# Introduction to Healthcare Supply Chain (HCSC) Analytics

# Industry 4.0



Data every where…
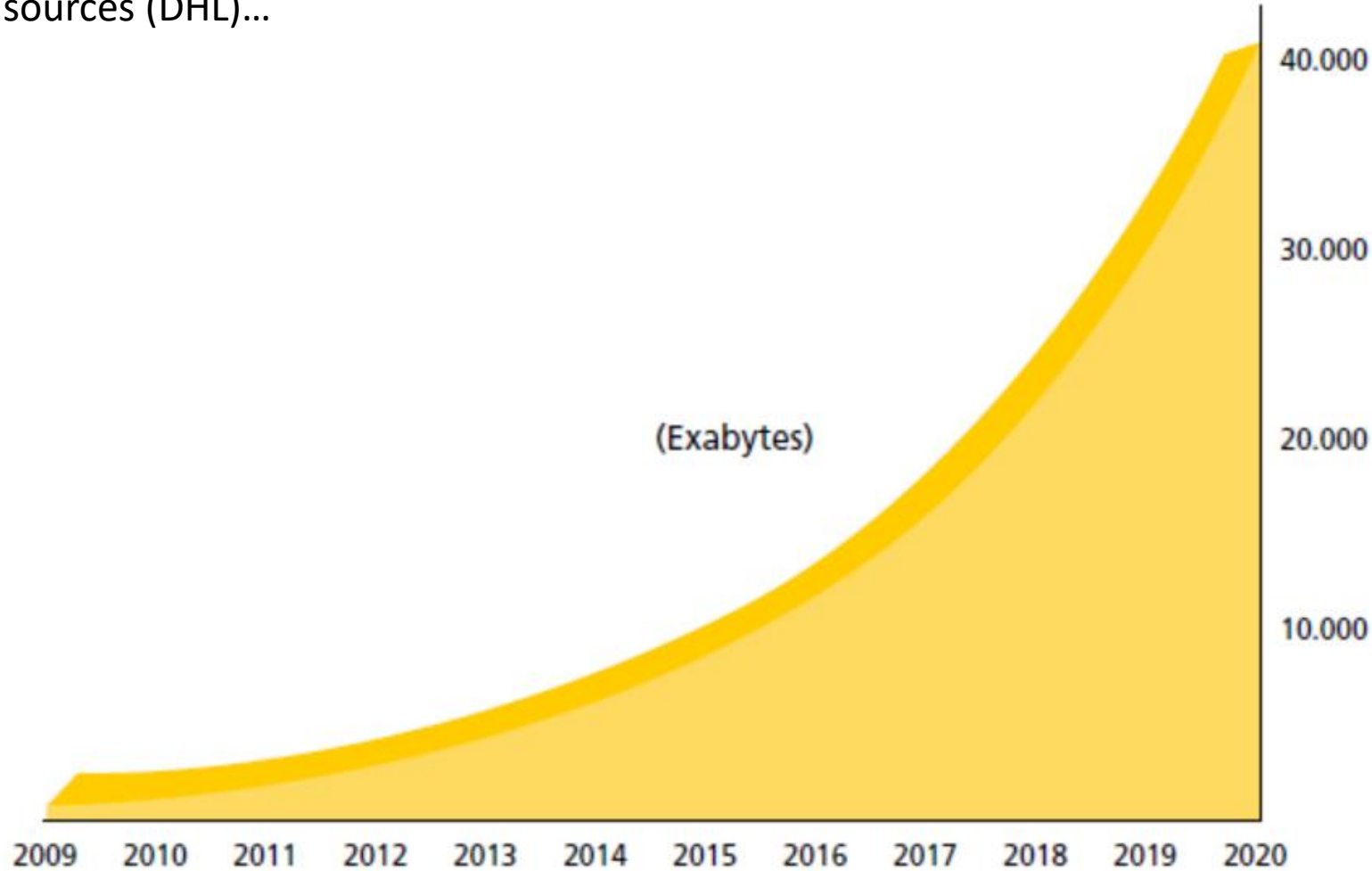
Gigamon Blog

The massive deployment of connected devices such as cars, smartphones, RFID readers, webcams, and sensor networks adds a huge number of autonomous data sources (DHL)...



(Exabytes)

40.000
30.000
20.000
10.000

2009  2010  2011  2012  2013  2014  2015  2016  2017  2018  2019  2020

Exponential data growth between 2010 and 2020; **Source:** IDC's Digital Universe Study, sponsored by EMC, December 2012

Gigamon Blog

# Data processing…(e.g. Artificial Intelligence)

**dataset - experiences**

**algorithms**

**human instructions**

**needed**

**computer…**

**capable to:**

**learning**

**decision making**

## Process efficiency

HR  FR  Time  MR

One purpose alone

# Data processing…(e.g. Artificial Intelligence)

dataset - experiences

**Is this sufficient to make decisions?**

algorithms

needed

human instructions

**computer…**

capable to:

**learning**

**decision making**

## Process efficiency

One purpose alone

HR    FR    Time    MR

# The great confusion…

KNOWLEDGE

INFORMATION
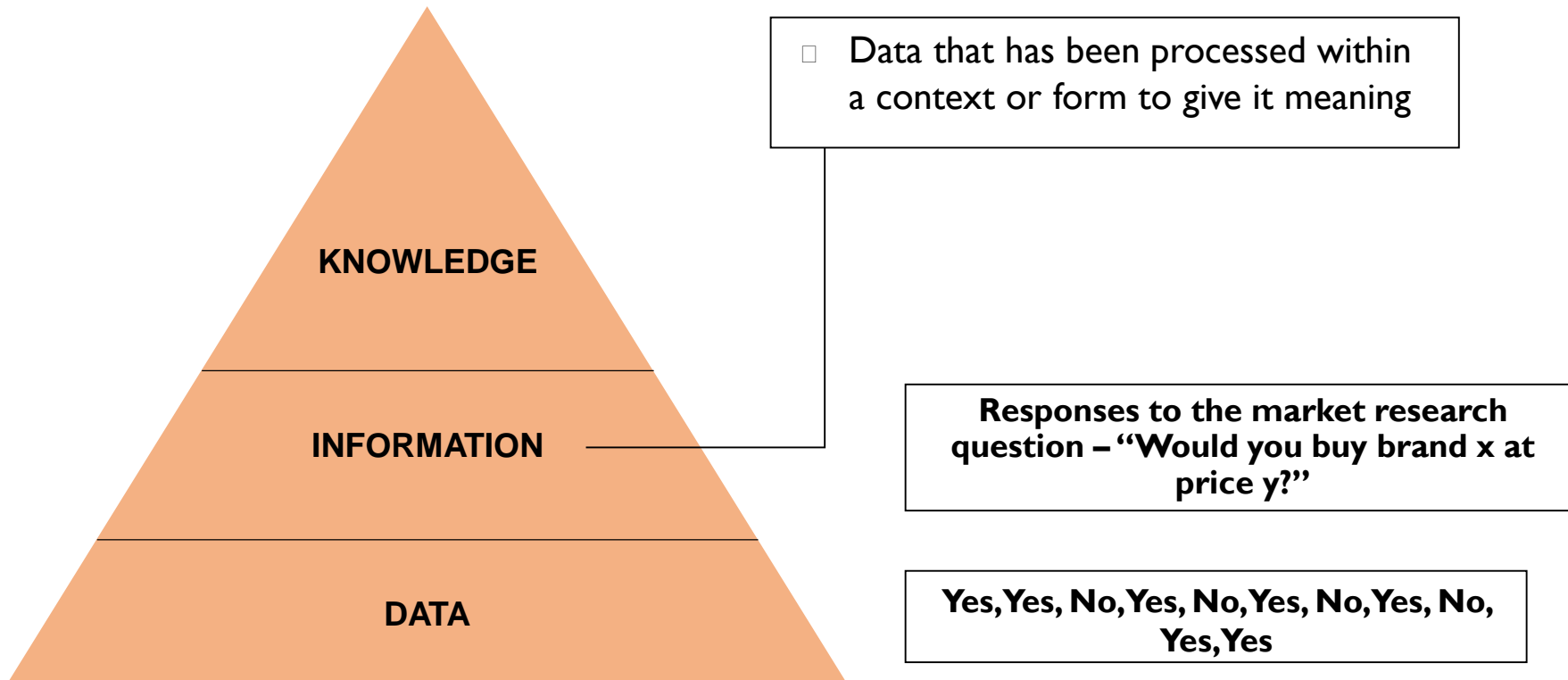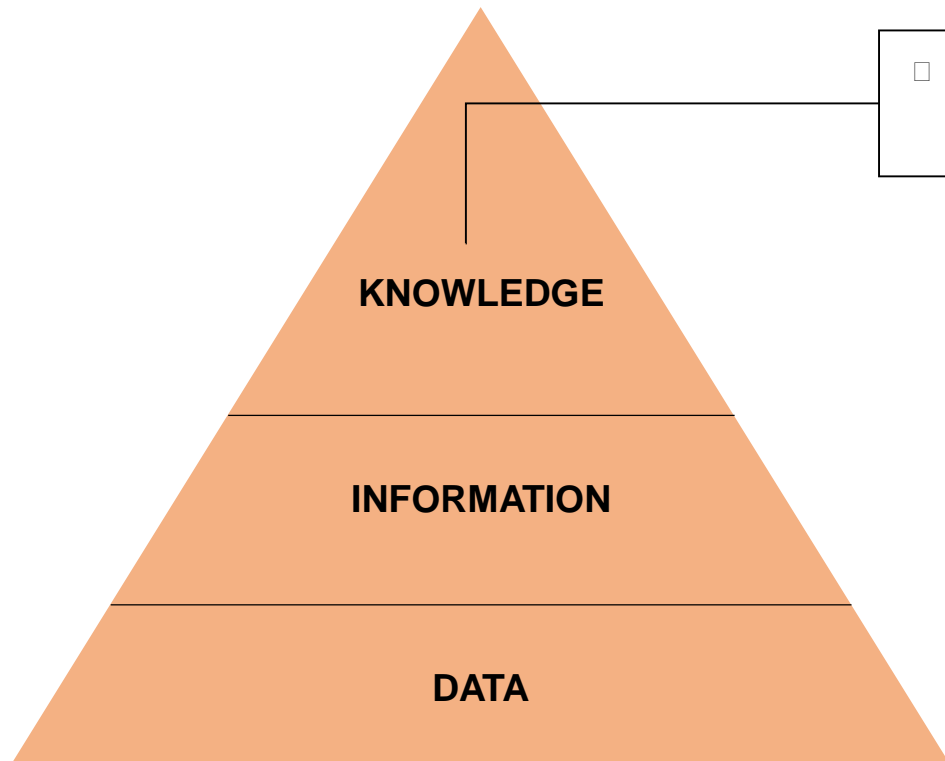
DATA

