# SF2930: Regresion analysis Lecture 1 Simple linear regression.

Tatjana Pavlenko

17 January 2018



3 D

< Ξ → </li>

Regression analysis is a statistical technique for investigating and modeling the relationship between variables. The reason why it is so widely applied is because it provides the answer to an everyday question, namely

how a response variable of special interest depends on several other, explanatory variables.



# WHAT IS REGRESSION? EXAMPLES

Applications of regression techniques are numerous; example include engineering, medical, biological and social sciences, physical and chemical sciences, economics, ... Regression is probably is the most widely used statistical methodology. Examples of applied problems and questions in which regression might be helpful:

- how the apartment prise depends on size, location, floor, closeness to subway, ...
- how growing of plants depends on fertilizer, soil quality, ...
- how home insurance premium depends on age of homeowner, value of the home and its contents, region, ...

In various quantitative settings, the regression techniques models the relationship between the *response variable* of special interests (Y) and a (set)  $x_1, \ldots x_k$  of *explanatory* or *predictor* variables.



回 と く ヨ と く ヨ と

## **REGRESSION MATHEMATICS**

Linking the response variable to the predictors:



 $\boldsymbol{\varepsilon}$  is the error term which is can neither be controlled or predicted.

- The goal is to learn about the function  $f(\cdot)$ .
- ► In full generality, finding f(·) without any conditions is very difficult: function space is infinite-dimensional!
- Solution: restrict the form of  $f(\cdot)$ .
- Linear modeling:  $Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$ .
- Finding f boils down to determining  $\beta_0, \beta_1, \dots, \beta_k$  from the data.



## **R**EGRESSION AND MODEL BUILDING

- Why modeling? It is stated once by George Box that All models are wrong, but some are useful.
- Certain degree of variability (uncertainty) is present in almost all processes. Most of this variability can not be modeled with deterministic methods of mathematics. Statistics (mathematics under uncertainty) provides a powerful toll of modeling phenomena under uncertainty.
- Larger data collected in real-world applications demand models in order to *extract big knowledge of big data*!
- In this course we start with simple but highly applicable, *linear* regression models. We then turn to more advance subjects.



In this course we start with highly applicable, *simple* linear regression models. It is rather simple modes but the tools constructed within the simple linear regression will be naturally extended to the multiple case.

Later in the course we turn to more advance subject such as

- multiple linear regression
- Iogistic regression
- Generalized linear models containing Poisson and negative binomial regression



## EXAMPLE



FIGURE: The wash water tank on SAS aircrafts. It is impossible to predict exactly the amount of wash water needed, therefore the tank is always filled to 100% at Arlanda airport. The project on minimizing the amount of water is performed with the main goal to lower the weight of the aircraft, and as a result reduce fuel consumption and cost. Goal: to investigate (to model!) the relationship between wash water consumption (target variable) and number of passengers, duration of a flight, time of the flight (night/day), ... (explanatory variables). Statistical approach: multiple linear regression analysis.



イロン イヨン イヨン イヨン

## EXAMPLE: LINEAR REGRESSION FOR ADVERTISING DATA



FIGURE: The Advertising data set. The plots displays sales for a particular product as a function of advertising budget for TV, radio and newspaper, for 200 different markes. Each blue line represents the simple model that can be used to predict sales from each feature variable.



## LINEAR REGRESSION FOR ADVERTISING DATA.

Advertising data, see Ch. 2 in ISL. The goal is using the data, to design a marketing plan (for the next year) that will result in high product sales. Specific questions are (see Intro to Ch. 3 in ISL):

- Q1: Is there a (linear) relationship between advertising budget and sales?
- Q2: How strong is the relationship between advertising budget and sales?
- Q3: Which media contribute to sales (Which predictors are associated with response)?
- Q4: How to estimate effect of each variable on sales? How accurate are these estimates?
- ▶ Q5: How can we predict future sales?
- Q6: Is the synergy among the advertising media?



#### LINEAR REGRESSION - AN APPROACH FOR MODELING AND PREDICTION

The goals of such regression analysis are usually two-fold:

- to model the relation between output Y and input vector
  x = (x<sub>1</sub>,..., x<sub>k</sub>), and to specify which of the predictor variables have effect on the output variable Y, as well as to quantify the strength of this relation (The regression equation, inference).
- ► Using the regression equation, to predict the expected response value of Y for an arbitrary (new) configuration of x<sub>1</sub>,..., x<sub>k</sub>.
- ► The dependence of Y on x = (x<sub>1</sub>,..., x<sub>k</sub>) is assumed to be *linear*, i.e. the *model* formula is

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \varepsilon,$$

where  $\beta_0, \ldots, \beta_k$  are *unknown* coefficients or parameters, and  $\varepsilon$  is the random error, which accounts for measurement error and effect of other variables not explicitly considered in the model.



