

# AUA CS 108, Statistics, Fall 2019

## Lecture 11

Michael Poghosyan

YSU, AUA

[michael@ysu.am](mailto:michael@ysu.am), [mpoghosyan@aua.am](mailto:mpoghosyan@aua.am)

18 Sep 2019

# Contents

- ▶ Q-Q Plots
- ▶ Sample Covariance and Correlation

## Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?

## Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?
- ▶ What is a Theoretical  $\alpha$ -Quantile?

## Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?
- ▶ What is a Theoretical  $\alpha$ -Quantile?
- ▶ When do we will write our next Quiz?

## Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?
- ▶ What is a Theoretical  $\alpha$ -Quantile?
- ▶ When do we will write our next Quiz?
- ▶ What is the 35% Quantile of the Distribution with the PDF

$$f(x) = \begin{cases} 5x^4, & x \in [0, 1] \\ 0, & \textit{otherwise} \end{cases} \quad ?$$

## Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?
- ▶ What is a Theoretical  $\alpha$ -Quantile?
- ▶ When do we will write our next Quiz?
- ▶ What is the 35% Quantile of the Distribution with the PDF

$$f(x) = \begin{cases} 5x^4, & x \in [0, 1] \\ 0, & \textit{otherwise} \end{cases} \quad ?$$

- ▶ What is a Q-Q Plot?

# Last Lecture ReCap

- ▶ What is a Sample  $\alpha$ -Quantile?
- ▶ What is a Theoretical  $\alpha$ -Quantile?
- ▶ When do we will write our next Quiz?
- ▶ What is the 35% Quantile of the Distribution with the PDF

$$f(x) = \begin{cases} 5x^4, & x \in [0, 1] \\ 0, & \textit{otherwise} \end{cases} \quad ?$$

- ▶ What is a Q-Q Plot?
- ▶ What is it for?



## Q-Q Plots, Data vs Data

**Problem:** we have two Datasets, not necessarily of the same size:

$$x : x_1, x_2, \dots, x_n \quad \text{and} \quad y : y_1, y_2, \dots, y_m,$$

and we want to see if  $x$  and  $y$  are coming from the same Distribution

## Q-Q Plots, Data vs Data

**Problem:** we have two Datasets, not necessarily of the same size:

$$x : x_1, x_2, \dots, x_n \quad \text{and} \quad y : y_1, y_2, \dots, y_m,$$

and we want to see if  $x$  and  $y$  are coming from the same Distribution

To draw a **Q-Q Plot** (Quantile-Quantile Plot), we take some levels of quantiles, say, for some  $p$ ,

$$\alpha = \frac{1}{p}, \frac{2}{p}, \dots, \frac{p-1}{p}$$

and then draw the points  $(q_\alpha^x, q_\alpha^y)$ .

## Q-Q Plots, Data vs Data

**Problem:** we have two Datasets, not necessarily of the same size:

$$x : x_1, x_2, \dots, x_n \quad \text{and} \quad y : y_1, y_2, \dots, y_m,$$

and we want to see if  $x$  and  $y$  are coming from the same Distribution

To draw a **Q-Q Plot** (Quantile-Quantile Plot), we take some levels of quantiles, say, for some  $p$ ,

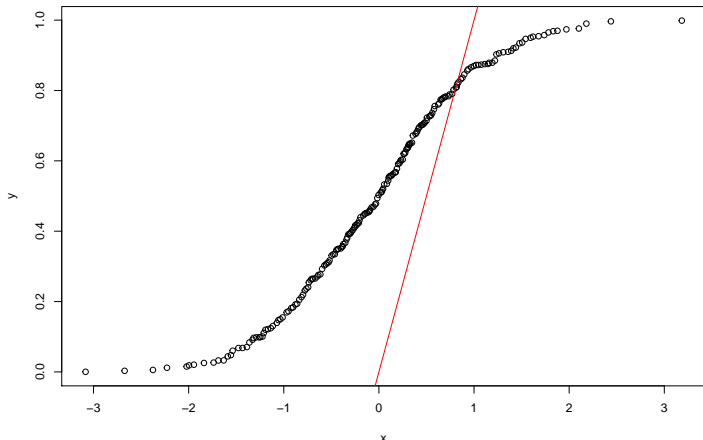
$$\alpha = \frac{1}{p}, \frac{2}{p}, \dots, \frac{p-1}{p}$$

and then draw the points  $(q_\alpha^x, q_\alpha^y)$ .

