MATH 204: Principles of Statistics 2

Dr David A. Stephens

Department of Mathematics & Statistics Room 1235, Burnside Hall

d.stephens@math.mcgill.ca www.math.mcgill.ca/~dstephens/204/ **Textbook:** McClave and Sincich (2008), *Statistics* (11th Edition), Chapters 10-15.

Note that the 10th Edition of McClave and Sincich contains essentially the same material.

Prerequisites: MATH 203 (or equivalent)

Some statistical computing knowledge useful.

Method of Assessment:

- Assignments
- ► Mid-Term
- ► Final

See syllabus handout for precise details.

Course Objectives	Three main sections
 Extensions of MATH 203 topics to other practical experimental contexts Introduction to statistical computation using standard software (SPSS) Practice in the use of statistical methods, in particular, hypothesis testing and linear modelling. 	I. THE ANALYSIS OF VARIANCE AND DESIGNED EXPERIMENTS II. LINEAR REGRESSION MODELLING III. NON-PARAMETRIC TESTING

Typical experimental scenario	Example: Pre-Natal Care
 two different groups of subjects single observation/measurement made on each subject scientific question of interest ARE THE TWO GROUPS OF SUBJECTS SIGNIFICANTLY DIFFERENT IN TERMS OF THEIR MEASURED VALUES ? 	 Objective: To compare the birthweights of babies in two groups of mothers. GROUP A: Received five or fewer pre-natal visits GROUP B: Received more than five pre-natal visits Do the GROUP A babies have significantly different birthweights from those from GROUP B ?





Statistical Testing Formally, we We adopt the following procedure to assess the "significance" of the difference between \overline{x}_A and \overline{x}_B . 1. Define a *test statistic*, *T*, that permits comparison of the two groups distribution. 2. Predict how T will behave assuming that the two groups are Similarly not significantly different. 3. Compare the prediction with what was actually observed.

▶ assume a Normal distribution for the data in the two groups

i.e. x_{A1}, \ldots, x_{An_A} are drawn from a population of birthweights that is well-modelled by a

Normal
$$(\mu_A, \sigma_A^2)$$

```
x_{B1},\ldots,x_{Bn_B} \sim Normal(\mu_B,\sigma_B^2)
```

We might initially assume that

$$\sigma_A^2 = \sigma_B^2$$

consider the two hypotheses

$$H_0: \quad \mu_A = \mu_B$$
$$H_a: \quad \mu_A \neq \mu_B$$

H_0 is the NULL HYPOTHESIS H_a is the ALTERNATIVE HYPOTHESIS

► define the test statistic

$$t = \frac{\overline{x}_A - \overline{x}_B}{s\sqrt{\frac{1}{n_A} + \frac{1}{n_B}}}$$

where

$$s^{2} = rac{(n_{A}-1)s_{A}^{2} + (n_{B}-1)s_{B}^{2}}{n_{A}+n_{B}-2}$$

 s^2 is the estimate of the common population variance

$$\sigma^2 = \sigma_A^2 = \sigma_B^2$$

Here

$$s^{2} = \frac{(10-1)61190.4 + (7-1)198679.6}{10+7-2} = 116186.1$$

so that

Thus

i.e. t should lie somewhere in the "high-probability region" of the Student-t(15) probability distribution



i.e. we are surprised to see <i>t</i> so far away from zero. The predicted behaviour of <i>t</i> , under the assumption that <i>H</i> ₀ is TRUE, DOES NOT MATCH THE OBSERVED BEHAVIOUR ! Therefore, the assumption that <i>H</i> ₀ is true MUST BE INCORRECT and we REJECT <i>H</i> ₀	 How do we quantify the "statistical significance" ? Two approaches: 1. Define the "high-probability" region, and reject H₀ if t does not lie in this region. 2. Compute the level of "surprise" at observing t under the assumption that H₀ is TRUE.
---	---



Equal Variances ?
Is the assumption of equal population variances

$$\sigma_A = \sigma_B$$

fair in this case ?
 $s_A^2 = 61190.4$
 $s_B^2 = 198679.6$
so that
 $\frac{s_A^2}{s_B^2} = 0.3080.$
Can we test $\sigma_A = \sigma_B$ formally ?

