① Introduction

1) derive ^whole popula. properties from random sample ← Point / interval
   ↖ stat tests
   Anova
   Survival

2) Rules for sample determination ← study design
   Power analysis

3) Interprete test results

② Probability triple (space)

1) sample space $\Omega$ : set of outcomes

2) set of events $F$ : collection of events $E$ (each: subset of outcomes)   $F = \{\emptyset, \Omega, \{1,2\}, \Omega/\{ \}$
   all, closed

3) probability measure $P$: each event $E$ of set $F$: probability $P(E)$   $P: F \to [0,1]$

Kolmogorov axioms for $P$

   1) $0 \leq P(E) \leq 1$

   2) $P(\Omega) = 1$

   3) $P(\cup_i E_i) = \sum_i P(E_i)$ · $E_i$ disjoint

Event-Algebra $F$
   1) $\Omega \in F$
   2) $A \in F \to \bar{A} \in F$
   3) $A_1,... \in F \to \cup_n A_n \in F$
   smalles $F = (\{\emptyset, \Omega\}, \{ \}) ...$ pote

Consequences

   1) $P(A \cup B) = P(A) + P(B) - P(A \cap B)$          (sum rule)

   2) $P(\Omega \backslash A) = 1 - P(A)$          (compl. prob)

   3) $P(B|A) = \dfrac{P(B \cap A)}{P(A)}$          (cond. prob.)     ... normieren, subset selection
   then count

Terms

   i) Independence $A, B$ if $P(A \cap B) = P(A) P(B)$

   ii) Disjoint $A, B$ if $P(A \cap B) = 0$    (mutually exclusive) ≡ dependant

Bayes theorem

N disjoint events $\Omega = A_1 \cup ... \cup A_N$, $B \notin \Omega$

$P(A_i | B) = \dfrac{P(A_i) \cdot P(B|A_i)}{\sum_k P(A_k) \cdot P(B|A_k)} \equiv \dfrac{\langle P(A_i) P(B|A_i) \rangle}{}$

disease        enzyme
1/2            test
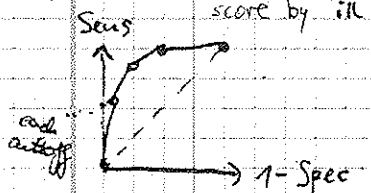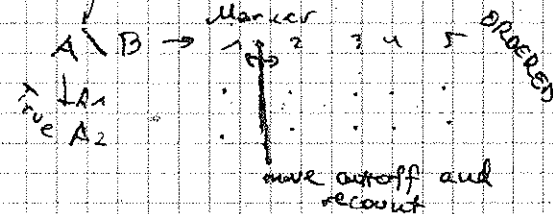
ROC curves   (receiver operating characteristics) → categorial biomarker

$Sens = \dfrac{TP}{TP + FN}$          $Spec = \dfrac{TN}{TN + FP}$

$P(B = pos | A_1)$                        $P(B = neg | A_2)$
score by ill                              score free by healthy



Sens
↑ → 1-Spec

$A = 0,5 \to$ no correl
$A \to 1 \to$ strong pos corr
$A \to 0 \to$ "  neg corr

Marker
$A \backslash B \to$ 1 2 3 4 5   ORDERED
true $A_1$
$A_2$
move cutoff and recount

Random variable

• function assigning real numbers to results of experi.

• Elements of $\Omega \longrightarrow$ Real numbers $\mathbb{R}$
   Random variable $X \in (d_1 + d_2)$

discrete: finite / countably infinite   $P(X \leq x)$
continuous: $P(X \leq x)$, $\forall x \in \mathbb{R}$ repr. by int. dens. func. $f \geq 0$, $P(X \leq x) = \int_{-\infty}^{x} f(x') dx'$
→ Probability mass (density) function; cumulative distribution function
   disc.    cont.

## Binomial distribution

$$f(k; n,p) = P(X=k) = \binom{n}{k} p^k (1-p)^k \qquad \binom{n}{k} = \frac{n!}{k!\,(n-k)!}$$

Prob. of $k$ successes in $n$ trials, prob $p$ success

$$E(x) = np \qquad\qquad = E(X) = \sum k \cdot f(k)$$

$$V(x) = np(1-p) = E((X-E(x))^2) = \sum f(k)(k-np)^2$$

## Poisson distribution

$$P(X=k) = \frac{\mu^k e^{-\mu}}{k!} \qquad\qquad k=0,1,2,\dots \qquad \mu: \text{expected nr of occurr. in interval}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad k: \text{nr of actual occurr of event}$$

$$E(X) = \mu$$

$$V(X) = \mu \qquad\qquad\qquad\qquad\qquad\qquad \text{FIXED INTERVAL}$$

## Gaussian distribution

$$f(x) = \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad\qquad \mu: \text{mean value}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \sigma: \text{standard deviat}$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \text{symmetric}$$

$$E(x) = \mu = \int x f(x)\, dx \qquad\qquad \mu \pm \sigma : 68\%$$

$$V(X) = \sigma^2 = \int (x-\mu)^2 f(x)\, dx \qquad \mu \pm 2\sigma : 95\%$$

$$\rightarrow \text{pattern when} \begin{cases} \text{large number} \\ \text{independent} \end{cases} \text{events}$$
$$\qquad\qquad\qquad\qquad \text{random} \\ \qquad\qquad\qquad\qquad \text{small effect}$$
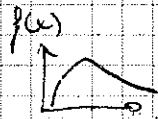
### Standard normal distribution

$$\mu = 0, \ \sigma = 1$$

$$f(x_{st}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x_{st}^2}{2}}$$

## Logarithmic Gaussian distribution

$$y = \log x \rightarrow x = e^y \qquad f(x)$$

$$N(\mu,\sigma^2) \ \text{log-normal}$$



## Uniform distribution

$$f(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases} \qquad\qquad E(x) = \frac{b+a}{2}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad V(x) = \frac{1}{12}(b-a)^2$$

## Exponential distribution

$$f(t) = \lambda e^{-\lambda t} \qquad\qquad \rightarrow \text{INTERVAL BETWEEN TWO EVENTS}$$

$$E(t) = \frac{1}{\lambda}$$

$$V(t) = \frac{1}{\lambda^2}$$

## $\chi^2$ distribution

$$X_1, \dots X_n \sim N(0,1) \qquad\qquad n: \text{degrees of freedom}$$

$$X_1^2 + \dots + X_n^2 \sim \chi_n^2 \qquad\qquad \rightarrow n \gg 1 \Rightarrow N(n,\sqrt{2n})$$

$$E(x) = n \qquad\qquad \Gamma\left(\frac{n}{2}, 2\right)$$

$$V(x) = 2n$$

## F - distribution (Fischer) homogeneity test

$$X \sim \chi_m^2, \ Y \sim \chi_n^2$$

$$F_{mn} = \frac{X/m}{Y/n} \qquad ; \qquad \frac{1}{F_{mn}} = F_{nm}$$

$$E(F) = \frac{n}{n-2} \qquad n > 2$$

$$V(F) = \frac{2n^2(m+n-2)}{m(n-2)^2(n-4)} \qquad n > 4$$

## Student's or t- distribution

$$f_n(x) = C_n \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}} \underset{n \gg 1}{\approx} \mathcal{N}(0,1) \qquad n: \text{degrees of freedom}$$

$$E(X) = \begin{cases} 0 & n > 1 \\ / & \text{oth} \end{cases}$$

$$V(X) = \begin{cases} \frac{n}{n-2} & n > 2 \\ \infty & 1 < n \le 2 \\ / & \text{oth} \end{cases}$$

$$X \sim N(0,1) \\ Q \sim \chi_n^2 \qquad : \quad T = \frac{X}{\sqrt{Q/n}} \sim t$$

