

ASDS Statistics, YSU, Fall 2020

Lecture 11

Michael Poghosyan

08 Oct 2020

Contents

- ▶ Sample and Theoretical Quantiles
- ▶ QQ Plot

Last Lecture Recap

- ▶ Give the definition and an interpretation of the Sample α -Quantile.

Last Lecture Recap

- ▶ Give the definition and an interpretation of the Sample α -Quantile.
- ▶ Give the definition and an interpretation of the Theoretical α -Quantile.

Examples

Example: Find the 70% quantile of the Distribution with the PDF

$$f(x) = \begin{cases} 3x^2, & x \in [0, 1] \\ 0, & \textit{otherwise} \end{cases}$$

Solution: OTB

Theoretical Quantiles, again

Now, if q_α is the α -quantile of some Distribution, and X is a r.v. from that Distribution, then

$$\mathbb{P}(X \leq q_\alpha) \geq \alpha \quad \text{and} \quad \mathbb{P}(X \geq q_\alpha) \geq 1 - \alpha.$$

Theoretical Quantiles, again

Now, if q_α is the α -quantile of some Distribution, and X is a r.v. from that Distribution, then

$$\mathbb{P}(X \leq q_\alpha) \geq \alpha \quad \text{and} \quad \mathbb{P}(X \geq q_\alpha) \geq 1 - \alpha.$$

Note: Here we are taking inequalities, and not, say, $\mathbb{P}(X \leq q_\alpha) = \alpha$, since, in the Discrete r.v. case, we can have no q_α with exact equality. Say, if $X \sim \text{Bernoulli}(0.2)$, and $\alpha = 0.4$, then no q_α exists with $\mathbb{P}(X \leq q_\alpha) = \alpha$.

Theoretical Quantiles, again

Now, if q_α is the α -quantile of some Distribution, and X is a r.v. from that Distribution, then

$$\mathbb{P}(X \leq q_\alpha) \geq \alpha \quad \text{and} \quad \mathbb{P}(X \geq q_\alpha) \geq 1 - \alpha.$$

Note: Here we are taking inequalities, and not, say, $\mathbb{P}(X \leq q_\alpha) = \alpha$, since, in the Discrete r.v. case, we can have no q_α with exact equality. Say, if $X \sim \text{Bernoulli}(0.2)$, and $\alpha = 0.4$, then no q_α exists with $\mathbb{P}(X \leq q_\alpha) = \alpha$.

Note: If $\alpha = 0.5$, we call $q_\alpha = q_{0.5}$ to be the **Median of the Distribution**.

Theoretical Quantiles, again

Now, if q_α is the α -quantile of some Distribution, and X is a r.v. from that Distribution, then

$$\mathbb{P}(X \leq q_\alpha) \geq \alpha \quad \text{and} \quad \mathbb{P}(X \geq q_\alpha) \geq 1 - \alpha.$$

Note: Here we are taking inequalities, and not, say, $\mathbb{P}(X \leq q_\alpha) = \alpha$, since, in the Discrete r.v. case, we can have no q_α with exact equality. Say, if $X \sim \text{Bernoulli}(0.2)$, and $\alpha = 0.4$, then no q_α exists with $\mathbb{P}(X \leq q_\alpha) = \alpha$.

Note: If $\alpha = 0.5$, we call $q_\alpha = q_{0.5}$ to be the **Median of the Distribution**. So if we consider a Continuous r.v. and draw the PDF of that r.v., then the Median is the (leftmost) point dividing the area under the PDF curve into 50%-50% portions.

Theoretical Quantiles, again

Later we will use a lot quantiles. When constructing Confidence Intervals or Hypothesis Testing, we will use Quantiles of the Normal Distribution, t -Distribution, χ^2 -Distribution.

Theoretical Quantiles, again

Later we will use a lot quantiles. When constructing Confidence Intervals or Hypothesis Testing, we will use Quantiles of the Normal Distribution, t -Distribution, χ^2 -Distribution.