Yes:

$$H_0: \quad \sigma_A = \sigma_E$$
$$H_2: \quad \sigma_A \neq \sigma_E$$

$$F = \frac{s_A^2}{s_B^2} = 0.3080$$

If H_0 is true, F should look like an observation from a

Fisher-F

distribution with

$$\left(n_{A}-1,n_{B}-1
ight)$$
 "degrees of freedom" .



$$C_{R_1} = 0.231$$
 $C_{R_2} = 5.523$

so the observed value of F does lie in the high probability region, and there is no reason to reject H_0 at $\alpha = 0.05$.

Can also compute a 95 % confidence interval for $\mu_A - \mu_B$

$$(\overline{x}_A - \overline{x}_B) \pm t_{n_A+n_B-2}(0.975)s\sqrt{\frac{1}{n_A} + \frac{1}{n_B}}$$

where

 $t_{n_A+n_B-2}(0.975) = 2.131$

that is, the 0.975 probability point of the Student - t(15)distribution.

Hence the 95 % confidence interval is

$$(-1144.59, -428.67)$$

- does not contain zero !



	1.1 DESIGNED EXPERIMENTS Data collection studies typically fall into one of two categories:
In this section	 (i) Observational studies: the experimenter has no control over the variables under study, and can only measure outcomes.
 introduction to the terminology of <i>designed experiments</i> extension of statistical testing theory to comparison of more 	 The IQ of MAC and PC users The relationship between environmental exposure to toxins and health status.
than two population means ► THE ANALYSIS OF VARIANCE (ANOVA) F-TEST	i.e. The experimenter does not control the exposure to variables that may cause changes in the outcome of interest.
	This type of study is common in medicine and epidemiology as it is relatively cheap to carry out.
	Common type of observational study:
	CASE-CONTROL STUDY

Example (Smoking and Lung Cancer)

A study (Doll and Hill, 1950) investigated 649 lung cancer cases and 649 matched healthy controls, both drawn from a population of men in the UK. They found out what proportion in each group were smokers.

Neither health status nor smoking status were controlled by the experimenter, but were merely observed.

	Smokers	Non-smokers	Total
Lung cancer	647	2	649
Controls	622	27	649

This type of study can be unreliable, and cannot uncover all the relationships of interest.

A preferred approach involves the experimenter controlling the variables that cause variation in the other variables.

Note that this may not be ethical in a smoking/lung cancer study.

(ii) Designed experiments: the experimenter can the levels of variables that may affect the variable of interest.

Example (Birthweight study)	Terminology
$\begin{array}{rll} & \mbox{GROUP A} & : & 5 \mbox{ or fewer visits} \\ & \mbox{GROUP B} & : & \mbox{More than 5 visits.} \end{array}$	 Response variable (dependent variable): the variable of interest in the study Factors : the variables that may have an effect of the response variable quantitative if measured on a numerical scale qualitative otherwise Factor Levels: the values of the factors utilized in the experiment Treatments: the factor-level combinations utilized. Experimental Units (subjects): the objects on which the factors are measured or observed.

Therefore:

- ► A *designed experiment* is one for which the analyst or experimenter **controls** the specification of treatments and the method of assigning units to treatments.
- ► An *observational experiment* or study is one for which the analyst simply **observes** the treatments and response on a sample of experimental units.

Example (Birthweight study)

- **Response:** Birthweight (g)
- ► Factor: Pre-natal treatment group
- ► Factor levels: GROUP A or GROUP B

that is, we have a single factor with two factor levels.

Example (SAT scores)

The SAT scores of female and male students in four schools are to be compared.

