BE
·DAC

# Introduction to Biostatistics

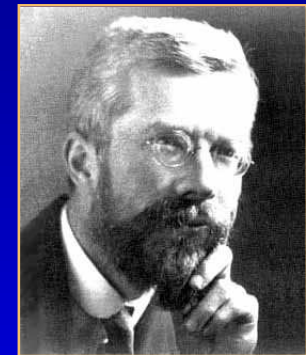## Part I: Basic Concepts

### Helmut Schütz
### BEBAC

π
ε
χ
ε
π Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

# Biometry, Biometrics, and Biostatistics

● Introduced in 1947 by R.A Fisher as *'Biometry'* and later *'Biometrics'*

*'Biometry, the active pursuit of biological knowledge by quantitative methods.'*

● The International Biometric Society

*'The terms "Biometrics" and "Biometry" have been used since early in the 20th century to refer to the field of development of statistical and mathematical methods applicable to data analysis problems in the biological sciences. Recently, the term "Biometrics" has also been used to refer to the emerging field of technology devoted to identification of individuals […]'*

● *'Biostatistics'* was introduced as a new term…

# Biometry, Biometrics, and Biostatistics

***Statistics.*** A subject which most statisticians find difficult but in which nearly all physicians are expert.
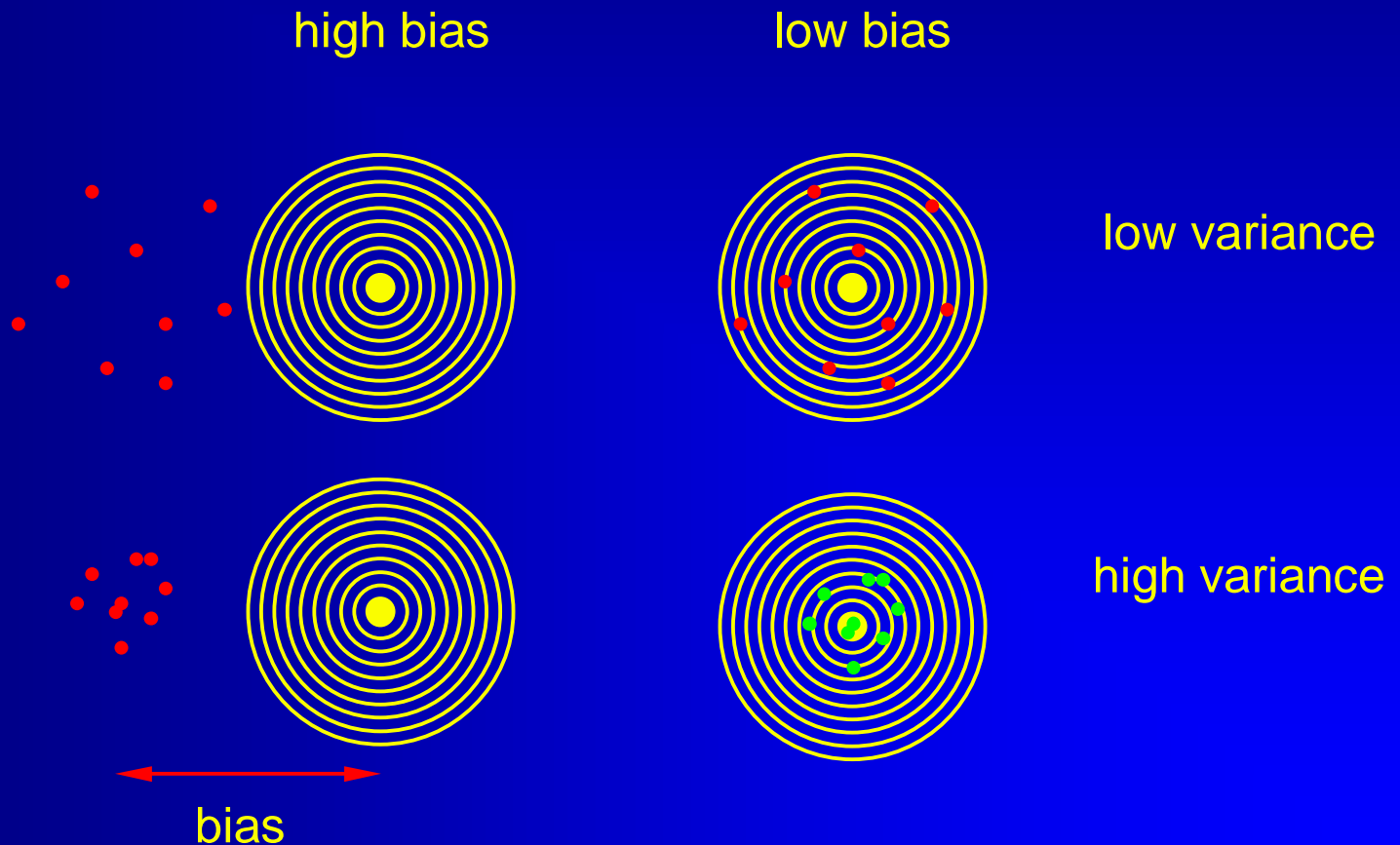
***Biostatistician.*** One who has neither the intellect for mathematics nor the commitment for medicine but likes to dabble in both.
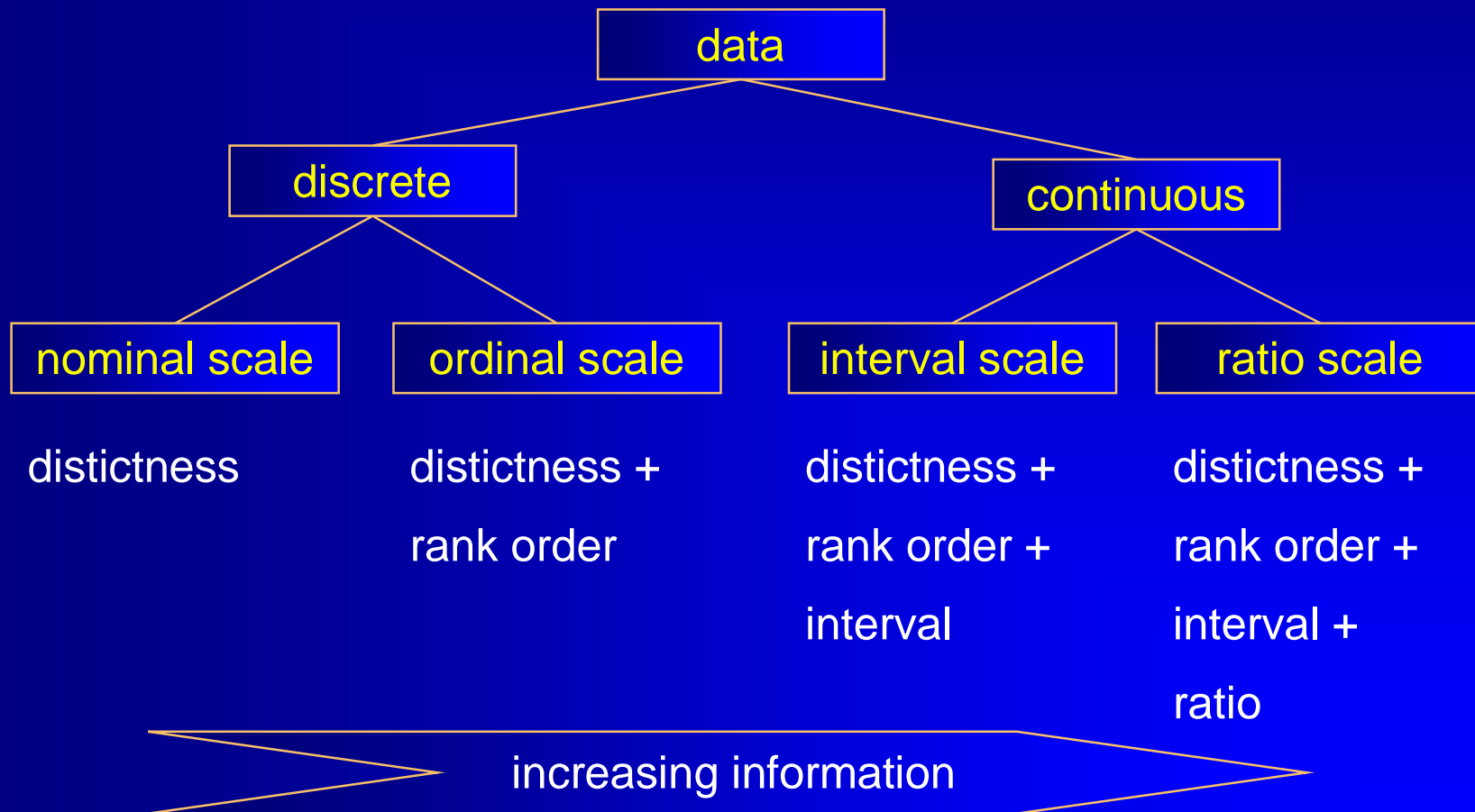
***Medical statistician.*** One who will not accept that Columbus discovered America… because he said he was looking for India in the trial plan.

