

## T4: Linear regression with gradient descent

Antonio Bahamonde  
Departamento de Informática

# Contents

- Linear regression
- Gradient descent
- Linear regression with gradient descent

## Linear regression

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
:	:	:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta$ 's are the **parameters** (also called **weights**) parameterizing the space of linear functions mappings

$$h : \mathcal{X} \mapsto \mathcal{Y}$$

## Linear regression

When there is no risk of confusion, we will drop the  $\theta$  subscript in

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

write it more simply as  $h(x)$ . To simplify our notation, we also introduce the convention of letting  $x_0 = 1$  (this is the **intercept** term), so that

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x,$$

# Linear regression

Given a training set, how do we pick, or learn, the parameters  $\theta$ ?

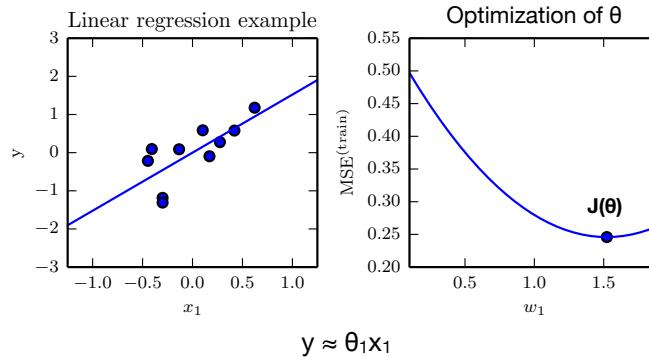
One reasonable method seems to be to make  $h(x)$  close to  $y$ , at least for the training examples we have.

To formalize this, we will define a function that measures, for each value of the  $\theta$ 's, how close the  $h(x)$ 's are to the corresponding  $y$ 's. We define the **cost function**:

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

**least-squares** cost function that gives rise to the ordinary least squares regression model

## Linear Regression



6

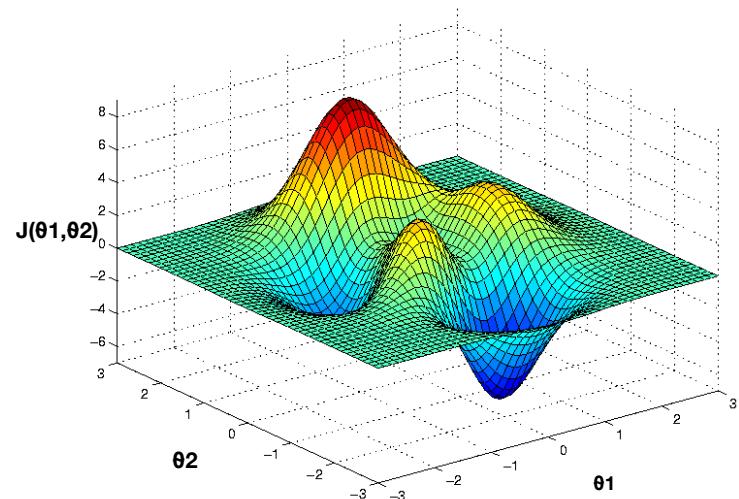
# Gradient descent

## Gradient descent

Given a function  $J(\theta_1, \theta_2)$  we look for the values of  $\theta_1$  and  $\theta_2$  that make the value of  $J$  to a minimum

Algorithm:

- Assign arbitrary values to  $\theta_1$  and  $\theta_2$
- Modify the values of  $\theta_1$  and  $\theta_2$  to reduce the value of  $J$   
Until we believe that we have reached a minimum



## Gradient Descent

To solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} J(\theta)$$

GD is the algorithm

```

 $\theta \leftarrow \theta_0$ 
repeat
 $\theta \leftarrow \theta - \gamma * \frac{\partial J}{\partial \theta}$ 
until no more improvement can be reached
return  $\theta$ 
```