## Moments

$$\mu'_{n(c)} = \int_{-\infty}^{+\infty} (x-c)^n f(x)\,dx \qquad // \text{ standardized: divide } \sigma: \left(\frac{x-c}{\sigma}\right)^n$$

- Mean $\mu = \mu'_1(0) = \int_{-\infty}^{+\infty} x\, f(x)\, dx$

- Variance $\sigma^2 = \mu'_2(\mu) = \int_{-\infty}^{+\infty} (x-\mu)^2 f(x)\, dx$

- Skewness $\gamma_1 = \mu'_3(\mu)/\sigma^3 = \int \left(\frac{x-\mu}{\sigma}\right)^3 f(x)\, dx \quad = \quad \frac{\mu_3(\mu)}{\sigma^3}$

- Kurtosis $\gamma_2 = \mu'_4(\mu)/\sigma^4 = \int \left(\frac{x-\mu}{\sigma}\right)^4 f(x)\, dx$

## Levels measure

| | | |
|---|---|---|
| Nominal | – label | $= \ne$ |
| Ordinal | – ordered | $< >$ |
| Interval | – difference | $+ \ -$ |
| Ratio | – zero | $* \ /$ |

## ③ Fundamentals

- Arithmetic mean $\quad \bar{x} = \frac{1}{n} \sum_i x_i \qquad \rightarrow$ minimizes $(x - \bar{x})^2$
  - ↳ linearity: $y_i = c_1 x_i + c_2 \rightarrow \bar{y} = c_1 \bar{x} + c_2$

- Geometric mean: $\bar{x}_g = \sqrt[n]{\prod x_i} \qquad \cdots$ growth processes

  $$\log(\bar{x}_g) = \frac{1}{n}(\log x_1 + \dots + \log x_n)$$

- Median of ordered sample $(n)$
  $$\tilde{x} = \begin{cases} x_{(n+1)/2} & n \text{ odd} \\ \frac{1}{2}(x_{n/2} + x_{n/2+1}) & n \text{ even} \end{cases}$$
  $\rightarrow$ robust against outliers
  $\rightarrow$ minimizes $|x - \tilde{x}|$

  $\bar{x} = \tilde{x} \qquad$ symmetric
  $\bar{x} > \tilde{x} \qquad$ pos. skewed $\quad \searrow \qquad \cdots$ paar verd. viel
  $\bar{x} < \tilde{x} \qquad$ neg. skewed $\quad \diagup$

- Mode: most occurring value $\quad \rightarrow$ not useful

**Spread**

● **Range**

$$r = x_{max} - x_{min}$$

→ sensitive outliers
→ depends on $n$

● **Percentiles** $p$ (sample $n$)  / Quantile

$$k = np/100$$

$$\begin{cases} V_p = X_{\text{ceil}(k)} & \text{if } k \text{ not integer} \\ V_p = \frac{1}{2}(X_k + X_{k+1}) & \text{if } k \text{ is integer} \end{cases}$$

$p$-th percentile is a value $V_p$ such that (at least) $p\%$ sample points $\leq V_p$

$$V_p = \Phi^{-1}(p)$$

$V_{25}, V_{50}, V_{75}$ → 1st, 2nd, 3rd quartile

Quantil distance $QD = V_{1-p} - V_p$

● **Variance**

$$s^2 = \frac{\sum_1 (x_i - \bar{x})^2}{n-1} \qquad \text{sample variance} = \frac{1}{n-1}(\sum x_i^2 - n\bar{x}^2)$$

$n-1$ → converge to true value ($n$ fact)

$$s = \sqrt{s^2} \qquad \text{sample standard devia}$$

(i) $y_i = x_i + c \longrightarrow s_x^2 = s_y^2$    $V(x) = E(x^2) - (EX)^2$

(ii) $y_i = cx_i \longrightarrow s_y^2 = c^2 s_x^2 \; ; \; s_y = c s_x$

● **Coefficient of variation**

$$C_v = \frac{s}{\bar{x}} \cdot 100\%$$

**Sample means**

assuming equal size, otherwise weighting

$$\overline{\bar{x}} = \frac{1}{m} \sum_i \bar{x}_i \qquad ; \qquad \overline{\bar{x}} = \overline{x}_{glob} \qquad \text{mean of sample means}$$

$$s_{\bar{x}}^2 = \frac{s^2}{n} \qquad \text{variance of sample means}$$

$$s_{\bar{x}} = \sqrt{s_{\bar{x}}^2} = \frac{s}{\sqrt{n}} \qquad \text{standard error of the mean} \\ \text{(standard devia of sample means)}$$

**Covariance**

$$Cov(X_1, X_2) = E[(X_1 - EX_1)(X_2 - EX_2)] = \frac{1}{n-1} \sum_i (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2)$$
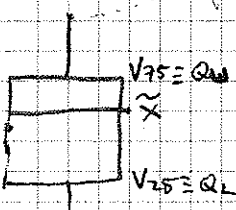
$$\rho(X_1, X_2) = Cov(X_1, X_2)/(\sigma_1 \sigma_2) \qquad \text{correlation coefficient}$$

If $\rho(X_1, X_2) = 0$ → $X_1$ and $X_2$ uncorrelated (independent)

**Box plot**  largest not outlying value

● learn width, skewness, median



$V_{75} = Q_u$
$\tilde{x}$
$V_{25} = Q_L$

(1) $Q_u - \tilde{x} \approx \tilde{x} - Q_L$ symm
(2) $>$  pos. skewed (right to the)
(3) $<$  negat. " (left)

● Outlier $x \gtrless Q_u + 1.5(Q_u - Q_L)$
    $\lessgtr Q_L$

● Extreme at $\gtrless \quad + 3$

○ outlier
+ extreme outlier

## Point - interval estimation

- sample of population unknown pdf → measure location ($\mu$) and spread ($\sigma$)

$\bar{X} \approx \mu$    (arith mean)

$s \approx \sigma$    (sample var)

- $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$     Distribution of arithmetic means

$$\Delta\mu = \frac{\sigma}{\sqrt{n}} \;\; ; \;\; \Delta\bar{x} = \frac{s}{\sqrt{n}}$$

(1) $X \sim N(\mu, \sigma^2)$     (4) $\frac{(\bar{X}-\mu)}{\sigma/\sqrt{n}} \sim N(0,1)$     (5) $\frac{\bar{X}-\mu}{s}\sqrt{n} \sim t_{n-1}$

(2) $\bar{X} \sim N(\mu, \sigma^2/n)$

(3) $X s \sigma = \frac{1}{\sigma}(x-\mu)$       $\sigma \Leftrightarrow s$     Student distr.