# The great confusion…

- Data **are** raw facts and figures that on their own have no meaning

- These can be any alphanumeric characters i.e. text, numbers, symbols

KNOWLEDGE

INFORMATION

DATA

HELP

# The great confusion…

KNOWLEDGE

INFORMATION

DATA

Data that has been processed within a context or form to give it meaning

**Responses to the market research question – "Would you buy brand x at price y?"**

**Yes, Yes, No, Yes, No, Yes, No, Yes, No, Yes, Yes**

HELP

# The great confusion…

KNOWLEDGE

INFORMATION

DATA

☐ Knowledge is the understanding of rules needed to interpret information

Ian H. Witten • Eibe Frank • Mark A. Hall

**DATA MINING**

Practical Machine Learning Tools and Techniques
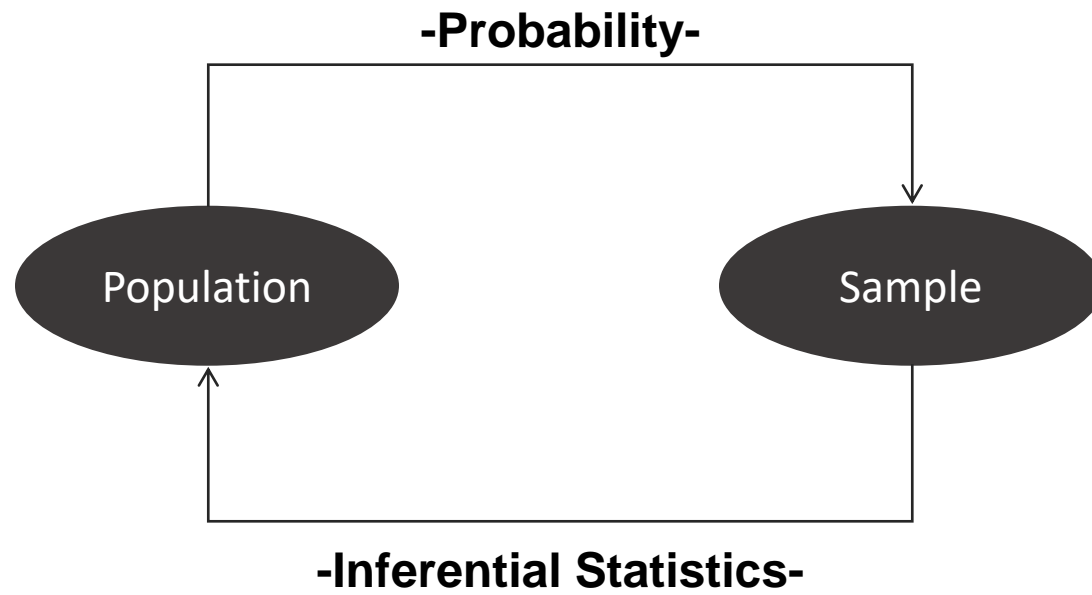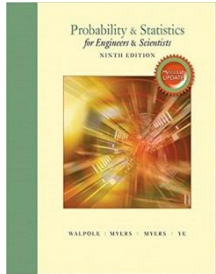
THIRD EDITION

MK

# HCSC Analytics…data processing for decision making aimed to:

| Hierarchical level | Facility | Supply and Inventory | Transportation | Customer service |
|---|---|---|---|---|
| **Strategic** | - Facility location<br>- Capacity setting<br>- Technology selection<br>- Process configuration<br>- Setting the IT system for planning and controlling | - Defining the inventory policy<br>- Identify the supplier list and select the best ones<br>- Product design<br>- Choose the IT system (S&I)<br>- Warehouse design<br>- Material handling system | - Transportation mode<br>- IT system (Trans) | - Define the service policy and strategy<br>- Portfolio indicators |
| **Tactical** | - Capacity planning during mid term | - Purchase planning (procurement)<br>- Definition of supplies<br>- Planning the inventory level<br>- Planning the safety stock | - Transportation system capacity<br>- Fleet routing<br>- Transportation planning during mid term | - Demand projection during mid term<br>- Advertisement planning |
| **Operative** | - Order scheduling<br>- Production execution<br>- Order control<br>- Maintenance planning | - Order dispatching and packing<br>- Material requirement planning<br>- Purchase control<br>- Stock control<br>- Discharge and loading operations | - Delivery planning<br>- Vehicle routing<br>- Control of transport operations | - Demand projection (short term)<br>- Tracking the customer service indicators<br>- Loyalty activities |

**Transforming data into predictive insights…**

...elements in probability allow us to draw conclusions about characteristics of hypothetical data taken from the population, based on known features of the population.

**-Probability-**

Population

Sample

**-Inferential Statistics-**

...elements in probability allow us to draw conclusions about characteristics of hypothetical data taken from the population, based on known features of the population.

**-Probability-**

Population

Sample

Most useful aspect of "theory of probabilities" in data analytics

**-Inferential Statistics-**

Probability distribution

| Probability distribution |

is a mathematical formula that describes the probabilities of occurrence of different possible outcomes in an experiment.

**Probability distribution**

is a mathematical formula that describes the probabilities of occurrence of different possible outcomes in an experiment.

finite number of outcomes

infinite number of outcomes

[predicting Success / Failure] **Discrete**

**Continuous** [predicting cost]

- ❑ Binomial Distribution
- ❑ Poisson Distribution
- ❑ Bernoulli Distribution
- ❑ Geometric Distribution
- ❑ Others.

