

# \*\*\* Applied Machine Learning Fundamentals \*\*\*

## Regression

Daniel Wehner

SAP SE

August 13, 2019



# Agenda August 13, 2019

## ① Introduction to Regression

What is Regression?

Least Squares Error Function

## ② Solutions to Regression

Closed-Form Solutions and Normal Equation

Gradient Descent

## ③ Probabilistic Regression

## ④ Basis Function Regression

## ⑤ Wrap-Up

Summary

Lecture Overview

Self-Test Questions

Recommended Literature and further Reading

Section:  
**Introduction to Regression**



# Regression

Type of target variable

Continuous

Type of training information

Supervised

Example Availability

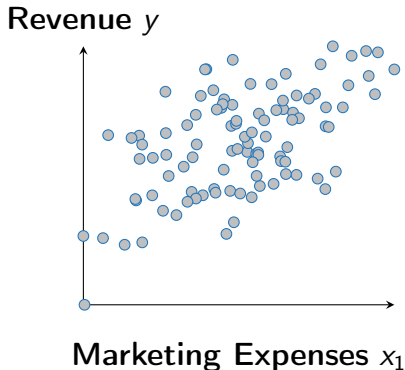
Batch learning

**Algorithm sketch:** Given the training data  $\mathcal{D}$  the algorithm derives a function of the type

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \cdots + \theta_{m+1} x_m \quad \mathbf{x} \in \mathbb{R}^m, \boldsymbol{\theta} \in \mathbb{R}^{m+1} \quad (1)$$

from the data.  $\boldsymbol{\theta}$  is the parameter vector containing the coefficients to be estimated by the regression algorithm. Once  $\boldsymbol{\theta}$  is learned it can be used for prediction.

## Example Data Set: Revenues



- Find a linear function:

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \cdots + \theta_{m+1} x_m$$

- Usually:  $x_0 = 1$ :

$$\hat{\mathbf{x}} \in \mathbb{R}^{m+1} = [1 \ \mathbf{x}]^T$$

$$h_{\theta}(\hat{\mathbf{x}}) = \sum_{j=0}^{m+1} \theta_j x_j = \boldsymbol{\theta}^T \hat{\mathbf{x}}$$

# Error Function for Regression

- In order to know how good the function fits we need an error function  $\mathcal{J}(\boldsymbol{\theta})$ :

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}^{(i)}) - y^{(i)})^2 \quad (2)$$

- We want to minimize  $\mathcal{J}(\boldsymbol{\theta})$ :

$$\min_{\boldsymbol{\theta}} \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\hat{\mathbf{x}}^{(i)}) - y^{(i)})^2$$

- This is **ordinary least squares (OLS)**

# Error Function Intuition

Section:  
**Solutions to Regression**





## Closed-Form Solutions

- Usual approach (for two unknowns): Calculate  $\theta_0$  and  $\theta_1$  according to

sample mean  $\bar{x}$

$$\theta_0 = \bar{y} - \theta_1 \bar{x} \quad \theta_1 = \frac{\sum_{i=1}^n (x^{(i)} - \bar{x}) \cdot (y^{(i)} - \bar{y})}{\sum_{i=1}^n (x^{(i)} - \bar{x})^2} \quad (3)$$

- 'Normal equation' (scales to arbitrary dimensions):

$$\theta = \underbrace{(\hat{\mathbf{X}}^\top \hat{\mathbf{X}})^{-1} \hat{\mathbf{X}}^\top}_{\text{Moore-Penrose pseudo-inverse}} \mathbf{y} \quad (4)$$

$\hat{\mathbf{X}}$  is called 'design matrix' or 'regressor matrix'

# Design Matrix / Regressor Matrix

- The design matrix  $\hat{\mathbf{X}} \in \mathbb{R}^{n \times (m+1)}$  looks as follows:

$$\hat{\mathbf{X}} = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_m^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_m^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \cdots & x_m^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \cdots & x_m^{(n)} \end{pmatrix} \quad (5)$$

In the following

$$\hat{\mathbf{X}} \equiv \mathbf{X}$$

- And the  $n \times 1$  label vector:

$$\mathbf{y} = (y^{(1)}, y^{(2)}, y^{(3)}, \dots, y^{(n)})^\top$$



# Derivation of the Normal Equation

- The derivation involves a bit of linear algebra
- Step **1**: Rewrite  $\mathcal{J}(\theta)$  in matrix-vector notation:

$$\begin{aligned}\mathcal{J}(\theta) &= \frac{1}{2}(\mathbf{X}\theta - \mathbf{y})^\top(\mathbf{X}\theta - \mathbf{y}) \\ &= ((\mathbf{X}\theta)^\top - \mathbf{y}^\top)(\mathbf{X}\theta - \mathbf{y}) \\ &= (\mathbf{X}\theta)^\top \mathbf{X}\theta - (\mathbf{X}\theta)^\top \mathbf{y} - \mathbf{y}^\top (\mathbf{X}\theta) + \mathbf{y}^\top \mathbf{y} \\ &= \theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2(\mathbf{X}\theta)^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

- To be continued...