10

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

**Simultaneously in all the components of  $\theta$**

$\alpha$  learning rate

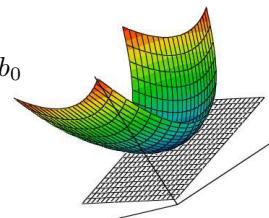
## Gradient descent

The Taylor polynomial of order 1 states that for values  $f(z)$  when  $z$  is close to  $a$  can be approximated using the derivative as follows

$$f(z) = f(a) + (z - a)^T \nabla_z f(a)$$

That is to say,

$$\begin{aligned} f(z) &= f(a) + z_1 b_1 + \dots + z_m b_m + b_0 \\ (b_1, \dots, b_m) &= \nabla_z f(a) \\ b_0 &= -a^T \nabla_z f(a) \end{aligned}$$



12

## Gradient descent

$$f(z) = f(a) + (z - a)^T \nabla_z f(a)$$

If we state that  $\mathbf{z}$  is  $\mathbf{a}$  modified by a small increase ( $\gamma$ ) in the address of a unit vector  $\mathbf{u}$  ( $\|\mathbf{u}\| = 1$ )

$$z = a + \gamma u$$

We have that

$$f(a + \gamma u) = f(a) + \gamma u^T \nabla_z f(a)$$

13

## Gradient descent

To find a value  $\mathbf{z}$  such that  $f(\mathbf{z}) < f(\mathbf{a})$ , we must find a direction  $(\mathbf{u})$  that

$$\begin{aligned} \operatorname{argmin}_{\mathbf{u}} \mathbf{u}^T \nabla_z f(\mathbf{a}) \\ = \operatorname{argmin}_{\mathbf{u}} \|\mathbf{u}\| \|\nabla_z f(\mathbf{a})\| \cos(\mathbf{u}, \nabla_z f(\mathbf{a})) \\ = \operatorname{argmin}_{\mathbf{u}} \cos(\mathbf{u}, \|\nabla_z f(\mathbf{a})\|) \end{aligned}$$

Therefore  $\mathbf{u}$  must be parallel to the opposite of the gradient ( $\cos = -1$ )

$$\mathbf{u} = -\nabla_z f(\mathbf{a})$$

14

## Gradient descent

Therefore (if  $\gamma$  is small)

$$\begin{aligned} f(\mathbf{a} - \gamma \nabla_z f(\mathbf{a})) &= f(\mathbf{a}) - \gamma \nabla_z f(\mathbf{a})^T \nabla_z f(\mathbf{a}) \\ &= f(\mathbf{a}) - \gamma \|\nabla_z f(\mathbf{a})\|^2 \\ &< f(\mathbf{a}) \end{aligned}$$

After the assignment

$$\mathbf{a} \leftarrow \mathbf{a} - \gamma \nabla_z f(\mathbf{a})$$

We have a point with lower  $f$ -value than  $f(\mathbf{a})$ , whenever  $\gamma$  is small

15

## Gradient descent

Tying up loose ends. The goal was to solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, h_{\theta}(\mathbf{x}))$$

We look for the  $\theta$  value that minimizes the mean of the loss function over the dataset D

$$\begin{aligned} J(\theta) &= \frac{1}{|D|} \sum_{(x,y) \in D} \text{loss}(y, h_{\theta}(\mathbf{x})) \\ &= \mathbb{E}_{(x,y) \sim p(D)} \text{loss}(y, h_{\theta}(\mathbf{x})) \end{aligned}$$

16

# Gradient descent

To solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} \mathbf{J}(\theta)$$

GD is the algorithm

```
θ ← θ₀
repeat
    θ ← θ - γ * ∂J / ∂θ
until no more improvement can be reached
return θ
```

17

# Mini Batch and SGD

GD has the problem to compute

$$\frac{\partial \mathbf{J}(\theta)}{\partial \theta} = \frac{1}{|D|} \sum_{(x,y) \in D} \frac{\partial}{\partial \theta} \text{loss}(\mathbf{y}, f_\theta(\mathbf{x}))$$

Thus, to make GD useful, we estimate this average with the mean of a subset  $D'$  of  $D$

$$\frac{\partial \mathbf{J}(\theta)}{\partial \theta} \cong \frac{1}{|D'|} \sum_{(x,y) \in D'} \frac{\partial}{\partial \theta} \text{loss}(\mathbf{y}, f_\theta(\mathbf{x}))$$

