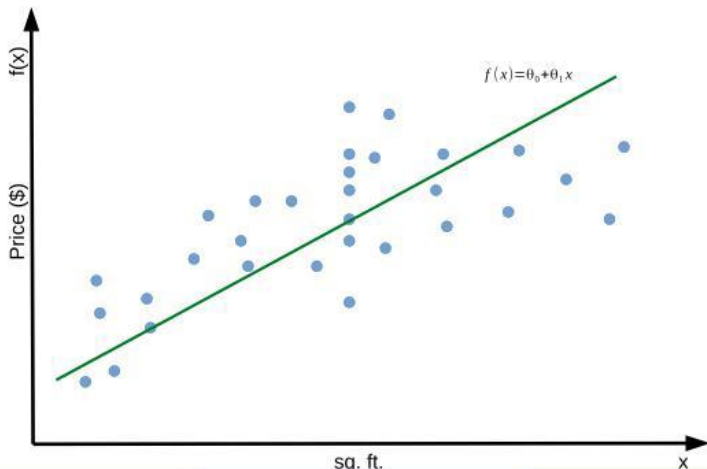
$$\frac{1}{N} \sum_{i=1}^N (-y^{(i)} \log(\theta x^{(i)}) - (1 - y^{(i)}) \log(1 - \theta x^{(i)}))$$

- Asymmetric:

$$\frac{1}{n} \sum_{i=1}^n \left| \alpha - \mathbb{1}_{(y^{(i)} - f(x^{(i)})) < 0} \right| \cdot \left( y^{(i)} - f(x^{(i)}) \right)^2$$

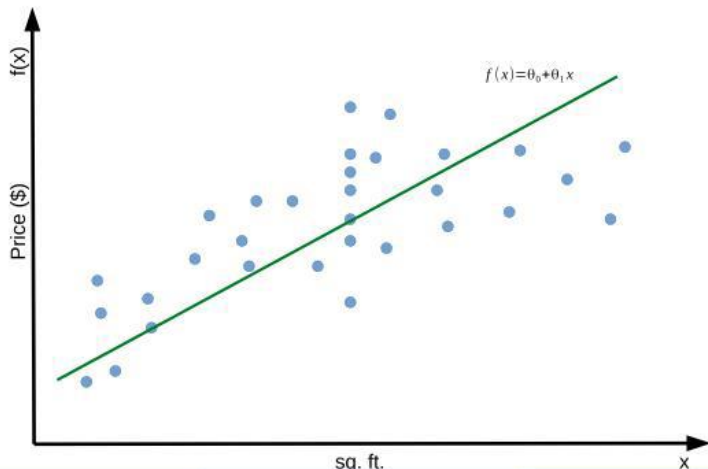
# Assymetric cost function

$$\frac{1}{n} \sum_{i=1}^n \left| \alpha - \mathbb{1}_{(y^{(i)} - f(x^{(i)})) < 0} \right| \cdot \left( y^{(i)} - f(x^{(i)}) \right)^2$$



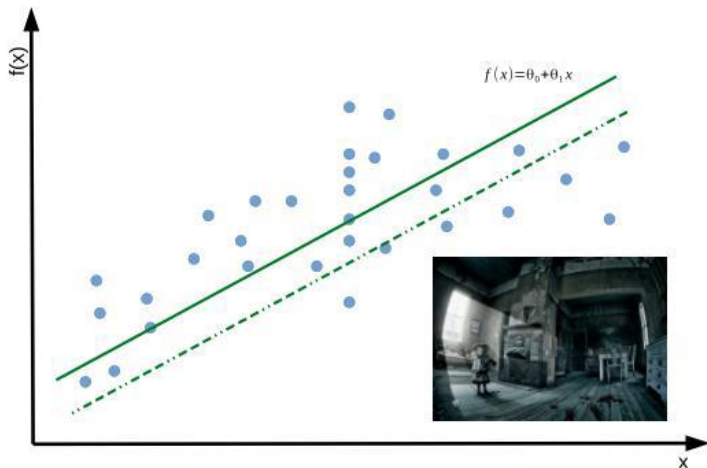
# Assymetric cost function

$$\frac{1}{n} \sum_{i=1}^n \left| \alpha - \mathbb{1}_{(y^{(i)} - f(x^{(i)})) < 0} \right| \cdot \left( y^{(i)} - f(x^{(i)}) \right)^2$$



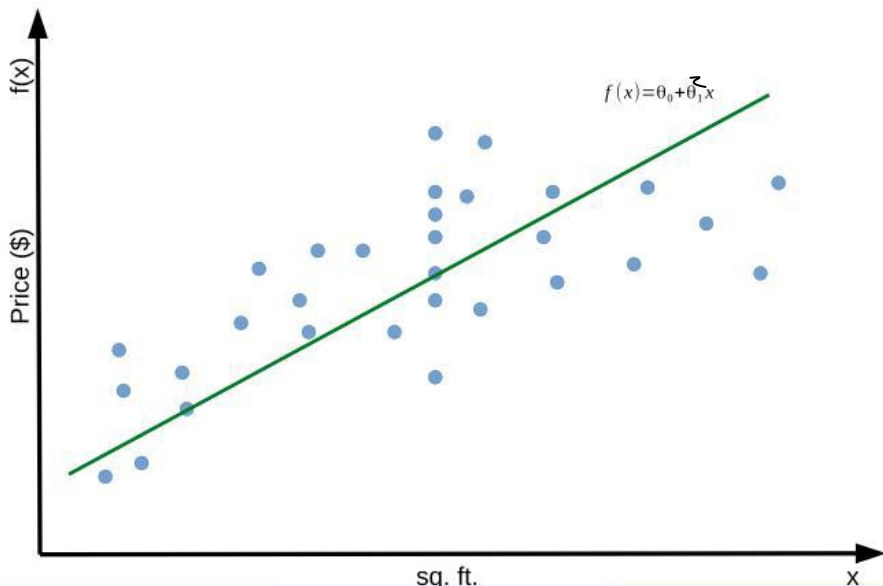
# Assymetric cost function

$$\frac{1}{n} \sum_{i=1}^n \left| \alpha - \mathbb{1}_{(y^{(i)} - f(x^{(i)})) < 0} \right| \cdot \left( y^{(i)} - f(x^{(i)}) \right)^2$$