・ロン ・回と ・ヨン・

#### GOALS IN REGRESSION ANALYSIS

- ► A "good" fit. Estimating a (hyper)-plane over x<sub>1</sub>, ..., x<sub>k</sub> to explain the response such that the errors are "small". The standard tool is Least Squares estimates, LS.
- Good parameter estimates are useful to describe the change of the response when varying some some predictor variable(s).
- Good prediction is useful to predict a new response as a function of new predictor variables.
- Uncertainties and significance with the three goals above.
- Development of a good model: In an interactive process, with the help of methods for the above goals, we may change parts of the model to come up with a better mode.



▲圖 ▶ ▲ 国 ▶ ▲ 国 ▶

## WHY CALL REGRESSION?

# History

 Developed by Sir Francis Galton (1822-1911) in his article "Regression towards mediocrity in hereditary structure"





## HOW REGRESSION GOT ITS NAME?

- In the late 1800's, a scientist Sir Francis Galton was working with large observational studies on humans, in part with the data on the heights of fathers and first sons.
- In terms of simple regression model x was the height of the father and Y was the height of the first, fully grown, son. Goal: to predict Y in terms of x.
- Galton concluded: For father's whose heights were taller than the average, the LS line predicts the son to be shorter than the father. Likewise, for father's whose heights were shorter than the average, the line predicts the son to be taller than the father.

#### There is a regression towards the mean effect!



回 と く ヨ と く ヨ と

## SIMPLE LINEAR REGRESSION

We begin with a mathematically most simple way of describing the relation between the variables: the *linear* relationship between a continuous response variable Y in dependence on

- a single explanatory variable x (simple linear regression model)
- several explanatory variables x<sub>1</sub>,..., x<sub>k</sub> (multiple regression model)
- Later in the course we will consider
  - modeling of the nonlinear relationship between a binary response Y and a set of explanatory variables x<sub>1</sub>,..., x<sub>k</sub> (logistic regression).
  - additional models when Y is discrete.



回 と く ヨ と く ヨ と

## SIMPLE LINEAR REGRESSION

The simple linear regression model is  $Y = \beta_0 + \beta_1 x + \varepsilon$ , where the *intercept*  $\beta_0$  and the *slope*  $\beta_1$  are unknown constants, and  $\varepsilon$  is a random error component.

Model, interpretation of variables and parameters:

- ➤ Y is a continuous dependent variable, assumed to be random and called for *response* or *outcome*.
- ➤ x is independent variable assumed to be **non-random**, i.e. we focus on a *fixed x-case*. (Random X-case will be considered later in the course).
- ► For the random error  $\varepsilon$ , assume that  $E(\varepsilon) = 0$  and  $V(\varepsilon) = \sigma^2$ ( $\sigma^2$  is unknown)
- β<sub>0</sub> and β<sub>1</sub> are called *regression coefficients*, assumed non-random.
- ►  $E(Y) = \beta_0 + \beta_1 x$  and  $V(Y) = \sigma^2$ , because  $\beta_0, \beta_1$  and x are non-random.

(ロ) (同) (E) (E) (E)

## BRING THE data INTO CONSIDERATION

Given is a set of paired data points  $(y_i, x_i)$  obtained for i = 1, ..., n observational units; each  $y_i$  is the observed value of a r.v.  $Y_i$ , i.e. a generic observation i is modeled as

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

Besides linearity of the model, we assume that

- The linear regression equation is **correct**:  $E(\varepsilon_i) = 0$  for all *i*.
- ▶ All *x*'s are **exact**, i.e we can observe *x<sub>i</sub>*'s perfectly.
- The variance of the errors is constant (homoscedasticity), i.e. V(ε<sub>i</sub>) = σ<sup>2</sup> for all i.
- ▶  $\varepsilon_i$  are pairwise **uncorrelated**, i.e  $Cov(\varepsilon_i, \varepsilon_j) = 0$  for all  $i \neq j$ .
- Since the only random element in the model is  $\varepsilon_i$ ,  $Y_i$ 's have also  $E(Y_i) = \beta_0 + \beta_1 x_i$  and **common variance**  $\sigma^2$ .
- ▶ For purpose of making tests of significance, we assume that  $\varepsilon_i$  are iid  $N(0, \sigma^2)$ . Then  $Y_i$  are iid and  $Y_i \in N(\beta_0 + \beta_1 x_i, \sigma^2)$ .

(ロ) (同) (E) (E) (E)

# A LINEAR FIT



FIGURE: Drawing the linear regression line for Advertising+ data (see Ex 1 in Ch 3 of ISL). The plot displays the regression of sales+ onto TV+. The linear fit is obtained by minimizing the sum of squared errors. Each gray vertical line represents an error and the fit make a *compromise* by averaging the squares of errors.



