# Normal Distribution

Andrey Chinnov, Sebastian Honermann, Carlos Zydorek

Case Studies
"Data Analytics"

## Outline

**1** Introduction
  - Normality as a requirement for statistical methods
  - Data Set Overview

**2** Normality Testing
  - Graphical Methods for Normality Testing
    - Q-Q-Plots
    - Chi-Square Plot
  - Quantitative Methods for Normality Testing
    - Shapiro-Wilk Test
    - Pearson's Chi-Squared Test
    - Kolmogorov-Smirnov Test

**3** Transformation to Normality
  - Box-Cox Transformation
  - Transformation Results Testing

**4** Summary

# Normality as a requirement for statistical methods

## Data Set Overview

# Outline

## Q-Q-Plots

**Sample:**

$$x = (x_1, x_2, \ldots, x_n)$$

**Empirical quantiles:**

$$x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)}$$
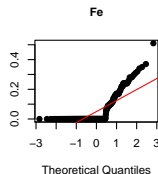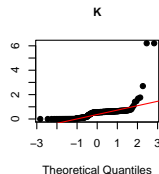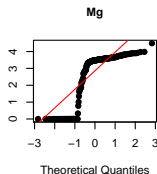
**Theoretical quantiles:**

$$q_{(j)} = \Phi^{-1}(p_{(j)}),$$

where

$$p_{(j)} = \frac{j - \frac{1}{2}}{n},$$

$\Phi$ - $N(0,1)$ c.d.f. .

$\implies$ Plot $x_{(i)}$ against $q_{(i)}$

# Q-Q-Plots

**QQ-Plots of the subdatasets:**

- Variables could be normally distributed within the subclasses

- For some cases there appear to be a linear relationships

- For other cases a linear relationship is questionable

- In some subdatasets a linear relationship seems plausible, however $n$ is very small

# Q-Q-Plots

**Results of the Transformation of the Full Dataset :**

- For some of the cases there seems to be a slight improvement

- For non-unimodal cases the transformation does not show significant improvements towards normality



Na Glass Type 1



Na Glass Type 1 transformed



Ca Glass Type 7



Ca Glass Type 7 transformed

## Q-Q-Plots

**Results of the Transformation of the Subdatasets :**

- For unimodal cases the transformation shapes the distribution closer to normality

- For non-unimodal cases the transformation does not show significant improvements towards normality



RI Glass Type 1



RI Glass Type 1 transformed



Si Glass Type 1



Si Glass Type 1 transformed

## Shapiro-Wilk Test

The test statistic $W$ indicates the deviation of the observed quantile values from the assumed cumulative distribution function quantiles

$$W = \frac{\sum\limits_{i=1}^{n} (a_i y_i)^2}{\sum\limits_{i=1}^{n} (y_i - \bar{y})^2},$$

where

- $a_i$ denotes the normalised "best linear unbiased" coefficients,
- $y_i$ denotes the observations.

The critical value for $W$ is obtained by the Monte Carlo Method
$\implies$ $p$-value is calculated

Important: If a variable contains only zeros the Shapiro-Wilk test is not applicable, since the term in the denominator sums up to zero.

## Shapiro-Wilk Test

**Testing the Full Dataset :**

Null hypothesis is rejected for all variables at a 1 % significance level

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | 0.87 | 0.01 | NA | 1.0766713449726e-12 | yes |
| Na | 0.95 | 0.01 | NA | 3.4655430546966e-07 | yes |
| Mg | 0.7 | 0.01 | NA | < 1.0e-15 | yes |
| Al | 0.94 | 0.01 | NA | 2.08315629600399e-07 | yes |
| Si | 0.92 | 0.01 | NA | 2.17503176825416e-09 | yes |
| K | 0.44 | 0.01 | NA | < 1.0e-15 | yes |
| Ca | 0.79 | 0.01 | NA | < 1.0e-15 | yes |
| Ba | 0.41 | 0.01 | NA | < 1.0e-15 | yes |
| Fe | 0.65 | 0.01 | NA | < 1.0e-15 | yes |

Test results of the Shapiro-Wilk test on the whole data sample

**After the Transformation :**

The null hypothesis can be rejected for the four transformed variables

⟹ Possible Explanation:

Combination of different distributions in the different glass types

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | NA | NA | NA | NA | NA |
| Na | 0.95 | 0.01 | NA | 8.75605777309153e-07 | yes |
| Mg | NA | NA | NA | NA | NA |
| Al | 0.97 | 0.01 | NA | 0.000244326513056066 | yes |
| Si | 0.93 | 0.01 | NA | 1.58998125691823e-08 | yes |
| K | NA | NA | NA | NA | NA |
| Ca | 0.89 | 0.01 | NA | 1.13880689831982e-11 | yes |
| Ba | NA | NA | NA | NA | NA |
| Fe | NA | NA | NA | NA | NA |

Test results of the Shapiro-Wilk test on the whole transformed data sample

# Shapiro-Wilk Test

**Testing the Full Dataset :**

Null hypothesis is rejected for all variables at a 1 % significance level

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | 0.87 | 0.01 | NA | 1.0766713449726e-12 | yes |
| Na | 0.95 | 0.01 | NA | 3.4655430546966e-07 | yes |
| Mg | 0.7 | 0.01 | NA | $< 1.0e-15$ | yes |
| Al | 0.94 | 0.01 | NA | 2.08315629600399e-07 | yes |
| Si | 0.92 | 0.01 | NA | 2.17503176825416e-09 | yes |
| K | 0.44 | 0.01 | NA | $< 1.0e-15$ | yes |
| Ca | 0.79 | 0.01 | NA | $< 1.0e-15$ | yes |
| Ba | 0.41 | 0.01 | NA | $< 1.0e-15$ | yes |
| Fe | 0.65 | 0.01 | NA | $< 1.0e-15$ | yes |

Test results of the Shapiro-Wilk test on the whole data sample

**After the Transformation :**

The null hypothesis can be rejected for the four transformed variables
$\implies$ Possible Explanation:
Combination of different distributions in the different glass types

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | NA | NA | NA | NA | NA |
| Na | 0.95 | 0.01 | NA | 8.75605777309153e-07 | yes |
| Mg | NA | NA | NA | NA | NA |
| Al | 0.97 | 0.01 | NA | 0.000244326513056066 | yes |
| Si | 0.93 | 0.01 | NA | 1.58998125691823e-08 | yes |
| K | NA | NA | NA | NA | NA |
| Ca | 0.89 | 0.01 | NA | 1.13880689831982e-11 | yes |
| Ba | NA | NA | NA | NA | NA |
| Fe | NA | NA | NA | NA | NA |

Test results of the Shapiro-Wilk test on the whole transformed data sample

# Graphical Methods for Normality Testing
## Chi-Square Plot

## Quantitative Methods for Normality Testing
Pearson's Chi-Squared Test

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.

### Definition

$F_n(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{x_i \leq x\}}(x)$
- empirical c. d. f. , where
$\mathbb{1}_{\{x_i \leq x\}}(x) = \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{otherwise.} \end{cases}$



Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.

### Definition

$F_n(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\{x_i \leq x\}}(x)$
- empirical c. d. f. , where
$\mathbb{1}_{\{x_i \leq x\}}(x) = \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{otherwise.} \end{cases}$

$F(x)$ - theoretical normal c. d. f.
with

$$\bar{x} = \frac{1}{n} \sum_i x_i, \quad \sigma_x^2 = \frac{1}{n}(x_i - \bar{x})^2$$



Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.

## Definition

$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x_i \leq x\}}(x)$
- empirical c. d. f. , where
$\mathbb{1}_{\{x_i \leq x\}}(x) = \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{otherwise.} \end{cases}$

$F(x)$ - theoretical normal c. d. f.
with

$$\bar{x} = \frac{1}{n} \sum_i x_i, \quad \sigma_x^2 = \frac{1}{n}(x_i - \bar{x})^2$$

$d = \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|$
- distance between them.