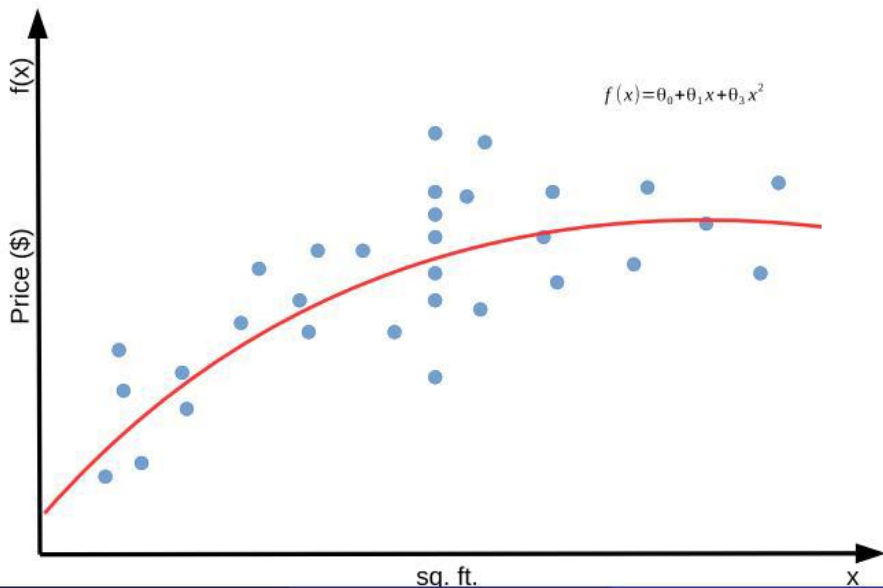


- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

# "Linear" is only regression

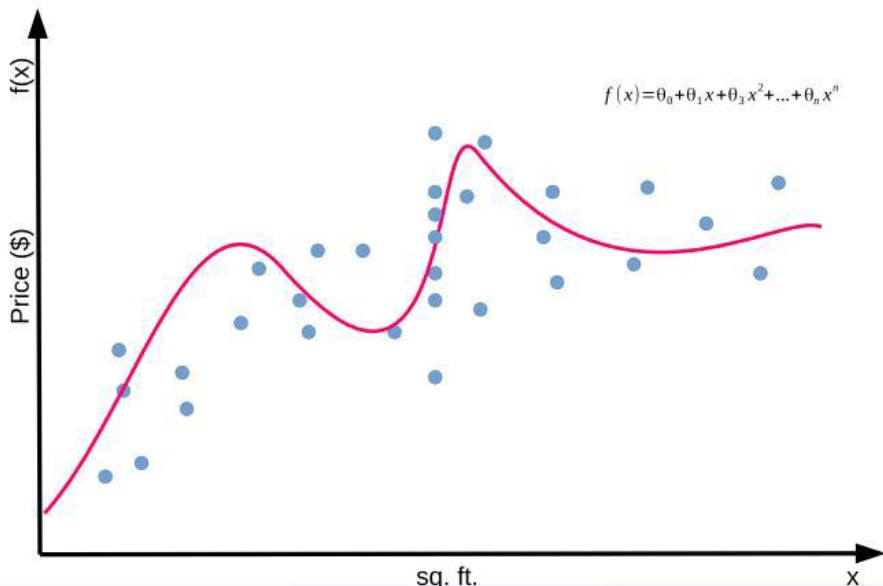


# "Linear" is only regression



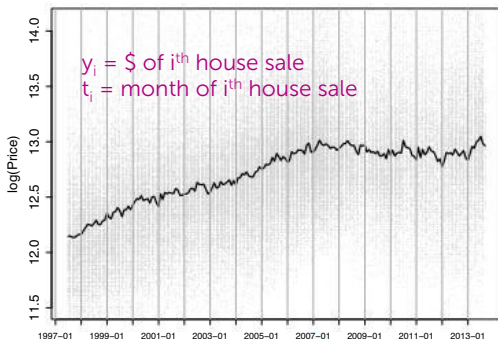


# "Linear" is only regression



# Even more complex features

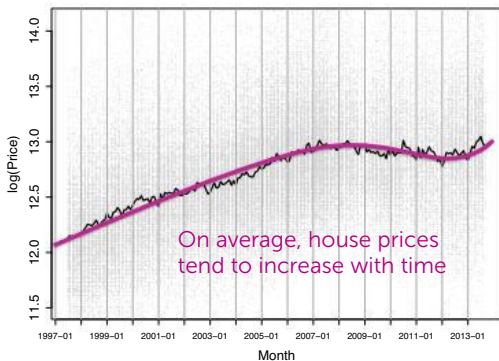
## Motivating application: Detrending time series



Month ← House sales recorded monthly

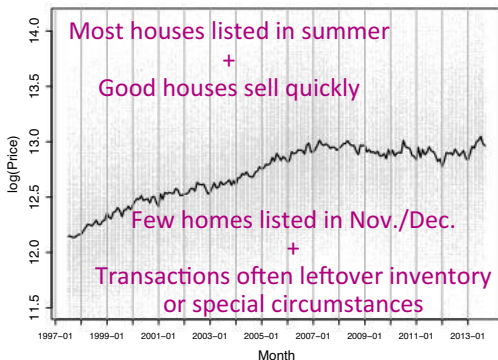
©2015 Emily Fox & Carlos Guestrin

## Trends over time



©2015 Emily Fox & Carlos Guestrin

## Seasonality



©2015 Emily Fox & Carlos Guestrin

## An example detrending

Model:  $X_0 = 1$

$$y_i = \theta_0 + \theta_1 t_i + \theta_2 \sin(2\pi t_i / 12 - \phi) + \varepsilon_i$$

Linear trend

Unknown phase/shift

Seasonal component = Sinusoid with period 12 (resets annually)

Trigonometric identity:  $\sin(a-b) = \sin(a)\cos(b) - \cos(a)\sin(b)$

$$\rightarrow \sin(2\pi t_i / 12 - \phi) = \sin(2\pi t_i / 12) \cos(\phi) - \cos(2\pi t_i / 12) \sin(\phi)$$

©2015 Emily Fox & Carlos Guestrin

## An example detrending

Equivalently,

$$y_i = \theta_0 + \theta_1 t_i + \theta_2 \sin(2\pi t_i / 12) + \theta_3 \cos(2\pi t_i / 12) + \varepsilon_i$$