*Stephen Senn*

# Terminology I

high bias          low bias

low variance

high variance

bias

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*4 • 57*

Pharma Edge

# Terminology II

```
                          ┌──────────────┐
                          │     data     │
                          └──────────────┘
                  ┌───────────────┴───────────────┐
          ┌──────────────┐                 ┌──────────────┐
          │   discrete   │                 │  continuous  │
          └──────────────┘                 └──────────────┘
          ┌───────┴───────┐                 ┌──────┴───────┐
┌──────────────┐  ┌──────────────┐  ┌──────────────┐  ┌──────────────┐
│ nominal scale│  │ ordinal scale│  │interval scale│  │  ratio scale │
└──────────────┘  └──────────────┘  └──────────────┘  └──────────────┘
```

| nominal scale | ordinal scale | interval scale | ratio scale |
|---|---|---|---|
| distictness | distictness + rank order | distictness + rank order + interval | distictness + rank order + interval + ratio |

increasing information

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*5 • 57*

Pharma Edge

# Data I

● Nominal scale (aka categorial)

■ Sex, ethnicity,…

■ Statistics:         mode, $\chi^2$ test

■ Transformations: equality

● Ordinal scale

■ School grades, disease states,…

■ Statistics:         median, percentile, sign test, Wilcoxon test

■ Transformations: monotonic increasing order

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

6 • 57

# Data II

- Interval scale
  - Calendar dates, temperature in °C, IQ,…
    - Statistics: mean, variance (standard deviation), correlation, regression, ANOVA
    - Transformations: linear

- Ratio scale
  - Measures with true zero point, temperature in K,…
    - Statistics: all of the above, geometric and harmonic mean, coefficient of variation
    - Transformations: multiplicative, logarithm

# Examples from PK

- Ordinal scale
  - $t_{max}$, $t_{lag}$
    - Statistics: median, percentile, sign test, Wilcoxon test
    - Transformations: monotonic increasing order
- Ratio scale
  - AUC, $C_{max}$, $\lambda_z$,…
    - Statistics: mean, variance (standard deviation), correlation, regression, ANOVA, geometric and harmonic mean, coefficient of variation
    - Transformations: multiplicative, logarithm

# Bell curve – and beyond

- Abraham de Moivre (1667–1754),
  Pierre-Simon Laplace (1749–1827)
  Central limit theorem 1733, 1812

- Carl F. Gauß (1777–1855)
  Normal distribution 1795
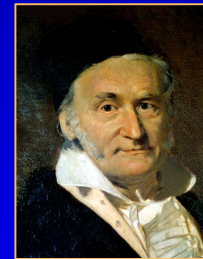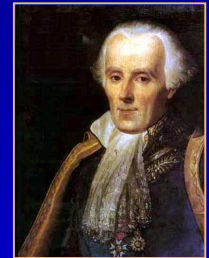
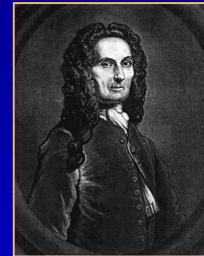- William S. Gosset, aka Student
  (1876–1937)
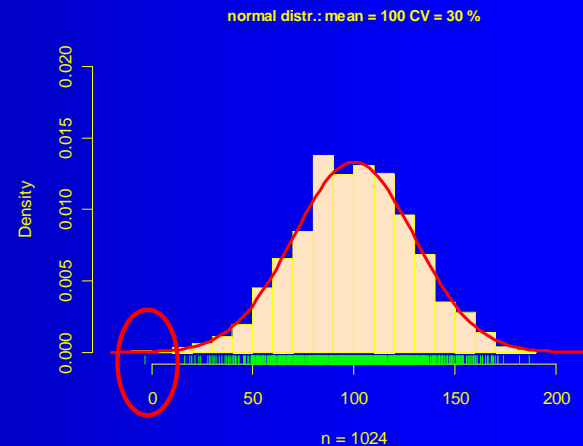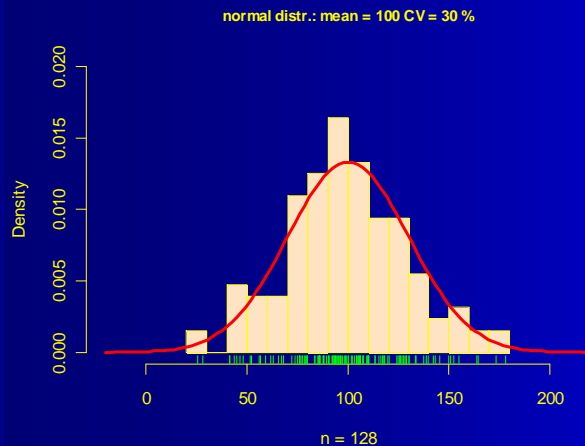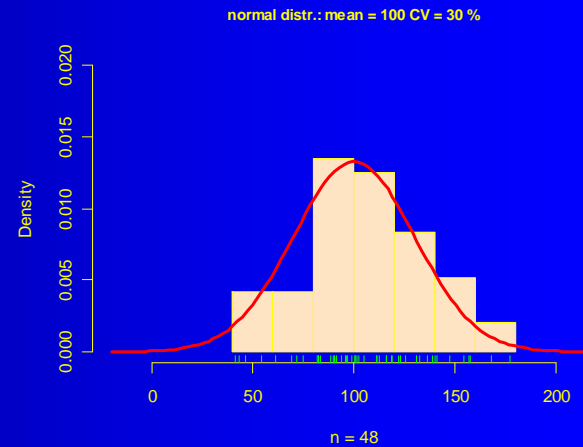  *t*-distribution 1908

