# AUA CS108, Statistics, Fall 2020
## Lecture 15

Michael Poghosyan

30 Sep 2020

# Contents

# Q-Q Plots, Theoretical vs Theoretical Distribution

Assume now we have two Theoretical Distributions (say, given by their CDFs $F$ and $G$).

## Q-Q Plots, Theoretical vs Theoretical Distribution

Assume now we have two Theoretical Distributions (say, given by their CDFs $F$ and $G$). The Problem is to estimate visually which Distribution has fatter tails.

# Q-Q Plots, Theoretical vs Theoretical Distribution

Assume now we have two Theoretical Distributions (say, given by their CDFs $F$ and $G$). The Problem is to estimate visually which Distribution has fatter tails.

To answer this question, we again take some levels of quantiles, say, for some $n$,

$$\alpha = \frac{1}{n}, \frac{2}{n}, ..., \frac{n-1}{n}$$

and then draw the points $(q_\alpha^F, q_\alpha^G)$, where $q_\alpha^F$ is the $\alpha$-quantile of the Theoretical Distribution with the CDF $F$, and $q_\alpha^G$ is the $\alpha$-quantile of the Theoretical Distribution with the CDF $G$.

# Q-Q Plots, Theoretical vs Theoretical Distribution

Assume now we have two Theoretical Distributions (say, given by their CDFs $F$ and $G$). The Problem is to estimate visually which Distribution has fatter tails.

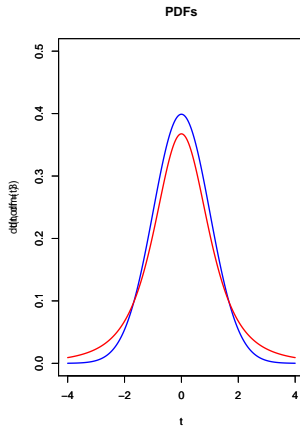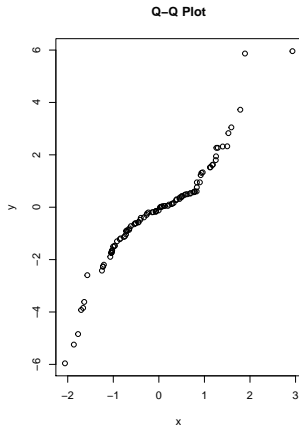To answer this question, we again take some levels of quantiles, say, for some $n$,

$$\alpha = \frac{1}{n}, \frac{2}{n}, ..., \frac{n-1}{n}$$

and then draw the points $(q_\alpha^F, q_\alpha^G)$, where $q_\alpha^F$ is the $\alpha$-quantile of the Theoretical Distribution with the CDF $F$, and $q_\alpha^G$ is the $\alpha$-quantile of the Theoretical Distribution with the CDF $G$.

**Idea:** If $G$ has fatter tails on both sides than $F$, then we will have graphically some cubic-function graph shape Quantiles.
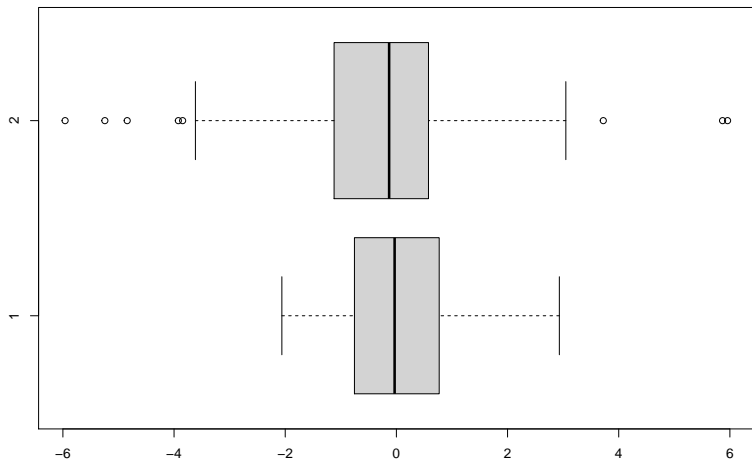
# Some Experiments

```r
par(mfrow = c(1,2))
x <- rnorm(100, mean=0, sd=1); y <- rt(100, df = 3)
qqplot(x,y, main = "Q-Q Plot")
t <- seq(-4,4,0.01)
plot(t, dnorm(t), type = "l", xlim = c(-4,4), ylim = c(0, 0.5), col ="blue", lwd = 2, main = "PDFs")
par(new = TRUE)
plot(t, dt(t, df = 3), type = "l", xlim = c(-4,4), ylim = c(0, 0.5), col ="red", lwd = 2)
```
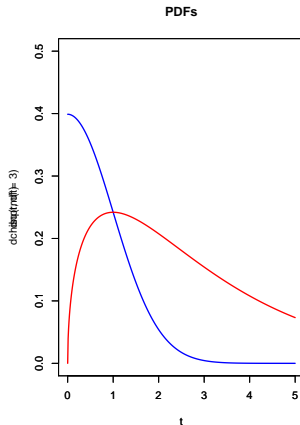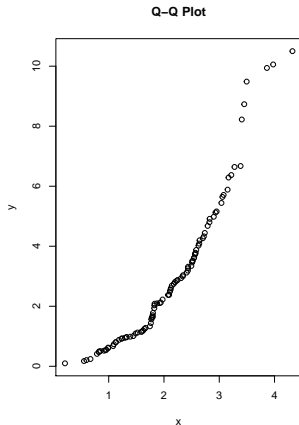
# Some Experiments