- ❑ Normal Distribution
- ❑ Uniform Distribution
- ❑ Chi-squared Distribution
- ❑ Exponential Distribution
- ❑ Others.

HELP

# Continuous Probability Distribution

$$P(a < X < b) = \int_a^b f(x)\,dx.$$

$$P(a < X < b) = \int_a^b f(x)\,dx$$

# Continuous Probability Distribution (Normal)

Highly probable…logical explanation

Density function…

…bell-shaped curve

$$F(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \; , \qquad -\infty < x < +\infty$$

$x$: random variable

$\mu$: mean

Deeper study later

$\sigma$: standard deviation

$\pi$: 3.14159 …

$e$: 2.71828 …

Normal Distribution: describes many phenomena that occur in nature, industry, and research

$$X \sim N(\mu, \sigma^2) \; \ldots \; \mu \pm \sigma$$

**Continuous Probability Distribution (Normal)**

Highly probable…logical explanation

Properties…

…bell-shaped curve

$$P(-\infty < X < +\infty) = \int_{-\infty}^{+\infty} f(x)\, dx = 1$$

$$P(-\infty < X < \mu) = P(\mu < X < +\infty)$$

$-3\sigma$ $-2\sigma$ $-1\sigma$ $\mu$ $1\sigma$ $2\sigma$ $3\sigma$ $X$

Normal Distribution: describes many phenomena that occur in nature, industry, and research

HELP

# Continuous Probability Distribution (Normal)



Normal curves with $\mu_1 < \mu_2$ and $\sigma_1 = \sigma_2$.

Normal curves with $\mu_1 < \mu_2$ and $\sigma_1 < \sigma_2$.

Normal curves with $\mu_1 = \mu_2$ and $\sigma_1 < \sigma_2$.

# Continuous Probability Distribution (Normal)

$$F(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad , \qquad -\infty < x < +\infty$$



Computing the probability values

# Continuous Probability Distribution (Normal)

$$F(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad , \qquad -\infty < x < +\infty$$



**Computing the probability values**

$$P(x_1 < X < x_2) = \frac{1}{\sqrt{2\pi}\sigma} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \, dx \quad \text{...hard to solve}$$

HELP

# Continuous Probability Distribution (Normal)

$$F(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad, \qquad -\infty < x < +\infty$$



**Computing the probability values**

$$P(x_1 < X < x_2) = \frac{1}{\sqrt{2\pi}\sigma} \int_{x_1}^{x_2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \quad \text{...hard to solve}$$

**Z-score**

| Z-score | benefits… |

- ❑ **Standardize** all the observations of any normal random variable *X* into a new set of observations;

- ❑ **Reduce the compl**exity of computing the probability;

- ❑ Make possible the **statistical comparison** between to random variables;

- ❑ The **tabulation** of Normal Distribution exists for Z-score only.

$$Z = \frac{x - \mu}{\sigma}$$

The distribution of a normal random variable with mean 0 and variance 1 is called a **standard normal distribution**.

$f(z)$

### Z Score Normal Distribution

Entire area under curve = 100% or 1.00

σ = 1

50% or .50 of values to left of mean

50% or .50 of the values to right of mean

μ = 0

Properties of the Z Score Normal Distribution:

1. Symmetrical

2. Mean = 0 and Standard Deviation = 1

z value: -3, -2, -1, 0, 1, 2, 3

HELP

**Practicing calculations of probabilities using the Normal Distribution…**

**practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

$N(\mu, \sigma^2)$    **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

$$ Z = \frac{x(15) - \mu(12)}{\sigma(\sqrt{4} = 2)} = \frac{15 - 12}{2} = \frac{3}{2} = 1.5 $$

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

$$Z = \frac{x(15) - \mu(12)}{\sigma(\sqrt{4} = 2)} = \frac{15 - 12}{2} = \frac{3}{2} = 1.5$$

Compute:

P (Z >= 1.5) = ?

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

Compute:

P (Z >= 1.5) = ?

Excel…

=NORM

| $f_x$ NORM |
| $f_x$ NORM.DIST |
| $f_x$ NORM.INV |
| $f_x$ NORM.S.DIST |
| $f_x$ NORM.S.INV |
| $f_x$ NORM_CONF |
| $f_x$ NORM_LOWER |
| $f_x$ NORM_UPPER |
| $f_x$ NORM1_POWER |
| $f_x$ NORM1_SIZE |
| $f_x$ NORM2_POWER |
| $f_x$ NORM2_SIZE |

Returns the standard normal distribution (has a mean of zero and a standard deviation of one)

HELP

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

Compute:

P (Z >= 1.5) = ?

Excel…

=NORM.S.DIST(

NORM.S.DIST(**z**, cumulative)

HELP

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

Compute:

P (Z >= 1.5) = ?

Excel…

=NORM.S.DIST(

NORM.S.DIST(**z**, cumulative)

Cumulative = 1

The area below the curve from the left asymptote (bell) to the define z-value

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

Cumulative = 1

Compute:

P (Z >= 1.5) = ?

Excel…

=NORM.S.DIST(

NORM.S.DIST(**z**, cumulative)

=NORM.S.DIST(1.5,TRUE)

| D | E | F | G |
|---|---|---|---|
|   |   |   |   |
|   |   |   | 0.933193 |

**H∃LP**

$N(\mu, \sigma^2)$ **practical examples…**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?
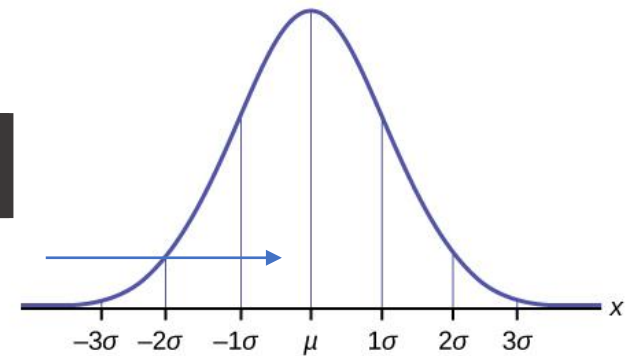
Compute:

P (Z >= 1.5) = ?

Excel…

| 0.933193 |

Prob: 0.066807

Prob: =1-G3

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

Compute:

P ( 250 ≤ X ≤ 325) = ?

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

Compute:

$$P(-2 \leq Z \leq +1) = ?$$

$$Z_{250} = \frac{x(250) - \mu(300)}{\sigma(25)} = \frac{-50}{25} = -2$$

$$Z_{325} = \frac{x(325) - \mu(300)}{\sigma(25)} = \frac{25}{25} = +1$$

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

$$Z_{250} = \frac{x(250) - \mu(300)}{\sigma(25)} = \frac{-50}{25} = -2$$

$$Z_{325} = \frac{x(325) - \mu(300)}{\sigma(25)} = \frac{25}{25} = +1$$

Compute:

$P ( -2 \leq Z \leq +1) =$

$P (Z >= -2) - P (Z >= +1)$

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

Compute:

$$P ( -2 \leq Z \leq +1) =$$

$$P (Z >= -2) - P (Z >= +1)$$

… 0.8186

$$Z_{250} = \frac{x(250) - \mu(300)}{\sigma(25)} = \frac{-50}{25} = -2$$

$$Z_{325} = \frac{x(325) - \mu(300)}{\sigma(25)} = \frac{25}{25} = +1$$

=NORM.S.DIST(1,TRUE)
NORM.S.DIST(z, cumulative)

=NORM.S.DIST(-2,TRUE)
NORM.S.DIST(z, cumulative)

$N(\mu, \sigma^2)$ **practical examples...**

Empirical evidences show that certain supplier can provide an important medical device within a normal distributed delivery time (with $\mu = 12$ and $\sigma^2 = 4$, days and squared-days, respectively). For the firm that receives the devices, more 15 days of lead time would make almost impossible to serve their customers. The main question is: how likely is that delivery time overcomes 15 days?