                                  Normal    t - dist      $n-1$ d.o.f.

## Central limit theorem

$x_1, \ldots, x_n$ ; population $\mu, \sigma^2$

$X \not\sim N(\mu, \sigma^2)$ not normal  (single not, but means are)

but for large $n$ ($>20$) the mean $\bar{X} \sim N(\mu, \sigma^2/n)$

### Interval of mean :

$Z = \frac{\bar{x}-\mu}{\sigma}\sqrt{n} \sim N(0,1)$    → $-1.96 < Z < 1.96 \Rightarrow 95\%$

                                                 $z_{97.5}$

$\mu \in (\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}})$ with $95\%$ probability, from repeated samples of size $n$

           $\underbrace{\quad}_{conf. interval}$                 $-1.96 < Z < 1.96$

$t = \frac{\bar{x}-\mu}{s}\sqrt{n} \sim t_{n-1}$

$\bar{\mu} \in (\bar{x} \pm t_{n-1, 97.5\%} \cdot \frac{s}{\sqrt{n}})$   w $95\%$ prob.

      $\underbrace{\quad}_{conf. interval}$      2.78 ($n=5$)

                       2 ($n=60$)        $-2 < t < 2$

                       1.96 ($n > 60$)

→ $P(t_{n-1, \frac{\alpha}{2}} < T < t_{n-1, 1-\frac{\alpha}{2}}) = 1-\alpha$

$1-\alpha$ : confidence level ~~error level~~

$\alpha$ : error probability    (that $\bar{x}$ outside $\mu$ interval)

- Higher $n$ : smaller (sharper) conf. interval
- $(1-\alpha) \cdot 100\%$ of all conf. int. will include true (unknown) mean

Only $\alpha \cdot 100\%$                       not    $\mu$        (error) !

- When $n > 200$

$$\mu \in (\bar{x} \pm Z_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}})$$

### For $\sigma^2$

If $X_i \sim N(\mu, \sigma^2)$ → $\frac{n-1}{\sigma^2} s^2 \sim \chi^2_{n-1}$

$P(\chi^2_{n-1, \frac{\alpha}{2}} < \frac{(n-1)s^2}{\sigma^2} < \chi^2_{n-1, 1-\frac{\alpha}{2}}) = 1-\alpha$

CI : $(n-1)s^2/\chi^2 < \sigma^2 < (n-1)s^2/\chi^2$

# Statistical tests

- Test hypothesis, unique decision making criterion
- $H_0$ : null hypothesis (of no effect), what you want to reject
- $H_1$: alternative hypoth.
  - × equality distribution parameters $(\mu, \sigma)$
  - × " pdf
  - × correld random variables

(0) $H_0, H_1$
(i) Test variable
(ii) Significance level / error probability $\alpha$ (5%)
(iii) Critical region $B$ such that $P(T \in B \mid H_0 \text{ true}) \leq \alpha$

## $t$ - test

- Comparison 2 mean values  • samples NOT < multimodal / too much skewed
- Normally distributed $X, Y$ , not very sens. to deviat
- $\sigma_X = \sigma_Y$ (reasonably) → F-test

ONE SAMPLE ONE-SIDED  [one arm vs popula]

- $H_0 : \mu = \mu_0$
- $H_1 : \mu < \mu_0$
- $t = \dfrac{\bar{x} - \mu_0}{s/\sqrt{n}}$

$H_1 : \mu > \mu_0$

$\equiv \dfrac{\text{Parameter estim - true value}}{\text{Param. error}}$   $\frac{s}{\sqrt{n}} \equiv \Delta \bar{x}$

→ deviat in units of std dev

- $\alpha = 0,05$
- Decision:   $t \gtrless t_{n-1, \alpha}$ : accept $H_0$        $t \lessgtr t_{n-1, \alpha}$
  $t < t_{n-1, \alpha}$ : reject $H_0$ w.s.l      $t > t_{n-1, \alpha}$

Due to $t \leq -$, $H_0$ is reject w $\alpha$% level, i.e. —

ONE SAMPLE TWO-SIDED (more conservative)
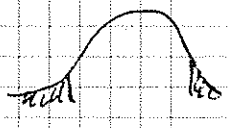
→ $t_{crit}$ larger

- $H_0 : \mu = \mu_0$
- $H_1 : \mu \neq \mu_0$
- $t = \dfrac{\bar{x} - \mu_0}{s/\sqrt{n}}$
- $\alpha = 0,05$
- $|t| \leq t_{n-1, 1-\frac{\alpha}{2}}$ : accept $H_0$ / cannot be rejected
  $|t| > t_{n-1, 1-\frac{\alpha}{2}}$ : reject $H_0$ w.s.l

TWO SAMPLES TWO-SIDED (indep.)      [double arm study]

- $H_0 : \mu_X = \mu_Y$     $\bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \frac{\sigma^2}{n_X} + \frac{\sigma^2}{n_Y}\right) \xrightarrow{\text{samp.} \ \sigma \to s}$ t distributed
- $H_1 : \mu_X \neq \mu_Y$
- $t = \dfrac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n_X} + \frac{1}{n_Y}}}$ , $s = \sqrt{\dfrac{(n_X - 1) s_X^2 + (n_Y - 1) s_Y^2}{n_X + n_Y - 2}}$    pooled empirical standard deviation
- $\alpha = 5\%$         $|t| \leq t_{K, 1-\frac{\alpha}{2}}$   accept $H_0$
- $K = n_X + n_Y - 2$     $|t| > t_{K, 1-\frac{\alpha}{2}}$   reject $H_0$

- two samples on same patient (paired), Individuum two diff states
  - better than two arms where $\bar{x} \neq \bar{y}$ diff pat.
  - e.g. diuretic/placebo ; before/after treat$\varnothing$

- Two states:
  $\{x_1, \ldots x_n\}$
  $\{y_1, \ldots y_n\}$

- $d_i = y_i - x_i$
  $\{d_1, \ldots d_n\}$

- Mean of diff.
  $$\bar{d} = \frac{1}{n} \sum_i d_i$$

- Sample st. dev of diff
  $$s_d = \sqrt{\frac{1}{n-1} \sum_i (d_i - \bar{d})^2}$$

Test
$\Longrightarrow$
$\mu_x \neq \mu_y$
$\mu = \mu_x - \mu_y$

- $H_0 : \mu = 0$
- $H_1 : \mu \neq 0$
- $\alpha = 5\%$
- $$t = \frac{\bar{d}}{s_d / \sqrt{n}}$$
- $K = n - 1$
- $|t| \leq t_{K, 1-\frac{\alpha}{2}}$ accept $H_0$
  $|t| > t_{K, 1-\frac{\alpha}{2}}$ reject $H_0$

