

### VIII Analysis of variance and regression

#### VIII.1 "One way" ANOVA

Idea: Consider data like

	<u>Treatments</u>				
	1	2	...	k	
<b>Observation</b> ↓	$y_{11}$	$y_{21}$		$y_{k1}$	We shall assume that the corresponding r.v. follow $Y_{ij} = \theta_i + \epsilon_{ij}$ $i = 1, \dots, k$ $j = 1, \dots, n_i$ with $\{\epsilon_{ij}\}$ independent. unknown noise or error
	$y_{12}$	$y_{22}$		$y_{k2}$	
	$\vdots$	$\vdots$		$\vdots$	
	$y_{1n_i}$	$y_{2n_i}$		$y_{kn_i}$	

Remark: we could consider

$$Y_{ij} = \mu + J_i + \epsilon_{ij}$$

but usually <sup>additional parameter</sup> impose that  $\sum_{i=1}^k J_i = 0$  which fixes one parameter.

Def. The model  $Y_{ij} = \theta_i + \epsilon_{ij}$  is called a **oneway ANOVA** if

- 1)  $\epsilon_{ij}$  is a r.v. following  $n(0, \sigma^2)$   $\sigma$  is independent of  $i, j$  for any  $i, j$ .
- 2)  $\epsilon_{ij}$  and  $\epsilon_{i'j'}$  are independent for any  $(i, j) \neq (i', j')$ .

Remark: These assumptions can also be weakened if necessary.

Def. **ANOVA null assumption** is  $H_0: \theta_1 = \theta_2 = \dots = \theta_k$ . <sup>very strong assumption</sup>

$$\Rightarrow H_1: \theta_j \neq \theta_k \text{ for } \geq 1 \text{ pair } (\theta_j, \theta_k) \text{ with } j \neq k$$

Let us set

$$A = \mathcal{A}_k = \{ \underline{a} = (a_1, a_2, \dots, a_k) \in \mathbb{R}^k \text{ with } \sum_{j=1}^k a_j = 0 \}$$

Examples:  $\underline{a} = (1, -1, 0, \dots, 0)$  or  $\underline{a} = (1, -\frac{1}{2}, -\frac{1}{2}, 0, \dots, 0)$ , etc

Def. Consider  $\underline{t} := (t_1, \dots, t_k)$  a set of  $k$  parameters or of  $k$  r.v.

If  $\underline{a} \in \mathcal{A}_k$ , then  $\sum_{j=1}^k a_j t_j \equiv \underline{a} \cdot \underline{t}$  is called a **contrast**.

Lemma:  $H_0 \Leftrightarrow \forall \underline{a} \in \mathcal{A}: \underline{a} \cdot \underline{\theta} = 0$   $\Rightarrow$  easy  $\Leftarrow$  as an exercise

Corollary:  $H_1 \Leftrightarrow \exists \underline{a} \in \mathcal{A}: \underline{a} \cdot \underline{\theta} \neq 0$

Under the ANOVA assumption one has

$$Y_{ij} \sim n(\theta_i, \sigma^2) \quad \forall i = 1, \dots, k, \quad \forall j = 1, \dots, n_i$$

$$\Rightarrow \bar{Y}_i := \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij} \sim n(\theta_i, \frac{\sigma^2}{n_i})$$

Then for any  $\underline{a} \in \mathbb{R}^k$  consider  $\sum_{i=1}^k a_i \bar{Y}_i =: \underline{a} \cdot \bar{\mathbf{Y}}$  one has

$$\underline{a} \cdot \bar{\mathbf{Y}} \sim n\left(\sum_{i=1}^n a_i \theta_i, \sigma^2 \sum_{i=1}^k \frac{a_i^2}{n_i}\right) \quad (\text{Ex. 11.8})$$

If  $\sigma^2$  is not known, then we can define the sample variance

$$s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 \sim \chi_{n_i-1}^2 \quad (\S II.2)$$

Since  $\sigma^2$  is the same for all experiments, one can set

$$s^2 = \frac{1}{N-k} \sum_{i=1}^k (n_i - 1) s_i^2 = \frac{1}{N-k} \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2 \text{ and}$$

$$(N-k) s^2 / \sigma^2 \sim \chi_{N-k}^2 \quad \text{sum of } \chi_{n_i-1}^2 \text{ (Lemma 5.3.2) with}$$

$$N := \sum_{i=1}^k n_i$$

Then 
$$\frac{\sum_{i=1}^k a_i \bar{Y}_i - \sum_{i=1}^k a_i \theta_i}{\sqrt{s^2 \sum_{i=1}^k \frac{a_i^2}{n_i}}} \sim t_{N-k}$$