- Ronald A. Fisher (1890–1962)
  Analysis of variance 1918

- Frank Wilcoxon (1892–1965)
  Nonparametric tests 1945

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*9 • 57*

Pharma Edge

# Statistical Distributions



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*10 • 57*

# Statistical Distributions



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*11 • 57*

# Statistical Distributions

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

● Normal Distribution

  ■ Defined by location (aka central tendency) and dispersion

  ■ Population

    ■ Location: population mean $\mu$

    ■ Dispersion: population variance $\sigma^2$

  ■ Sample

    ■ Location: sample mean $\overline{x}$

    ■ Dispersion: sample variance $s^2$

  ■ Probability = 1 within <u>-∞</u> and <u>+∞</u>

$$F(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{x} e^{-\frac{1}{2}\left(\frac{t-\mu}{\sigma}\right)^2} dt$$

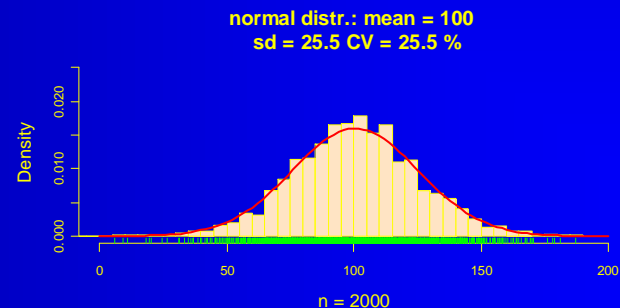**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

π ε χ ε π Pharma Edge

12 • 57

# Statistical Distributions

● Lognormal Distribution

$$f(x) = \begin{cases} \dfrac{1}{\sigma x \sqrt{2\pi}}\, e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}} & x > 0 \\ 0 & x \leq 0 \end{cases}$$

  ■ Defined by location and dispersion

  ■ Population

    ■ Location:    population mean $\mu$

    ■ Dispersion: population variance $\sigma^2$

  ■ Sample

    ■ Location:    sample mean $\overline{x}$

    ■ Dispersion: sample variance $s^2$

  ■ Probability = 1 within <u>0</u> and <u>+∞</u>

$$F(x) = \frac{1}{\sigma \sqrt{2\pi}} \int_0^x \frac{1}{t}\, e^{-\frac{(\ln t - \mu)^2}{2\sigma^2}}\, dt$$

Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011

Pharma Edge

13 • 57

# Statistical Distributions



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*14 • 57*

# Statistical Distributions



population: mean = 100
sd = 20.08 CV = 20.08 %
N = 1e+06

sample 1: mean = 101
sd = 15.94 CV = 15.78 %
n = 36

sample 2: mean = 100.1
sd = 20.41 CV = 20.38 %
n = 36

sample 3: mean = 100.2
sd = 19.61 CV = 19.57 %
n = 36

sample 4: mean = 96.69
sd = 19.31 CV = 19.97 %
n = 36

sample 5: mean = 102.5
sd = 19.31 CV = 18.85 %
n = 36

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*15 • 57*

Pharma Edge

# Statistical Distributions



population: mean = 100
sd = 20.03 CV = 20.03 %

20 samples drawn
from population

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

π Pharma Edge

*16 • 57*

# Central Limit Theorem

- If samples are drawn by a random process from a population with a normal distribution, distribution of sample means is also normal.

- The mean of the distribution of sample means is identical to the mean of the 'parent population' – the population from which the samples are drawn.

- The higher the sample size that is drawn, the 'narrower' will be the dispersion of the distribution of sample means.

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*17 • 57*

# Normal Distribution I

Standard normal distribution: $\mu = 0$, $\sigma = 1$

±1σ p 68.27%

Standard normal distribution: $\mu = 0$, $\sigma = 1$

±2σ p 95.45%

Standard normal distribution: $\mu = 0$, $\sigma = 1$

±3σ p 99.73%

Standard normal distribution: $\mu = 0$, $\sigma = 1$

±4σ p 99.99%

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*18 • 57*

# Normal Distribution II



**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*19 • 57*

# Normal Distribution III



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*20 • 57*

Pharma Edge

# Confidence Interval I

● If we have drawn a sample from a population, we get the sample mean $\overline{x}$ and the sample standard deviation $s$.

● Can we make a prediction about the population mean?

● Yes. That's called a Confidence Interval (CI).

  ■ If $\sigma$ is known:

$$[\mu] = \overline{x} \pm z \frac{\sigma}{\sqrt{n}}$$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*21 • 57*

# Confidence Interval II

● Example from previous slides: $\mu$ 100, $\sigma$ 20
Sample sizes 36, $z_{0.05}$ 1.960

| Samples' means | Confidence Interval | |
|---|---|---|
| 101.0 | 94.47 | 107.5 |
| 100.1 | 93.57 | 106.6 |
| 100.2 | 93.67 | 106.7 |
| 96.69 | 90.16 | 103.2 |
| 102.5 | 95.97 | 109.0 |

● But generally we don't know $\sigma$!

● Help is on the way…

# Student's *t* Distribution

$$f(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\,\Gamma\left(\frac{\nu}{2}\right)}\left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

$$\Gamma(x) = \int\limits_0^{+\infty} t^{x-1}e^{-t}dt$$

- Depends on one parameter, the 'degrees of freedom $\nu$'. In the most simple case df = n – 1.

- The *t* Distribution is 'heavy tailed' compared to the normal distribution. Small sample sizes are penalized.

- Approaches quickly the normal distribution for df >≈30.

- Allows calculation of a CI of the sample mean based on the sample standard deviation $s$.

$\pi$ Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*23 • 57*

# Student's *t* Distribution



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*24 • 57*

Pharma Edge

# Confidence Interval III

- Example from previous slides: $\mu$ 100, $\sigma$ 20
  Sample sizes 36, $z_{0.05}$ 1.960, $t_{36-1,0.05}$ 2.030

| Samples' mean | stand. dev. | Confidence Intervals based on $z$ | | based on $t$ | |
|---|---|---|---|---|---|
| 101.0 | 15.94 | 94.47 | 107.5 | 95.61 | 106.4 |
| 100.1 | 20.41 | 93.57 | 106.6 | 93.19 | 107.0 |
| 100.2 | 19.61 | 93.67 | 106.7 | 93.57 | 106.8 |
| 96.69 | 19.31 | 90.16 | 103.2 | 90.16 | 103.2 |
| 102.5 | 19.31 | 95.97 | 109.0 | 95.97 | 109.0 |

$$[\mu] = \bar{x} \pm z \frac{\sigma}{\sqrt{n}} \qquad [\mu] = \bar{x} \pm t \frac{s}{\sqrt{n}}$$

Pharma Edge
**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**
25 • 57

# Location parameters

- x = [91,72,141,119,92,124,92,101,90,145]
  ranks [3, 1, 9, 7, 4.5, 8, 4.5, 6, 2, 10]
  ordered [72,90,91,92,92,101,119,124,141,145]

- Mode: 92 (most frequent number)

- Median: 96.5 (middle value)
  - If n=odd:   value at $x_{n/2}$
  - If n=even: value $(x_{n/2}+x_{n/2+1})/2 = (92+101)/2$