The above Datasets, using BoxPlots:

```
boxplot(x,y, horizontal = T)
```

# Some Experiments

```
par(mfrow = c(1,2))
x <- rnorm(100, mean=2, sd=1); y <- rchisq(200, df = 3)
qqplot(x,y, main = "Q-Q Plot")
t <- seq(0,5,0.01)
plot(t, dnorm(t), type = "l", xlim = c(0,5), ylim = c(0, 0.5), col ="blue", lwd = 2,  main = "PDFs")
par(new = TRUE)
plot(t, dchisq(t, df = 3), type = "l", xlim = c(0,5), ylim = c(0, 0.5), col ="red", lwd = 2)
```

# Addition, Q-Q Plot

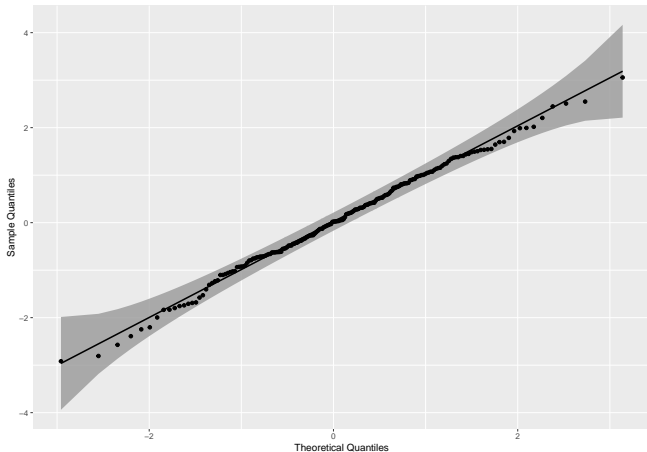here you can find some interpretaitons of different shapes of Q-Q Plots: StackExchange Page.

# Addition, Q-Q Plot with a Confidence Band

```
require(qqplotr)
x <- data.frame(variable = rnorm(200))
ggplot(data = x, mapping = aes(sample = variable)) + stat_qq_band() +
  stat_qq_line() + stat_qq_point() + labs(x = "Theoretical Quantiles", y = "Sample Quantiles")
```

# Numerical Summaries for Bivariate Data

# Sample Covariance and the Correlation Coefficient

Assume now we have a bivariate Dataset

$$(x_1, y_1), ..., (x_n, y_n),$$

or just two 1D Datasets of the same size:

$$x : x_1, ..., x_n \quad and \quad y : y_1, ..., y_n.$$

# Sample Covariance and the Correlation Coefficient

Assume now we have a bivariate Dataset

$$(x_1, y_1), ..., (x_n, y_n),$$

or just two 1D Datasets of the same size:

$$x : \quad x_1, ..., x_n \qquad and \qquad y : \quad y_1, ..., y_n.$$

Our aim is to see if some linear relationship, association exists between $x$ and $y$.

# Sample Covariance and the Correlation Coefficient

Assume now we have a bivariate Dataset

$$(x_1, y_1), ..., (x_n, y_n),$$

or just two 1D Datasets of the same size:

$$x : x_1, ..., x_n \quad and \quad y : y_1, ..., y_n.$$

Our aim is to see if some linear relationship, association exists between $x$ and $y$. Of course, the best way is to visualize our Dataset by a ScatterPlot.