Say, later, by z_α we will denote the α -quantile of the Standard Normal Distribution, $\mathcal{N}(0, 1)$.

Theoretical Quantiles, again

Later we will use a lot quantiles. When constructing Confidence Intervals or Hypothesis Testing, we will use Quantiles of the Normal Distribution, t -Distribution, χ^2 -Distribution.

Say, later, by z_α we will denote the α -quantile of the Standard Normal Distribution, $\mathcal{N}(0, 1)$.

Say, we will take $\alpha \in (0, 1)$ and find two points $a, b \in \mathbb{R}$ such that for $X \sim \mathcal{N}(0, 1)$

$$\mathbb{P}(X \leq a) = \mathbb{P}(X \geq b) = \frac{\alpha}{2}.$$

Theoretical Quantiles, again

Later we will use a lot quantiles. When constructing Confidence Intervals or Hypothesis Testing, we will use Quantiles of the Normal Distribution, t -Distribution, χ^2 -Distribution.

Say, later, by z_α we will denote the α -quantile of the Standard Normal Distribution, $\mathcal{N}(0, 1)$.

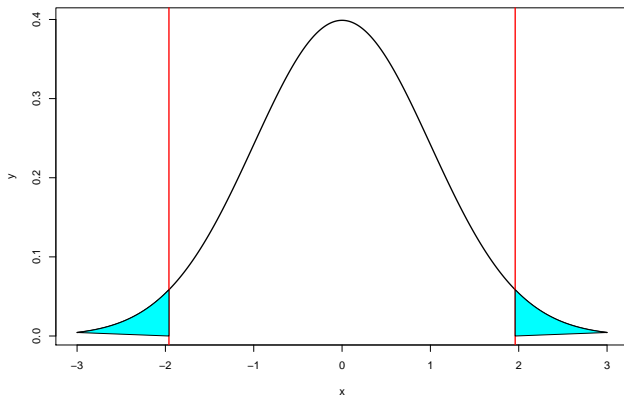
Say, we will take $\alpha \in (0, 1)$ and find two points $a, b \in \mathbb{R}$ such that for $X \sim \mathcal{N}(0, 1)$

$$\mathbb{P}(X \leq a) = \mathbb{P}(X \geq b) = \frac{\alpha}{2}.$$

The idea is to find a symmetric (in fact, the smallest length) interval $[a, b]$ such that for a Standard Normal r.v. X , the chances of $X \notin [a, b]$ are small, are exactly α .

Graphically

```
alpha <- 0.05; z.alpha <- qnorm(alpha/2, mean = 0, sd = 1)
x <- seq(-3,3, by = 0.01)
y <- dnorm(x, mean = 0, sd = 1)
plot(x,y, type = "l", xlim = c(-3,3), lwd = 2)
abline(v = z.alpha, lwd = 2, col = "red")
abline(v = -z.alpha, lwd = 2, col = "red")
polygon(c(x[x<=z.alpha], z.alpha),c(y[x<=z.alpha],0),col="cyan")
polygon(c(x[x>=-z.alpha], -z.alpha),c(y[x>=-z.alpha],0),col="cyan")
```



Theoretical Quantiles, again

Then, it is easy to see, if $\alpha \in (0, 0.5)$ because of the symmetry, that $b = -a$, and

$$a = z_{\alpha/2}.$$

Theoretical Quantiles, again

Then, it is easy to see, if $\alpha \in (0, 0.5)$ because of the symmetry, that $b = -a$, and

$$a = z_{\alpha/2}.$$

So

$$b = -z_{\alpha/2} = z_{1-\alpha/2}$$

Theoretical Quantiles, again

Then, it is easy to see, if $\alpha \in (0, 0.5)$ because of the symmetry, that $b = -a$, and

$$a = z_{\alpha/2}.$$

So

$$b = -z_{\alpha/2} = z_{1-\alpha/2}$$

Note: Please be careful when using Normal Tables. Usually, there is a picture above the table, on which you can find the explanation of the process. Just search “Normal tables” in Google Images.