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*26 • 57*

# Location parameters

● Harmonic mean: 101.9516

$$\overline{x}_{harm} = \frac{n}{\displaystyle\sum_{i=1}^{i=n} \frac{1}{x_i}} = \frac{n}{\dfrac{1}{x_1} + \dfrac{1}{x_2} \cdots \dfrac{1}{x_i}}$$

● Geometric mean: 104.2814

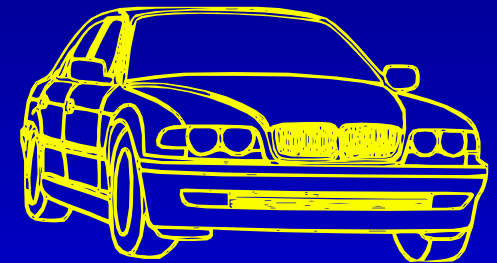$$\overline{x}_{geom} = \sqrt[n]{\prod_{i=1}^{i=n} x_i} = \sqrt[n]{x_1 \cdot x_2 \cdots x_n} = e^{\frac{1}{n}\sum_{i=1}^{i=n} \ln x_i}$$

● Arithmetic mean: 106.7

$$\overline{x}_{arithm} = \frac{1}{n} \sum_{i=1}^{i=n} x_i = \frac{x_1 + x_2 \cdots x_i}{n}$$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

27 • 57

# A note on the harmonic mean

● Driving at 100 km/h from A to B; distance is 100 km.

● Driving back at 50 km/h.

● What is the average speed for the round-trip?

❑ 75 km/h ❑ 70.71 km/h ☑ 66.67 km/h

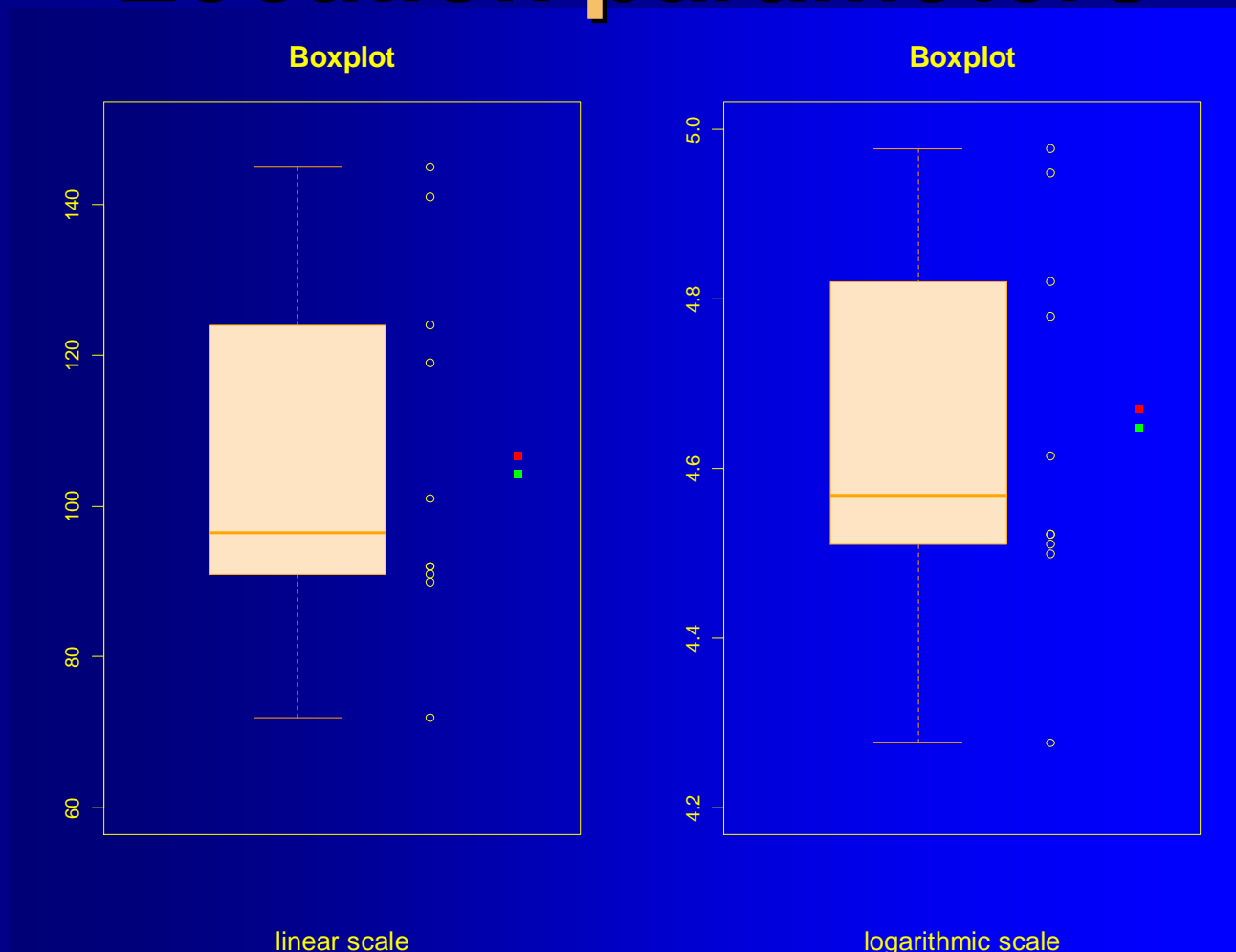● 1 h for 100 km (A→B) and 2 h for 100 km (A←B); 200 km/3 h = 66.67 km/h.

● Harmonic mean for rates!

$$\bar{x}_{harm} = \frac{2}{\frac{1}{100} + \frac{1}{50}} = \frac{2}{0.01 + 0.02} = 66.\dot{6}$$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

π Pharma Edge

*28 • 57*

# Location parameters

- Application of *any* location parameter *always* (!) implies an underlying distributional assumption.
  - Median: discrete (or unknown)
  - Arithmetic mean: normal distribution
  - Geometric mean: lognormal distribution
  - Harmonic mean: rates
- Example from above sampled from a lognormal distribution
  - Arithmetic mean: 106.7 (too high!)
  - Geometric mean: 104.3 (correct)

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

29 • 57

# Location parameters



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*30 • 57*

Pharma Edge

# Nitpicking terminology

- If we are estimating parameters of a distribution, we are using **Estimators**; *e.g.*, the arithmetic mean is the unbiased estimator of the central tendency of the normal distribution.

- The numerical outcomes (*i.e.*, values one give in the report) are **Estimates**.

- Don't write something like
    *'The point estimator was 95.34 %.'*

  … when it was actually a maximum likelihood estimator based on least squares means in log-scale ☺

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge                                                                                            *31 • 57*

# Dispersion parameters

- ordered [72,90,91,92,92,101,119,124,141,145]
- Quartiles (25%, 75%): Be cautious! Different methods implemented in software…
  - 90.00, 124.00:   SAS
  - 91.25, 122.75:   S, $R$, M$-Excel
  - 90.75, 128.25:   Minitab, SPSS, Phoenix/WinNonlin
  - 90.92, 125.42:   Hyndman & Fan (1996)
  - … and many others!