How do I know this?

Historical dataset provide sufficient evidence to assume our oxygenated water demand is normally distributed, with mean 300 Kgs and standard deviation of 25 Kgs. After a discussion with the financial department, we realize that for overcoming the breaking-even point our sales should be between 250 and 325 kilograms. How probable it is that our sales are between 250 and 325 kilograms?

HELP

The *goodness of fit test...*

The *goodness of fit test* is used to *test* if sample data fits a distribution from a certain population (i.e. a population with a normal distribution or one with a Weibull distribution).

Professional software

...carry it out in Excel...

**Real Statistics Using Excel**

*Everything you need to do real statistical analysis using Excel*

Home   Free Download   Basics   Distributions   ANOVA   Miscellaneous   Regression   Multivariate   Appendix   Blogs   Tools   Contact Us

Charles Zaiontz

| Surgical gloves |
|---|
| 12.8 |
| 11.51 |
| 19.32 |
| 18.4 |
| 14.34 |
| 17.2 |
| 18.78 |
| 12.69 |
| 16.09 |
| 14.06 |
| 14.64 |
| 14 |
| 17.88 |
| 15.49 |
| 18.61 |
| 11.48 |
| 12.58 |

Real Statistics

Desc | Reg | Anova | Time S | Multivar | Corr | Misc

T Test and Non-parametric Equivalents
Non-parametric Tests
Chi-square Test for Independence
Cochran-Mantel-Haenszel Test
Goodness of Fit
Distribution Fitting
Resampling
Multiple Imputation (MI)
Full Information Max Log-Likelihood (FIML)
EM for Missing Multivariate Normal Data
Multiple Tests
Statistical Power and Sample Size

OK
Cancel
Help
Config

For more info see www.real-statistics.com

The *goodness of fit test* is used to *test* if sample data fits a distribution from a certain population (i.e. a population with a normal distribution or one with a Weibull distribution).
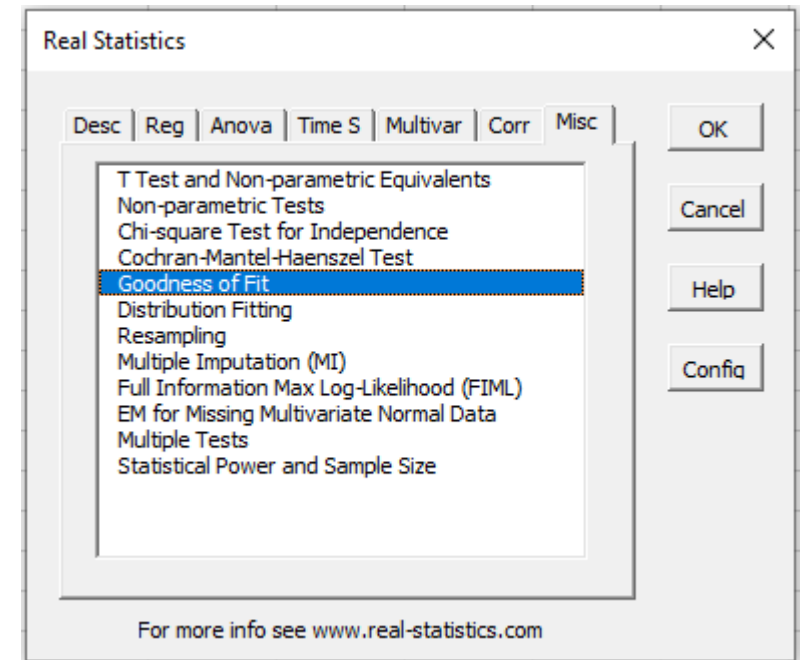
Professional software

…carry it out in Excel…

| Surgical gloves |
|---|
| 12.8 |
| 11.51 |
| 19.32 |
| 18.4 |
| 14.34 |
| 17.2 |
| 18.78 |
| 12.69 |
| 16.09 |
| 14.06 |
| 14.64 |
| 14 |
| 17.88 |
| 15.49 |
| 18.61 |
| 11.48 |
| 12.58 |

Goodness of Fit                    ×

Input Range    Sheet2!$B$4:$B$21    _  Fill    OK

Alpha          0.05                        Cancel

☑ Column headings included with data    Help

Test
○ Two sample KS (freq data)
○ Two sample KS (raw data)
● One sample Anderson-Darling
○ One sample Chi-square

Estimation
● MLE
○ Moments
○ Pure Moments
○ Regression
○ Specify parameters

Distribution
○ Generic    ● Normal    ○ Gamma    ○ Weibull
○ Exponential    ○ Beta    ○ Uniform

Specify parameters
Param 1    N/A         Param 2    N/A

Output Range    Sheet2!$E$4    _  New

The *goodness of fit test* is used to *test* if sample data *fits* a distribution from a certain population (i.e. a population with a normal distribution or one with a Weibull distribution).

Professional software

…carry it out in Excel…

| Surgical gloves |
| --- |
| 12.8 |
| 11.51 |
| 19.32 |
| 18.4 |
| 14.34 |
| 17.2 |
| 18.78 |
| 12.69 |
| 16.09 |
| 14.06 |
| 14.64 |
| 14 |
| 17.88 |
| 15.49 |
| 18.61 |
| 11.48 |
| 12.58 |

**Anderson-Darling Test**

| | | | |
| --- | --- | --- | --- |
| Alpha | 0.05 | mean | 15.28647 |
| Distrib | Normal | std dev | 2.586869 |
| Method | MLE | | |
| | | | |
| AD stat | 0.509352 | | |
| p-value | 0.19793 | | |
| crit value | 0.713947 | | |

**Goodness of Fit** ×

Input Range Sheet2!$B$4:$B$21 — Fill OK

Alpha 0.05 Cancel

☑ Column headings included with data Help

Test
- ○ Two sample KS (freq data)
- ○ Two sample KS (raw data)
- ◉ One sample Anderson-Darling
- ○ One sample Chi-square

Estimation
- ◉ MLE
- ○ Moments
- ○ Pure Moments
- ○ Regression
- ○ Specify parameters

Distribution
- ○ Generic   ◉ Normal   ○ Gamma   ○ Weibull
- ○ Exponential   ○ Beta   ○ Uniform

Specify parameters
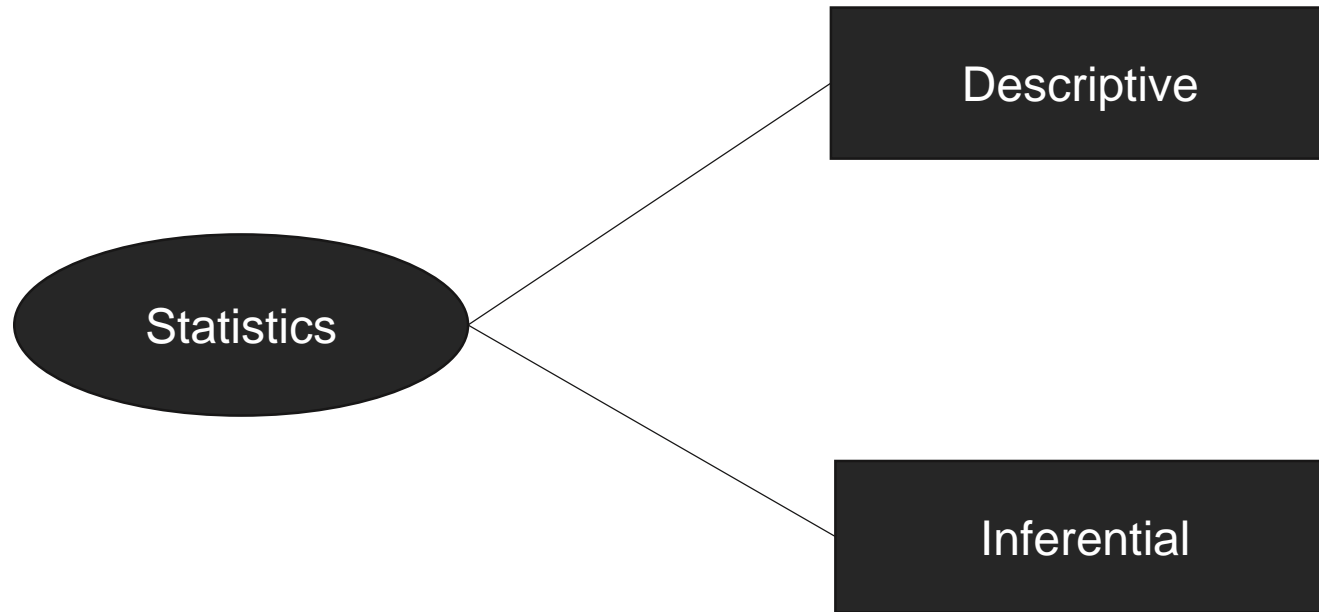Param 1 N/A     Param 2 N/A