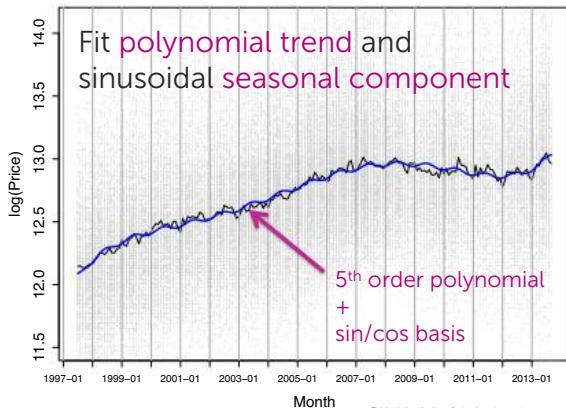
*feature 1 = 1 (constant)*

*feature 2 =  $t$*

*feature 3 =  $\sin(2\pi t/12)$*

*feature 4 =  $\cos(2\pi t/12)$*

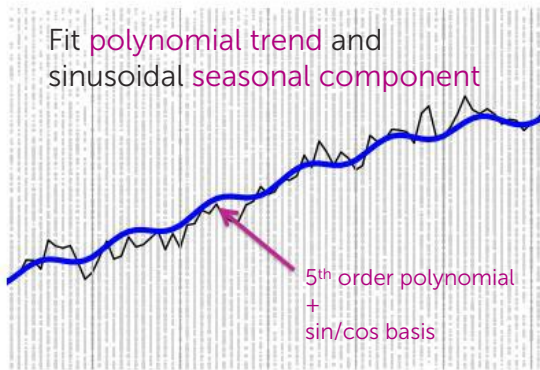
## Detrended housing data



©2015 Emily Fox & Carlos Guestrin

## Zoom in...

Fit **polynomial trend** and  
sinusoidal **seasonal component**



©2015 Emily Fox & Carlos Guestrin

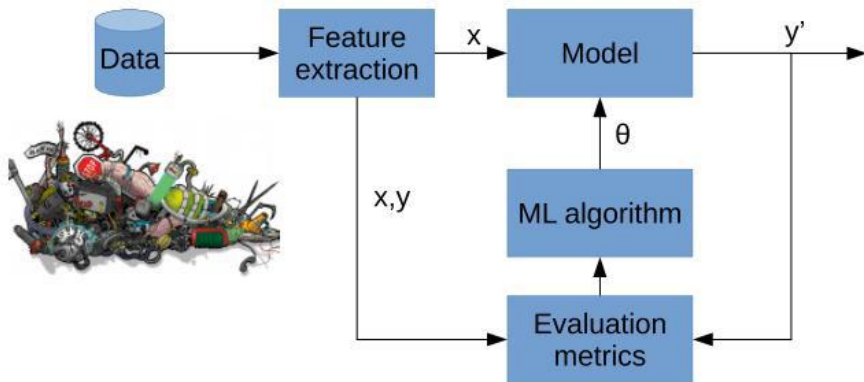


# Even more complex features

Adding more dimensions to the problem can be done mainly in two ways:

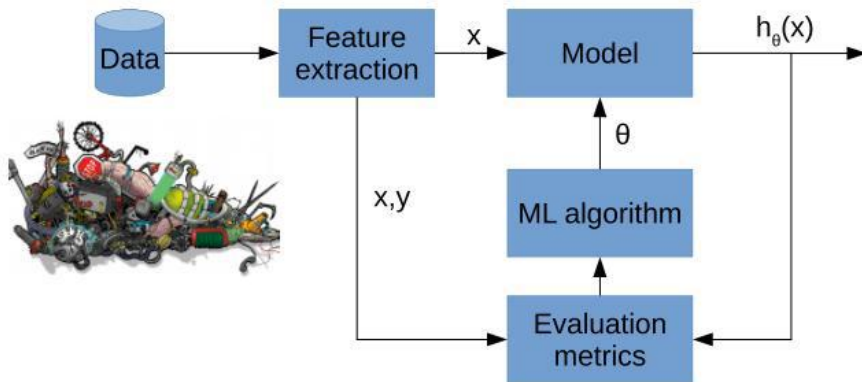
- By adding additional features (like number of bathrooms)
- By adding artificial features being just a combinations or functions of existing ones (squared, log, etc.)
- To allow more compact notation, usually, we add fake *zero-th* feature that equals one, and is multiplied by  $\theta_0$ .