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*32 • 57*

# Dispersion parameters

- Standard deviation (SD) of arithmetic mean: 24.2400

$$SD_{arithm} = \sqrt{\frac{1}{n-1}\sum_{i=1}^{i=n}\left(x_i - \overline{x}\right)^2}$$

- SD of geometric mean: 23.8000

$$SD_{geom} = e^{\sqrt{\frac{1}{n-1}\sum_{i=1}^{i=n}\left(\ln x_i - \ln \overline{x}_{geom}\right)^2}}$$

Pharma Edge
**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**
*33 • 57*

# Dispersion parameters

● SD of harmonic mean: 22.8519

$$SD_{harm} = \sqrt{(n-1)\sum_{i=1}^{i=n}\left(\bar{\bar{H}}_i - \bar{\bar{H}}\right)^2}$$

$$\bar{\bar{H}} = \frac{1}{n}\sum_{i=1}^{i=n}\bar{\bar{H}}_i$$

$$\bar{\bar{H}}_i = \frac{n-1}{\left(\displaystyle\sum_{j=1}^{i=n}\frac{1}{x_j}\right) - \frac{1}{x_i}}$$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*34 • 57*

# Dispersion parameters

- Coefficient of Variation
  (sometimes given in percent of mean):
  $$CV\% = 100 \cdot SD/\overline{x}$$

| Population (N=10$^6$), parameters | | Sample (n=36), parameters | | | |
|---|---|---|---|---|---|
| $\mu$ | 100.00 | $\overline{x}_{arithm}$ | 106.70 | $\overline{x}_{geom}$ | 104.28 |
| $\sigma$ | 20.00 | $SD_{arithm}$ | 24.24 | $SD_{geom}$ | 23.80 |
| $CV\%$ | 20.00 | $CV\%$ | 22.72 | $CV\%$ | 22.83 |

π ε χ ε π Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

35 • 57

# A remark on Variances

● Whilst means and variances are additive, standard deviations (and CVs as well) *are not!*

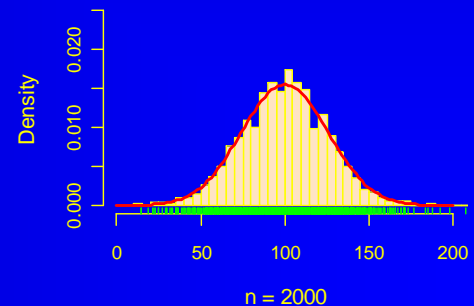| sample | mean | s | | $s^2$ | |
|--------|------|---|---|-------|---|
| 1 | 100 | 20 | 400 | | |
| 2 | 100 | 30 | 900 | | |
| 1+2 | (100+100)/2 | 25?? | $\sum s^2/2$ | $\sqrt{650}=25.5!!$ | |

**mean = 100 sd = 20 CV = 20 %**     **mean = 100 sd = 30 CV = 30 %**     **mean = 100 sd = 25.5 CV = 25.5 %**



n = 1000          n = 1000          n = 2000

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*36 • 57*

# Arithm. *vs.* geom. means



**Biostatistics: Basic concepts & applicable principles for various designs
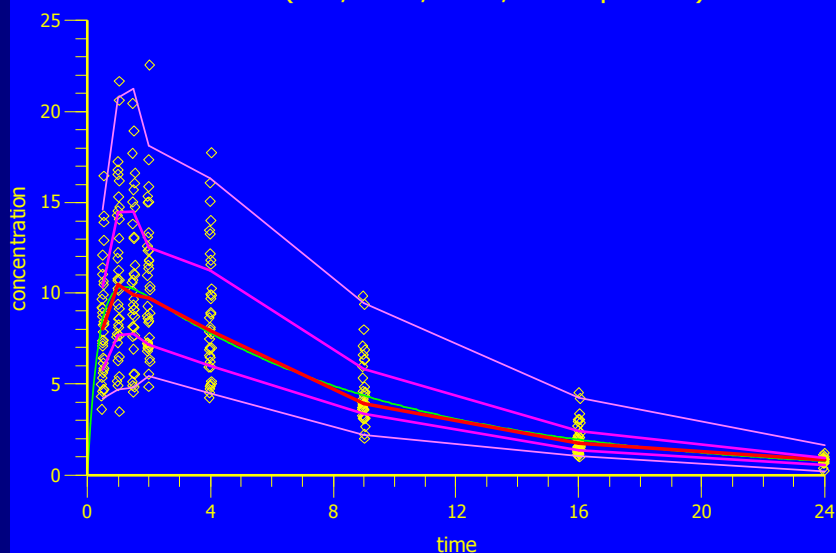in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*37 • 57*

π Pharma Edge

# Median and quantiles



Median (5%, 25%, 75%, 95% quantile)

Median (5%, 25%, 75%, 95% quantile)