## $t$ - test $\quad \sigma_x \neq \sigma_y$

$X \sim \mathcal{N}(\mu_x, \sigma_x^2)$
$Y \sim \mathcal{N}(\mu_y, \sigma_y^2)$

- $H_0 : \mu_x = \mu_y$
- $H_1 : \mu_x \neq \mu_y$

$\bar{x} - \bar{y} \sim \mathcal{N}(0, \frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y})$

- $t \approx \dfrac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}}$ (not $t$-distrib except if $n$ large)

- $d' = \left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}\right)^2 \Big/ \left[\frac{\left(s_x^2 / n_x\right)^2}{n_x - 1} + \frac{\left(\frac{s_y^2}{n_y}\right)^2}{n_y - 1}\right]$

- $d'' = \text{floor}(d')$

- $|t| > t_{d'', 1-\alpha/2}$ : reject $H_0$

## Nonparametric methods
$t$ - test not applicable if
  - not normal distri
    $(\mu, \sigma)$ not enough
  - $\sigma_x \neq \sigma_y$
  - central limit theorem not applic.

Assess choice based on:

- F - test eq. of variances
- $\bar{x}, \tilde{x}, Q_v, Q_2 \rightarrow$ box plot
- histogram

Opinions
(1) Parametric if no evidence of no-normal
  - more powerful than non-parametric
  - use nonp. only if positive evidence no-normal

(2) Nonparametric always, except pos. evidence that par. are applic.
  - 95% of power of param.
  - should assume as little as possible on data

# RANK TESTS

- group/rank → ordered serial
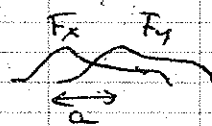- based on relative sizes of observations
- not on the size itself

Parametric                  Nonparametric

1) Two samples indep. t-test ⟶ Mann-Whitney rank sum test

2) Two samples paired t-test ⟶ Wilcoxon signed-rank test

## Conditions for applicability

- 2 indep. random variables $X, Y$ distributed $F_X, F_Y$
- $F_Y(x) = F_X(x-a)$ ⟶ differ only by a shift         $F_X$    $F_Y$
  ↳ $\sigma_x = \sigma_y$
- Independent samples $X_1, ... X_m$; $Y_1, ... Y_n$
- $H_0: a = 0$ ; $H_1: a \neq 0$

## Mann-Whitney rank sum test    $H_0$: same popula, drug no effect

(1) Combine data both samples

(2) Order values lowest to highest    Sampl   A    B     A     B

(3) Assign ranks to indiv. values     $x_i$ | 100   101    103    105

                                        $r$ | 1     2     3     4

(4) Group same value ⟶ $r' = r + \frac{1+g}{2}$    $(X_{r+1}, ..., X_{r+g})$

(5) Compute rank sum $R_1$ for first sample (lowest $n$)

(6) Calculate critical value for $T(R_1)$

    (i) $T = R_1$ → histogram of all rank sum possibilities    $n_1, n_2 < 5$

    (ii) $T = R_1$ → lookup table: if $T_{crit}^l < T < T_{crit}^h$ :    $n_1, n_2 < 10$
                         ↳ accept $H_0$

    (iii) $T(R_1) \sim N(0,1)$ if $T \leq z_{1-\frac{\alpha}{2}}$   ↗    $n_1, n_2 > 10$

## Wilcoxon signed-rank test    <sup>extra</sup> condit: symmetric (often neglected)
                                    $H_0: d = 0$, $H_1: d \neq 0$ ; $\hat{a} = \bar{x} - \bar{y}$

(1) $d_i = x_i - y_i$ ; arrange order abs. values

(2) Ignore $d_i = 0$ ; rank rest 1 to $n$ with sign!

(3) Group same abs. value $r' = r + \frac{1+g}{2}$   $(|d_{r+1}| ... |d_{r+g}|)$

(4) Calculate signed rank sum $W$

(5) Compute critical value $W_{crit}$

    (i) histogram of all possib. + − to $n$ differences          $n \leq 6$

    (ii) lookup table   $|W| \leq W_{crit}$ → accept $H_0$         $n < 16$

    (iii) $T(W) \sim N(0,1)$   $T \leq z_{1+\frac{\alpha}{2}}$   ↗

# p - Value

- significance level $\alpha$ at which $t = t_{crit,\alpha}$
- $p = \alpha_{crit}$ ($t = t_{crit,\alpha}$) $= P(t_{n-1} \leq t)$

$\quad p < 0.001$    very highly significant

$0.001 \leq p < 0.01$    highly significant

$0.01 \leq p < 0.05$       significant

$0.05 \leq p < 0.1$   trend towards significance

$p > 0.1$     not statistically significant

- Two methods for determining stat. significance

1) Critical value method $t \Leftrightarrow t_{crit}$ : rejection/acceptance

2) p-value method $p$ (exact) $< 0.05$ reject; but gives more info

$\quad p = 2 \times (1 - \Phi(t))$

# ⑥ Analysis of variance

- Means of more than two groups have to be compared
  ↳ pair-wise two-samples cannot be applied directly!

⇒ ANOVA method

## $\chi^2$ test one sample

- $X \sim N(\mu, \sigma^2)$ ; empir. var $s^2$; compare with known $\sigma_0^2$
- $X_0^2 = \dfrac{(n-1)s^2}{\sigma_0^2}$ ⟶ $\chi^2_{n-1}$ distributed, $n-1$ d.o.f.

- $H_0 : \sigma = \sigma_0$ ; $H_1 : \sigma \neq \sigma_0$, $\alpha$
- If $X_0^2 > \chi^2_{n-1, 1-\frac{\alpha}{2}}$ or $X_0^2 < \chi^2_{n-1, \frac{\alpha}{2}}$ reject $H_0$

## F - test

- tests if $\sigma_x = \sigma_y$ (precondition for t-test)
- Two samples $\{X\}$ and $\{Y\}$, independent, $\{n_x, n_y\}$, normally distributed

- $F = \begin{cases} \dfrac{s_x^2}{s_y^2} & \text{for } s_x > s_y \\[2mm] \dfrac{s_y^2}{s_x^2} & \text{for } s_y > s_x \end{cases}$ follows F distribution with $\begin{cases} m_1 = n_x - 1, \; m_2 = n_y - 1 \\[2mm] m_1 = n_y - 1, \; m_2 = n_x - 1 \end{cases}$ d.o.f

$\qquad \dfrac{s_{max}^2}{s_{min}^2} \quad \begin{array}{l} m_1 = n_{max} - 1 \\ m_2 = n_{min} - 1 \end{array}$

- $H_0 : \sigma_x = \sigma_y$, $\alpha$
- $H_1$ one-sided : $H_1 : \sigma_x < \sigma_y$ or $H_1 : \sigma_x > \sigma_y$ → reject $H_0$ if $F > F_{1-\alpha, m_1, m_2}$
- $H_1$ two-sided : $H_1 : \sigma_x \neq \sigma_y$       $F > F_{1-\frac{\alpha}{2}, m_1, m_2}$