# Sample Covariance and the Correlation Coefficient

Assume now we have a bivariate Dataset

$$(x_1, y_1), ..., (x_n, y_n),$$

or just two 1D Datasets of the same size:

$$x : x_1, ..., x_n \quad and \quad y : y_1, ..., y_n.$$

Our aim is to see if some linear relationship, association exists between $x$ and $y$. Of course, the best way is to visualize our Dataset by a ScatterPlot.

Now we want to answer, numerically, how strong/week is the linear relationship between our variables $x$ and $y$.

# Sample Covariance

The **Sample Covariance** of Variables (1D Datasets) $x$ and $y$ is

$$cov(x, y) = s_{xy} = \frac{\sum\limits_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n}$$

# Sample Covariance

The **Sample Covariance** of Variables (1D Datasets) $x$ and $y$ is

$$cov(x, y) = s_{xy} = \frac{\sum\limits_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n}$$

or

$$cov(x, y) = s_{xy} = \frac{\sum\limits_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n - 1}$$

# Sample Covariance

The **Sample Covariance** of Variables (1D Datasets) $x$ and $y$ is

$$cov(x, y) = s_{xy} = \frac{\sum\limits_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n}$$

or

$$cov(x, y) = s_{xy} = \frac{\sum\limits_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n - 1}$$

Here $\overline{x}$ and $\overline{y}$ are the Sample Means for the Datasets $x$ and $y$.

# Sample Covariance

The **Sample Covariance** of Variables (1D Datasets) $x$ and $y$ is

$$cov(x, y) = s_{xy} = \frac{\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n}$$

or

$$cov(x, y) = s_{xy} = \frac{\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n - 1}$$

Here $\overline{x}$ and $\overline{y}$ are the Sample Means for the Datasets $x$ and $y$.

**Note:** Recall that for a r.v. $X$, $Cov(X, X) = Var(X)$.

# Sample Covariance

The **Sample Covariance** of Variables (1D Datasets) $x$ and $y$ is

$$cov(x, y) = s_{xy} = \frac{\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n}$$

or

$$cov(x, y) = s_{xy} = \frac{\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{n - 1}$$

Here $\overline{x}$ and $\overline{y}$ are the Sample Means for the Datasets $x$ and $y$.

**Note:** Recall that for a r.v. $X$, $Cov(X, X) = Var(X)$. Here, for Datasets, we have two definitions for the Sample Variance $var(x)$. And we give two definitions of the Sample Covariance, so the property $cov(x, x) = var(x)$ will hold in both cases.

# Sample Covariance

**Definition:** We say that the Variables (Datasets) $x$ and $y$ are **uncorrelated**, if $cov(x, y) = 0$.

# Sample Covariance

**Definition:** We say that the Variables (Datasets) $x$ and $y$ are **uncorrelated**, if $cov(x, y) = 0$.

**Remark:** In Probability, we have 2 notions: *Independence* and *Correlation*. Here, in the case of Datasets, we do not have the notion of *Independence*.

# Sample Covariance

**Definition:** We say that the Variables (Datasets) $x$ and $y$ are **uncorrelated**, if $cov(x, y) = 0$.

**Remark:** In Probability, we have 2 notions: *Independence* and *Correlation*. Here, in the case of Datasets, we do not have the notion of *Independence*.

**Remark:** For almost all numerical summaries for 1D data, first step was sorting the Dataset to obtain Order Statistics. But please note that for calculating Covariance or Correlation Coefficient (as well as for ScatterPlotting), sorting the Datasets will give incorrect results.

# Sample Covariance

**Definition:** We say that the Variables (Datasets) $x$ and $y$ are **uncorrelated**, if $cov(x, y) = 0$.

**Remark:** In Probability, we have 2 notions: *Independence* and *Correlation*. Here, in the case of Datasets, we do not have the notion of *Independence*.

**Remark:** For almost all numerical summaries for 1D data, first step was sorting the Dataset to obtain Order Statistics. But please note that for calculating Covariance or Correlation Coefficient (as well as for ScatterPlotting), sorting the Datasets will give incorrect results. This is because we want to find a relationship between $x_1$ and $y_1$, $x_2$ and $y_2$, ..., not the relationship between the minimal elements of Datasets etc.