# LEAST SQUARES FIT

- Example with Advertising+ above.
- Q: How to fit a straight line that fits the data well?
- Our strategy is to fit the line in such a way such that the squared errors are minimal. This is called Least Squares (LS) fitting.
- ► Residuals vs. errors. The residual e<sub>i</sub> = y<sub>i</sub> ŷ<sub>i</sub> is the the difference between the observed and the fitted y-value for the *i*th observation. Residuals are numerical realizations of random errors ε<sub>i</sub>.
- Illustration of residuals: see white board.

LS fitting strategy: to fit the regression line for which the sum of squared residuals is minimized, i. e.  $\sum_{i=1}^{n} e_i^2 \rightarrow \min$ .



# LEAST SQUARES FIT (CONT.)

More precisely, given a set of observations  $(y_i, x_i)$ , i = 1, ..., n the goal is to obtain estimators of  $\beta_0$  and  $\beta_1$ , (say  $\hat{\beta}_0$  and  $\hat{\beta}_1$ ) which minimize the *LS objective function* 

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2.$$

Solution strategy: taking partial derivatives on  $S(\beta_0, \beta_1)$  with respect to  $\beta_0$  and  $\beta_1$ , and setting them to zero.

$$\frac{\partial S(\beta_0,\beta_1)}{\partial \beta_0}|_{\hat{\beta}_0,\hat{\beta}_1}=0, \quad \frac{\partial S(\beta_0,\beta_1)}{\partial \beta_1}|_{\hat{\beta}_0,\hat{\beta}_1}=0.$$

This results in a system of linear equations, called for *normal* equations. These can be solved simultaneously to yield computing formulas for the *LS* estimates,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  in terms of the observed data  $(y_i, x_i)$ , i = 1, ..., n.

回 と く ヨ と く ヨ と

### LEAST SQUARES SOLUTION

**The LS estimates** of  $\beta_0$  and  $\beta_1$  are computed as

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}.$$

The fitted model:  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ Residuals:  $e_i = y_i - \hat{y}_i = \text{observed} - \text{predicted}.$ 

$$SS_{\text{Res}} = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2.$$

A convenient working formula:  $SS_{\text{Res}} = SS_{\text{T}} - \hat{\beta}_1 S_{xy}$ , where  $SS_{\text{T}} = \sum_{i=1}^{n} (y_i - \bar{y})^2$ . See mathematics on the board.

We also obtain the **estimate of the variance of the error term** (residual mean square – measure of residuals variability):

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y})^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{SS_{\text{Res}}}{n-2} = MS_{\text{Res}}$$

## PROPERTIES OF LS ESTIMATORS

The LS estimators of  $\beta_0$  and  $\beta_1$  have several important properties:

•  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are **linear combinations** of the observations  $y_i$ :

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \sum_{i=1}^n c_i y_i, \ c_i = \frac{x_i - \bar{x}}{S_{xx}}, \ i = 1, \dots, n.$$

• Assuming that the model is correct i.e.  $E(Y_i) = \beta_0 + \beta_1 x_i$ ,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are **unbiased**:

$$E(\hat{eta}_0)=eta_0,\quad E(\hat{eta}_1)=eta_1$$

The variances of β̂<sub>0</sub> and β̂<sub>1</sub> are found as (see math and interpretation of V(β̂<sub>1</sub>) on the board)

$$V(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right), \quad V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}$$

.

(신문) (신문)

# **OPTIMALITY OF LS ESTIMATORS**

Gauss-Markov theorem states the result on the mathematical optimality of the LS: assuming that the following conditions are met,

•  $E(\varepsilon_i) = 0$ , (i.e. the relation is a straight line),

• 
$$V(\varepsilon_i) = \sigma$$
 for  $i = 1, ..., n$ ,

Cov(ε<sub>i</sub>, ε<sub>j</sub>) = 0 for i ≠ j = 1,..., n, (the errors are uncorrelated),

the LS estimators are unbiased and have minimum variance when compared with all other unbiased estimators that are linear combinations of  $y_i$ .

- By Gauss-Markov thm, LS estimators are Best Linear Unbiased Estimators, (BLUE), where *best* means minimum variance.
- Later on in the course: Gauss-Markov theorem for the more general multiple regression, of which simple linear regression is a special cause.

# BENEFITS OF LS APPROACH

Mathematically ...

- ► The fitted line goes through the center of gravity (x̄, ȳ) of the data.
- The LS technique is simple in the sense that the solutions, β̂<sub>0</sub> and β̂<sub>1</sub>, are obtained in closed form as functions of (x<sub>i</sub>, y<sub>i</sub>).
- $\hat{\beta}_0$  and  $\hat{\beta}_1$  are unbiased estimators
- Large sample properties of the LS estimates have some deeper mathematical advantages (approximate normality of coefficient estimators), and exact properties assuming Gaussian errors. Will be presented later in the course.



(4) (3) (4) (3) (4)