## Derivation of the Normal Equation (Ctd.)

- Step ②: Calculate the derivative of  $\mathcal{J}(\theta)$  and set it to zero:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) &= 2\mathbf{X}^T \mathbf{X} \theta - 2\mathbf{X}^T \mathbf{y} \stackrel{!}{=} 0 \\ \Leftrightarrow \mathbf{X}^T \mathbf{X} \theta &= \mathbf{X}^T \mathbf{y}\end{aligned}$$

- If  $\mathbf{X}^T \mathbf{X}$  is invertible, we can multiply both sides by  $(\mathbf{X}^T \mathbf{X})^{-1}$ :

Normal equation:

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

# Problems with Matrix Inversion?

- What if  $(\mathbf{X}^\top \mathbf{X})^{-1}$  does not exist?
- Problems and solutions:
  - ① Linearly dependent (redundant) features or design matrix does not have full rank? (E. g. size in  $\text{m}^2$  and size in  $\text{feet}^2$ )  
 $\Rightarrow$  **Delete correlated features**
  - ② Too many features ( $m > n$ )?  
 $\Rightarrow$  **Delete features (e. g. using PCA) / add training examples**
  - ③ Other numerical instabilities?  
 $\Rightarrow$  **Add a regularization term** (later)
  - ④ Computationally too expensive?  
 $\Rightarrow$  **Use gradient descent**

# Gradient Descent

- We want to minimize a smooth function  $\mathcal{J} : \mathbb{R}^{m+1} \rightarrow \mathbb{R}$ :

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{m+1}} \mathcal{J}(\boldsymbol{\theta})$$

- Update the parameters iteratively:

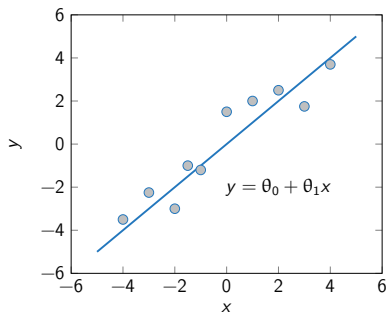
$$\boldsymbol{\theta}^{(t+1)} \longleftarrow \boldsymbol{\theta}^{(t)} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}^{(t)}) \quad (6)$$

- where  $\alpha > 0$  (**learning rate**) and  $\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta})$  is the gradient of  $\mathcal{J}(\boldsymbol{\theta})$  w. r. t.  $\boldsymbol{\theta}$ :

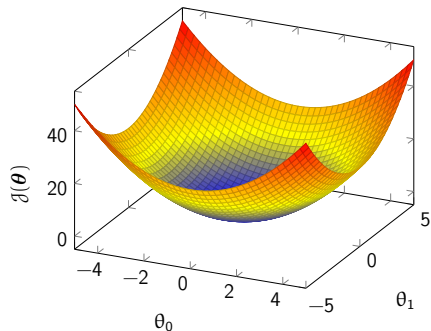
$$\nabla_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) = \left( \frac{\partial \mathcal{J}(\boldsymbol{\theta})}{\partial \theta_0}, \frac{\partial \mathcal{J}(\boldsymbol{\theta})}{\partial \theta_1}, \dots, \frac{\partial \mathcal{J}(\boldsymbol{\theta})}{\partial \theta_{m+1}} \right)^{\top}$$

# Data Input Space vs. Hypothesis Space

## Data input space



## Hypothesis space $\mathcal{H}$

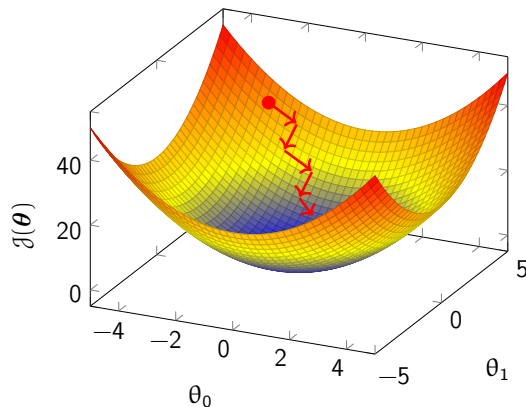


## Data Input Space vs. Hypothesis Space (Ctd.)

- **Data input space**
  - Determined by the  **$m$  attributes** of the data set  $x_1, x_2, \dots, x_m$
  - Often high-dimensional
- **Hypothesis space  $\mathcal{H}$** 
  - Determined by the **number of parameters** of the model
  - Each point in the hypothesis space corresponds to a **specific assignment of model parameters**
  - The error function gives information about how good this assignment is
  - **Gradient descent is applied in the hypothesis space  $\mathcal{H}$**



