

# Multivariate Analysis of Variance

---

**Jiří Kléma**

Department of Computer Science,  
Czech Technical University in Prague



<http://cw.felk.cvut.cz/wiki/courses/b4m36san/start>













## Independence test for two categorical variables

---

- clarification why the Pearson's test statistic follows  $\chi^2$  distribution,
- for simplicity, concern a simple goodness of fit test with only 2 categories
  - $N$  trials,  $X$  observations in cat 1,  $N - X$  observations in cat 2,
  - $p_1 = p = \frac{X}{N}$ ,  $p_2 = 1 - p = \frac{N-X}{N}$
  - $H_0 : p = p_0$  (compare to a statistical model, a single number only here),
  - $X$  follows binomial distribution

$$\Pr(X = k) = \binom{N}{k} p^k (1 - p)^{N-k}$$

- the probability of getting exactly  $k$  successes in  $N$  trials, each trial successful with probability  $p$ ,
- for large  $N$ s can be approximated by  $N(Np, Np(1 - p))$ ,
- we can standardize  $X$  as  $z = \frac{X - Np}{\sqrt{Np(1-p)}}$ .





















# Analysis of variance (ANOVA)

---

## ■ assumptions

- the subjects are **independently sampled**
  - \* employ repeated measures ANOVA otherwise,
- the data are **normally distributed** in each group
  - \*  $E(Y_{i.}) = \mu_i$ , e.g., no group sub-populations with different means,
  - \* residuals of the model below show the normal distribution
$$y_{ij} = \mu + \alpha_i + \epsilon_{ij} = \mu_i + \epsilon_{ij}$$
  - \* employ non-parametric Kruskal-Wallis test otherwise,
- the data are **homoscedastic**
  - \* the variability in the data does not depend on group membership,
  - \* there is a common variance  $var(Y_{ij}) = \sigma^2$ ,
- multiple comparisons problem for more groups,

## ■ the hypotheses of interest

- $H_0 : \mu_1 = \mu_2 = \dots = \mu_g$ ,
- $H_a : \mu_i \neq \mu_j$  for at least one  $i \neq j$ .

# Analysis of variance (ANOVA)

---

■ method

- partition  $SS_{total}$ , the total variation in a response variable,
- distinguish **within groups variability**  $SS_{error}$ ,
- and **between groups variability**  $SS_{treat}$ ,

$$\begin{aligned}SS_{total} &= \sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{..})^2 = \\&= \sum_{i=1}^g \sum_{j=1}^{n_i} ((y_{ij} - \bar{y}_{i.}) + (\bar{y}_{i.} - \bar{y}_{..}))^2 = \\&= \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i.})^2}_{SS_{error}} + \underbrace{\sum_{i=1}^g n_i (\bar{y}_{i.} - \bar{y}_{..})^2}_{SS_{treat}}\end{aligned}$$

- \*  $\bar{y}_{i.} = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}$  ... group  $i$  sample mean,
- \*  $\bar{y}_{..} = \frac{1}{N} \sum_{i=1}^g \sum_{j=1}^{n_i} y_{ij}$  ... grand mean.











# Post-hoc ANOVA tests

---

- after performing ANOVA (and rejecting the null hypothesis)
  - we only assume that there is some difference in group means,
- a post-hoc test identifies which particular groups stand behind the test outcome,
- Tukey's HSD (honest significant difference) test
  - a t-test that controls for family-wise error rate (FWER),
  - compares all pairs of group means,
  - identifies all pairs whose difference is larger than expected standard error,
  - observed test statistics related to the studentized range distribution,

$$q_{obs} = \frac{\bar{y}_i - \bar{y}_j}{\sqrt{\frac{MS_{error}}{n^*}}} \sim q_{g, N-g}$$

- $n^*$  ... number of observations per group (their harmonic mean if not equal),
- always positive, sort the means before its application.

# ANOVA extensions/alternatives

---

- up to now we talked about ANOVA that
  - is parametric,
  - deals with independent measurements,
  - is one-way (with a single factor),
  - concerns a single target variable only,
- other options
  - non-parametric analysis (Wilcoxon test → Kruskal-Wallis analysis),
  - compares all possible group means (repeated measures ANOVA, Friedman test if non-parametric too),
  - main effects ANOVA and factorial ANOVA,
  - multivariate ANOVA (MANOVA).





# Analysis of variance (ANOVA)

- method

- the analogy of  $SS_{total}$  in ANOVA is a  $p \times p$  **cross products matrix**  $\mathbf{T}$ ,
- similarly to ANOVA, it can be decomposed into the **Error Sum of Squares and Cross Products**  $\mathbf{E}$ , and the **Hypothesis Sum of Squares and Cross Products**  $\mathbf{H}$ .

$$\begin{aligned}
 \mathbf{T} &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_{..})(\mathbf{y}_{ij} - \bar{\mathbf{y}}_{..})' = \\
 &= \sum_{i=1}^g \sum_{j=1}^{n_i} \{(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i) + (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_{..})\} \{(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i) + (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_{..})\}' = \\
 &= \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)'}_{\mathbf{E}} + \underbrace{\sum_{i=1}^g n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_{..})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}_{..})'}_{\mathbf{H}}
 \end{aligned}$$

\*  $\bar{\mathbf{y}}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{y}_{ij}$  ... sample mean vector for group  $i$ ,

\*  $\bar{\mathbf{y}}_{..} = \frac{1}{N} \sum_{i=1}^g \sum_{j=1}^{n_i} \mathbf{y}_{ij}$  ... grand mean vector of length  $p$ .







