

STAC67: Regression Analysis

Lecture 1

Sohee Kang

Jan. 13, 2021

Regression: What is it?

- Simply: The **most widely used** statistical tool for understanding relationships among variables.
- A conceptually simple method for **investigating relationships between one or more factors of outcome of interest.**
- The relationship is expressed in the form of an equation or a model connecting the outcome to the factors.
- Examples of application:
 - **Epidemiology**: what are the social factors to contribute the death rate of Covid-19 in Canada?
 - **Business**: determining price and marketing strategy:
 - **Estimate** the effect of price and advertisement on sales
 - **Decide** what is optimal price and campaign?
 - Straight prediction questions:
 - What will be the interest rates be next month?
 - What will be the length of hospital stay of a surgical patient?
 - Explanation and understanding:
 - Much more . . .

Origin of Regression (History)

- In 1886, Sir Francis Galton (1822-1911) invented the term, **"regression"** after he analyzed the data on the heights of parents and their children.

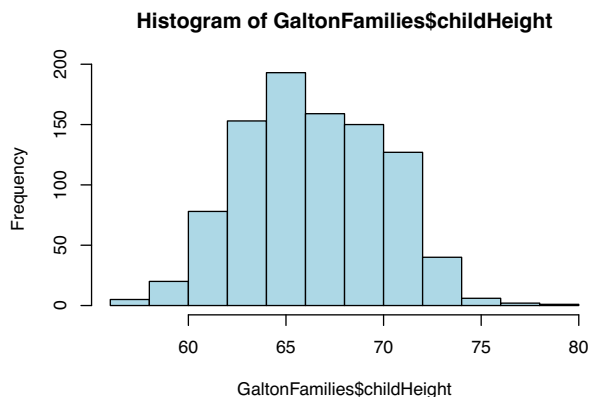
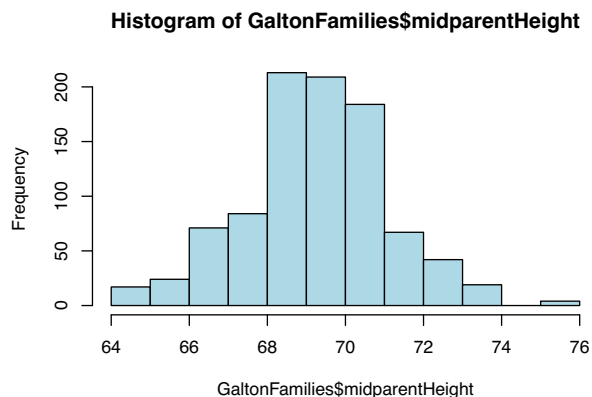


"Regression towards mediocrity in hereditary stature". The Journal of the Anthropological Institute of Great Britain and Ireland, Vol 15, pages 246-263 (or Wikipedia!)