- ► Response: SAT score
- ► Factors: SEX and SCHOOL (both qualitative)
- Factor levels:
 - ► SEX: Female and Male
 - ► SCHOOL: A,B,C,D

that is, we have a two factors, SEX with two factor levels and SCHOOL with four factor levels. There are 8 possible treatments:

(F, A), (F, B), (F, C), (F, D), (M, A), (M, B), (M, C), (M, D)

Example (Pain Relief)

Different pain relief remedies are to be compared : factors are

- REMEDY (quantitative/qualitative, 3 levels)
 - Dose level 0
 - Dose level 1
 - Dose level 2
- ► AGE GROUP (quantitative/qualitative, 4 levels)
 - ▶ 0-16 years
 - ► 17-40 years
 - ► 41-65 years
 - 66 years and over
- SEX (qualitative, 2 levels)
 - ► Female
 - ► Male

A total of $3 \times 4 \times 2 = 24$ possible treatment combinations; REMEDY is the only factor that can be assigned by the analyst.

Completely Randomized Design Statistical Objectives The experimental units assigned to different treatments A completely randomized design (CRD) is a design for which (factor-level combinations) form treatments are randomly assigned to experimental units, or in which random samples of experimental units are selected for each independent samples treatment. from The term can be applied to both experimental and observational different populations studies. For example, in a CRD. ▶ if the treatments are FEMALE/MALE for the factor SEX, a CRD draws independent samples of FEMALES and MALES We wish to **compare** treatments: specifically, we wish to compare for the two treatment groups. the treatment means. ▶ if the treatments are DOSE 0/DOSE 1, a CRD assigns experimental units independently to the two treatment groups A Multiple Group Comparison of Means ! at random.

Suppose that there are k treatments:	Comparing <i>k</i> Treatments
TMT 1Mean μ_1 TMT 2Mean μ_2 \vdots \vdots	Suppose
TMT k Mean μ_k We wish to test the hypotheses	TMT 1has n_1 experimental unitsTMT 2has n_2 experimental units \vdots \vdots TMT khas n_k experimental units
$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ $H_a : \text{At least two of the } k \text{ treatment means are different}$ How do we do this ? What is the relevant test statistic ? (41)	Denote by x_{ij} the response for unit j in treatment group i , for $j = 1,, n_i$ and $i = 1,, k$.

Let

$$\overline{x}_i = \frac{1}{n_i} \sum_{i=1}^{n_i} x_{ij}$$

denote the sample mean for treatment i, and

$$s_i^2 = rac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - \overline{x}_i)^2$$

denote the sample variance for treatment i.

Now we consider pooling, that is, combining all units into a single group.

Define

► the total sample

 $n = n_1 + \dots + n_k = \sum_{i=1}^k n_i$

► the overall sample mean

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{k} \sum_{j=1}^{n_i} x_{ij}$$

the overall sample variance

$$s^{2} = rac{1}{n-1} \sum_{i=1}^{k} \sum_{j=1}^{n_{i}} (x_{ij} - \overline{x})^{2}$$

Finally, consider the pooled sample variance

$$s_P^2 = rac{1}{n-k} \sum_{i=1}^k (n_i - 1) s_i^2$$

- the extension of the pooled estimate of the population variance in a two-sample t-test.

Using these quantities, we can derive a test statistic for multiple group comparison.

We wish to compare how much variation is due to the

A DIFFERENCE BETWEEN TREATMENTS

and how much is due to

We measure A using the statistic

$$SST = \sum_{i=1}^{k} n_i (\overline{x}_i - \overline{x})^2$$

SST - \underline{S} um of \underline{S} quares for \underline{T} reatments

We measure B using the statistic

$$SSE = \sum_{i=1}^{k} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$
$$= \sum_{i=1}^{k} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$
$$= (n-k)s_P^2$$

SSE - $\underline{S}um$ of $\underline{S}quares$ for $\underline{E}rror$

NOTE: This measure of random or error variability implicitly assumes that the variability **within** the treatment groups is the **same for each group**. That is, population variances

$$\sigma_1^2,\ldots,\sigma_n^2$$

are equal.