Glass Type 1, Natrium (Na)

## Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

$$
\begin{aligned}
H_0 &: \quad \mathbb{P} = \mathbb{P}_0, \\
H_1 &: \quad \mathbb{P} \neq \mathbb{P}_0.
\end{aligned}
$$

## Quantitative Methods for Normality Testing
### Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

$$H_0 \; : \; \mathbb{P} = \mathbb{P}_0,$$
$$H_1 \; : \; \mathbb{P} \neq \mathbb{P}_0.$$

KS test statistics:

$$D_n = \sqrt{n} \cdot \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|.$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

$$H_0 \quad : \quad \mathbb{P} = \mathbb{P}_0,$$
$$H_1 \quad : \quad \mathbb{P} \neq \mathbb{P}_0.$$

KS test statistics:

$$D_n = \sqrt{n} \cdot \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|.$$

Properties of $D_n$ in case $H_0$ is TRUE:

- Distribution of $\hat{D}_n := (D_1, D_2, \ldots, D_n)$ does not depend on $F$

## Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

$$H_0 \quad : \quad \mathbb{P} = \mathbb{P}_0,$$
$$H_1 \quad : \quad \mathbb{P} \neq \mathbb{P}_0.$$

KS test statistics:

$$D_n = \sqrt{n} \cdot \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|.$$

Properties of $D_n$ in case $H_0$ is TRUE:

- Distribution of $\hat{D}_n := (D_1, D_2, \ldots, D_n)$ does not depend on $F$

$$\implies \text{tabulated}$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

Let $x = (x_1, x_2, \ldots, x_n)$ be a sample of unknown distribution $\mathbb{P}$.
Theoretical c. d. f. $F$ defines a distribution $\mathbb{P}_0$.

$$H_0 \quad : \quad \mathbb{P} = \mathbb{P}_0,$$
$$H_1 \quad : \quad \mathbb{P} \neq \mathbb{P}_0.$$

KS test statistics:

$$D_n = \sqrt{n} \cdot \sup_{x \in \mathbb{R}} |F_n(x) - F(x)|.$$

Properties of $D_n$ in case $H_0$ is TRUE:

- Distribution of $\hat{D}_n := (D_1, D_2, \ldots, D_n)$ does not depend on $F$

$$\implies \text{tabulated}$$

- $\forall t > 0 :$

$$P(D_n \leq t) \xrightarrow[n \to \infty]{} H(t) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} \exp^{-2i^2 t^2}$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value

## Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that
depends on a significance level $\alpha$:

## Quantitative Methods for Normality Testing
### Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{ccc} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that
depends on a significance level $\alpha$:

$\alpha = P(\delta \neq H_0 | H_0)$

# Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that
depends on a significance level $\alpha$:

$\alpha = P(\delta \neq H_0 | H_0) = P(D_n > c | H_0)$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that

depends on a significance level $\alpha$:

$$\alpha = P(\delta \neq H_0 | H_0) = P(D_n > c | H_0) = 1 - P(D_n \leq c | H_0)$$

## Quantitative Methods for Normality Testing
Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that
depends on a significance level $\alpha$:

$$\alpha = P(\delta \neq H_0 | H_0) = P(D_n > c | H_0) = 1 - P(D_n \leq c | H_0) \approx 1 - H(c).$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq c \\ H_1 & : & D_n > c \end{array} \right. ,$$

where $c$ - critical value that
depends on a significance level $\alpha$:

$$\alpha = P(\delta \neq H_0 | H_0) = P(D_n > c | H_0) = 1 - P(D_n \leq c | H_0) \approx 1 - H(c).$$

$$\implies c \approx H_{1-\alpha}$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty}(-1)^{i-1}\exp^{-2i^2t^2}$$

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty}(-1)^{i-1}\exp^{-2i^2 t^2}$$
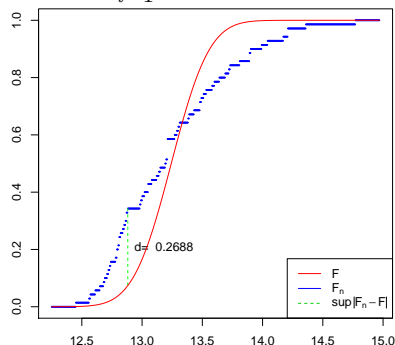
**Example:**



Glass Type 1, Natrium (Na)
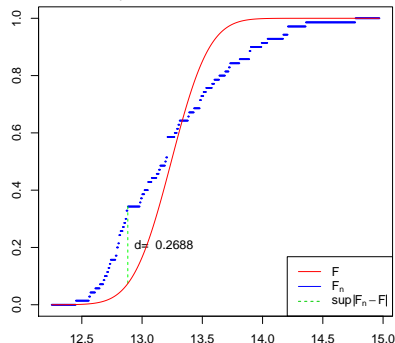
# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2 \sum_{i=1}^{\infty} (-1)^{i-1} \exp^{-2i^2 t^2}$$