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*38 • 57*

Pharma Edge
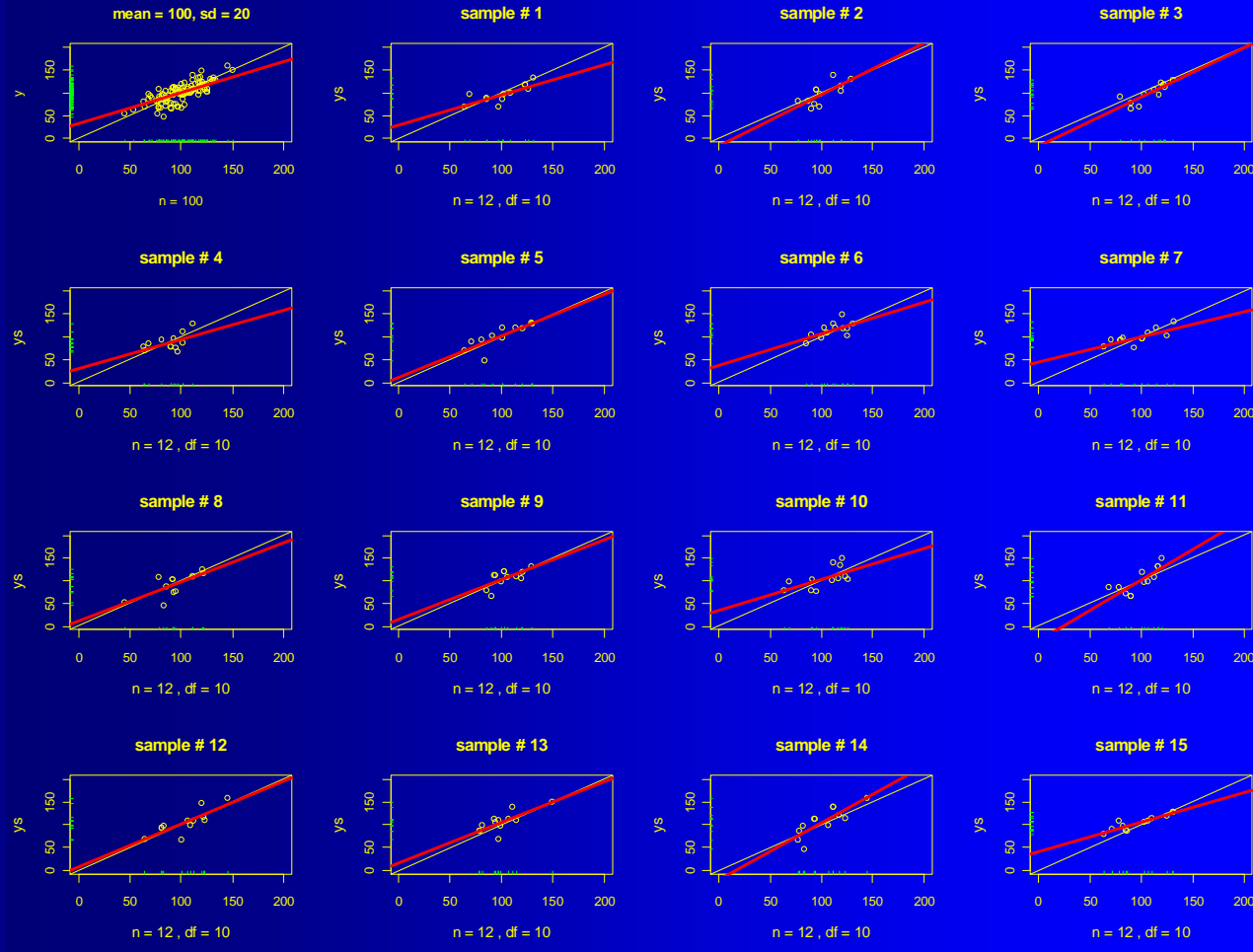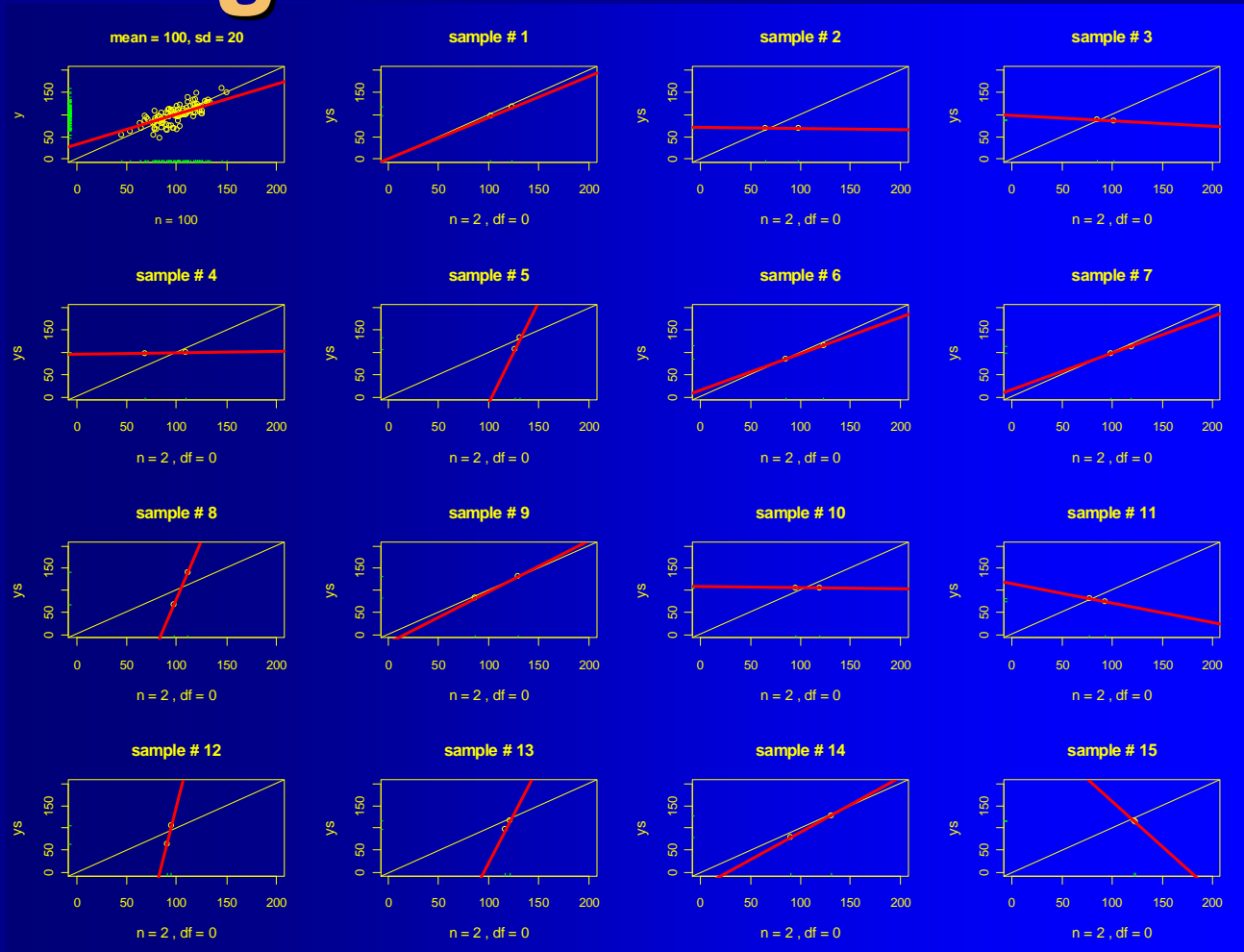
# Degrees of freedom…

- For every estimated parameter in a statistical model one degree of freedom is 'lost' from the number of samples (df = n – p).
  - Any model becomes useless if df=0, and impossible to fit if df<0 (p>n). Example:
    - Linear regression: Two parameters are fit (slope, intercept; since any line is defined by two points $(x_1/y_1|x_2,y_2)$ at least three data points are needed (df=1).

Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*39 • 57*

# Degrees of freedom…



**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*40 • 57*

Pharma Edge

# Degrees of freedom…



**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**
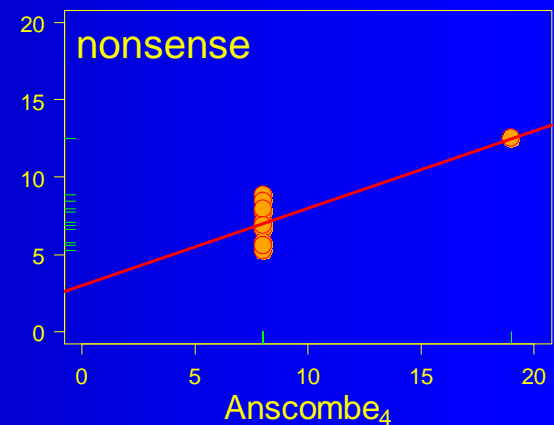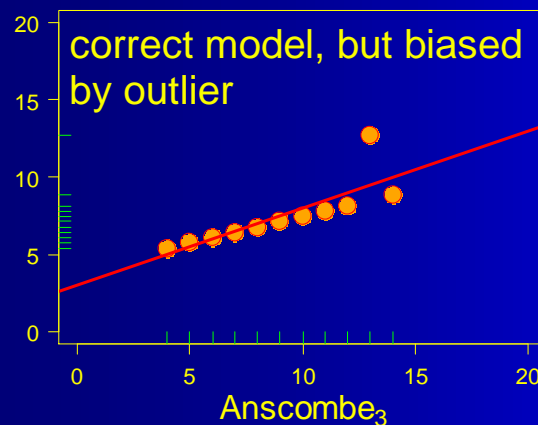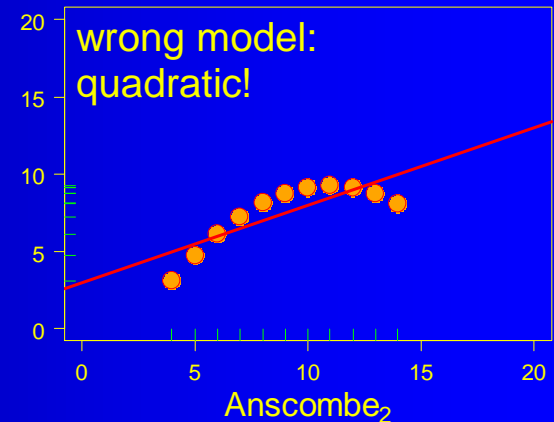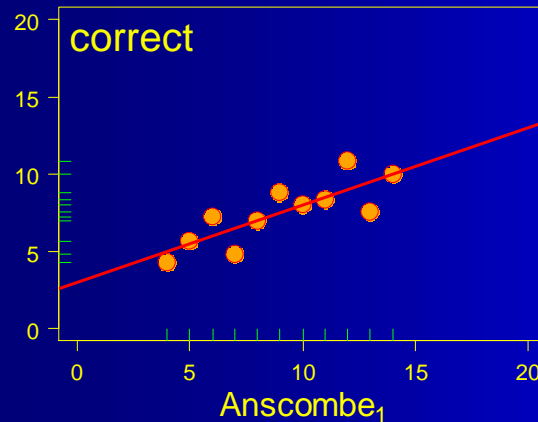
*41 • 57*