Data Visualization

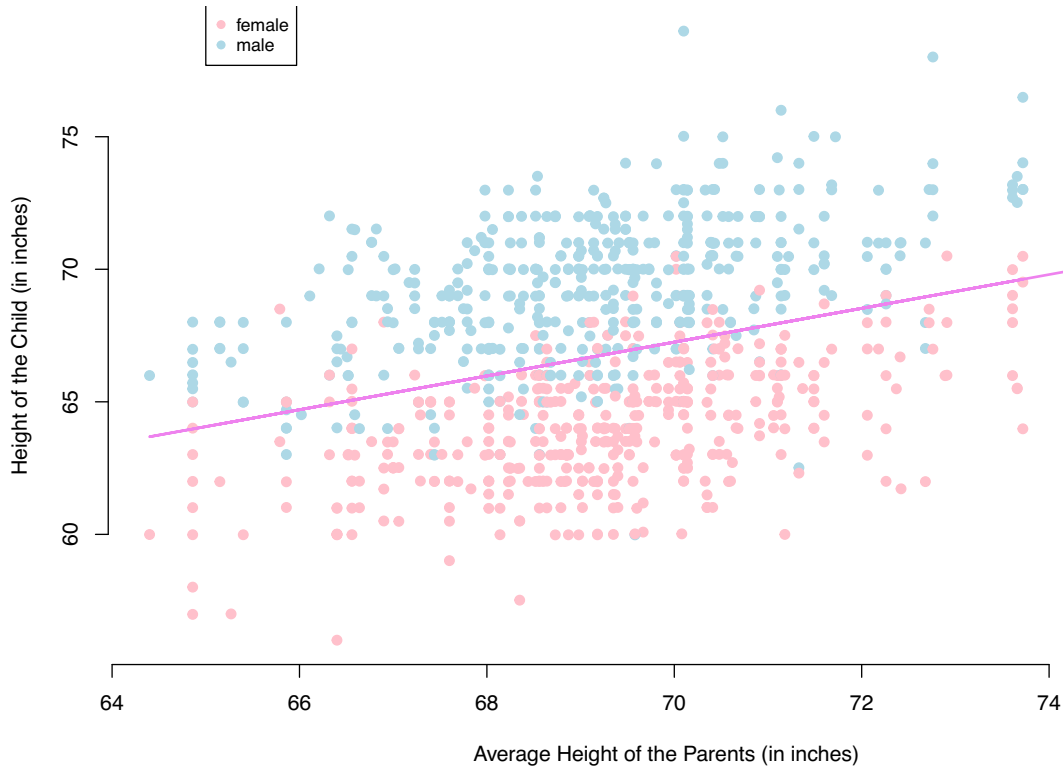
Example: Galton family data

- The GaltonFamilies(HistData) dataset lists the individual observations for 934 adult children born to 205 fathers and mothers. He wrote that, “the average regression of the offspring is a constant fraction of their respective mid-parental deviations.” For height, Galton estimated this regression coefficient to be about two-thirds ($2/3$).



Data Visualization

Figure 1. Scatterplot of Galton Family Data with Fitted Values



Review

- Covariance and Correlation Coefficient

Suppose we have observations on n subjects consisting of a **dependent** or **response variable** Y and an **independent** or **explanatory variable** X .

- Measure both **direction** and **strength** of the relationship between Y and X .

Obs	Y	X
1	y_1	x_1
2	y_2	x_2
\vdots	\vdots	\vdots
n	y_n	x_n

Covariance and Correlation

calculate this $\int \int (x - \mu_x)(y - \mu_y) f_{xy} dx dy$

Def. $\text{Cov}(X, Y) = E((X - \mu_x)(Y - \mu_y))$, where $\mu_x = E(X)$, $\mu_y = E(Y)$

$$Z_x = \frac{X - \mu_x}{\sqrt{\text{Var}(X)}}, \quad Z_y = \frac{Y - \mu_y}{\sqrt{\text{Var}(Y)}}$$

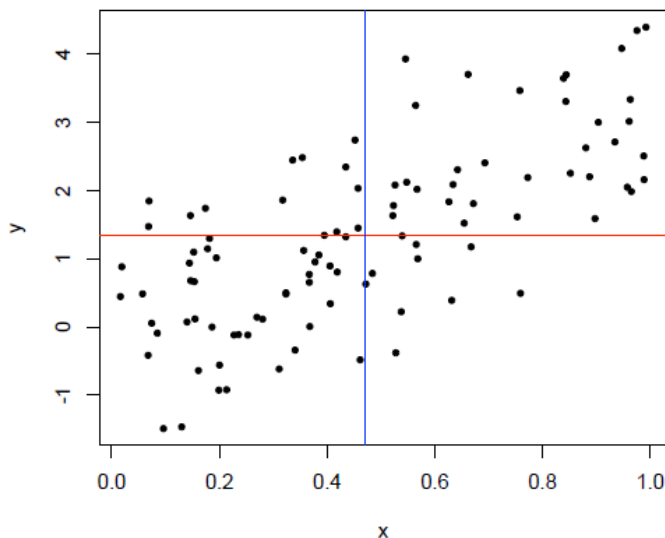
standardized X standardized Y

$$\text{Cov}(Z_x, Z_y) = \rho_{xy} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \text{Corr}(X, Y)$$

$$\begin{aligned} &= E((Z_x - \mu_{Z_x})(Z_y - \mu_{Z_y})) \\ &= E(Z_x Z_y) = \frac{1}{\sigma_x \sigma_y} E((X - \mu_x)(Y - \mu_y)) \end{aligned}$$