# Note on notation

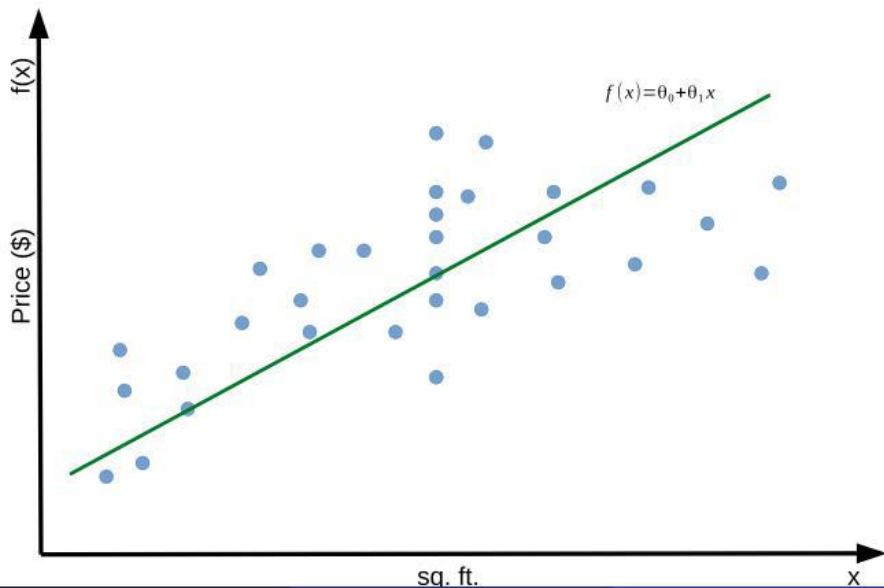


# Note on notation

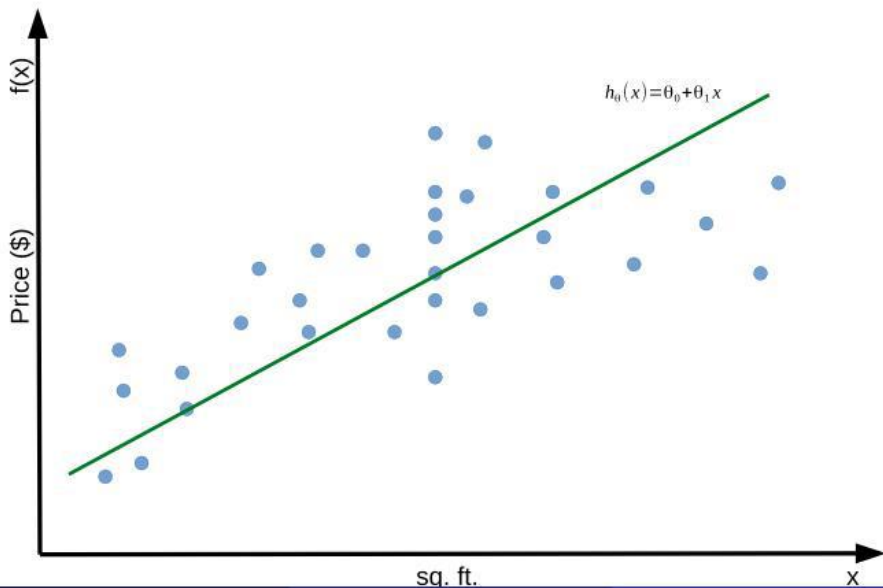
$\theta, \theta_k?$   
Sim



# Note on notation



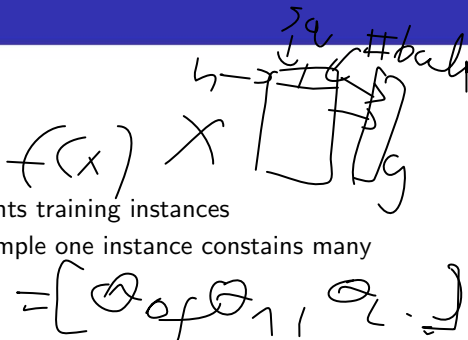
# Note on notation



# Note on notation

- $N$  – number of observations
  - $D$  – number of dimensions
  - $h_{\theta}(x)$  – a hypothesis to be learned
  - $X$  – training set, where rows represents training instances
  - $x^{(j)}$  –  $j$ -th training instance. For example one instance contains many feature values:
    - $x_1^{(j)}$  – sq. ft.
    - $x_2^{(j)}$  – no. of bathrooms
    - $x^{(j)}_i$  – some uber cool feature
  - $\theta$  – vector of coefficients to learn (each  $\theta_i$  corresponds to each  $x_i$ )
  - $y$  – ground truth values for training set.  $y^{(j)}$  value corresponds to  $j$ -th row from  $X$ , thus to  $x^{(j)}$ .
  - We will denote  $J(\theta)$  as a cost function which in our case will be RSS.
- Therefore

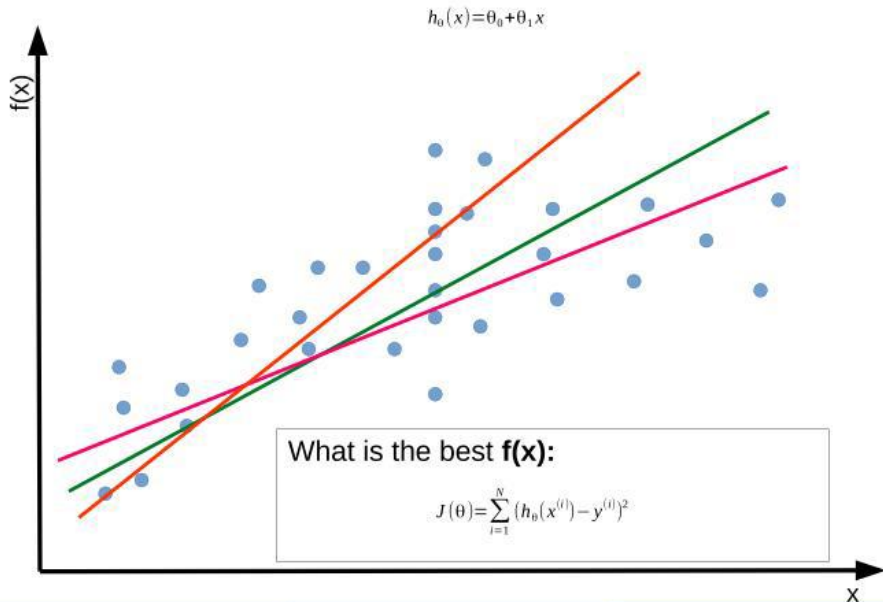
$$J(\theta) = \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



# Presentation Outline

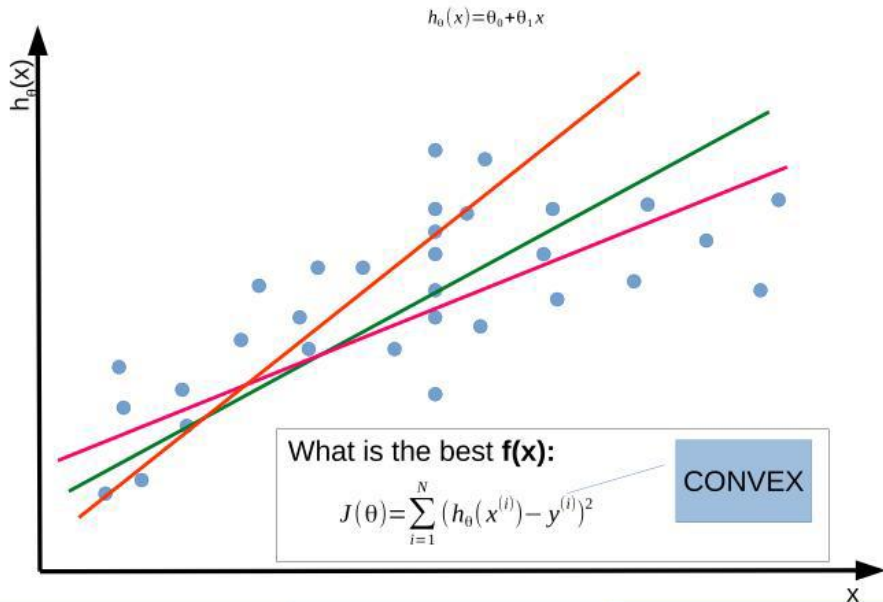
- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

# Minimizing the error

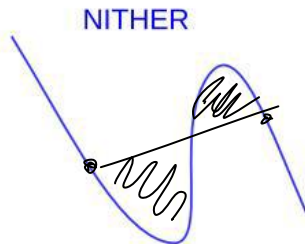
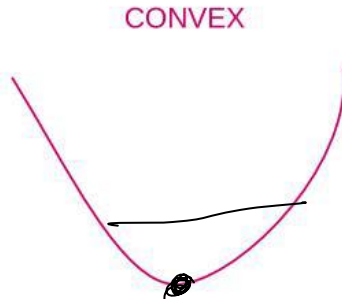
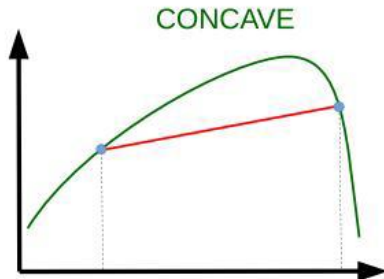