# Visualize your data!

Anscombe's
Quartet (1973)

All datasets:
$\text{mean}_x$ 9.0, $s^2_x$ 10
$\text{mean}_y$ 7.5, $s^2_y$ 3.75
$\text{Corr}_{yx}$ 0.898
$\text{Regr}_{yx}$ $y = 3 + 0.5x$

Don't rely *solely* on numerical results.



correct — Anscombe$_1$

wrong model: quadratic! — Anscombe$_2$

correct model, but biased by outlier — Anscombe$_3$

nonsense — Anscombe$_4$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*42 • 57*

Pharma Edge

# Data Transformation?



Clearly in favor of a lognormal distribution. Shapiro-Wilk test highly significant for normal distribution (rejected).

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*43 • 57*

# Data Transformation!

**MPH, 12 subjects**

Density

AUC [ng×h/mL]
Shapiro-Wilk p= 0.29668

**MPH, 12 subjects**

Density

ln(AUC [ng×h/mL])
Shapiro-Wilk p= 0.85764

**Normal Q-Q Plot**

Sample Quantiles

Theoretical Quantiles

**Normal Q-Q Plot**

Sample Quantiles

Theoretical Quantiles

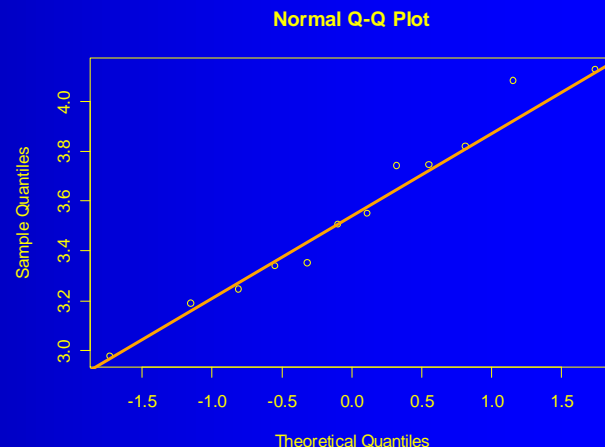Data set from a real study. Both tests *not* significant (assumed distributions not reject-ed).

Tests not acceptable according to GLs; log-transforma-tion based on prior knowledge (PK)!

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*44 • 57*

Pharma Edge

# Data Transformation

- BE testing started in the early 1980s with an acceptance range of 80% – 120% of the reference based on the normal distribution.

- Was questioned in the mid 1980s
  - Like many biological variables AUC and $C_{max}$ *do not* follow a normal distribution
    - Negative values are impossible
    - The distribution is skewed to the right
    - Might follow a lognormal distribution
  - Serial dilutions in bioanalytics lead to multiplicative errors

**Pharma Edge**

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*45 • 57*

# Data Transformation: PK

$$F_T = \frac{AUC_T \cdot \cancel{CL_T}}{\cancel{D_T}}, F_R = \frac{AUC_R \cdot \cancel{CL_R}}{\cancel{D_R}}$$

$$F_{rel}(BA) = \frac{AUC_T}{AUC_R}$$

Assumption 1: $D_1 = D_2$ $(D_1/D_2 = 1^*)$

Assumption 2: $CL_1 = CL_2$

# Data Transformation

- 'Problems' with logtransformation
  - If we transform the *'old'* acceptance limits of 80% − 120%, we get −0.2231, +0.1823.
  - These limits are *not symetrical* around 100% any more, the maximum power is obtained at $e^{0.1823–0.2231}$ = 96%…
  - Solution:
    lower limit = 1 − 0.20, upper limit = 1/lower limit
    $ln$(0.80) = −0.2231 and $ln$(1.25) = +0.2231.
    Symetrical around 0 in the log-domain and around 100% in the backtransformed domain ($e^0$=1).

Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*47 • 57*

# Data Transformation

- 'Problems' with logtransformation
  - Discussion, whether more bioinequivalent formula-
    tions will pass due to *'5% wider'* limits
    lower limit = 1 – 0.20, upper limit = 1/lower
    <span style="color:yellow">80.00% – 125.00%</span> (width <span style="color:yellow">45.00%</span>)
    instead of keeping the *'old'* width
    lower limit = 1 – 0.1802, upper limit = 1/lower
    <span style="color:yellow">81.98% – 121.98%</span> (width <span style="color:yellow">40.00%</span>)
    or even become more strict by setting
    upper limit = 1 + 0.20, lower limit = 1/upper
    <span style="color:yellow">83.33% – 120.00%</span> (width <span style="color:yellow">36.67%</span>)
    **80% – 125% was chosen for convenience (!)**

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge
*48 • 57*

# *F* Distribution

● Allows comparison of variances (depending on $\nu$) of two distributions. We will need that in ANOVA.

$$F(x \mid \nu_1, \nu_2) = \begin{cases} \nu_1^{\frac{\nu_1}{2}} \nu_2^{\frac{\nu_2}{2}} \cdot \dfrac{\Gamma\left(\dfrac{\nu_1}{2} + \dfrac{\nu_2}{2}\right)}{\Gamma\left(\dfrac{\nu_1}{2}\right)\Gamma\left(\dfrac{\nu_1}{2}\right)} \cdot \dfrac{x^{\frac{\nu_1}{2}-1}}{\left(\nu_1 x + \nu_2\right)^{\frac{\nu_1+\nu_2}{2}}} & x \geq 0 \\[2em] 0 & x < 0 \end{cases}$$

$$\Gamma(x) = \int\limits_0^{+\infty} t^{x-1} e^{-t} dt$$

● Note that if one of the degrees of freedom = 1, there is a relationshop to the *t* distribution:

$$F(\nu_1 = 1, \nu_2 = \nu) = \left(t(\nu)\right)^2$$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*49 • 57*

# Significance tests

- In statistics (as well as in science in general) it is <u>not possible to *prove* something</u>.