In practice this assumption must be checked.

Finally, we define the test statistic using the mean levels of variability

► MST - <u>Mean Square for Treatments</u>

$$MST = \frac{SST}{k-1} = \frac{1}{k-1} \sum_{i=1}^{k} n_i (\overline{x}_i - \overline{x})^2$$

► MSE - <u>Mean Square for Error</u>

$$MSE = \frac{SSE}{n-k} = \frac{1}{n-k} \sum_{i=1}^{k} \sum_{j=1}^{n_i} (x_{ij} - \overline{x}_i)^2 = s_P^2$$

Then the test statistic is

$$F = \frac{MST}{MSE} = \frac{\text{Average Variation due to Treatments}}{\text{Average Variation due to Treatments}}$$

$$F = \frac{MST}{MSE} = \frac{\text{Average Variation due to Treatments}}{\text{Average Variation due to Errors}}$$

$$F = \frac{MST}{MSE} = \frac{MST}{\text{Average Variation due to Errors}}$$

$$F = \frac{MST}{MSE}$$

NOTE: If

 $SS = \sum_{i=1}^{k} \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2$ 1. The samples are randomly selected in an independent manner from the k treatment populations. [Satisfied in a CRD] is the overall or total sum of squares, then 2. All k populations have distributions that are approximately normal. SS = SST + SSE3. The k population variances are equal. so we can decompose the overall variation (SS) into the variation $\sigma_1^2 = \sigma_2^2 = \cdots \sigma_k^2.$ due to treatments (SST) and the variation due to the errors (SSE).

Assumptions behind the ANOVA F-test

n - k = 1334

Therefore

$$MST = \frac{SST}{k-1} = \frac{10.606}{2} = 5.303$$

$$MSE = \frac{SSE}{n-k} = \frac{136.432}{1334} = 0.102$$
and

$$F = \frac{MST}{MSE} = 51.851$$
If H_0 is true, that is,

$$\mu_1 = \mu_2 = \mu_3$$
then F should look like an observation from a
Fisher-F($k-1, n-k$)
distribution.

$$55$$
Here we are dealing with the
Fisher-F(2,1334)
distribution. From tables, we discover that if $\alpha = 0.05$, then

$$F_{\alpha}(2,1334) = 3.002$$
and thus we
Reject H_0
and conclude that there is a significant impact on milk protein
level due to diet.

Are the assumptions met ? 1. Independent samples : Not possible to tell with current information. In fact, data comprise repeated measurements on 79 cows - potentially not independent, as observations on the same cow are likely to be more similar. 2. Normal Distributions : Visual inspection of boxplots so we cannot look up $F_{0.05}(2, 1334)$. However, we know that indicates that this may be valid. 3. Equal variances : $s_1^2 = 0.102$ $s_2^2 = 0.091$ $s_3^2 = 0.114$ so assumption appears to be valid - can we test this formally ?



Note: Tables in McClave and Sincich only give

and here the test statistic is F = 51.851.

 $F_{0.05}(2, 120) = 3.07$

 $F_{0.05}(2,\infty) = 3.00$

 $3.00 < F_{0.05}(2, 1334) < 3.07$

0 1 2 2 27.0 22.8 21.9 23.5 26.2 23.1 23.4 19.6 28.8 27.7 20.1 23.7 33.5 27.6 27.8 20.8 28.8 24.0 19.3 23.9	For $\alpha = 0.05$, from McClave and Sincich tables $F_{0.05}(3,16) = 3.24$ and so we
vve find that	Reject H ₀
SST = 140.094 $SSE = 116.324$ $SS = 256.418$	at $\alpha = 0.05$ and conclude that there is a significant difference between treatment groups
MST = 46.698 $MSE = 7.270$	p-value is 0.0046.
and	
F = 6.423	
which we need to compare with the Fisher- $F(3, 16)$ distribution.	
61	62