This is called the Stochastic Gradient Descent with mini batch. Whenever  $D'$  has only one element, we have SGD

18

# Optimization results

Note that, while gradient descent can be susceptible to [local minima in general](#), the optimization problem we have posed here

for linear regression

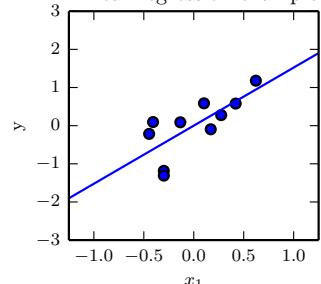
- has only one global, and no other local, optima; gradient descent always converges (assuming the learning rate ( $\gamma$  or  $\alpha$ ) is not too large) to the global minimum.
- $J$  is a convex quadratic function.
- Stochastic gradient descent gets  $\theta$  “close” to the minimum much faster than batch gradient descent.

(Note however that it may never “converge” to the minimum, and the parameters  $\theta$  will keep oscillating around the minimum of  $J(\theta)$ ; but in practice most of the values near the minimum will be reasonably good approximations to the true minimum.)

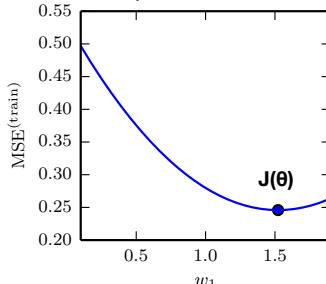
# Linear regression algorithm

## Linear Regression

Linear regression example



Optimization of  $\theta$



$$y \approx \theta_1 x_1$$

21

## Gradient descent

To solve

$$\theta \leftarrow \operatorname{argmin}_{\theta} J(\theta)$$

GD is the algorithm

```

 $\theta \leftarrow \theta_0$ 
repeat
   $\theta \leftarrow \theta - \gamma * \frac{\partial J}{\partial \theta}$ 
until no more improvement can be reached
return  $\theta$ 

```

22

## Linear regression



$$\begin{aligned} \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta). \\ \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\ &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^d \theta_i x_i - y \right) \\ &= (h_\theta(x) - y) x_j \end{aligned}$$

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

A 3D surface plot of the cost function  $J(\theta)$  showing a parabolic shape. The vertical axis is labeled  $J(\theta)$ , and the horizontal axes are labeled  $\theta_0$  and  $\theta_1$ .

## Linear regression



For a single training example, this gives the LMS ("least mean squares") update rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}.$$

the magnitude of the update is proportional to the error term ( $y^{(i)} - h_\theta(x^{(i)})$ ):

- if we are encountering a training example on which our prediction nearly matches the actual value of  $y^{(i)}$ , then we find that there is little need to change the parameters;
- a larger change to the parameters will be made if our prediction has a large error

## Probabilistic interpretation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}, \quad p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right).$$

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

## Probabilistic interpretation

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2. \end{aligned}$$

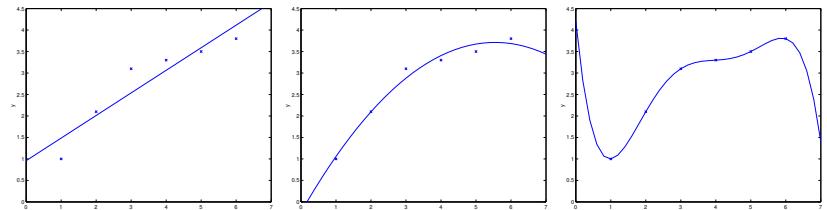
## Probabilistic interpretation

Least-squares regression corresponds to finding the maximum likelihood estimate of  $\theta$ .

$\theta$  did not depend on what was  $\sigma$ , even if  $\sigma$  were unknown.

$$\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2,$$

## Overfitting underfitting



## Más información

- Linear regression (Main notes, chapter 1)
- Linear Algebra Review and Reference. Zico Kolter (updated by Chuong Do and Tengyu Ma) June 20, 2020

[https://cs229.stanford.edu/lectures-spring2022/cs229-linear\\_algebra\\_review.pdf](https://cs229.stanford.edu/lectures-spring2022/cs229-linear_algebra_review.pdf)