- We can only state a hypothesis and try to *reject* this so called null hypothesis by evaluating data from an experiment.

- Example:

  - $H_0$: $\mu_1 = \mu_2$ (no difference in means, null hypothesis)
      *vs.*

  - $H_a$: $\mu_1 \neq \mu_2$ (different means; alternative hypothesis)

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*50 • 57*

# $\alpha$- vs. $\beta$-Error

- All formal decisions are subjected to two types of error:
  - Error Type I ($\alpha$-Error, Risk Type I)
  - Error Type II ($\beta$-Error, Risk Type II)
    Example from the justice system:

| Verdict | Defendant innocent | Defendant guilty |
|---|---|---|
| Presumption of innocence not accepted (guilty) | **Error type I** | **Correct** |
| Presumption of innocence accepted (not guilty) | **Correct** | **Error type II** |

Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*51 • 57*

# $\alpha$- vs. $\beta$-Error

- … in more statistical terms:

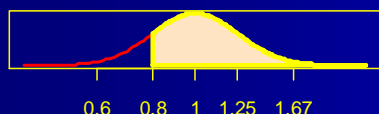| Decision | Null hypothesis true | Null hypothesis false |
|---|---|---|
| Null hypothesis rejected | Error type I | Correct ($H_a$) |
| Failed to reject null hypothesis | Correct ($H_0$) | Error type II |

- In BE-testing the null hypothesis is bioinequivalence ($\mu_1 \neq \mu_2$)!

| Decision | Null hypothesis true | Null hypothesis false |
|---|---|---|
| Null hypothesis rejected | Patients' risk | Correct (BE) |
| Failed to reject null hypothesis | Correct (not BE) | Producer's risk |

Pharma Edge
**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**
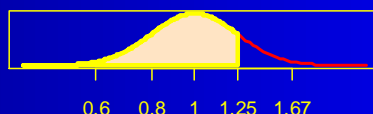*52 • 57*

# $\alpha$- vs. $\beta$-Error

- $\alpha$-Error: Patients' Risk to be treated with a bioinequivalent formulation (H$_0$ falsely rejected)
  - BA of the test compared to reference in a *particular* patient is risky *either* below 80% *or* above 125%.
  - If we keep the risk of particular patients at 0.05 (5%), the risk of the entire population of patients ($<$80% *and* $>$125%) is 2$\times\alpha$ (10%) is:
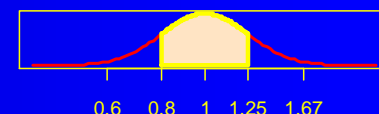  90% CI = 1 $-$ 2$\times\alpha$ = 0.90

| 95% one-sided CI | 95% one-sided CI | 90% two-sided CI = two 95% one-sided |
|:---:|:---:|:---:|
| 0.6  0.8  1  1.25  1.67 | 0.6  0.8  1  1.25  1.67 | 0.6  0.8  1  1.25  1.67 |
| particular patient | particular patient | population of patients |

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge

*53 • 57*

# $\alpha$- vs. $\beta$-Error

- $\beta$-Error: Producer's Risk to get no approval for a bioequivalent formulation (H$_0$ falsely not rejected)
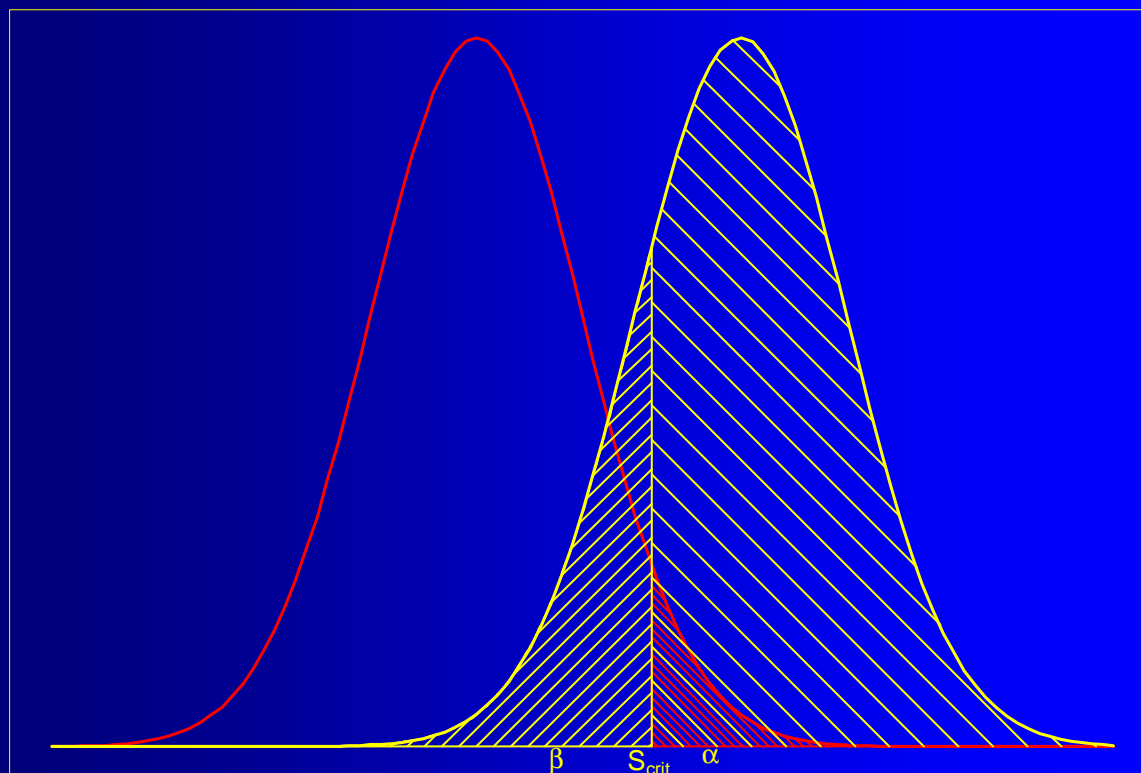  - *Set* in study planning to $\leq 0.2$, where power $= 1 - \beta = \geq 80\%$
  - If power is set to 80 %

    **One out of five studies will fail just by chance!**

| $\alpha$ 0.05 | BE |
|---|---|
| not BE | $\beta$ 0.20 |

# Significance test ($\alpha$- vs. $\beta$)



Significance test: $\alpha$, $\beta$

**Biostatistics: Basic concepts & applicable principles for various designs in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

$55 \cdot 57$

Pharma Edge

# Part I: Basic Concepts



Helmut Schütz

**BEBAC**

Consultancy Services for
Bioequivalence and Bioavailability Studies
1070 Vienna, Austria
helmut.schuetz@bebac.at

π Pharma Edge

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

*56 • 57*

# To bear in Remembrance...

In these matters the only certainty is
that nothing is certain.
*Gaius Plinius Secundus (Pliny the Elder)*

The theory of probabilities is at bottom
nothing but common sense reduced to calculus.
*Pierre-Simon Laplace*

It is a good morning exercise for a research scientist
to discard a pet hypothesis every day before
breakfast.
It keeps him young.                              *Konrad Lorenz*

**Biostatistics: Basic concepts & applicable principles for various designs
in bioequivalence studies and data analysis | Mumbai, 29 – 30 January 2011**

Pharma Edge                                                                                                 *57 • 57*