ר



The ANOVA TableFor a completely randomized design, we may report the results of
the ANOVA F-test in a stylized form, the ANOVA TableSOURCE DF SS MS FTREATMENTS
$$k-1$$
 SST $MST = \frac{SST}{(k-1)}$ $F = \frac{MST}{MSE}$ ERROR $n-k$ SSE $MSE = \frac{SSE}{(n-k)}$ TOTAL $n-1$ SS

SST

 $\overline{(k-1)}$

SSE (n-k)

MST F =

MSE

(i)
$$(k-1) + (n-k) = (n-1)$$

(ii)
$$SST + SSE = SS$$

i.e. we can fill in missing values if they are not given.

Sometimes an extra column is added at the right of the table to give the *p*-value of the ANOVA F-test.

SOURCEDFSSMSFpTMT
$$k-1$$
SSTMST $F = \frac{MST}{MSE}$ p-valERROR $n-k$ SSEMSETOTAL $n-1$ SS

where *p*-val solves

$$\frac{MST}{MSE} = F_{p-val}(k-1, n-k)$$

and $F_{\alpha}(\nu_1, \nu_2)$ is the $(1 - \alpha)$ probability point of the Fisher-F distribution.



Levene's Test Example (PTSD Example (see handout)) To test n = 45, k = 4. $H_0 = \sigma_1^2 = \sigma_2^2 = \cdots = \sigma_k^2$ H_1 = At least one pair of σ^2 different. F-statistic F = 3.046Critical Value $F_{0.05}(3,41) \simeq 2.84$ $F_{0.025}(3,41) \simeq 3.46$ Test statistic $W = \frac{(n-k)}{(k-1)} \frac{SST_Z}{SSE_Z} = \frac{MST_Z}{MSE_Z}$ $F_{0.01}(3, 41) \simeq 4.31$ where SST_7 and SSE_7 are the usual sums of squares evaluated for Tables give $F_{\alpha}(3, 40)$. the new data z_{ii} where \implies Reject H_0 at $\alpha = 0.05$ (p = 0.039). $z_{ij} = |x_{ij} - \overline{x}_j|.$ BUT Levene's Test suggests that the assumption of equal If H_0 is true variances is NOT valid. $W \sim \text{Fisher-F}(k-1, n-k).$

Why do we need the three assumptions ?

- ► independence
- Normality
- equal variances
- so that we can predict (under H_0) that

 $F \sim \text{Fisher-F}(k-1, n-k)$

and complete the test (compute *p*-values and the rejection region).

But our hypothesis of interest is

 H_0 : No difference between treatments

Under this hypothesis, the treatment labels

SHOULD NOT MATTER !

i.e. we should be able to exchange the labels, and not notice any major difference in the test statistic.

This leads us to consider permutation or randomization tests.

i.e. we compute the test statistic for all possible relabellings consistent with H_0 , retaining the group sample sizes, and use these values to compute the rejection region.

Randomization/Permutation Tests

Suppose that there are ${\it N}$ possible relabellings that give rise to test statistics

 F_1, F_2, \ldots, F_N

Then the rejection region for significance level $\boldsymbol{\alpha}$ is the interval to the right of

 $N(1-\alpha)$ th largest of the values F_1, F_2, \ldots, F_N

and the *p*-value is

 $\frac{\text{Number of } F_1, F_2, \dots, F_N \ge F}{N}$

where

 $F = \frac{MST}{MSE}$

is the true test statistic.

If the group sample sizes are n_1, n_2, \ldots, n_k then

$$N=\frac{n!}{n_1!n_2!\dots n_k!}$$

where

 $n! = n(n-1)(n-2)\dots 3.2.1$

("*n* factorial") - potentially very large.