**Idea:** If  $x$  and  $y$  are coming from the same Distribution, then the Quantiles of  $x$  and  $y$  need to be approximately the same,  $q_\alpha^x \approx q_\alpha^y$ , so geometrically, the points  $(q_\alpha^x, q_\alpha^y)$  need to be close to the bisector line.

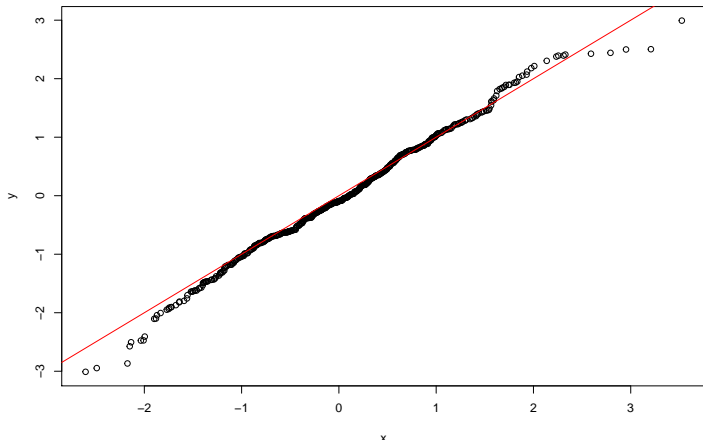
## Example, Q-Q Plots, Data vs Data

```
x <- rnorm(1000)
y <- runif(200)
qqplot(x,y)
abline(0,1, col="red")
```



## Example, Q-Q Plots, Data vs Data

```
x <- rnorm(1000)
y <- rnorm(500)
qqplot(x,y)
abline(0,1, col="red")
```



## Example, Q-Q Plot by Hands, Data vs Data

**Example:** Assume

$$x : -1, 2, 1, 2, 3, 2, 1 \quad y : 0, 3, 4, 1, 1, 1, 1, 2$$

Draw the Q-Q Plot for  $x$  and  $y$ .

## Q-Q Plots, Data vs Theoretical Distribution

Assume now we have a Dataset  $x$  and a Theoretical Distribution (say, given by its CDF  $F$  or PDF  $f$ ).

## Q-Q Plots, Data vs Theoretical Distribution

Assume now we have a Dataset  $x$  and a Theoretical Distribution (say, given by its CDF  $F$  or PDF  $f$ ). The Problem is to estimate visually if the Dataset comes from that Distribution.



## Q-Q Plots, Data vs Theoretical Distribution

Assume now we have a Dataset  $x$  and a Theoretical Distribution (say, given by its CDF  $F$  or PDF  $f$ ). The Problem is to estimate visually if the Dataset comes from that Distribution.

**Example:** Say, is the following Dataset

```
## [1] -0.83 -0.84 -0.61 -0.40  0.56 -0.96  0.58  0.41 -0.01
## [13]  0.45 -0.63  0.94  0.22 -0.54 -0.28 -0.12 -0.87
```

from a Normal Distribution?

## Q-Q Plots, Data vs Theoretical Distribution

Assume now we have a Dataset  $x$  and a Theoretical Distribution (say, given by its CDF  $F$  or PDF  $f$ ). The Problem is to estimate visually if the Dataset comes from that Distribution.

**Example:** Say, is the following Dataset

```
## [1] -0.83 -0.84 -0.61 -0.40  0.56 -0.96  0.58  0.41 -0.
## [13]  0.45 -0.63  0.94  0.22 -0.54 -0.28 -0.12 -0.87
```

from a Normal Distribution?

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

$$\alpha = \frac{1}{p}, \frac{2}{p}, \dots, \frac{p-1}{p}$$

and then draw the points  $(q_{\alpha}^F, q_{\alpha}^x)$ , where  $q_{\alpha}^F$  is the  $\alpha$ -quantile of the Theoretical Distribution, and  $q_{\alpha}^x$  is the  $\alpha$ -quantile of  $x$ .

## Q-Q Plots, Data vs Theoretical Distribution

Assume now we have a Dataset  $x$  and a Theoretical Distribution (say, given by its CDF  $F$  or PDF  $f$ ). The Problem is to estimate visually if the Dataset comes from that Distribution.