- $-1 \leq \rho_{xy} \leq 1$
- When the relationship is perfectly linear then $|\rho| = 1$.
- if two variables are independent then $\rho = 0$. (Note: the inverse does not hold)

Sample Covariance and Correlation



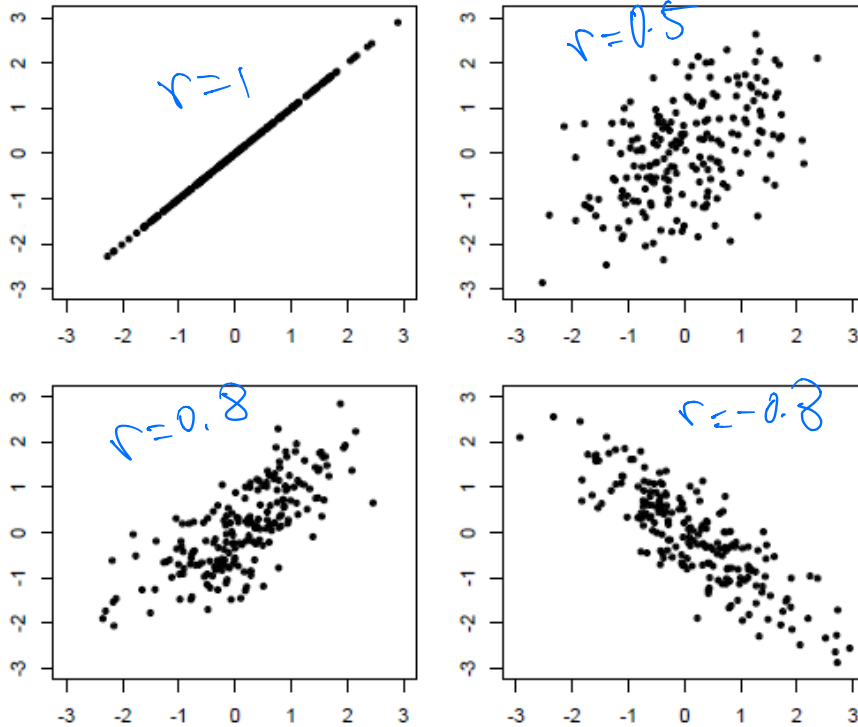
$$\text{Cov}(X, Y) = \frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{n-1}$$

= Sample Covariance

$$\text{Let } Z_y = \frac{y - \bar{y}}{s_y} \quad \text{Let } Z_x = \frac{x - \bar{x}}{s_x}$$

$$\text{Cor}(Y, X) = \text{Cov}(Z_y, Z_x) = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{y_i - \bar{y}}{s_y} \right) \left(\frac{x_i - \bar{x}}{s_x} \right) = \text{Sample Correlation Coefficient}$$

Correlation



$$r_{xy} = r_{yx}$$

Question: what are main differences between correlation and regression model?
Corr: measures linear relation between x and y

Test for Population correlation

~ population correlation

When $\rho = 0$, and the joint distribution of (X, Y) is bivariate normal, and it can be shown that:

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

has a student's t distribution with $n - 2$ degrees of freedom

$$H_0 : \rho = 0 \quad \text{vs} \quad H_1 : \rho \neq 0$$

- Testing Procedure

- Calculating the observed value of t (call this t_{obs})
- Compute the p-value for the test