Output Range Sheet2!$E$4 — New
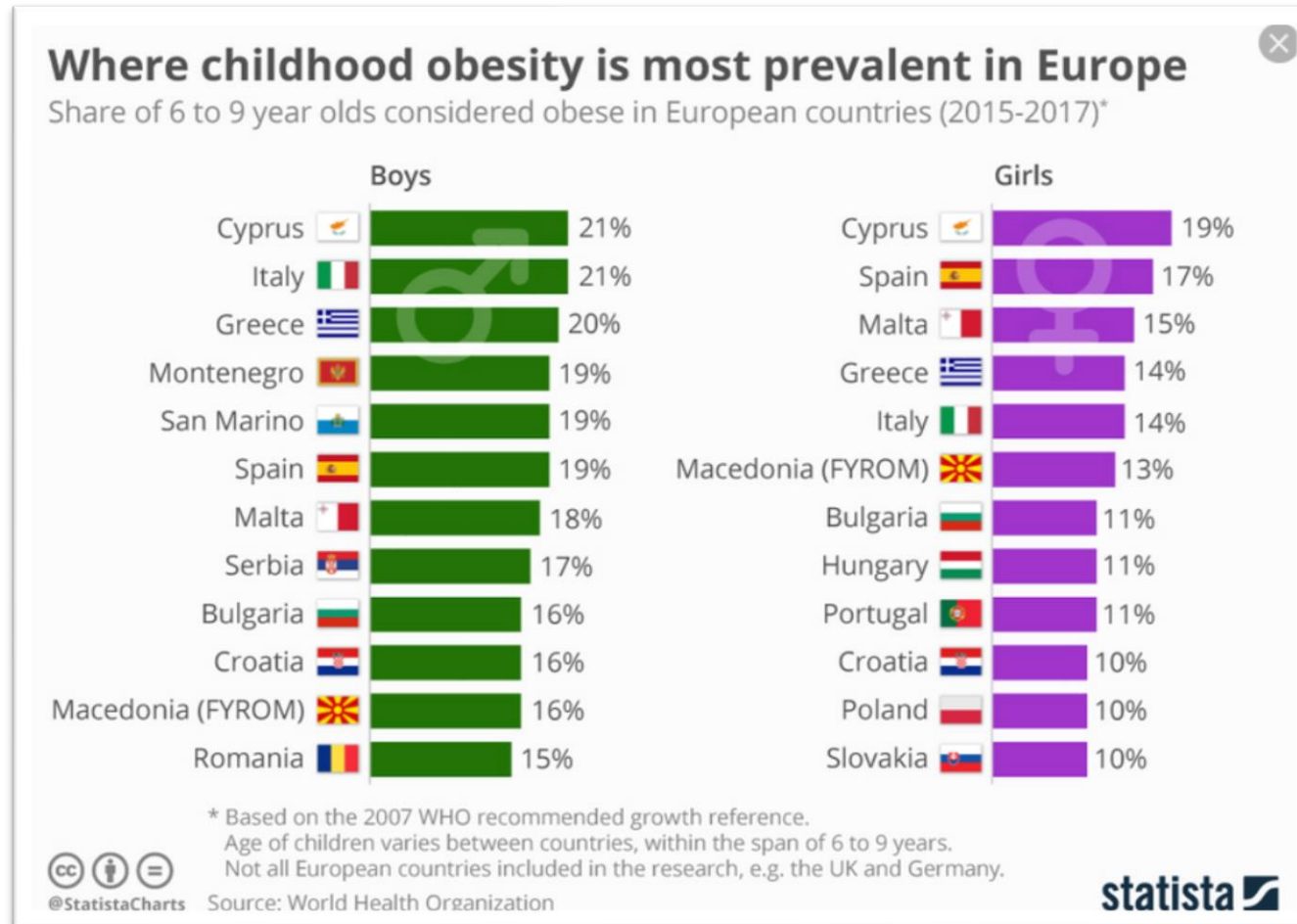
# Refreshing statistics…

## Descriptive

Descriptive statistics is the term given to the analysis of data that helps to describe, show or summarize data in a meaningful way such that, for example, patterns might emerge from the data.

HELP

# Example of descriptive results

e.g.



Where childhood obesity is most prevalent in Europe
Share of 6 to 9 year olds considered obese in European countries (2015-2017)*

**Boys**

| Country | % |
|---|---|
| Cyprus | 21% |
| Italy | 21% |
| Greece | 20% |
| Montenegro | 19% |
| San Marino | 19% |
| Spain | 19% |
| Malta | 18% |
| Serbia | 17% |
| Bulgaria | 16% |
| Croatia | 16% |
| Macedonia (FYROM) | 16% |
| Romania | 15% |

**Girls**

| Country | % |
|---|---|
| Cyprus | 19% |
| Spain | 17% |
| Malta | 15% |
| Greece | 14% |
| Italy | 14% |
| Macedonia (FYROM) | 13% |
| Bulgaria | 11% |
| Hungary | 11% |
| Portugal | 11% |
| Croatia | 10% |
| Poland | 10% |
| Slovakia | 10% |

* Based on the 2007 WHO recommended growth reference.
Age of children varies between countries, within the span of 6 to 9 years.
Not all European countries included in the research, e.g. the UK and Germany.

@StatistaCharts    Source: World Health Organization

statista

# Elements

**Measures of central tendency:** these are ways of describing the central position of a frequency distribution for a group of data

Mean (average, geometric, harmonic)

Median

Mode

Later in Excel

HELP

# Elements

**Measures of central tendency:** these are ways of describing the central position of a frequency distribution for a group of data

Mean (average, geometric, harmonic)

Median

Mode

Later in Excel as well

working with samples…no compensation

$$G = \sqrt[n]{x_i \cdot x_{i+1} \cdot x_{i+2} \cdot \cdots \cdot x_n}$$

# Elements

**Measures of central tendency:** these are ways of describing the central position of a frequency distribution for a group of data

Mean (average, geometric, harmonic)

Median

Mode

Examine in details

$$H = \frac{N}{\sum_{i=1}^{n} 1/X_i}$$

working with samples…less important the positive extreme values

**Using the measures of central tendency in forecasting…**

| Date | Surgical gloves | | Mean(average) | Errors |
|---|---|---|---|---|
| Jan-17 | 12.8 | | 16.74529412 | =ABS(C3-E3 |
| Feb-17 | 11.51 | | 16.74529412 | ABS(number) |
| Mar-17 | 19.32 | | 16.74529412 | |
| Apr-17 | 18.4 | | 16.74529412 | |
| May-17 | 14.34 | | 16.74529412 | |
| Jun-17 | 31 | | 16.74529412 | |
| Jul-17 | 18.78 | | 16.74529412 | |
| Aug-17 | 12.69 | | 16.74529412 | |
| Sep-17 | 16.09 | | 16.74529412 | |
| Oct-17 | 14.06 | | 16.74529412 | |
| Nov-17 | 14.64 | | 16.74529412 | |
| Dec-17 | 25 | | 16.74529412 | |
| Jan-18 | 17.88 | | 16.74529412 | |
| Feb-18 | 15.49 | | 16.74529412 | |
| Mar-18 | 18.61 | | 16.74529412 | |
| Apr-18 | 11.48 | | 16.74529412 | |
| May-18 | 12.58 | | 16.74529412 | |