**Example:** Say, is the following Dataset

```
## [1] -0.83 -0.84 -0.61 -0.40  0.56 -0.96  0.58  0.41 -0.
## [13]  0.45 -0.63  0.94  0.22 -0.54 -0.28 -0.12 -0.87
```

from a Normal Distribution?

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

$$\alpha = \frac{1}{p}, \frac{2}{p}, \dots, \frac{p-1}{p}$$

and then draw the points  $(q_{\alpha}^F, q_{\alpha}^x)$ , where  $q_{\alpha}^F$  is the  $\alpha$ -quantile of the Theoretical Distribution, and  $q_{\alpha}^x$  is the  $\alpha$ -quantile of  $x$ .

**Idea:** If  $x$  is from the Distribution given by  $F$ , then we need to have  $q_{\alpha}^F \approx q_{\alpha}^x$ , so, graphically, the point will be close to the bisector.

## Normal Q-Q Plot

In **R**, we have a function `qqnorm` which plots the Q-Q Plot for the Dataset  $x$  vs the Normal Distribution.

## Normal Q-Q Plot

In **R**, we have a function `qqnorm` which plots the Q-Q Plot for the Dataset  $x$  vs the Normal Distribution. Unfortunately, we do not have this kind of function for other standard distributions, say, Uniform.

---

<sup>1</sup>or one can write his/her own function `qqunif` or `qqexp`, say

## Normal Q-Q Plot

In **R**, we have a function `qqnorm` which plots the Q-Q Plot for the Dataset  $x$  vs the Normal Distribution. Unfortunately, we do not have this kind of function for other standard distributions, say, Uniform. But one can use the `qqplot(x,y)` command, by generating  $y$  from the given Distribution<sup>1</sup>.

---

<sup>1</sup>or one can write his/her own function `qqunif` or `qqexp`, say

## Normal Q-Q Plot

In **R**, we have a function `qqnorm` which plots the Q-Q Plot for the Dataset  $x$  vs the Normal Distribution. Unfortunately, we do not have this kind of function for other standard distributions, say, Uniform. But one can use the `qqplot(x,y)` command, by generating  $y$  from the given Distribution<sup>1</sup>.

Another **R** command is `qqline` which adds a line passing (by default) through the first and third Quartiles,

$$(q_{0.25}^F, q_{0.25}^x) \quad \text{and} \quad (q_{0.75}^F, q_{0.75}^x).$$

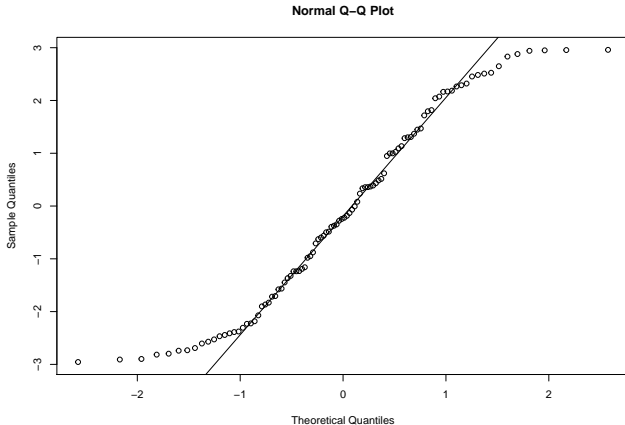
---

<sup>1</sup>or one can write his/her own function `qqunif` or `qqexp`, say

# Some Experiments

Here are some experiments with `qqnorm`

```
x <- runif(100, -3, 3)
qqnorm(x)
qqline(x)
```