# Visualization of Gradient Descent in 3 Dimensions



# Versions of Gradient Descent

- Assume some training data  $\mathcal{D}$ :  $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^n$
- Squared error for a **single** example:  $\ell(y_{pred}, y_{true}) = (y_{pred} - y_{true})^2$
- Our objective is to minimize the **total** error:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{m+1}} \mathcal{J}(\boldsymbol{\theta}) = \min_{\boldsymbol{\theta} \in \mathbb{R}^{m+1}} \sum_{i=1}^n \ell(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)})$$

- Three versions of gradient descent:
  - ① Batch gradient descent
  - ② Stochastic gradient descent
  - ③ Mini-batch gradient descent

## Versions of Gradient Descent (Ctd.)

- **Batch gradient descent**: Compute gradient based on ALL data points

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \alpha \sum_{i=1}^n \nabla \ell(h_{\boldsymbol{\theta}^{(t)}}(\mathbf{x}^{(i)}), y^{(i)}) \quad (7)$$

- **Stochastic gradient descent**: Compute gradient based on a SINGLE data point (**pick training example randomly and not sequentially!**)
- For  $i \in \{1, \dots, n\}$  do:

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \alpha \nabla \ell(h_{\boldsymbol{\theta}^{(t)}}(\mathbf{x}^{(i)}), y^{(i)}) \quad (8)$$

# Solving linear Regression using Gradient Descent

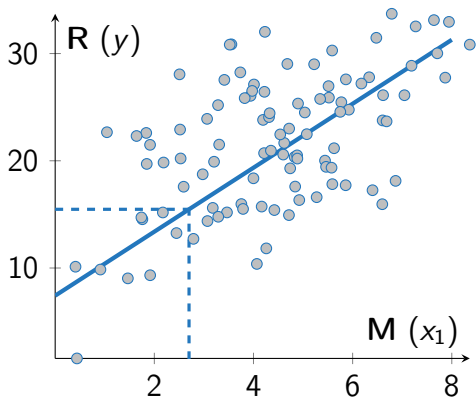
- Randomly initialize  $\theta$
- To minimize the error, keep changing  $\theta$  according to:

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \alpha \nabla_{\theta} \mathcal{J}(\theta^{(t)}) \quad (9)$$

- We need to calculate  $\nabla_{\theta_j} \mathcal{J}(\theta^{(t)})$ : (based on a single example)

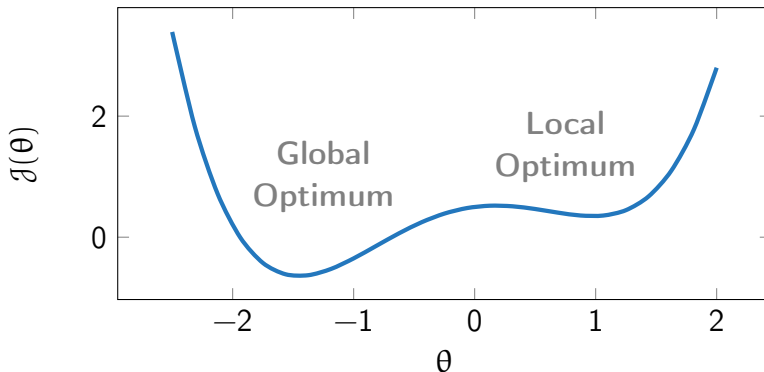
$$\begin{aligned} \nabla_{\theta_j} \mathcal{J}(\theta^{(t)}) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(\mathbf{x}) - y)^2 = 2 \cdot \frac{1}{2} (h_{\theta}(\mathbf{x}) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(\mathbf{x}) - y) \\ &= (h_{\theta}(\mathbf{x}) - y) \cdot \frac{\partial}{\partial \theta_j} (\theta_0 x_0 + \cdots + \theta_{m+1} x_{m+1} - y) = \boxed{(h_{\theta}(\mathbf{x}) - y) x_j} \end{aligned}$$

# Solving the introductory Example



- $\theta_0 \approx 7.4218$
- $\theta_1 \approx 2.9827$
- $\mathcal{J}(\boldsymbol{\theta}) \approx 446.9584$
- $h_{\boldsymbol{\theta}}(\mathbf{x}) = 7.4218 + 2.9827 \cdot x_1$
- $R = h_{\boldsymbol{\theta}}(2.7) = \underline{\underline{15.4750}}$

# Disadvantage of Gradient Descent



Section:  
**Probabilistic Regression**







Section:  
**Basis Function Regression**





# Summary

# Lecture Overview

## Unit I: Machine Learning Introduction

# Self-Test Questions

# Recommended Literature and further Reading

# Thank you very much for the attention!

**Topic:** \*\*\* Applied Machine Learning Fundamentals \*\*\* Regression

**Date:** August 13, 2019

**Contact:**

Daniel Wehner (D062271)

SAP SE

[daniel.wehner@sap.com](mailto:daniel.wehner@sap.com)

## Do you have any questions?