Example (PTSD Example) $k = 4, n = 45$ $(n_1 = 14, n_2 = 10, n_3 = 11, n_4 = 10)$	Example (PTSD Example (continued)) Using this approach, we compute for $\alpha = 0.05$ CRITICAL VALUE : $C_R = 2.844$ p-VALUE : $p = 0.040$
There are $\frac{45!}{14!10!11!10!} = 2.610 \times 10^{24}$ possible relabellings: a very big number. We compute $F = \frac{MST}{MSE}$ for each relabelling. For the real data, $F = 3.046$.	Compare this with the ANOVA F-test values $CRITICAL VALUE : C_R = 2.833$ $p-VALUE : p = 0.039$ (using the Fisher-F(3,41) distribution. Thus we obtain virtually identical results; but the randomization test does not need the assumptions of normality or equal variances.
7	75



If
$$k = 2$$
, consider $F = MST/MSE$;

$$MST = \frac{1}{k-1} \sum_{i=1}^{k} n_i (\overline{x}_i - \overline{x})^2 = n_1 (\overline{x}_1 - \overline{x})^2 + n_2 (\overline{x}_2 - \overline{x})^2$$

$$= \frac{n_1 n_2}{n_1 + n_2} (\overline{x}_1 - \overline{x}_2)^2$$

$$MSE = \frac{1}{n-k} \sum_{i=1}^{k} \sum_{j=1}^{n_i} (x_{ij} - \overline{x}_i)^2 = s_P^2$$

$$= \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

Therefore

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

Thus $F = t^2$, where t is the two-sample t-test statistic.

Thus if k = 2, the ANOVA F-test and the two sample *t*-test are **EQUIVALENT**

$$t \sim \text{Student-t}(n-2)$$

 $F \sim \text{Fisher-F}(1, n-2)$

and we must get the same conclusion (to reject H_0 or otherwise) using either statistic.

Summary	1.3 Multiple Comparison of Means
If the assumptions • independence (holds by design in a CRD) • Normal populations • equal variances hold, use ANOVA F-test If the assumptions do not hold • use Randomization/Permutation test • use Non-parametric test (see Section 3)	If the ANOVA F-test null hypothesis $H_0: \mu_1 = \dots = \mu_k$ is rejected , then it is of interest to discover which of the means are different. For k groups, there are $c = k(k-1)/2$ pairs of group means that can be compared. Consider a "family" of hypothesis tests - a collection of tests of different hypotheses carried out independently on different data sets. For each test in the family, we consider testing the hypothesis at significance level α .
18	8

Notation

Label the tests $i = 1, \ldots, c$, and for each i, label

- the null hypotheses H_{0i}
- the test statistics T_i
- the rejection regions \mathcal{R}_i

that are potentially different for each i.

We specify for each *i*,

 $\alpha = P[T_i \in \mathcal{R}_i | H_{0i} \text{ is } \mathbf{TRUE}]$

which implicitly defines \mathcal{R}_i . Note that α is the

"Test Type-I Error Rate" or "Comparisonwise Error Rate"

Now consider the results of all tests in the family; what is the "Familywise" Type-I error rate ?

Using the laws of probability

$$P[T_i \in \mathcal{R}_i | H_{0i} \text{ is } \mathbf{TRUE}] = \alpha$$

means that

$$P[T_i \notin \mathcal{R}_i | H_{0i} \text{ is } \mathbf{TRUE}] = 1 - \alpha$$

giving the probability that the test **does not reject** H_{0i} , if H_{0i} is in fact true, is $1 - \alpha$.

Now we consider all tests together;

 $P[\text{Each } T_i \notin \mathcal{R}_i | \text{Each } H_{0i} \text{ is } \mathbf{TRUE}] = (1 - \alpha)^c$

This is the probability that each test results in the null hypothesis **not** being rejected, that is, the probability that we **never** commit a Type-I error.

Therefore the probability of at least one Type-I error is

 $\alpha_F = 1 - (1 - \alpha)^c$

 α_F is the Familywise Error Rate.