HELP

# Using the measures of central tendency in forecasting…

| Date | Surgical gloves | | Mean(average) | Errors | | Geometric-Mean | Errors |
|---|---|---|---|---|---|---|---|
| Jan-17 | 12.8 | | 16.74529412 | 3.945294 | | =GEOMEAN(C3:C19) | |
| Feb-17 | 11.51 | | 16.74529412 | 5.235294 | | GEOMEAN(**number1**, [number2], …) | |
| Mar-17 | 19.32 | | 16.74529412 | 2.574706 | | 16.13918802 | |
| Apr-17 | 18.4 | | 16.74529412 | 1.654706 | | 16.13918802 | |
| May-17 | 14.34 | | 16.74529412 | 2.405294 | | 16.13918802 | |
| Jun-17 | 31 | | 16.74529412 | 14.25471 | | 16.13918802 | |
| Jul-17 | 18.78 | | 16.74529412 | 2.034706 | | 16.13918802 | |
| Aug-17 | 12.69 | | 16.74529412 | 4.055294 | | 16.13918802 | |
| Sep-17 | 16.09 | | 16.74529412 | 0.655294 | | 16.13918802 | |
| Oct-17 | 14.06 | | 16.74529412 | 2.685294 | | 16.13918802 | |
| Nov-17 | 14.64 | | 16.74529412 | 2.105294 | | 16.13918802 | |
| Dec-17 | 25 | | 16.74529412 | 8.254706 | | 16.13918802 | |
| Jan-18 | 17.88 | | 16.74529412 | 1.134706 | | 16.13918802 | |
| Feb-18 | 15.49 | | 16.74529412 | 1.255294 | | 16.13918802 | |
| Mar-18 | 18.61 | | 16.74529412 | 1.864706 | | 16.13918802 | |
| Apr-18 | 11.48 | | 16.74529412 | 5.265294 | | 16.13918802 | |
| May-18 | 12.58 | | 16.74529412 | 4.165294 | | 16.13918802 | |

$$G = \sqrt[n]{x_i \cdot x_{i+1} \cdot x_{i+2} \cdot \cdots \cdot x_n}$$

# Using the measures of central tendency in forecasting…

$$H = \frac{N}{\sum_{i=1}^{n} 1/X_i}$$

| Date | Surgical gloves | Mean(average) | Errors | Geometric-Mean | Errors | Harmonic-Mean | Errors |
|---|---|---|---|---|---|---|---|
| Jan-17 | 12.8 | 16.74529412 | 3.945294 | 16.13918802 | 3.339188 | =HARMEAN(C3:C19) | |
| Feb-17 | 11.51 | 16.74529412 | 5.235294 | 16.13918802 | 4.629188 | HARMEAN(**number1**, [number2], ...) | |
| Mar-17 | 19.32 | 16.74529412 | 2.574706 | 16.13918802 | 3.180812 | 15.6313613 | |
| Apr-17 | 18.4 | 16.74529412 | 1.654706 | 16.13918802 | 2.260812 | 15.6313613 | |
| May-17 | 14.34 | 16.74529412 | 2.405294 | 16.13918802 | 1.799188 | 15.6313613 | |
| Jun-17 | 31 | 16.74529412 | 14.25471 | 16.13918802 | 14.86081 | 15.6313613 | |
| Jul-17 | 18.78 | 16.74529412 | 2.034706 | 16.13918802 | 2.640812 | 15.6313613 | |
| Aug-17 | 12.69 | 16.74529412 | 4.055294 | 16.13918802 | 3.449188 | 15.6313613 | |
| Sep-17 | 16.09 | 16.74529412 | 0.655294 | 16.13918802 | 0.049188 | 15.6313613 | |
| Oct-17 | 14.06 | 16.74529412 | 2.685294 | 16.13918802 | 2.079188 | 15.6313613 | |
| Nov-17 | 14.64 | 16.74529412 | 2.105294 | 16.13918802 | 1.499188 | 15.6313613 | |
| Dec-17 | 25 | 16.74529412 | 8.254706 | 16.13918802 | 8.860812 | 15.6313613 | |
| Jan-18 | 17.88 | 16.74529412 | 1.134706 | 16.13918802 | 1.740812 | 15.6313613 | |
| Feb-18 | 15.49 | 16.74529412 | 1.255294 | 16.13918802 | 0.649188 | 15.6313613 | |
| Mar-18 | 18.61 | 16.74529412 | 1.864706 | 16.13918802 | 2.470812 | 15.6313613 | |
| Apr-18 | 11.48 | 16.74529412 | 5.265294 | 16.13918802 | 4.659188 | 15.6313613 | |
| May-18 | 12.58 | 16.74529412 | 4.165294 | 16.13918802 | 3.559188 | 15.6313613 | |

HELP

# Using the measures of central tendency in forecasting…

| Date | Surgical gloves | | Mean(average) | Errors | | Geometric-Mean | Errors | | Harmonic-Mean | Errors |
|---|---|---|---|---|---|---|---|---|---|---|
| Jan-17 | 12.8 | | 16.74529412 | 3.945294 | | 16.13918802 | 3.339188 | | 15.6313613 | 2.831361 |
| Feb-17 | 11.51 | | 16.74529412 | 5.235294 | | 16.13918802 | 4.629188 | | 15.6313613 | 4.121361 |
| Mar-17 | 19.32 | | 16.74529412 | 2.574706 | | 16.13918802 | 3.180812 | | 15.6313613 | 3.688639 |
| Apr-17 | 18.4 | | 16.74529412 | 1.654706 | | 16.13918802 | 2.260812 | | 15.6313613 | 2.768639 |
| May-17 | 14.34 | | 16.74529412 | 2.405294 | | 16.13918802 | 1.799188 | | 15.6313613 | 1.291361 |
| Jun-17 | 31 | | 16.74529412 | 14.25471 | | 16.13918802 | 14.86081 | | 15.6313613 | 15.36864 |
| Jul-17 | 18.78 | | 16.74529412 | 2.034706 | | 16.13918802 | 2.640812 | | 15.6313613 | 3.148639 |
| Aug-17 | 12.69 | | 16.74529412 | 4.055294 | | 16.13918802 | 3.449188 | | 15.6313613 | 2.941361 |
| Sep-17 | 16.09 | | 16.74529412 | 0.655294 | | 16.13918802 | 0.049188 | | 15.6313613 | 0.458639 |
| Oct-17 | 14.06 | | 16.74529412 | 2.685294 | | 16.13918802 | 2.079188 | | 15.6313613 | 1.571361 |
| Nov-17 | 14.64 | | 16.74529412 | 2.105294 | | 16.13918802 | 1.499188 | | 15.6313613 | 0.991361 |
| Dec-17 | 25 | | 16.74529412 | 8.254706 | | 16.13918802 | 8.860812 | | 15.6313613 | 9.368639 |
| Jan-18 | 17.88 | | 16.74529412 | 1.134706 | | 16.13918802 | 1.740812 | | 15.6313613 | 2.248639 |
| Feb-18 | 15.49 | | 16.74529412 | 1.255294 | | 16.13918802 | 0.649188 | | 15.6313613 | 0.141361 |
| Mar-18 | 18.61 | | 16.74529412 | 1.864706 | | 16.13918802 | 2.470812 | | 15.6313613 | 2.978639 |
| Apr-18 | 11.48 | | 16.74529412 | 5.265294 | | 16.13918802 | 4.659188 | | 15.6313613 | 4.151361 |
| May-18 | 12.58 | | 16.74529412 | 4.165294 | | 16.13918802 | 3.559188 | | 15.6313613 | 3.051361 |
| | | | | | | | | | | |
| | | | | | | | | | | |
| | | | Average-all-errors: | 3.737993 | | | 3.631033 | | | 3.595374 |

Outliers

Small errors in predictions…

# Elements

**Measures of spread:** these are ways of summarizing a group of data by describing how spread out the scores are.

$$S = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}}$$

Standard deviation

Variance

Range

CV

Later in Excel…and
Real Stat Add-Ins

The coefficient of variation (CV) represents the ratio of the standard deviation to the mean, and it is a useful statistic for comparing the degree of variation from one data series to another, even if the means are drastically different from each other.