Q-Q Plots

Next, we consider three important statistical problems: Check visually if

Q-Q Plots

Next, we consider three important statistical problems: Check visually if

- ▶ two given Datasets (possibly, of different sizes) are from the same Distribution;

Q-Q Plots

Next, we consider three important statistical problems: Check visually if

- ▶ two given Datasets (possibly, of different sizes) are from the same Distribution;
- ▶ a given Dataset comes from a given Distribution;

Q-Q Plots

Next, we consider three important statistical problems: Check visually if

- ▶ two given Datasets (possibly, of different sizes) are from the same Distribution;
- ▶ a given Dataset comes from a given Distribution;
- ▶ given two theoretical Distributions, check if one of them is a shifted-scaled version of the other one, or check if one has *fatter tails* than the other one

Q-Q Plots, Data vs Data

Now, assume we have two Datasets, not necessarily of the same size:

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

Q-Q Plots, Data vs Data

Now, assume we have two Datasets, not necessarily of the same size:

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

Question: Are x and y coming from the same Distribution?

Q-Q Plots, Data vs Data

Now, assume we have two Datasets, not necessarily of the same size:

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

Question: Are x and y coming from the same Distribution?

Q-Q Plot helps to answer to this question visually.

Q-Q Plots, Data vs Data

Now, assume we have two Datasets, not necessarily of the same size:

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

Question: Are x and y coming from the same Distribution?

Q-Q Plot helps to answer to this question visually. To draw the Q-Q Plot for Datasets, we take some levels of quantiles, say, for some n ,

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

and then draw the points (q_α^x, q_α^y) .

Q-Q Plots, Data vs Data

Now, assume we have two Datasets, not necessarily of the same size:

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

Question: Are x and y coming from the same Distribution?

Q-Q Plot helps to answer to this question visually. To draw the Q-Q Plot for Datasets, we take some levels of quantiles, say, for some n ,

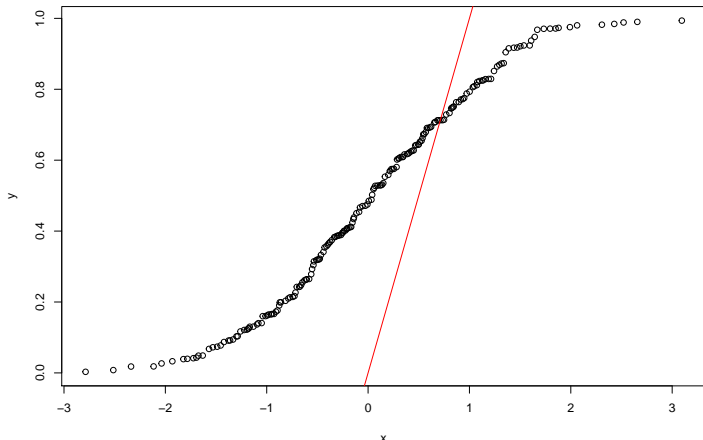
$$\alpha = \frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}$$

and then draw the points (q_α^x, q_α^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_α^x, q_α^y) need to be close to the bisector line.

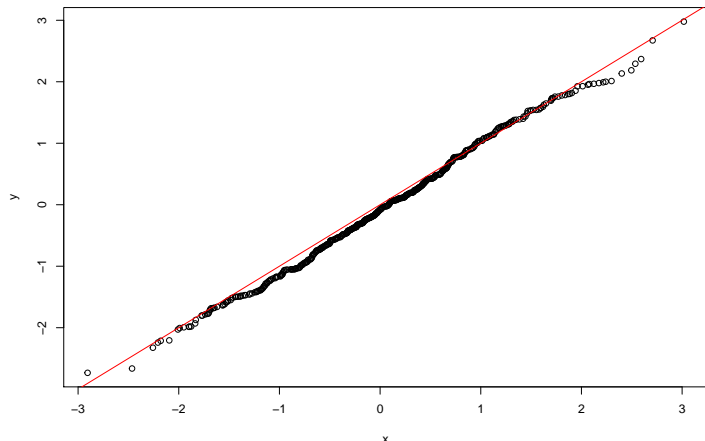
Example, Q-Q Plots, Data vs Data

```
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] -0.063  0.224  0.441  0.300  0.357  0.301  0.182 -0.063
## [11]  0.974 -0.322  0.288 -0.946  0.425  0.854 -0.820 -0.063
```