## The one way ANOVA

- Comparison of means of arbitrary number of groups
- Each group $N(?, \sigma^2)$, same $\sigma$
- Determine whether variability dominated by $\Big\langle$ spread within groups / spread between groups $(\mu_1, \mu_2)$

- $k$ groups, $n_i$, $s_i$, $n = \sum n_i$

- $s^2_{within} = \frac{1}{n-k} \sum_{i=1}^{k} (n_i - 1) s_i^2$     pooled variance

- $s^2_{between} = \frac{1}{k-1} \left[ \sum_{i=1}^{k} n_i \bar{x}_i^2 - \frac{1}{n} \left( \sum_{i=1}^{k} n_i \bar{x}_i \right)^2 \right] = n \, s_{\bar{x}}^2$

  $\hookrightarrow = \frac{1}{n} \sum_{i=1}^{k} n_i (\bar{x}_i - \bar{\bar{x}})^2$

## Recipe

(1) $H_0$: all groups have same mean ($s_{between} = s_{within}$)

  $H_1$: at least one group diff mean

$\to s_{between} \overset{?}{>} s_{within}$

$\Rightarrow F > 1$

(2) $F = \dfrac{s^2_{between}}{s^2_{within}}$

(3) If $F \leq F_{k-1, n-k, 1-\alpha}$ accept $H_0$

  $F >$                       reject $H_0$

(4) $p$-value

rejection region
p-value
$F_{k-1, n-k, 1-\alpha}$

## If difference, compare specific groups

- Two specific (of $k$) groups
- $H_0: \bar{x}_1 = \bar{x}_2$ ;   $H_1: \bar{x}_1 \neq \bar{x}_2$ ; $\alpha$

(1) $s^2 = s^2_{within}$

(2) $t = \dfrac{\bar{x}_1 - \bar{x}_2}{\sqrt{s^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$

(3) If $t > t_{n-k, 1-\frac{\alpha}{2}}$ or $t < t_{n-k, \frac{\alpha}{2}}$ reject $H_0$

  otherwise accept

## Multiple comparisons: Bonferroni approach

- ensure overall prob. any signif. differences all possible group not $> \alpha^*$

  $\alpha^* = \dfrac{\alpha}{\binom{k}{2}}$   $\to$ test at $\alpha^*$ instead $\alpha$ $\to$ more demanding

  $\alpha^* < \alpha$

- $p$ value stays the same

[1] ANOVA showed at least one group diff $\mu$

[2] Bonferroni mult comparisons specific groups

## Bonferroni - Holm

Groups: $p_1$      $p_2$      $p_3$      $\to$ exit after 1st non-significant

$\tilde{p}_1 = 4 \cdot p_1$   $\tilde{p}_2 = 3 \cdot p_2$   $\tilde{p}_3 = 2 \cdot p_3$

$\therefore$ tests

# Survival analysis

- final state (result) of some individuals is unknown ← *dif start time*
  (binary)   *patient not follow*   *#ody dura*
  ↳ clinical studies / survival ← 5y surr
    5y local tm control
    5y progr - free

- Censored observation → no follow-up / study end → at least survived X y
  ↳ X +
  ↳ affects only di
- Endpoints types — death
    — disease recurrence
    — explosion

  $S(t)$



## Survival function / probability $S(t)$

$$S(t) = \frac{Nr.\ indiv.\ surviving > t}{Total\ nr.\ indiv.} = P(T > t)$$   $T$: time until death

$t_{50}$ : median survival time $= S^{-1}(0.5) \approx V_{50}$

$\hat{S}(t)$ of sample → observe until all die

## Kaplan - Meier estimator

(1) $n_i$ probands being observable at beginning of time interval $i$ → $[t_{i-1}, t_i]$; $t_0 = 0$

(2) $d_i$ individuals die ; $l_i$ censored at end of interval $i = t_i$
  ↳ $n_{i+1} = n_i - d_i - l_i$ at beginning of time int. $i+1$

$$\hat{S}(t_i) = \prod_{j=1}^{i} \left(1 - \frac{d_j}{n_j}\right)$$   → only for $t_i$ where death, not censoring
  → do not count $l_i$ as $d_i$

## Greenwood formula

$$S_{\hat{S}(t_i)} = \hat{S}(t_i) \sqrt{\sum_{j=1}^{i} \frac{d_j}{n_j(n_j - d_j)}}$$   → all previous time intervals
  → standard deviat of $\hat{S}$

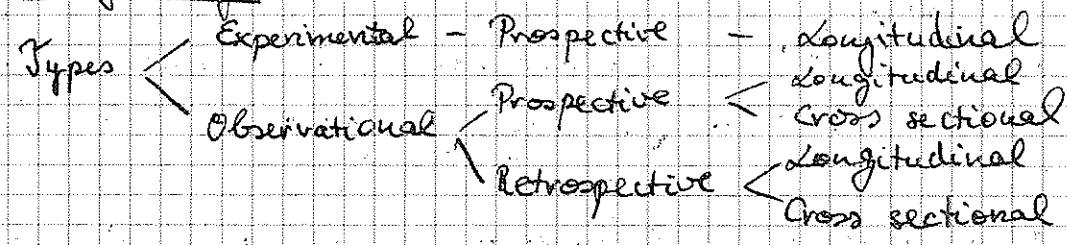Confidence interval $\hat{S}(t_i) \pm z_\alpha \cdot S_{\hat{S}(t_i)}$   $z_\alpha$ : $N(0,1)$
  $+$ truncate $[0,1]$   two - sided
  $z_{5\%} → 1.96$



## Log - rank test — Comparison of survival curves (two samples)

- nonparametric test, combine all, and suppress censored but reduce $n_{i+1}$
  times two cohorts
- $t_i$ : times of death end int. , censored not shown
- $d_i$ : deaths end of interval $i$
- $n_i$ : persons alive & observ. at begin t int. $i$
- $f_i = \frac{d_i^{(tot)}}{n_i^{(tot)}}$ : total (conditional) prob. to die at $t_i$ at combined gr
- $e_i = n_i^{(a)} \cdot f_i$ : expected death group (a) → to be compared w. combined (tot)
  from meas. combined
- $U_{L,i} = d_i^{(a)} - e_i^{(a)}$ : diff exp. - meas. (a)
- $S_i(u_L)^2$ : contrib. to empirical deviat of $U_L = \frac{n_i^{(b)} n_i^{(a)} d_i^{(tot)}(n_i^{(tot)} - d_i^{(tot)})}{n_i^{(tot)2}(n_i^{(tot)} - 1)}$
- $U_L = \sum U_{L,i}$ ; $S_{u_L}^2 = \sum S_{u_L,i}^2$
- $Z = \frac{|U_L| - \frac{1}{2}}{S_{u_L}} \sim N(0,1)$ distributed
- $H_0$: survival curves same, $H_1$: diff ; $\alpha$
- $Z > z_{crit} →$ reject $H_0$

× only two groups
× cannot test if other factors like age have influence

# Study design

Types
- Experimental — Prospective — Longitudinal
- Observational
  - Prospective
    - Longitudinal
    - Cross sectional
  - Retrospective
    - Longitudinal
    - Cross sectional

- **Observational** : DACQ without intervention        (winter study e.g.)
- **Experimental** :  (1) set hypothesis
  - (2) define intervention        (regular sauna)
  - (3) measure effect

- **Prospective** : intervention and later DACQ (exp. alw prosp)
  - (1) sauna prev. flu
  - (2) establish two groups < no sauna / 1 /week    + select sample (age, gender) problem
  - (3) diagnosis flu for each group

- **Retrospective** : data related to past
  - Interview flu / sauna
  - ⊖ Large samples
  - ⊖ Degree of truth?