$$CV = \frac{S}{\bar{X}}$$

Useful in risk analysis…

**Supplier 1(Kgs)** *vs* **Supplier 2 (units)**

**Measures of Distribution**

**Skewness**: distribution (aggregations of observations) can be spread around both sides of the central tendency.



**Kurtosis:** is the measure of the peak of a distribution, and indicates how high is around the mean.

**Analyzing the weekly patient arrivals…**

Comprehensive module for descriptive statistics

# Analyzing the weekly patient arrivals…

| Descriptive Statistics | |
|---|---|
| | *Patients arrivals* |
| Mean | 102.1730769 |
| Standard Error | 2.122692368 |
| Median | 103.5 |
| Mode | 110 |
| Standard Deviation | 15.30695235 |
| Sample Variance | 234.3027903 |
| Kurtosis | 1.534557486 |
| Skewness | 0.397113782 |
| Range | 87 |
| Maximum | 151 |
| Minimum | 64 |
| Sum | 5313 |
| Count | 52 |
| Geometric Mean | 101.0481308 |
| Harmonic Mean | 99.90476696 |
| AAD | 11.62795858 |
| MAD | 9.5 |
| IQR | 19.25 |

102 patients arrive…on average

$$SE = \frac{S}{\sqrt{n}}$$

↓(data better distributed)

The most repeated value

A little high…

Symmetric respect to the mean

HELP

# Analyzing the weekly patient arrivals…

| Descriptive Statistics | |
| --- | --- |
| | *Patients arrivals* |
| Mean | 102.1730769 |
| Standard Error | 2.122692368 |
| Median | 103.5 |
| Mode | 110 |
| Standard Deviation | 15.30695235 |
| Sample Variance | 234.3027903 |
| Kurtosis | 1.534557486 |
| Skewness | 0.397113782 |
| Range | 87 |
| Maximum | 151 |
| Minimum | 64 |
| Sum | 5313 |
| Count | 52 |
| Geometric Mean | 101.0481308 |
| Harmonic Mean | 99.90476696 |
| AAD | 11.62795858 |
| MAD | 9.5 |
| IQR | 19.25 |

Average of the Absolute Deviation…

$$AAD = \frac{1}{n}\sum|x_i - \bar{x}|$$

# Analyzing the weekly patient arrivals…

| Descriptive Statistics | |
|---|---|
| | |
| | *Patients arrivals* |
| Mean | 102.1730769 |
| Standard Error | 2.122692368 |
| Median | 103.5 |
| Mode | 110 |
| Standard Deviation | 15.30695235 |
| Sample Variance | 234.3027903 |
| Kurtosis | 1.534557486 |
| Skewness | 0.397113782 |
| Range | 87 |
| Maximum | 151 |
| Minimum | 64 |
| Sum | 5313 |
| Count | 52 |
| Geometric Mean | 101.0481308 |
| Harmonic Mean | 99.90476696 |
| AAD | 11.62795858 |
| MAD | 9.5 |
| IQR | 19.25 |

Median Absolute Deviation…

$$\text{Median } \{|x_i - \tilde{x}| : x_i \text{ in } S\}$$

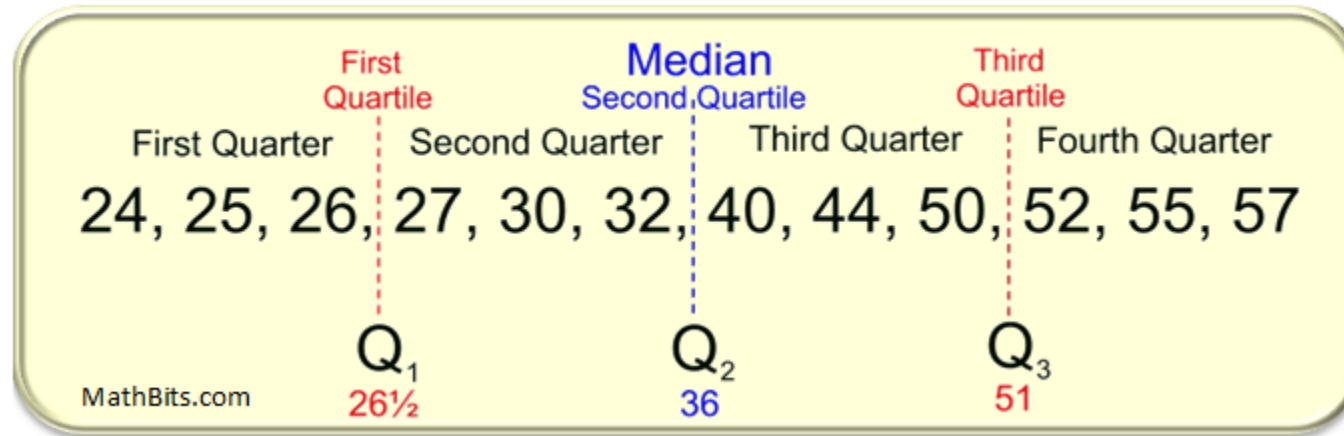where $\tilde{x}$ = median of the data elements in $S$.

**Analyzing the weekly patient arrivals…**

| Descriptive Statistics | |
|---|---|
| | |
| | *Patients arrivals* |
| Mean | 102.1730769 |
| Standard Error | 2.122692368 |
| Median | 103.5 |
| Mode | 110 |
| Standard Deviation | 15.30695235 |
| Sample Variance | 234.3027903 |
| Kurtosis | 1.534557486 |
| Skewness | 0.397113782 |
| Range | 87 |
| Maximum | 151 |
| Minimum | 64 |
| Sum | 5313 |
| Count | 52 |
| Geometric Mean | 101.0481308 |
| Harmonic Mean | 99.90476696 |
| AAD | 11.62795858 |
| MAD | 9.5 |
| IQR | 19.25 |

Inter-quartile Range…

# Descriptive Stats

Inter-quartile Range…



First Quartile
First Quarter

Median
Second Quartile
Second Quarter

Third Quartile
Third Quarter
Fourth Quarter

24, 25, 26, 27, 30, 32, 40, 44, 50, 52, 55, 57

$Q_1$ 26½  $Q_2$ 36  $Q_3$ 51

MathBits.com

Sorting data

Descriptive Stats

Inter-quartile Range…

Median | Median
Q1 | Q2 | Q3

25% | 25% | 25% | 25%

Interquartile Range
= Q3 - Q1

Example 1: (even number)

4  7  9 | 11  12  20

IQR = 12 - 7 = 5

Example 2: (odd number)

5  8  10  10  15  18  23

IQR = 18 - 8 = 10

https://www.pinterest.com/pin/2039286893534078866/

# Descriptive Stats

**Box plot**



**Outliers** are:

greater than $Q_3 + (1.5 \cdot IQR)$

(referred to as the upper fence)

or less than $Q_1 - (1.5 \cdot IQR)$

(referred to as the lower fence)