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] -0.063  0.224  0.441  0.300  0.357  0.301  0.182 -0.063
## [11]  0.974 -0.322  0.288 -0.946  0.425  0.854 -0.820 -0.946
```

from a Normal Distribution?

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

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

and then draw the points $(q_{\alpha}^F, q_{\alpha}^x)$, where q_{α}^F is the α -quantile of the Theoretical Distribution, and q_{α}^x is the α -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] -0.063  0.224  0.441  0.300  0.357  0.301  0.182 -0.063
## [11]  0.974 -0.322  0.288 -0.946  0.425  0.854 -0.820 -0.063
```

from a Normal Distribution?

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

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

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

¹or 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 Distribution¹.

¹or 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 Distribution¹.

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

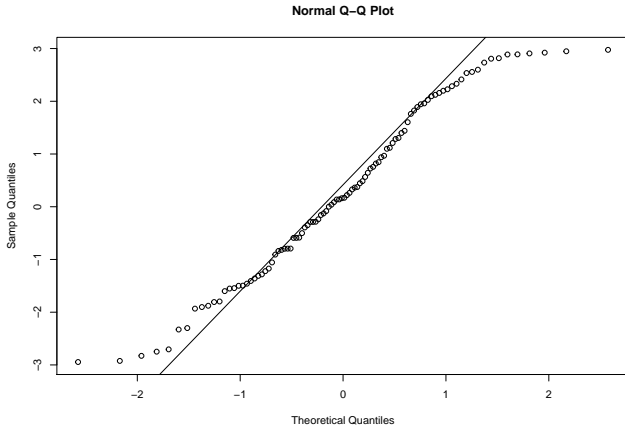
$$(q_{0.25}^F, q_{0.25}^x) \quad \text{and} \quad (q_{0.75}^F, q_{0.75}^x).$$

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

Some Experiments

Here are some experiments with `qqnorm`

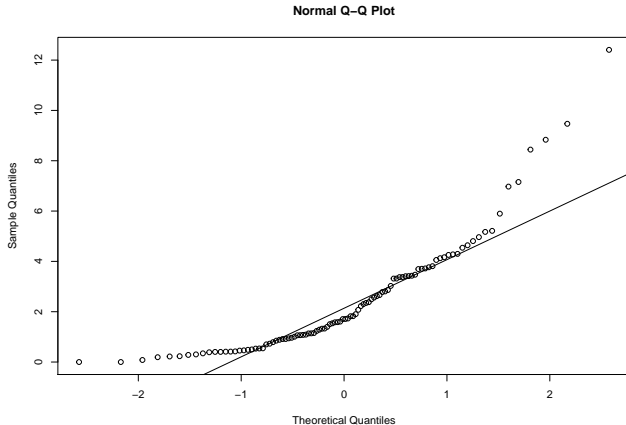
```
x <- runif(100, -3, 3)
qqnorm(x)
qqline(x)
```



Some Experiments

Here are some experiments with `qqnorm`

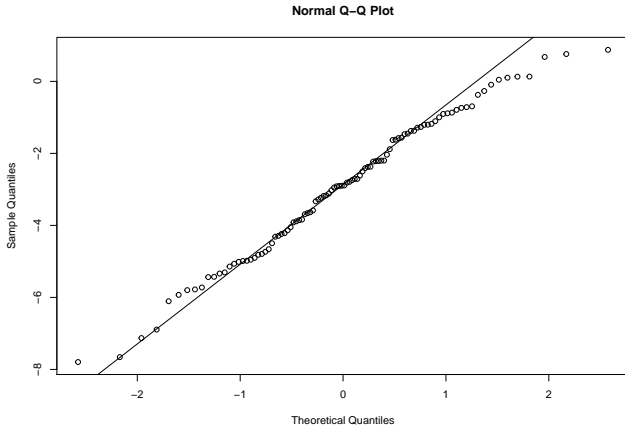
```
x <- rexp(100,0.4)
qqnorm(x)
qqline(x)
```