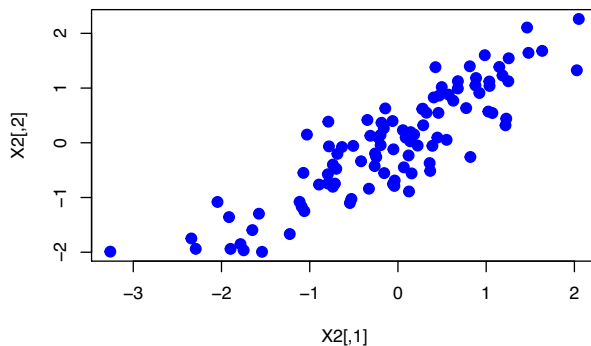
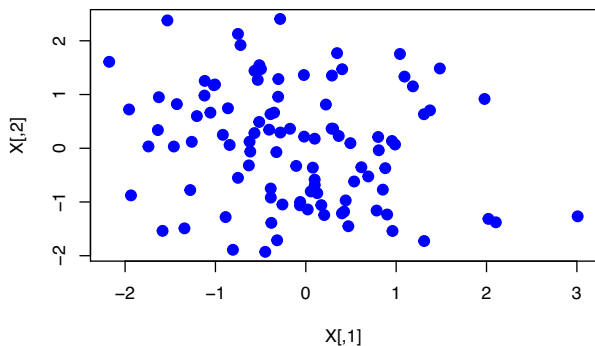
Simulation

```
par(mfrow=c(2,2))  
library(mvtnorm) — multivariate normal package
```

Warning: package 'mvtnorm' was built under R version 3.5.3

```
sigma.1 = matrix(c(1, 0, 0, 1), ncol=2)  
sigma.2 = matrix(c(1, 0.9, 0.9, 1), ncol=2)  
X = rmvnorm(100, mean=c(0,0), sigma.1)  
plot(X, pch=20, cex=2, col="blue")  
X2 = x = rmvnorm(100, mean=c(0,0), sigma.2)  
plot(X2, pch=20, cex=2, col="blue")
```

Handwritten notes:
- σ_1 (under $\sigma.1$)
- σ_2 (under $\sigma.2$)
- μ (under mean)
- generate sample (under rmvnorm)



Simulation

```
x = X2[,1]
y = X2[, 2]
cor.test(x, y)
```

```
##
## Pearson's product-moment correlation
##
## data: x and y
## t = 18.084, df = 98, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.822453 0.915802
## sample estimates:
##      cor
## 0.8771691
```

```
r= cor(x, y)
t = r*sqrt(98)/sqrt(1-r^2)
t
```

```
## [1] 18.08385
```

For $H_0: \rho = 0$ $H_1: \rho \neq 0$

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t(98) \\ = 20.73$$

$$\alpha = 5\%$$

since $p\text{-value} < \alpha$
Thus we reject the null hypothesis

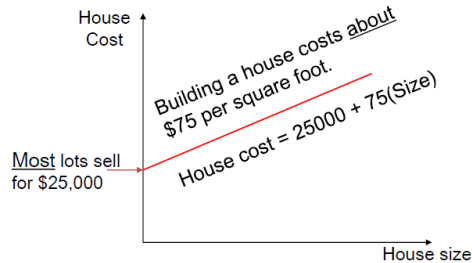
Relationship between variables

What factor or variable affects the price of house?

- Relation of the form

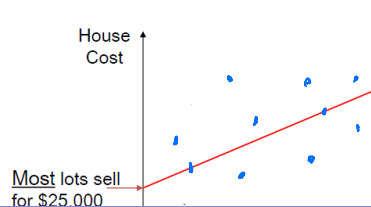
① Mathematical Relation: $Y = f(X)$

where, X , Y are variables and f is a function



② Statistical Relation:

$$Y = f(X) + \epsilon$$



Data Collection for regression analysis

- **Observational study**

- Investigator has no control over the explanatory variables (X)
- Limitation: not adequate for cause-and-effect **A strong association does not necessarily mean a cause-and-effect relationship**

- **Experiment**

- Investigator exercises control over the explanatory variables (X) through random assignment
- Random assignment balances out effect of other variables that might affect Y
- **Gold standard for cause-and-effect conclusions**

• Always understand how data is collected

The Regression Process

- ① The researcher must clearly define the question(s) of interest in the study
- ② The response variable Y must be decided on, based on the question of interest.
- ③ A set of potentially relevant covariates, which can be measured, needs to be defined.
- ④ Data is collected.
- ⑤ Model Specification.
- ⑥ Decide on a method for fitting the specified model
- ⑦ Fit the model - typically using software such as R
- ⑧ Examine the fitted model for violations of assumptions.
- ⑨ Conduct hypothesis testing for questions of interest.
- ⑩ Report the results from statistical inference.