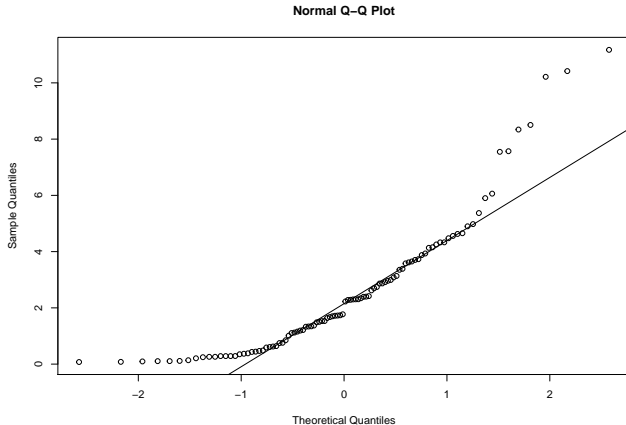




# Some Experiments

Here are some experiments with `qqnorm`

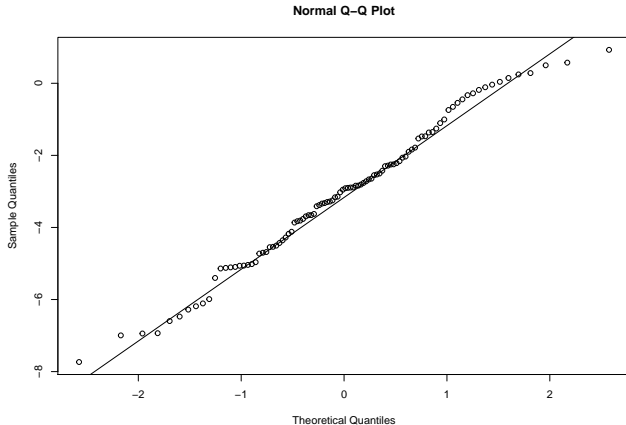
```
x <- rexp(100,0.4)
qqnorm(x)
qqline(x)
```



# Some Experiments

Here are some experiments with `qqnorm`

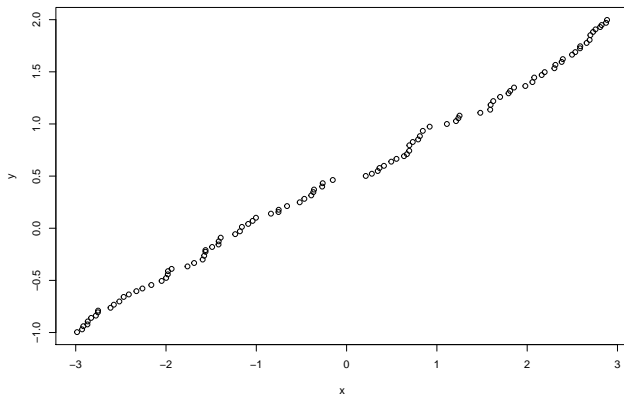
```
x <- rnorm(100, mean = -3, sd = 2)
qqnorm(x)
qqline(x)
```



## Some Experiments

Now, assume we want to see if our Dataset  $x$  is from  $Unif[-1, 2]$ :

```
x <- runif(100, -3, 3)
y <- runif(1000, -1, 2)
qqplot(x, y)
```



## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*.

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

---

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting.

---

<sup>2</sup>Can you state rigorously and prove this?

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting. This means that the Quantiles of any Normal Distribution can be obtained by a linear transform from the Standard Normal Quantiles<sup>2</sup>.

---

<sup>2</sup>Can you state rigorously and prove this?

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting. This means that the Quantiles of any Normal Distribution can be obtained by a linear transform from the Standard Normal Quantiles<sup>2</sup>.

So if, say,  $x$  is a sample from  $\mathcal{N}(2, 3^2)$ , then

- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(2, 3^2)$ , the Quantiles will be

---

<sup>2</sup>Can you state rigorously and prove this?



## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting. This means that the Quantiles of any Normal Distribution can be obtained by a linear transform from the Standard Normal Quantiles<sup>2</sup>.

So if, say,  $x$  is a sample from  $\mathcal{N}(2, 3^2)$ , then

- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(2, 3^2)$ , the Quantiles will be on the bisector;

---

<sup>2</sup>Can you state rigorously and prove this?

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting. This means that the Quantiles of any Normal Distribution can be obtained by a linear transform from the Standard Normal Quantiles<sup>2</sup>.

So if, say,  $x$  is a sample from  $\mathcal{N}(2, 3^2)$ , then

- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(2, 3^2)$ , the Quantiles will be on the bisector;
- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(0, 1)$ , the Quantiles will be

---

<sup>2</sup>Can you state rigorously and prove this?

## Important Note

It is important, that, using `qqnorm`, we can check if our Dataset comes from a Normal Distribution, *with some mean and variance*. I mean, the above idea was, say, to check if given Dataset  $x$  comes from given Distribution, say,  $\mathcal{N}(2, 3^2)$ .

But, for the Normal Distribution, we can use the fact that all Normal Distributions can be obtained from the Standard Normal, by scaling and shifting. This means that the Quantiles of any Normal Distribution can be obtained by a linear transform from the Standard Normal Quantiles<sup>2</sup>.