Some Experiments

Here are some experiments with `qqnorm`

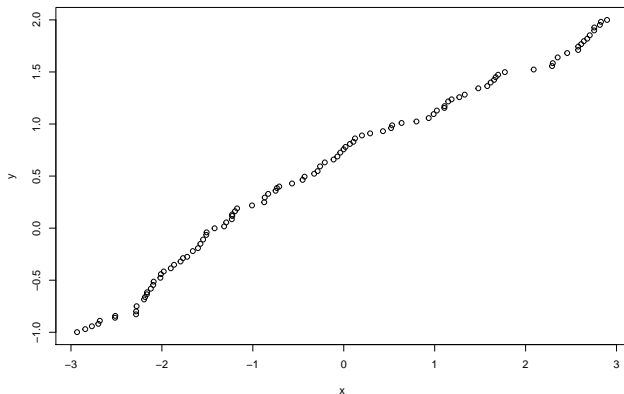
```
x <- rnorm(100, mean = -3, sd = 2)
qqnorm(x)
qqline(x)
```



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.

²Can 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 Quantiles².

²Can 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 Quantiles².

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

²Can 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 Quantiles².

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;

²Can 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 Quantiles².

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

²Can 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 Quantiles².

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?);

²Can 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.

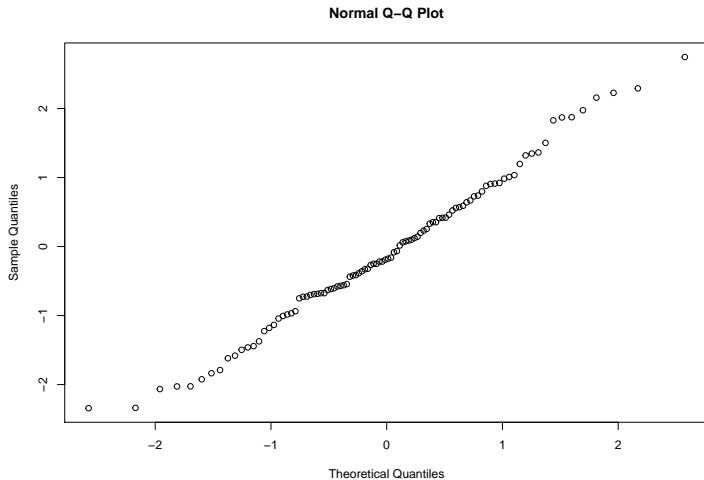
Important Note

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.

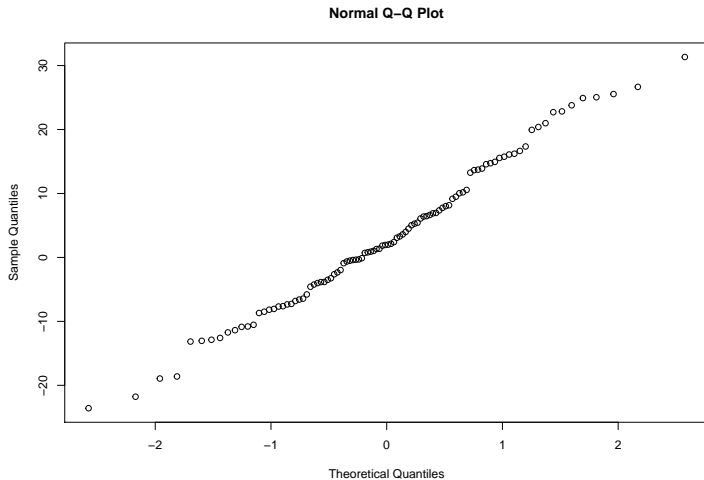
Some Experiments

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



Some Experiments

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



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.