# Example

Here is the **R** code to calculate the Covariance between the Speed and Dist variables in the cars Dataset:

```r
cov(cars$speed, cars$dist)
```

```
## [1] 109.9469
```

# Sample Correlation Coefficient

Another measure of the linear relationship between the Variables $x$ and $y$ of Bivariate Dataset is the *Pearson's Correlation Coefficient*:

# Sample Correlation Coefficient

Another measure of the linear relationship between the Variables $x$ and $y$ of Bivariate Dataset is the *Pearson's Correlation Coefficient*:

**Definition:** The **Sample Correlation Coefficient** of $x$ and $y$ is

$$cor(x, y) = \rho_{xy} = \frac{cov(x, y)}{\sqrt{Var(x) \cdot Var(y)}} = \frac{cov(x, y)}{sd(x) \cdot sd(y)} = \frac{s_{xy}}{s_x \cdot s_y},$$

where $s_x$ and $s_y$ are the standard deviations for $x$ and $y$, respectively.

# Sample Correlation Coefficient

Another measure of the linear relationship between the Variables $x$ and $y$ of Bivariate Dataset is the *Pearson's Correlation Coefficient*:

**Definition:** The **Sample Correlation Coefficient** of $x$ and $y$ is

$$cor(x, y) = \rho_{xy} = \frac{cov(x, y)}{\sqrt{Var(x) \cdot Var(y)}} = \frac{cov(x, y)}{sd(x) \cdot sd(y)} = \frac{s_{xy}}{s_x \cdot s_y},$$

where $s_x$ and $s_y$ are the standard deviations for $x$ and $y$, respectively.

If $s_x = 0$ or $s_y = 0$, then we take $cor(x, y) = 0$ by definition.

# Sample Correlation Coefficient

Another measure of the linear relationship between the Variables $x$ and $y$ of Bivariate Dataset is the *Pearson's Correlation Coefficient*:

**Definition:** The **Sample Correlation Coefficient** of $x$ and $y$ is

$$cor(x, y) = \rho_{xy} = \frac{cov(x, y)}{\sqrt{Var(x) \cdot Var(y)}} = \frac{cov(x, y)}{sd(x) \cdot sd(y)} = \frac{s_{xy}}{s_x \cdot s_y},$$

where $s_x$ and $s_y$ are the standard deviations for $x$ and $y$, respectively.

If $s_x = 0$ or $s_y = 0$, then we take $cor(x, y) = 0$ by definition.

**Note:** Please note that we need to calculate the Standard Deviations and Covariance by using the same denominator: either everywhere take $n$, or take everywhere $n - 1$.

## Sample Correlation Coefficient

In both cases, when one calculates Standard Deviations and Covariance by using $n$ simultaneously or $n - 1$ simultaneously in the denominator, we will obtain

$$cor(x, y) = \rho_{xy} = \frac{\displaystyle\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{\sqrt{\displaystyle\sum_{k=1}^{n}(x_k - \overline{x})^2 \cdot \sum_{k=1}^{n}(y_k - \overline{y})^2}}$$

## Sample Correlation Coefficient

In both cases, when one calculates Standard Deviations and Covariance by using $n$ simultaneously or $n-1$ simultaneously in the denominator, we will obtain

$$cor(x,y) = \rho_{xy} = \frac{\sum_{k=1}^{n}(x_k - \overline{x}) \cdot (y_k - \overline{y})}{\sqrt{\sum_{k=1}^{n}(x_k - \overline{x})^2 \cdot \sum_{k=1}^{n}(y_k - \overline{y})^2}}$$
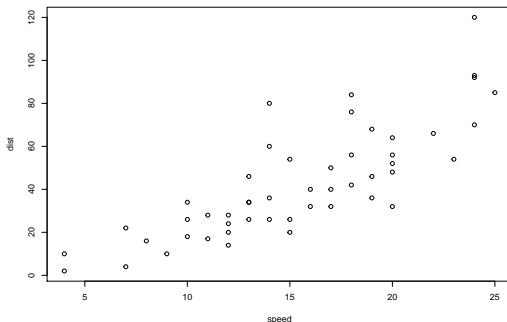
Another formula to calc the correlation coefficient is

$$cor(x,y) = \rho_{xy} = \frac{\sum_{k=1}^{n}x_k y_k - n \cdot \overline{x} \cdot \overline{y}}{\sqrt{\sum_{k=1}^{n}x_k^2 - n \cdot (\overline{x})^2} \cdot \sqrt{\sum_{k=1}^{n}y_k^2 - n \cdot (\overline{y})^2}}.$$

## Examples:

Now, the **R** code to calculate the Covariance between the Speed and Dist variables in the cars Dataset:

```r
plot(dist~speed, data = cars)
```



```r
cor(cars$speed, cars$dist)
```

```
## [1] 0.8068949
```