So if, say,  $x$  is a sample from  $\mathcal{N}(2, 3^2)$ , then

- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(2, 3^2)$ , the Quantiles will be on the bisector;
- ▶ when doing a Q-Q Plot of  $x$  vs  $\mathcal{N}(0, 1)$ , the Quantiles will be on some line (can you find the line equation?);

---

<sup>2</sup>Can you state rigorously and prove this?

## Important Note

So if `qqnorm` shows that the quantiles are close to a line, that means that the Dataset is possibly from a Normal Distribution.

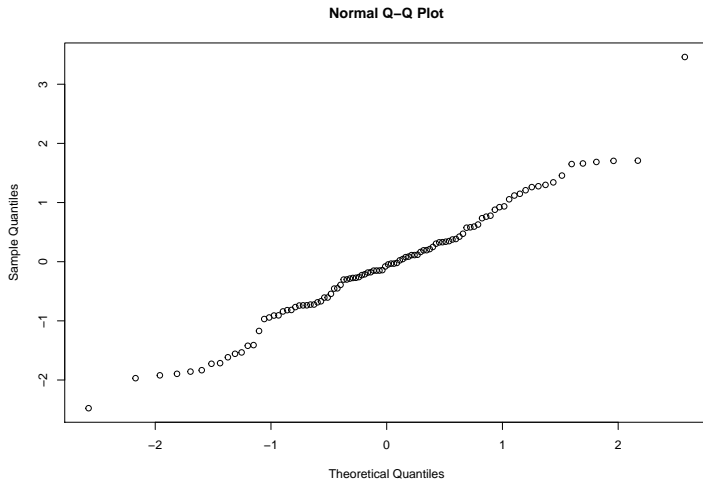
## Important Note

So if `qqnorm` shows that the quantiles are close to a line, that means that the Dataset is possibly from a Normal Distribution.

And if `qqnorm` shows that the quantiles are close to the bisector, that means that the Dataset is possibly from the Standard Normal Distribution.

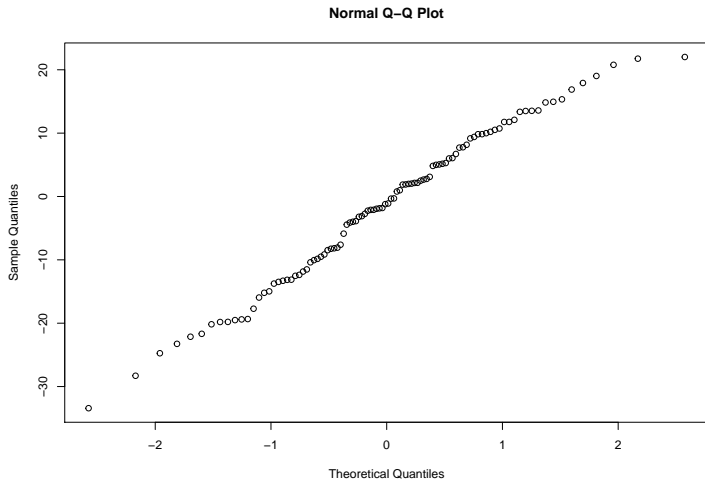
# Some Experiments

```
x <- rnorm(100, mean=0, sd=1)
qqnorm(x)
```



# Some Experiments

```
x <- rnorm(100, mean=2, sd=12)
qqnorm(x)
```



**Exercise:** Express the Quantiles of  $\mathcal{N}(\mu, \sigma^2)$  in terms of the quantiles of  $\mathcal{N}(0, 1)$ .



## Important Note, v2

The above important note works also for the Uniform Distribution. This is again because all Uniform Distributions are the scaled-translated versions of the Standard Uniform  $Unif[0, 1]$ .

## Important Note, v2

The above important note works also for the Uniform Distribution. This is again because all Uniform Distributions are the scaled-translated versions of the Standard Uniform  $Unif[0, 1]$ .

So if you will compare your Dataset with  $Unif[0, 1]$ , and Q-Q Plot will show that the Quantiles are close to a line, that means that probably your Dataset is from a Uniform Distribution, with some parameters.

**Exercise:** Express the Quantiles of  $Unif[a, b]$  in terms of the quantiles of  $Unif[0, 1]$ .