# Minimizing the error



# Minimizing the error



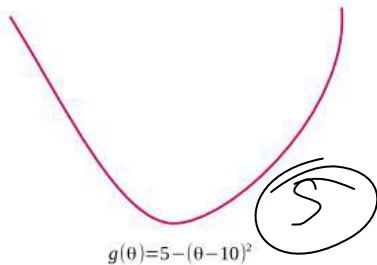
# Outline

- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

# Calculate gradient

CONVEX

$$J(\theta) = \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



$\theta = 10$

$$0 = \frac{d}{d\theta} = 0 - 2(\theta - 10)$$
$$0 = -2(\theta - 10)$$

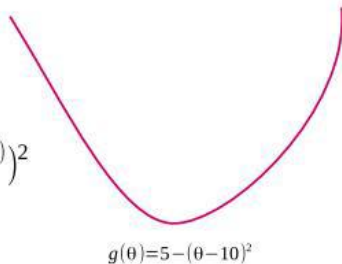
# Calculate gradient

CONVEX

$$J(\theta) = \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \sum_{i=1}^N (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} - y^{(i)})^2$$

↪ (7)

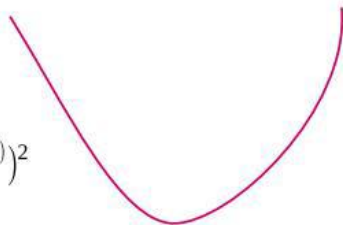


# Calculate gradient

CONVEX

$$J(\theta) = \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \sum_{i=1}^N (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} - y^{(i)})^2$$



$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} - y^{(i)}) x_1^{(i)} \end{bmatrix} \approx \mathbf{0}$$

# Calculate gradient

$$J(\theta) = \sum_{i=1}^N (\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} - y^{(i)})^2$$

$$J(\theta) = \sum_{i=1}^N (x_0^{(i)} \theta_0 + x_1^{(i)} \theta_1 - y^{(i)})^2$$

$$J(\theta) = \sum_{i=1}^N ( \begin{bmatrix} \cancel{x_0} & \cancel{x_1} & \dots & \boxed{x_0} & \boxed{x_1} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_0 \\ \theta_1 \end{bmatrix} - \boxed{y} )^2$$

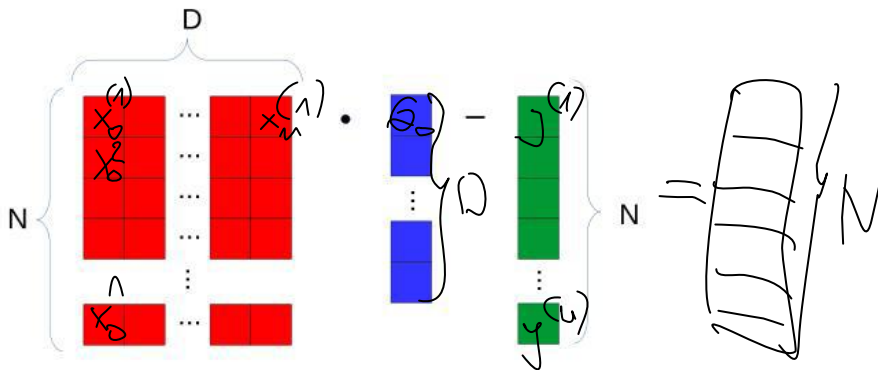
Handwritten diagram illustrating the dot product of the feature vector and the parameter vector. The feature vector is shown as a row of red boxes, with the first two boxes crossed out and the last two boxes highlighted. The parameter vector is shown as a column of blue boxes, with the first two boxes crossed out and the last two boxes highlighted. The result is a green box representing the target value  $y$ .

$$J(\theta) = \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)})^2$$



# Calculate gradient

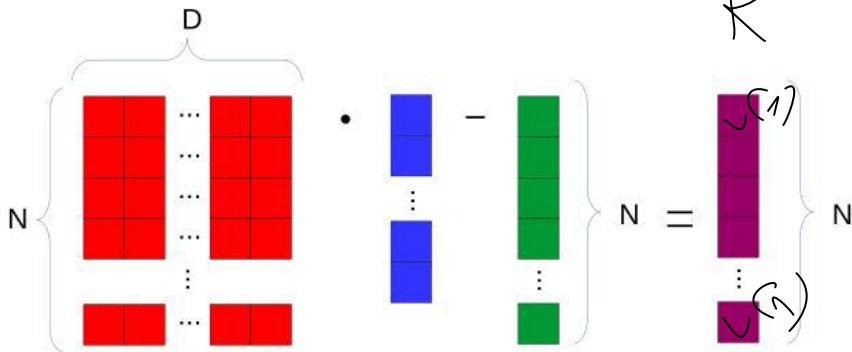
$$J(\theta) = \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \boldsymbol{\theta} - y^{(i)})^2$$





# Calculate gradient

$$J(\theta) = \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)})^2$$



# Calculate gradient

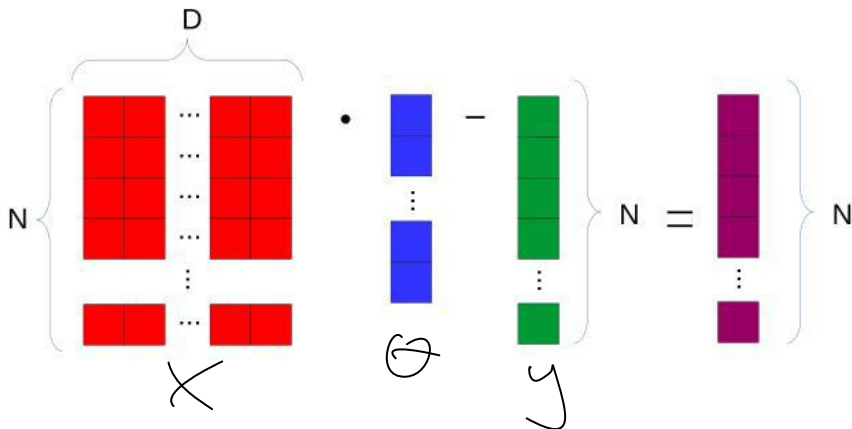
$$J(\theta) = \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \boldsymbol{\theta} - y^{(i)})^2 = \sum_{i=1}^N (r^{(i)})^2$$

Handwritten notes above the equation:  $\begin{pmatrix} 1 & 1 & 1 & 1 \end{pmatrix}$  and  $\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$  with a circled  $\mathbf{R}$  and a circled  $\mathbf{R}^T$ .