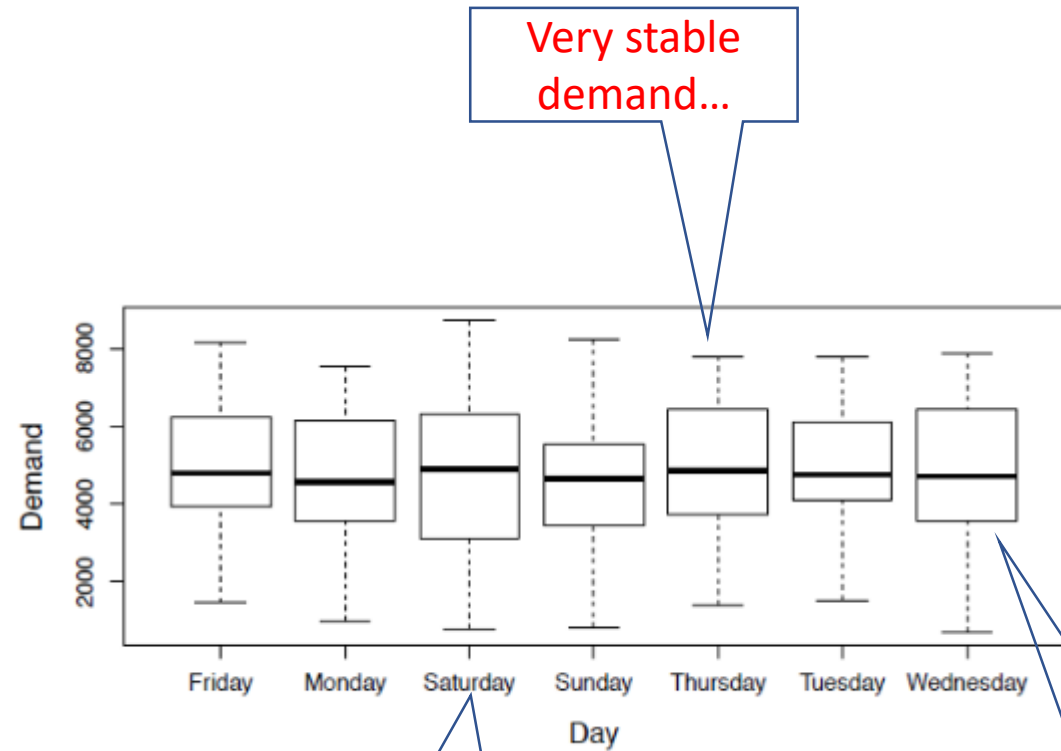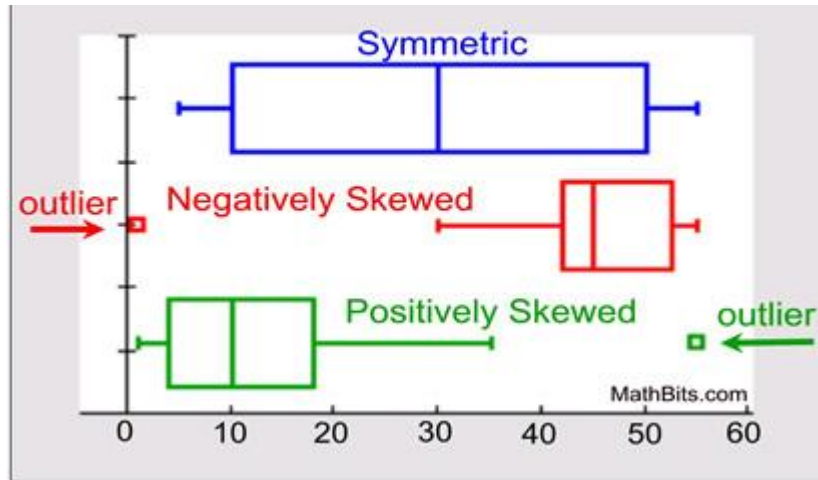| 25 | 26½ | 36 | 51 | 57 |
|----|-----|-----|-----|-----|
| Minimum | First Quartile | Median | Third Quartile | Maximum |

MathBits.com

0   10   20   30   40   50   60

# Descriptive Stats

# Descriptive Stats

# Analyzing the weekly patient arrivals…

| Multiplier | 1.5 |
|---|---|

| Patients arrivals | |
|---|---|
| Min | 64 |
| Q1-Min | 27 |
| Med-Q1 | 12.5 |
| Q3-Med | 6.75 |
| Max-Q3 | 27.75 |
| Mean | 102.1731 |

| Min | 64 |
|---|---|
| Q1 | 91 |
| Median | 103.5 |
| Q3 | 110.25 |
| Max | 138 |
| Mean | 102.1731 |

| Grand Min | 0 |
|---|---|

| Outliers | 151 |
|---|---|

**Box Plot**

Patients arrivals

One outlier…the data could be removed from the dataset

**Analyzing the weekly patient arrivals…**

| Shapiro-Wilk Test | |
|---|---|
| | |
| *Patients arrivals* | |
| W-stat | 0.96806 |
| p-value | 0.174667 |
| alpha | 0.05 |
| normal | yes |
| | |
| d'Agostino-Pearson | |
| | |
| DA-stat | 5.134768 |
| p-value | 0.076736 |
| alpha | 0.05 |
| normal | yes |

**There is normality**

1) Define the hypothesis

2) Identify the proper statistical test

3) Compute the p-value

4) Compare p-value against an "acceptable significance value(α)"…then make a decision…

If $p$-value $\leq$ α Then
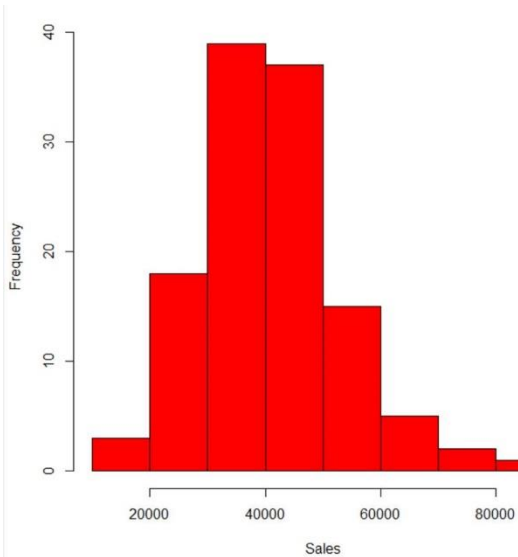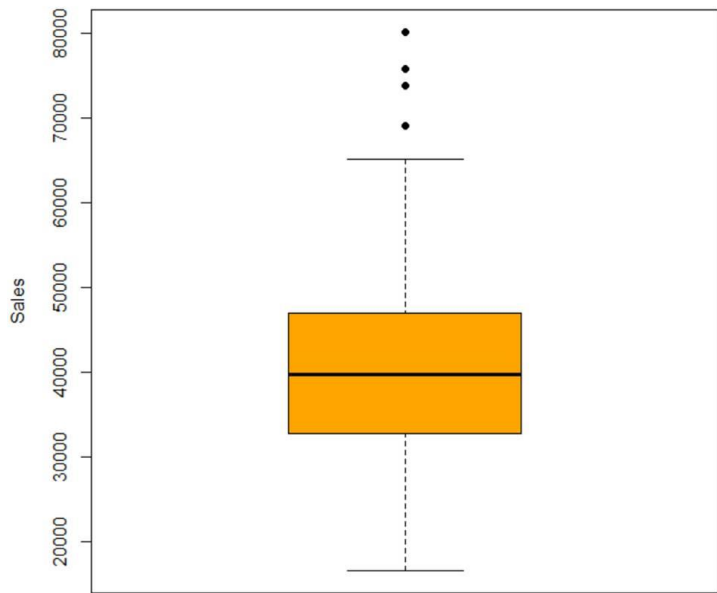the null hypothesis is ruled out, and the alternative hypothesis is valid.
Else
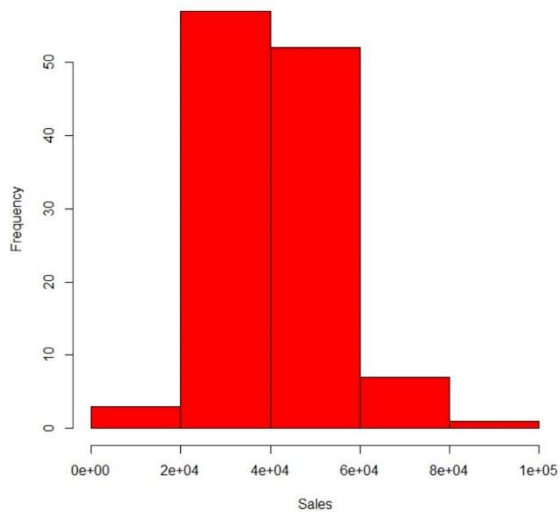The null hypothesis is valid

H0: The data follow a Normal Distribution
H1: The data do not follow a Normal Distribution

HELP

# Descriptive-Stats in R…





Histogram with too little categories



Histogram with too many categories