Now for any  $\underline{a} \in \mathbb{R}^k$  fixed, a hypothesis test could be

$$\text{Reject } H_0: \sum_{i=1}^k a_i \theta_i = 0 \text{ if } \left| \frac{\underline{a} \cdot \bar{\mathbf{Y}} - \underline{a} \cdot \underline{\theta}}{\sqrt{s^2 \sum_{i=1}^k \frac{a_i^2}{n_i}}} \right| > t_{N-k, \frac{\alpha}{2}} \text{ for a given } \alpha.$$

Equivalently,

$$\alpha \text{ confident interval with confidence coefficient } 1-\alpha \text{ is given by}$$

$$\sum_{i=1}^k a_i \bar{y}_i - t_{N-k, \frac{\alpha}{2}} \sqrt{s^2 \sum_{i=1}^k \frac{a_i^2}{n_i}} \leq \sum_{i=1}^k a_i \theta_i \leq \sum_{i=1}^k a_i \bar{y}_i + t_{N-k, \frac{\alpha}{2}} \sqrt{s^2 \sum_{i=1}^k \frac{a_i^2}{n_i}}$$

Key fact: If we choose  $\underline{a} = (1, -1, 0, \dots, 0)$  then we are testing  $H_0: \theta_1 - \theta_2 = 0$  or for  $\underline{a} = (1, -\frac{1}{2}, -\frac{1}{2}, 0, \dots, 0)$  then we are testing  $H_0: \theta_1 = \frac{1}{2}(\theta_2 + \theta_3)$ .

Observe that all data are used for the computation of  $s^2$ .

⚠ If we choose 2 different  $\underline{a}, \underline{a}' \in \mathbb{R}^k$ , the 2 computations for  $H_0: \underline{a} \cdot \underline{\theta} = 0$  and  $H_0': \underline{a}' \cdot \underline{\theta} = 0$  are not independent, which implies that the 2 confidence coefficients are not both  $1-\alpha$ .

Remark: The test for the ANOVA null assumption  $H_0: \theta_1 = \theta_2 = \dots = \theta_k$  can be obtained by the hypothesis union-intersection

$$H_0 \Leftrightarrow \underline{\theta} \in \bigcap_{\underline{a} \in \mathcal{A}} \{ \underline{\xi} \in \mathbb{R}^k \mid \underline{a} \cdot \underline{\xi} = 0 \}$$

To obtain a criterion for this intersection corresponds to a maximization problem. See 11.2.4

An  $\alpha$ -level test for rejecting  $H_0$  is given by

$$\frac{\sum_{i=1}^k n_i (\bar{Y}_i - \bar{\bar{Y}})^2}{s^2} > (k-1) F_{k-1, N-k; \alpha} \text{ with } \bar{\bar{Y}} := \frac{1}{N} \sum_{i,j} Y_{ij}$$

↑ F distribution

from which one can obtain a  $(1-\alpha)$  confidence interval. Called ANOVA F-test

## VIII.2 Simple Linear Regression

Idea: Consider relation of the form:  $f(x_i)$

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

$\uparrow$  response       $\uparrow$  predictor       $\leftarrow$  noise or error  
 2 unknown coeff.

Remark: It is called **linear regression** because it is linear in the coefficients  $\alpha, \beta$  (and not because of  $x_i$ )

Remark: Once some data are collected, just evaluating  $\alpha$  and  $\beta$  corresponds to "data fitting" but there is no statistical inference.

For any random variables  $\{(X_i, Y_i)\}_{i=1}^n$ , let us set

			$i$	$x$	$Y$
			1	$x_1$	$y_1$
			2	$x_2$	$y_2$
			$\vdots$	$\vdots$	$\vdots$
			$n$	$x_n$	$y_n$

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad S_{XX} = \sum_{i=1}^n (X_i - \bar{X})^2$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \quad S_{YY} = \sum_{i=1}^n (Y_i - \bar{Y})^2 \quad S_{XY} = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

About data fitting

Lemma: (least squares) One has

$$\min_{\alpha, \beta} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 = \sum_{i=1}^n (y_i - (a + bx_i))^2 \text{ for}$$

$$b = S_{XY} / S_{XX} \text{ and } a = \bar{y} - b\bar{x}$$

Remark: The previous computation is quite natural if we think that  $y_i$  is a function of  $x_i$  with relation  $y = \alpha + \beta x$ , but if we just collect  $(x_i, y_i)$ , it is not clear that this is the best choice.

In order to do some statistical inference let us assume that

$$Y_i = \alpha + \beta x_i + \varepsilon_i$$

with  $\{\varepsilon_i\}$  indep. r. v., with  $E(\varepsilon_i) = 0$  and  $\text{var}(\varepsilon_i) = \sigma^2$  independent of  $i$

Then  $E(Y_i) = \alpha + \beta x_i$  and  $\text{var}(Y_i) = \sigma^2$ .