The diagram illustrates the matrix representation of the cost function  $J(\theta)$ . It shows a matrix of size  $N \times D$  (red blocks) multiplied by a vector of size  $D$  (blue block), minus a vector of size  $N$  (green block), resulting in a vector of size  $N$  (purple block). The resulting vector is labeled  $\mathbf{r}$  and has handwritten checkmarks and a circled  $\mathbf{R}$  next to it.

# Calculate gradient

$$J(\theta) = \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \boldsymbol{\theta} - y^{(i)})^2 = (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

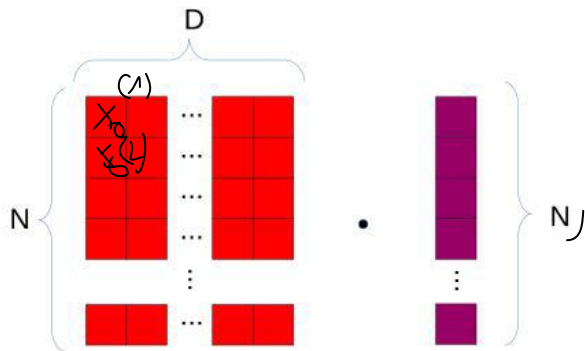


# Calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (\mathbf{x}_0^{(i)} r^{(i)}) \\ \sum_{i=1}^N (\mathbf{x}_1^{(i)} r^{(i)}) \end{bmatrix}$$

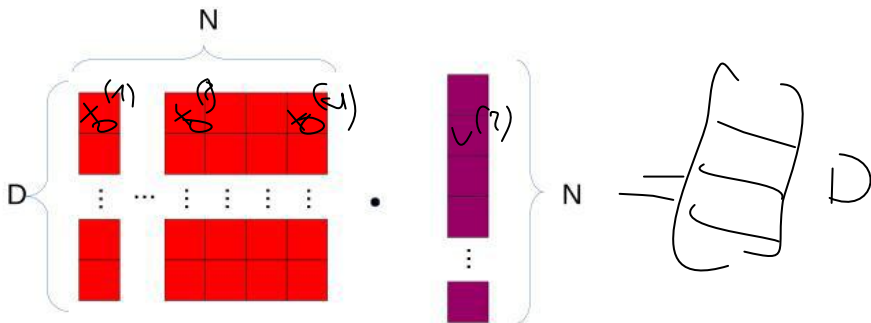
# Calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (x_0^{(i)} r^{(i)}) \\ \sum_{i=1}^N (x_1^{(i)} r^{(i)}) \end{bmatrix}$$



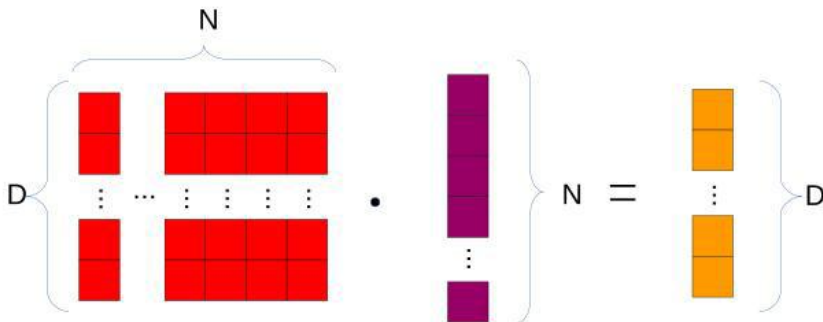
# Calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (x_0^{(i)} r^{(i)}) \\ \sum_{i=1}^N (x_1^{(i)} r^{(i)}) \end{bmatrix}$$



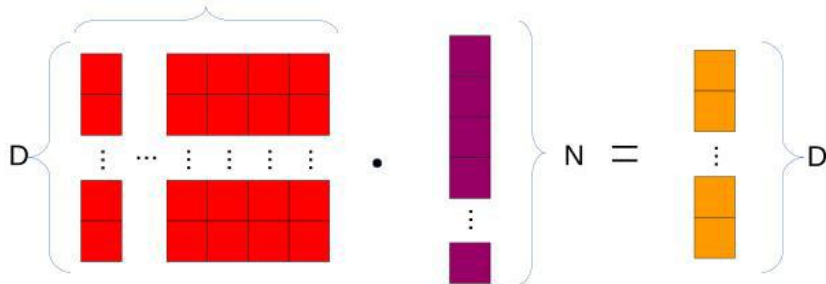
# Calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N (\mathbf{x}_0^{(i)} r^{(i)}) \\ \sum_{i=1}^N (\mathbf{x}_1^{(i)} r^{(i)}) \end{bmatrix} = \mathbf{X}^T \mathbf{R}$$



# Calculate gradient

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((x^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = 2 \underbrace{X^T (X \theta - y)}_{\mathcal{R}}$$





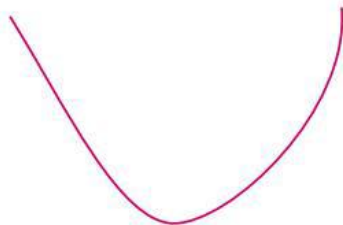
# Calculate gradient

CONVEX

$$J(\theta) = (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y})$$

$$\nabla J(\theta) = 2\mathbf{X}^T (\mathbf{X}\theta - \mathbf{y})$$

$$\nabla J(\theta) = 0$$



$$g(\theta) = 5 - (\theta - 10)^2$$

$$\begin{aligned} 2\mathbf{X}^T (\mathbf{X}\theta - \mathbf{y}) &= 0 \\ 2\mathbf{X}^T \mathbf{X} \theta - 2\mathbf{X}^T \mathbf{y} &= 0 \\ 2\mathbf{X}^T \mathbf{X} \theta &= 2\mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

$\mathbf{A} \mathbf{A}^T$   
 $(\mathbf{X}^T \mathbf{X})$

# Calculate gradient

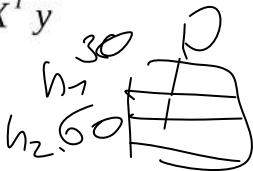
CONVEX

$$J(\theta) = (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y})$$

$$\nabla J(\theta) = 2\mathbf{X}^T (\mathbf{X}\theta - \mathbf{y})$$

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

Invertible if:



$$g(\theta) = 5 - (\theta - 10)^2$$

Complexity of inverse:



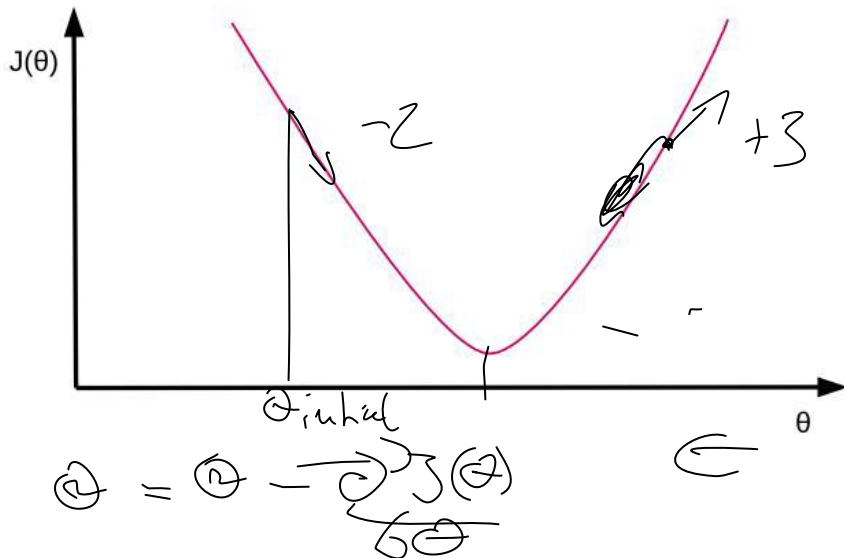
$$O(D^3)$$

$$h_2 \geq h_1$$

# Outline

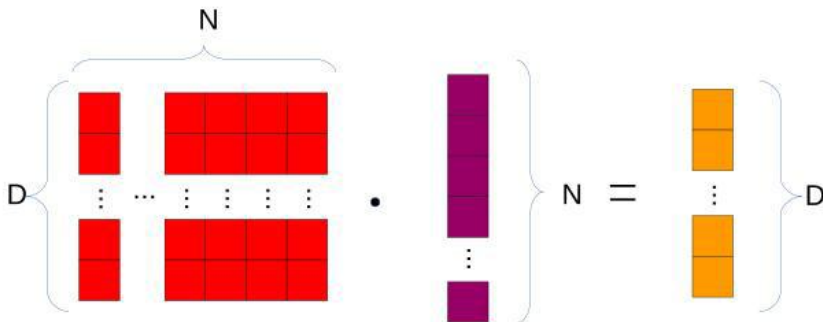
- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - **Gradient descent**
  - Features normalization

# Gradient descent



# Gradient descent

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = 2 \mathbf{X}^T (\mathbf{X} \theta - \mathbf{y})$$



# Gradient descent

$$\nabla J(\theta) = \begin{bmatrix} 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_0^{(i)} \\ 2 \sum_{i=1}^N ((\mathbf{x}^{(i)})^T \theta - y^{(i)}) x_1^{(i)} \end{bmatrix} = 2 \mathbf{X}^T (\mathbf{X} \theta - \mathbf{y})$$

Q



=

Q



$-\alpha$

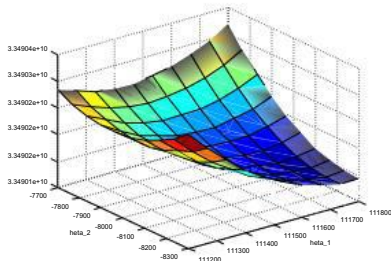
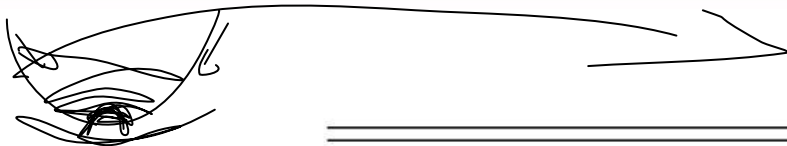


} D

0 - 1

0 -

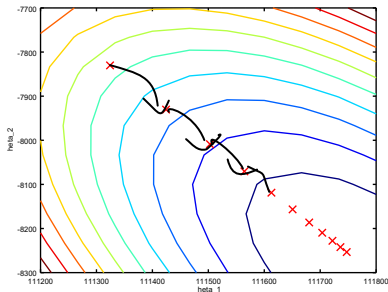
# Gradient descent



```
1 while  $\left| \frac{\partial J(\theta)}{\partial \theta} \right| > \epsilon$  do
2   for  $i \in \{1, \dots, D\}$  do
3      $\frac{\partial J(\theta)}{\partial \theta_i} = 2 \sum_j^N (\theta x^{(j)} - y^{(j)}) x_i^{(j)}$ 
4   end
5   for  $i \in \{1, \dots, D\}$  do
6      $\theta_i = \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$ 
7   end
8 end
```

# Gradient descent

1000



---

```
1 while  $\left| \frac{\partial J(\theta)}{\partial \theta} \right| > \epsilon$  do
2   for  $i \in \{1, \dots, D\}$  do
3      $\frac{\partial J(\theta)}{\partial \theta_i} = 2 \sum_j^N (\theta x^{(j)} - y^{(j)}) x_i^{(j)}$ 
4   end
5   for  $i \in \{1, \dots, D\}$  do
6      $\theta_i = \theta_i - \alpha \frac{\partial J(\theta)}{\partial \theta_i}$ 
7   end
8 end
```