- **Longitudinal** : × consecutive intervention and observation of events
  - × multiple observations in time    e.g. 5 years survival prosp. sauna study

- **Cross sectional** : single data taking in sample   e.g. screening, surveys retrosp. sauna study
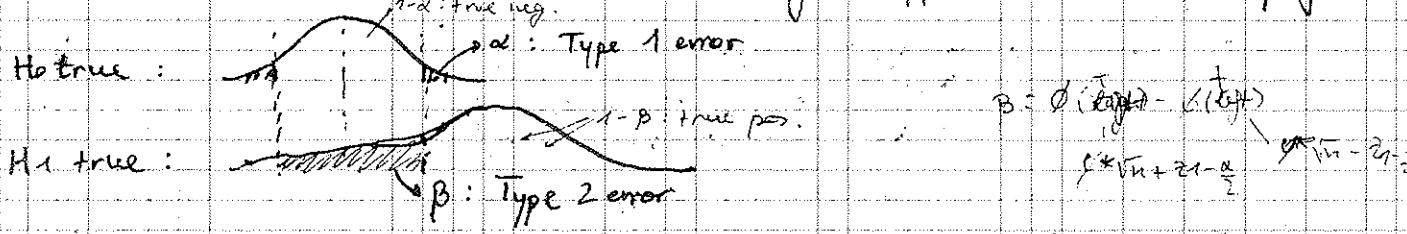
# Randomisation

- For prosp. exp. studies w. alternative interventions < phot ther. / prot ther.
  - → random selection patients two arms
  - → eliminate subject bias : blinded study
  - → " inv + " " : double blinded study
  - → equal dist of known/unknown bias factors on both arms

- **Sample randomisd** : select patients via random numbers
  - └ bad : two arms diff size

- **Block randomisd** : — equal distrib. of subjects onto the study arms
  - A A B B    — define block size : 4
  - A B A B    — both arms equal weight each block
  - → r number 1-6 : block allocation sequence

- **Stratified randomisd** : — balance of imp. features (age) in each arm
  - — block randomisd for each stratum

- → 3 subgroups per age

# Power of statistical tests

- methods decide whether data compatible with hypothesis $H_0$
- $F, t, z, W, U$ test statistics
- $H_0$ rejected if test value out of 95% acceptance region (assuming $H_0$)
  - × 2 samples same popula  $t$-test ; Mann-Whitney rank test
  - × 2 populor same varia  $F$-test
  - × 2 survival curves equal  log-rank
- $p < 0.05$  stat significant → test value out of 95% region
- $p \gtrless 0.05$  not stat. sig → $H_0$ cannot be rejected $\neq$ $H_0$ is valid
  - could not be proved that is not valid

| Reality \ Decision | $H_0$ rejected ($\equiv$ positive) | $H_0$ accepted ($\equiv$ negative) |
|---|---|---|
| $H_0$ false | True positive $p = 1-\beta$ | False negative $p = \beta$ <br> Type II error |
| $H_0$ true | False positive $p = \alpha$ <br> Type I error | True negative $p = 1-\alpha$ |

## Power of test $\boxed{1-\beta}$
→ probability of true positive decision
= correctly rejecting $H_0$
→ detect stat sign. diff when $H_0$ really false

$H_0$ true :

$H_1$ true :



$1-\alpha$ : true neg.
$\alpha$ : Type 1 error
$1-\beta$ : true pos.
$\beta$ : Type 2 error

$\beta = \phi(z_{1-\alpha} - \phi(z_{1-\alpha}))$
$\zeta * \sqrt{n} + z_{1-\frac{\alpha}{2}}$ $\sqrt{n} - z_{1-\frac{\alpha}{2}}$

Power depends on: (1) chosen $\alpha$, (signif. level) $P \propto \alpha$ ; $n$  $P \propto t_{crit; n, \alpha}$
$n_1 = n_2 = n$  (2) ratio between diff to be detected and SEM  $P \propto \frac{\Delta \mu}{\sigma \sqrt{\frac{2}{n}}}$

$t' = \frac{\mu_1 - \mu_2}{\sigma \sqrt{2/n}} \longrightarrow \Phi = \frac{\delta}{\sigma}$ noncentrality parameter
$\sigma \cdots$ of sample, not of mean
"... prop to ther. effect in units of stdev

$P \uparrow \begin{cases} n \uparrow \ (2\times) \\ \alpha \uparrow \\ \Delta \mu \uparrow \\ \sigma \downarrow \end{cases} \zeta = \frac{\delta}{\sigma} = \Phi \uparrow$

Only size of 1 group



$n=40$
$n=5$

## $t$-test power function
- families of curves $[n]$, fixed $\alpha$, depending on $\Phi$
- obtain $n$ ($P = 80\%$, $\Phi = 1$, $\alpha = 5\%$) inverting
- study should designed / $P \geqslant 80\%$

## Sample size estim

$n \approx \frac{2 s^2}{\delta^2} \left( z_{1-\beta} + z_{1-\frac{\alpha}{2}} \right)^2$  two-sided

$z_{1-\beta}$ $z_{1-\frac{\alpha}{2}}$



$\boxed{z_{1-\frac{\alpha}{2}} + z_{1-\beta} = \delta'}$

$\delta' = \frac{\delta}{s}\sqrt{\frac{n}{2}}$

## Hazard function
$F(t) = 1 - S(t)$  cdf $\longrightarrow h(t) = -\frac{S'(t)}{S(t)}$
$f(t) = S'(t)$  pdf

# Regression

- stochastic dependency between two random variables $X, Y$    (phenomena in nature)
  - ↳ stoch. factors influencing $\langle$ either $X$ or $Y$ / both

$$X = X(U_1, \dots U_M, V_1, \dots V_2)$$
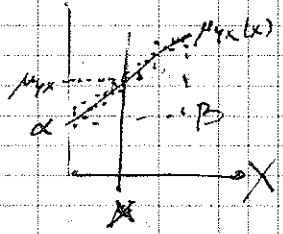$$Y = Y(U_1, \dots U_M, W_1, \dots W_j)$$    are stoch. dependent

## Regression theory

- predict random variable $Y$ if $X$ is known / fixed
- regression curves in $xy$ plane

$$\hat{y}(x) = E(Y \mid X = x), \quad \hat{x}(y) = E(X \mid Y = y)$$
   ↳ condit. expect.