**Example:**

- $n = 70$



Glass Type 1, Natrium (Na)
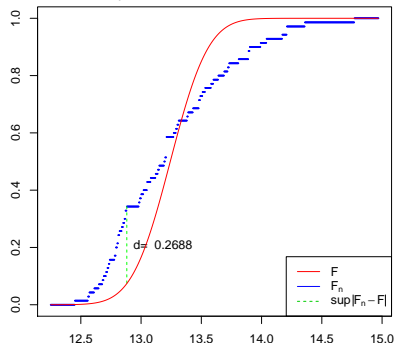
# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty} (-1)^{i-1} \exp^{-2i^2 t^2}$$

**Example:**

- $n = 70$
- $D_n = \sqrt{n} \sup |F_n - F| = 2.2493$



d= 0.2688

| | F |
| | $F_n$ |
| | $\sup |F_n - F|$ |

Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lcl} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty}(-1)^{i-1}\exp^{-2i^2t^2}$$

**Example:**

- $n = 70$
- $D_n = \sqrt{n}\sup|F_n - F| = 2.2493$
- $\alpha = 0.01$
  $$\implies c = H_{1-\alpha} = 1.6276$$



Glass Type 1, Natrium (Na)

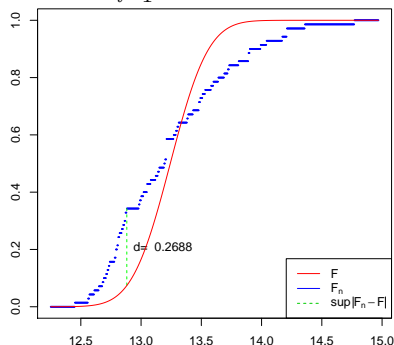# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty}(-1)^{i-1}\exp^{-2i^2t^2}$$

**Example:**

- $n = 70$
- $D_n = \sqrt{n}\sup|F_n - F| = 2.2493$
- $\alpha = 0.01$
  $$\implies c = H_{1-\alpha} = 1.6276$$
- $D_n > c \implies H_0$ rejected



Glass Type 1, Natrium (Na)
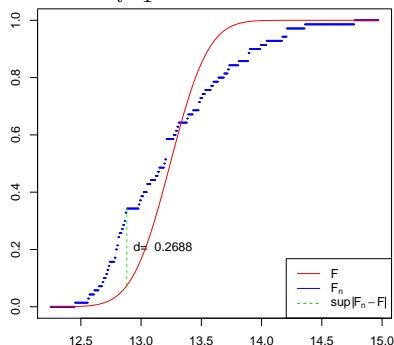
# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \left\{ \begin{array}{lll} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{array} \right. , \quad H(t) = 1 - 2\sum_{i=1}^{\infty} (-1)^{i-1} \exp^{-2i^2 t^2}$$

**Example:**

- $n = 70$
- $D_n = \sqrt{n} \sup |F_n - F| = 2.2493$
- $\alpha = 0.01$
    $$\implies c = H_{1-\alpha} = 1.6276$$
- $D_n > c \implies H_0$ rejected
- $\implies \mathbb{P} \neq \mathbb{P}_0$



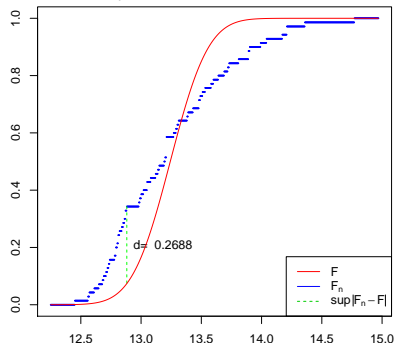Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
## Kolmogorov-Smirnov Test

The KS test uses the decision rule for a given significance level $\alpha$

$$\delta = \begin{cases} H_0 & : & D_n \leq H_{1-\alpha} \\ H_1 & : & D_n > H_{1-\alpha} \end{cases}, \quad H(t) = 1 - 2\sum_{i=1}^{\infty}(-1)^{i-1}\exp^{-2i^2t^2}$$