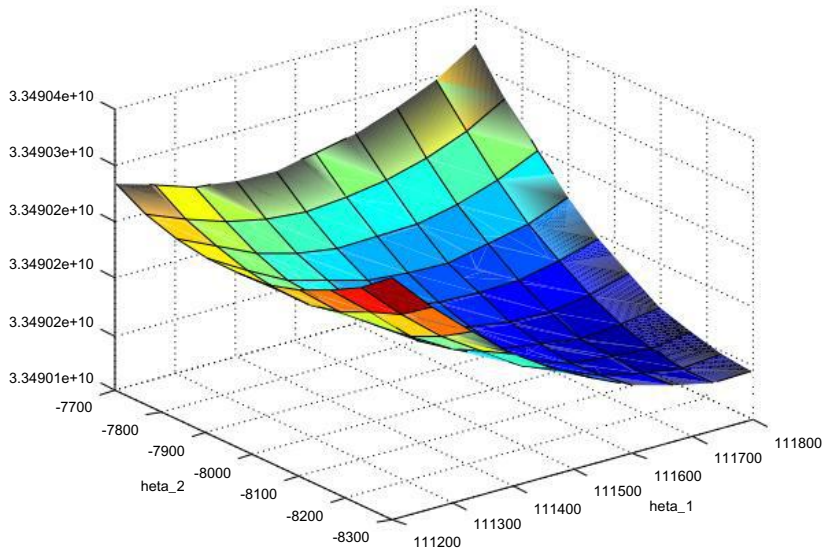
---



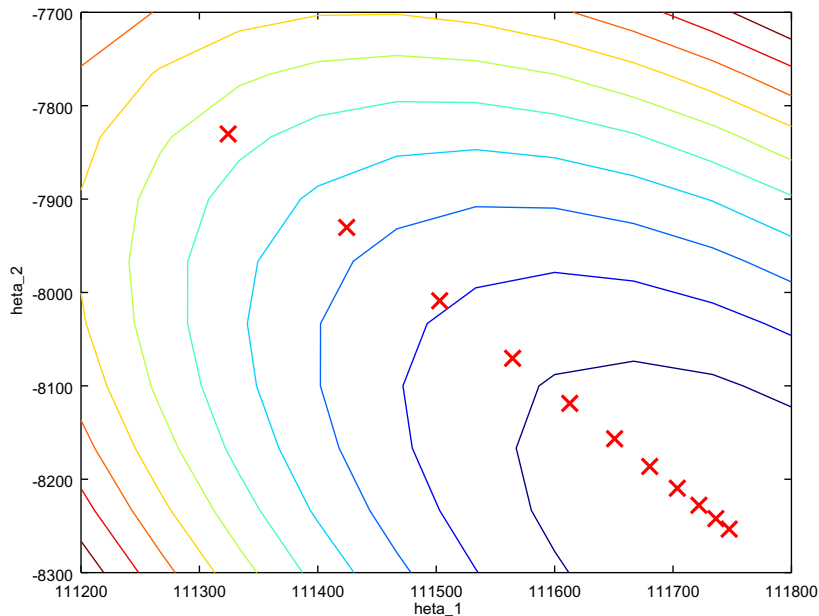
# Outline

- 1 Roadmap
- 2 Use case
- 3 Model
- 4 Simple linear regression
  - General idea
  - Why RSS?
  - Other cost functions
  - Linearity in linear regression
- 5 Solution to linear regression problem
  - The overall idea
  - Closed form solution
  - Gradient descent
  - Features normalization

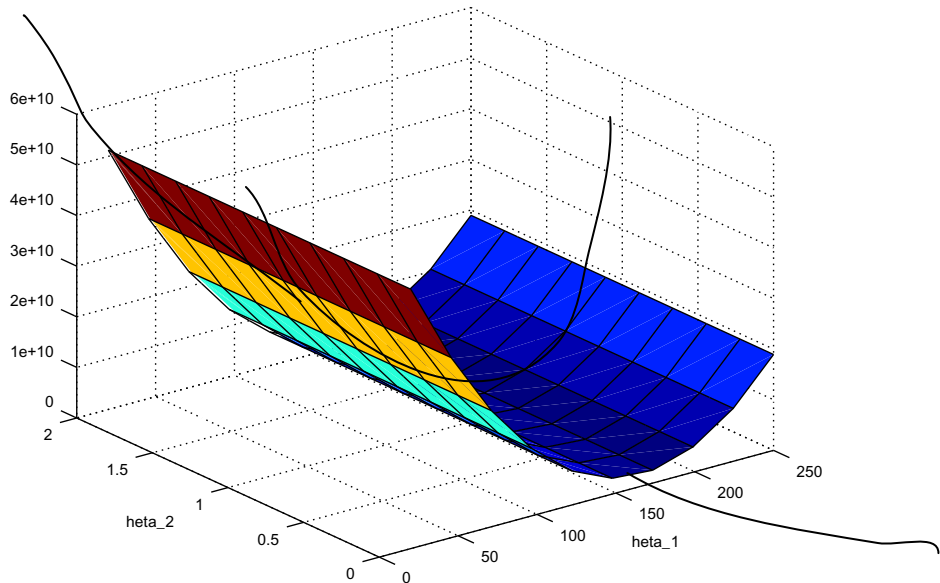
# Convergence of gradient descent



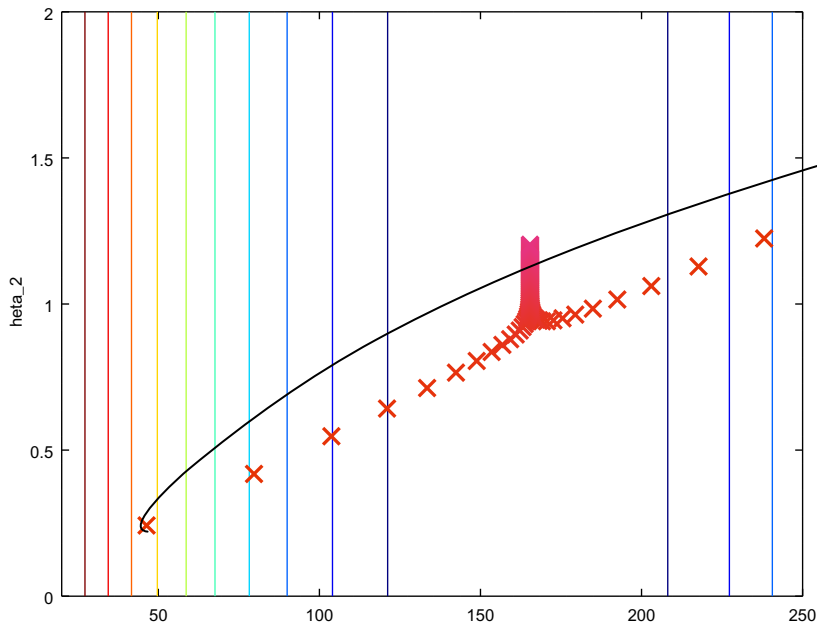
# Convergence of gradient descent



# Convergence of gradient descent



# Convergence of gradient descent



## Min-Max scaling

- Values after scaling are within fixed range. Usually  $[0; 1]$ .
- We use following equation:

$$x_i = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

## Mean standardization

- Values are oriented around zero-mean and standard deviation 1
- We use the following equation:

$$x_i = \frac{x_i - \mu}{\sigma}$$

# Thank you!

**Szymon Bobek**

Institute of Applied Computer Science

AGH University of Science and Technology

21 March 2017

<http://geist.agh.edu.pl>