   ↳ location of most accurate prediction for $Y$ if $X$ has value $x$
   ↳ minimising $E[Y - \hat{y}(x)]^2$ mean square error

## Linear regression

$\hat{y}(x)$ straight → $X$ and $Y$ linearly correlated

- $\mu_{yx}(x) = \alpha + \beta x = E[Y(x)]$ : mean of all $y$ at certain $x$
- $\sigma_{yx}(x)$ : stdev of all $y$ at certain $x$
- Preconditions $\left\{ \begin{array}{l} \mu_{yx} = \alpha + \beta x \\ \text{For all } x, \ Y \sim \mathcal{N}(\mu_{yx}, \sigma_{yx}) \\ \sigma_{yx} \text{ constant for all } x \end{array} \right.$



population
$\alpha, \beta$ → $a, b$ estimators sample

$$\sum_k (y_k - a - b x_k)^2 \longrightarrow \frac{\partial}{\partial a} \stackrel{!}{=} 0 \ ; \ \frac{\partial}{\partial b} \stackrel{!}{=} 0$$

$$b = \frac{n \sum XY - \sum X \sum Y}{n(\sum x^2) - (\sum x)^2} \quad ; \quad a = \overline{Y} - b\overline{X}$$

$$S_{yx} = \sqrt{\frac{\sum [Y - (a + bx)]^2}{n - 2}} = \sqrt{\frac{n-1}{n-2}(s_y^2 - b^2 s_x^2)}$$

$$S_a = S_{yx} \sqrt{\frac{1}{n} + \frac{\overline{x}^2}{(n-1)s_x^2}} \quad ; \quad S_b = \frac{1}{\sqrt{n-1}} \frac{S_{yx}}{S_x}$$

## Confidence intervals

$$t = \frac{b - \beta}{S_b} \longrightarrow \beta \in \left( b \pm t_{k, 1-\frac{\alpha}{2}} S_b \right) \qquad k = n - 2 \ \hat{=} \ \text{d.o.f.}$$
   ↳ $\beta > 0$ : sign.

$$t = \frac{a - \alpha}{S_a} \longrightarrow \alpha \in \left( \alpha \pm t_{k, 1-\frac{\alpha}{2}} S_a \right) \longrightarrow \text{contains } 0 : \text{trend}$$

$$\underset{\text{popul. line}}{\mu_y = \alpha + \beta x} \ne \underset{\text{regr. line}}{\overset{\text{curve}}{\hat{y} = a + bx}} \longrightarrow S_{\hat{y}(x)} = S_{yx} \sqrt{\frac{1}{n} + \frac{(x - \overline{x})^2}{(n-1)s_x^2}}$$

wider towards end

$\left[ \begin{array}{l} \text{error of mean} \\ \text{(when } x = \overline{x} \text{)} \\ \text{at regr. line} \end{array} \right] = \frac{S_{yx}}{\sqrt{n}}$

$$\hat{y} \in \left( y \pm t_{k, 1-\frac{\alpha}{2}} S_{\hat{y}} \right)$$

- For individual observation (instead of $\overline{y}$)
  - 1) variability determined by $S_{yx}$
  - 2) given by uncertainty line of means $S_{\hat{y}}$

$$S_{yN} = \sqrt{S_{yx}^2 + S_{\hat{y}}^2}$$

$$\longrightarrow \hat{y} \in \left( y \pm t_{k, 1-\frac{\alpha}{2}} S_{yN} \right)$$

# Comparison two regression lines

1) test differences slopes just
2) test differences intercepts just ........ based on $t$-test
3) test equality whole line

| SLOPES | $K = n_1 + n_2 - 4$ | INTERCEPTS |
|---|---|---|

$$t = (b_1 - b_2)/s_{b_1 - b_2}$$

$H_0: b_1 = b_2$

$$s_{b_1 - b_2} = \sqrt{s_{b_1}^2 + s_{b_2}^2} \qquad (n_1 = n_2)$$

$$t = (a_1 - a_2)/s_{a_1 - a_2}$$

$H_0: a_1 = a_2$

$$s_{a_1 - a_2} = \sqrt{s_{a_1}^2 + s_{a_2}^2}$$

$(n_1 \neq n_2)$ : pooled variance estimator var. around reg. lines

$$s_{yxp}^2 = \frac{(n_1 - 2) s_{yx_1}^2 + (n_2 - 2) s_{yx_2}^2}{K}$$

$$s_{b_1 b_2} = s_{yxp} \sqrt{\frac{1}{(n_1 - 1) s_{x_1}^2} + \frac{1}{(n_2 - 1) s_{x_2}^2}}$$

$$s_{a_1 a_2} = s_{yxp} \sqrt{\frac{1}{n_1} + \frac{1}{n_2} + \frac{\overline{x_1}^2}{(n_1 - 1) s_{x_1}^2} + \frac{\overline{x_2}^2}{(n_2 - 1) s_{x_2}^2}}$$
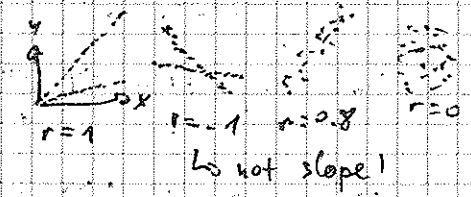
# EQUALITY

$H_0$: Two regression lines are equal

(1) Calculate regr. lines $\{x_1, y_1\}$ ; $\{x_2, y_2\}$

(2) Pooled variance estimate both lines $s_{yxp}$ (combined)

(3) Calculate common regression line and $s_{yx_s}^2 =$ ........ $(n = n_1 + n_2)$

(4) Estimate variance reduction when separate fits $s_{yx_\Delta}^2 = (n_1 + n_2 - 2) s_{yx_s}^2 - (n_1 + n_2 - 4) s_{yxp}^2$ over $2$

(5) Perform F-test $F = \dfrac{s_{yx_\Delta}^2}{s_{yxp}^2}$ $\quad m_1 = 2$
$\quad m_2 = n_1 + n_2 - 4$

(6) If $F > F_{m_1, m_2, 1-\alpha}$ : reject $H_0$

# Correlation

- Regr. analysis $\begin{cases} \text{change of dep. variable following change of indep} \\ \text{conf. int. for predicting dep. var. at fixed value of indep.} \end{cases}$

- Correlation $\begin{cases} \text{causality unknown (which dep/indep) not defined} \\ \text{describes strength of relationship between two variables} \end{cases}$

○ Pearson product-moment correlation coefficient

$$r = \frac{\sum (x - \overline{x})(y - \overline{y})}{\sqrt{\sum (x - \overline{x})^2 \sum (y - \overline{y})^2}} \qquad -1 \leq r \leq 1$$

$r = 1$ $\quad r = -1$ $\quad r = 0.8$ $\quad r = 0$
↳ not slope !

multiv. regr. missing

# Hazard function $h(t)$ ; $S(t)$ : survival function

$F(t) = 1 - S(t)$ cdf

$f(t) = - S'(t)$ pdf

$\to h(t) = - \dfrac{S'(t)}{S(t)} = \lim_{\Delta t \to 0} P(T \in (t, t + \Delta t) \mid T > t)/\Delta t$

Rate event occurs if not happened until $t$

$\to S(t) = \exp\left( - \int_0^t h(t') dt' \right) \qquad h(t) = \lambda \to S(t) = e^{\lambda t}$

$h(t)$ Hazard

Incr. Weibull = leukemia

Lognormal = tuberculose

Constant = healthy

Falling Weibull = operat.