**Example:**

- $n = 70$
- $D_n = \sqrt{n}\sup|F_n - F| = 2.2493$
- $\alpha = 0.01$
  $$\implies c = H_{1-\alpha} = 1.6276$$
- $D_n > c \implies H_0$ rejected
- $\implies \mathbb{P} \neq \mathbb{P}_0$
- $\not\implies$ data not normally distributed!!!



d = 0.2688

Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

$R$ code used:

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu,\sigma) = \sup_{x\in\mathbb{R}} |F_n(x) - F(x,\mu,\sigma)| \to \min.$$

$R$ code used:

```
c(mean(dat),var(dat))
```

[1] 13.2422857   0.2493019

```
#optim is a predifined R function in stats package
#defalut method of optimization is Nelder and Mead
result = optim(c(mean(dat),var(dat)),KS)
result$par
```

[1] 13.1769501   0.4682486
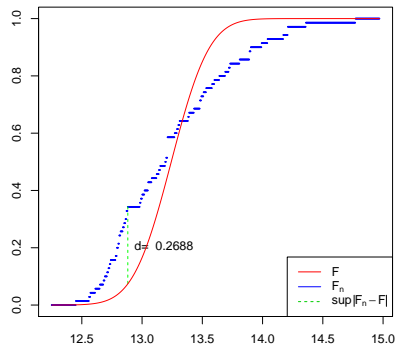
```
result$value
```

[1] 0.07870673

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$
- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$

## Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$

- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$
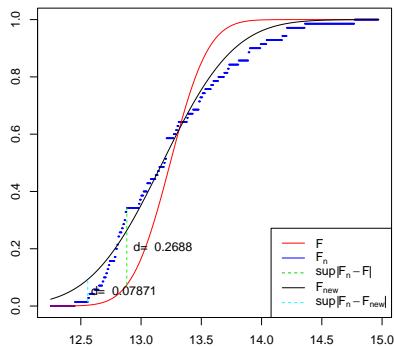


Glass Type 1, Natrium (Na)

## Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$

- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$

- $D_n = \sqrt{n} \sup |F_n - F_{new}| = 0.6585$
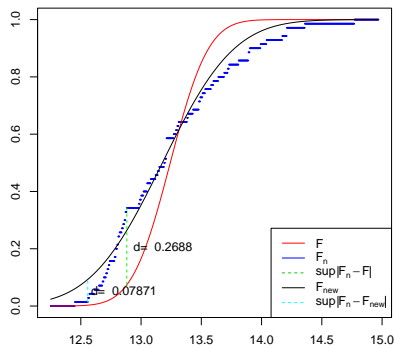


Glass Type 1, Natrium (Na)

## Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$

- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$

- $D_n = \sqrt{n} \sup |F_n - F_{new}| = 0.6585$
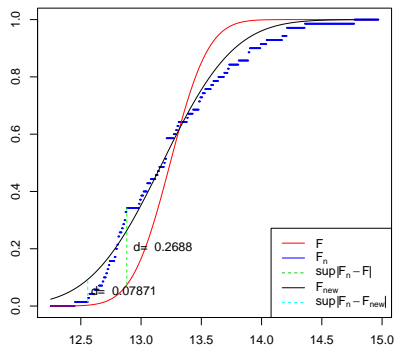
- $c = 1.6276$



Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$

- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$

- $D_n = \sqrt{n} \sup |F_n - F_{new}| = 0.6585$

- $c = 1.6276$

- $D_n < c \implies H_0$ accepted



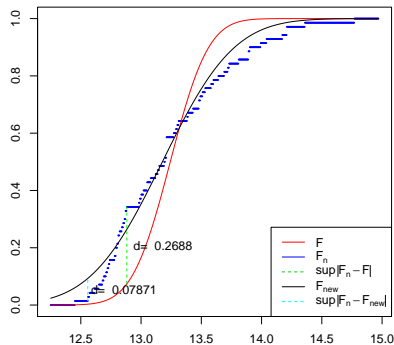Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$
- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$
- $D_n = \sqrt{n} \sup |F_n - F_{new}| = 0.6585$
- $c = 1.6276$
- $D_n < c \implies H_0$ accepted
- $\implies \mathbb{P} = \mathbb{P}_0$
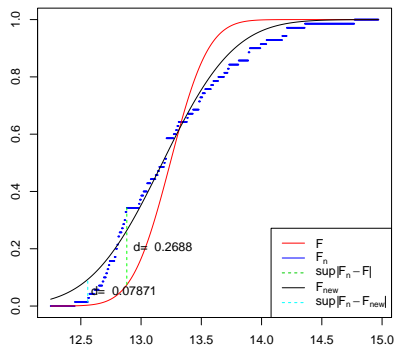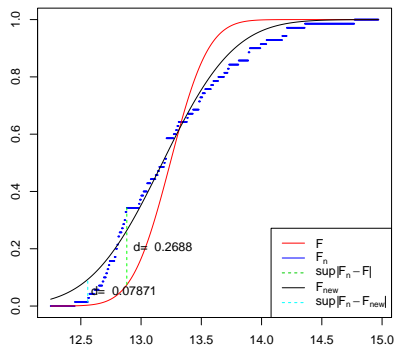


Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test

KS test is improved by solving the following optimization problem

$$KS(\mu, \sigma) = \sup_{x \in \mathbb{R}} |F_n(x) - F(x, \mu, \sigma)| \to \min.$$

- Initial vector of parameters
  $\mu = 13.2423, \quad \sigma^2 = 0.2493$
- Optimized vector of parameters
  $\hat{\mu} = 13.1770, \quad \hat{\sigma}^2 = 0.4682$
- $D_n = \sqrt{n} \sup |F_n - F_{new}| = 0.6585$
- $c = 1.6276$
- $D_n < c \implies H_0$ accepted
- $\implies \mathbb{P} = \mathbb{P}_0$
- $\implies$ data normally distributed!



Glass Type 1, Natrium (Na)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test Results

Results of Improved KS test on the whole data set:

| variable | test statistic | sig. level | critical value | p-value | rejected |
|----------|----------------|------------|----------------|---------|----------|
| RI | 1.34 | 0.01 | 1.63 | 0.0561963016778131 | no |
| Na | 0.87 | 0.01 | 1.63 | 0.43825271603342 | no |
| Mg | 2.94 | 0.01 | 1.63 | 6.18457917100912e-08 | yes |
| Al | 0.84 | 0.01 | 1.63 | 0.474757887353829 | no |
| Si | 0.96 | 0.01 | 1.63 | 0.314710019077325 | no |
| K | 2.14 | 0.01 | 1.63 | 0.000212776619708754 | yes |
| Ca | 1.33 | 0.01 | 1.63 | 0.057710602872685 | no |
| Ba | 2.60 | 0.01 | 1.63 | 2.75476085742632e-06 | yes |
| Fe | 4.68 | 0.01 | 1.63 | $< 1.0e-15$ | yes |

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test Results

Results of Improved KS test on the whole data set:

| variable | test statistic | sig. level | critical value | p-value | rejected |
|----------|----------------|------------|----------------|---------|----------|
| RI | 1.34 | 0.01 | 1.63 | 0.0561963016778131 | no |
| Na | 0.87 | 0.01 | 1.63 | 0.43825271603342 | no |
| Mg | 2.94 | 0.01 | 1.63 | 6.18457917100912e-08 | yes |
| Al | 0.84 | 0.01 | 1.63 | 0.474757887353829 | no |
| Si | 0.96 | 0.01 | 1.63 | 0.314710019077325 | no |
| K | 2.14 | 0.01 | 1.63 | 0.000212776619708754 | yes |
| Ca | 1.33 | 0.01 | 1.63 | 0.057710602872685 | no |
| Ba | 2.60 | 0.01 | 1.63 | 2.75476085742632e-06 | yes |
| Fe | 4.68 | 0.01 | 1.63 | $< 1.0e-15$ | yes |

- 5 variables are normaly distributed (RI,Na,Al,Si,Ca)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test Results

Results of Improved KS test on the whole data set:

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | 1.34 | 0.01 | 1.63 | 0.0561963016778131 | no |
| Na | 0.87 | 0.01 | 1.63 | 0.43825271603342 | no |
| Mg | 2.94 | 0.01 | 1.63 | 6.18457917100912e-08 | yes |
| Al | 0.84 | 0.01 | 1.63 | 0.474757887353829 | no |
| Si | 0.96 | 0.01 | 1.63 | 0.314710019077325 | no |
| K | 2.14 | 0.01 | 1.63 | 0.000212776619708754 | yes |
| Ca | 1.33 | 0.01 | 1.63 | 0.057710602872685 | no |
| Ba | 2.60 | 0.01 | 1.63 | 2.75476085742632e-06 | yes |
| Fe | 4.68 | 0.01 | 1.63 | $< 1.0e-15$ | yes |

- 5 variables are normaly distributed (RI,Na,Al,Si,Ca)
- 4 variables are not (Mg,K,Ba,Fe)

# Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test Results

Results of Improved KS test on the whole data set:

| variable | test statistic | sig. level | critical value | p-value | rejected |
|---|---|---|---|---|---|
| RI | 1.34 | 0.01 | 1.63 | 0.0561963016778131 | no |
| Na | 0.87 | 0.01 | 1.63 | 0.43825271603342 | no |
| Mg | 2.94 | 0.01 | 1.63 | 6.18457917100912e-08 | yes |
| Al | 0.84 | 0.01 | 1.63 | 0.474757887353829 | no |
| Si | 0.96 | 0.01 | 1.63 | 0.314710019077325 | no |
| K | 2.14 | 0.01 | 1.63 | 0.000212776619708754 | yes |
| Ca | 1.33 | 0.01 | 1.63 | 0.057710602872685 | no |
| Ba | 2.60 | 0.01 | 1.63 | 2.75476085742632e-06 | yes |
| Fe | 4.68 | 0.01 | 1.63 | $< 1.0e{-}15$ | yes |

- 5 variables are normaly distributed (RI,Na,Al,Si,Ca)
- 4 variables are not (Mg,K,Ba,Fe)
- The best statistics test value for Al

## Quantitative Methods for Normality Testing
Improved Kolmogorov-Smirnov Test Results

Results of Improved KS test on the whole data set:

| variable | test statistic | sig. level | critical value | p-value | rejected |
|----------|----------------|------------|----------------|---------|----------|
| RI | 1.34 | 0.01 | 1.63 | 0.0561963016778131 | no |
| Na | 0.87 | 0.01 | 1.63 | 0.43825271603342 | no |
| Mg | 2.94 | 0.01 | 1.63 | 6.18457917100912e-08 | yes |
| Al | 0.84 | 0.01 | 1.63 | 0.474757887353829 | no |
| Si | 0.96 | 0.01 | 1.63 | 0.314710019077325 | no |
| K | 2.14 | 0.01 | 1.63 | 0.000212776619708754 | yes |
| Ca | 1.33 | 0.01 | 1.63 | 0.057710602872685 | no |
| Ba | 2.60 | 0.01 | 1.63 | 2.75476085742632e-06 | yes |
| Fe | 4.68 | 0.01 | 1.63 | $< 1.0e-15$ | yes |

- 5 variables are normaly distributed (RI,Na,Al,Si,Ca)
- 4 variables are not (Mg,K,Ba,Fe)
- The best statistics test value for Al
- The worst statistic test value for Fe

## Quantitative Methods for Normality Testing

**Test Results:**

| variable | rejected |
|:--------:|:--------:|
| RI | no |
| Na | no |
| Mg | yes |
| Al | no |
| Si | no |
| K | yes |
| Ca | no |
| Ba | yes |
| Fe | yes |

# Quantitative Methods for Normality Testing

**Test Results:**

| variable | rejected |
|----------|----------|
| RI | no |
| Na | no |
| Mg | yes |
| Al | no |
| Si | no |
| K | yes |
| Ca | no |
| Ba | yes |
| Fe | yes |

# Outline

**1** Introduction
- Normality as a requirement for statistical methods
- Data Set Overview

**2** Normality Testing
- Graphical Methods for Normality Testing
  - Q-Q-Plots
  - Chi-Square Plot
- Quantitative Methods for Normality Testing
  - Shapiro-Wilk Test
  - Pearson's Chi-Squared Test
  - Kolmogorov-Smirnov Test

**3** Transformation to Normality
- Box-Cox Transformation
- Transformation Results Testing

**4** Summary

# Box-Cox Transformation

Transformation Results Testing

## Outline

# Summary