$S(t)$ Survival

## Cox regression

- Estimate $S(t)$ after account for covariates (age, gender, smoker...)
  ↳ not possible Kaplan-Meier
1) Proportional Hazard models (Cox) < Covariates multipl. Hazard; $\log(HR) \propto \sum x_i$
2) Accelerated Failure Time models (AFT) → Cov multipl. Survival; $\log(S) \propto \sum x_i$

M=1 univariate

Cox: $i$ in N patients, M covariates $X_{ij}$ ⟹ $S(t)$ wrt baseline group ($X_{ij}=0$ all)

MODEL exp fact lasar.
$$h(t, \vec{x_i}) = h_0(t) \cdot e^{\beta_1 x_{i1}} e^{\beta_2 x_{i2}} \dots e^{\beta_M x_{iM}} = h_0(t) \cdot f(x_i) \text{ proportional hazard assumption}$$

- $X_{ij}$ not $X_{ij}(t)$; $X_j(x_k)$
- $\beta_j$ fit parameter; not $\beta_j(t)$

- $h_0(t)$ Baseline Hazard-function
  ↳ not explicitly needed (semiparametric model)
  ↳ robust against $\Delta$ model

Goal: find best $\beta_j$ of model to fit data; get HR

$$\frac{h(t, \vec{x_i})}{\sum_{k \in n_i, t_k \geq t_i} h(t, \vec{x_k})} = L_i = \frac{h_0(t_i) \prod_j^M \exp(\beta_j x_{ij})}{\sum_{k \in n_i, t_k \geq t_i} h_0(t) \prod_j^M \exp(\beta_j x_{kj})} \to L(\vec{\beta}) = \prod_i^N L_i^{\delta_i}$$

$h_0$ disappear, not $\delta$

$\delta_i = 0$ othw $\delta_i = 1$
- censored ignored BUT count for sum other $L_i$
- current $t_i$ irrelevant, just order $v_i$ (idem to wilcoxon rank sum)

- $LL(\vec{\beta})$ → easier to maximize; Newton-Raphson iterative $\frac{\partial LL}{\partial \beta}$
- $Cov(\vec{\beta}) = -I^{-1}$ (2nd deriv of LL)
- $\hat{\beta_j} \pm z_{1-\frac{\alpha}{2}} \sqrt{Cov(\hat{\beta})_{jj}}$ → Test Wald if $\hat{\beta_j} \neq 0$ or Likel. Quotient Test = Loglik Diff Test
  ↳ $W_j = (\hat{\beta_j}/s_j)^2 \sim \chi^2_1$ → p-value } $W_j > \chi^2_{1,0.95}$
  ↳ $LL - LL\{\beta_j=0\} \sim \chi^2$ → p-value } ↳ survival signif. affected by covariate $x_j$

### Hazard ratio (two groups w. diff. covariates)

$$HR_{ij} = \exp(\hat{\beta_j}) = h(t) x_{ij}=1 / h(t; x_{ij}=0)$$

$$HR_i = h(t; \vec{x_i^*}) / h(t; \vec{x_i}) = \exp(\beta_1(x_{i1}^* - x_{i1})) \dots \exp(\beta_M(x_{iM}^* - x_{iM}))$$

↳ how much risk increase if cov. increased by 1
↳ factor to multiply treat group wrt to control, rest constant
$$HR_j \begin{cases} < 1 \text{ longer survival} \\ > 1 \text{ shorter survival} \\ = 1 \text{ no diff} \end{cases} \text{ wrt control group } ; \quad (W_{ij} > W_{crit})$$

### Logistic regression   binary criterion y; covariates $x_i$ (dose, age...)
probab. (TCP) ↳ $P(y=1) \in [0,1]$   predictor

univariate $$P(y=1) = \frac{\exp(b_0 + b_1 x)}{1 + \exp(b_0 + b_1 x)}$$
$b_0$: constant
$b_1$: slope
$x$: dose
$$x_{50} = -\frac{b_0}{b_1} \to y(x_{50}) = 0.5$$

multivariate $$P(y=1) = \frac{\exp(\vec{b} \cdot \vec{x})}{1 + \exp(\vec{b} \cdot \vec{x})}$$
$L_i(\vec{b})$
$\vec{b} = (b_0, b_1, \dots b_M)$   $y_i$ N points $= \{0,1\}$
$\vec{x} = (1, x_1, \dots x_M)$   $x_j$ M predictors

P Logit


Goal: find optimum $\vec{b}$ to fit data; get $x_{50}$
- $L_i(\vec{b}) = e^{\vec{b} \cdot \vec{x_i}} / (1 + e^{\vec{b} \cdot \vec{x_i}}) \Rightarrow L(\vec{b}) := P(y_1 \dots y_n) = \prod_{i=1}^N L_i(b)^{y_i} (1 - L_i(b))^{1-y_i}$
- $LL(\vec{b})$ → easier max → NR, iterat
- $Cov(\vec{b}) \dots$ 2nd deriv. LL < $\hat{b_m} \pm z_{\alpha/2} \sqrt{Cov(\hat{b})_{mm}}$ ; $\hat{b} \cdot \vec{x} \pm z_{\alpha/2} \sqrt{\vec{x} \cdot Cov(\hat{b}) \vec{x}^T}$
- $H_0$: $\hat{b_m} = 0$ → Wald stat $W_m = (\hat{\beta_m}/s_m)^2 \sim \chi^2_1$ or LL quotient test

### Two group comparison
- Two dummy predictors < $x_i \cdot e_i$; $(x_i, e_i)$ (0/1 group no); $x_i$: dose; & combine all data
  1) log regr. $\{x_i; e_i; x e_i\}$; if $x e_i = 0 \to b_1 = b_2$; otherwise $b_1 \neq b_2$; $b_0$ no info
  2) if $b_1 = b_2$, log regr $\{x; e_i\}$ if $e_i \neq 0 \to b_{01} \neq b_{02} \Rightarrow$ implies $TCD50_1 \neq TCD50_2$
OR *) Likel. Quotient Test $\{x_i, e_i, x e_i\}$ vs $\{x\} \to$ signif $\equiv$ regr. line different but don't know which parameter