	$\alpha = 0.05$	$\alpha = 0.01$
С	α_F	α_F
5	0.226	0.049
10	0.401	0.096
50	0.923	0.395
100	0.994	0.634

Therefore, whenever we carry out a "family" of tests, we should not use the traditional choices of $\alpha=0.05$ or 0.01 on each test.

The Bonferroni Method

To fix $\alpha_F = 0.05$, say, we need to use α on each test where

$$\alpha_F = 1 - (1 - \alpha)^c \iff \alpha = 1 - (1 - \alpha_F)^{1/c}$$

For example, if $\alpha_F = 0.05$ and c = 10, use

$$\alpha = 1 - (1 - 0.05)^{1/10} = 0.0051$$

It can be shown that

$$1 - (1 - \alpha)^c \approx c\alpha$$

86

Therefore, if α_F is the required **familywise error rate**, we must set the **comparisonwise error** rate to be $\alpha = \alpha_F/c$.

 α_F/c is known as the Bonferroni Correction.

Confidence Intervals SPSS gives twelve different methods for correcting the confidence interval for use in different experimental situations. For example For the k = 2 group comparison of means, a $100(1 - \alpha)\%$ • planned comparisons $\mu_1 = \mu_3$, $\mu_7 = \mu_{10}$ etc. confidence interval for $\mu_1 - \mu_2$ is ► all comparisons $(\overline{x}_1 - \overline{x}_2) \pm t_{\alpha/2}(n_1 + n_2 - 2)s_P \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$ Three methods are recommended: ► Tukey's Method Bonferroni's Method where $t_{\alpha}(\nu_1)$ is the $1 - \alpha$ probability point of the Student-t Scheffé's Method distribution with u_1 degrees of freedom (under the assumptions of independence, Normality and equal group variances). Having selected a multiple comparison correction method, we compute simultaneous confidence intervals for each comparison of If we move to a family of *c* tests, to get simultaneous confidence means, and identify intervals for the differences in means $\mu_i - \mu_i$ for all pairs of *i* and *j*, we should adjust α to α_F when computing the $100(1-\alpha)\%$ which means are significantly different confidence interval. • the ranking of differences $\mu_i - \mu_i$ in terms of magnitude.

1.4 Randomized Block Designs We wish to compare treatments whilst acknowledging that there may be differences between the blocks. A randomized block design used matched experimental units organized into sets known as **blocks** and assigns one member from That is, the observed variation is due to the set to each treatment. TREATMENTS and BLOCKS and ERROR For k treatments rather than merely 1. Compile b blocks of k experimental units, with each block comprising units that are similar. TREATMENTS and ERROR 2. Assign one unit from each block to each treatment at random. as in the CRD. Then there are a total of n = bk measured responses.

 Response : Measured SAT Score Eactor : Sex 	Example (SAT Scores (continued))
► Factor-levels : $k = 2$ (Female/Male) ► Blocks : $b = 5$ (Previous GPA, within same school) i.e. $k = 2, b = 5 \therefore n = 10$. Block Female SAT Male SAT 1 A: 2.75 540 530 2 B: 3.00 570 550 3 C: 3.25 590 580 4 D: 3.50 640 620 5 E: 3.75 690 690	 This design recognizes that GPA score and school are likely to explain some variation in SAT Score, but that neither is directly related to the "treatment" of interest (SEX - Female/Male). i.e. the blocking variable removes systematic variation in response that is not of primary interest. We pick one Female and one Male in each school/GPA category, and pair them.
	91



Let

$$SST = \sum_{i=1}^{k} b(\overline{x}_i - \overline{x})^2$$
$$SSB = \sum_{j=1}^{b} k(\overline{x}_j^{(B)} - \overline{x})^2$$
$$SS = \sum_{i=1}^{k} \sum_{j=1}^{b} (x_{ij} - \overline{x})^2$$

SST: Sum of Squares for **Treatments** SSB: Sum of Squares for **Blocks** SS: **Total** Sum of Squares

Finally

$$SS = SST + SSB + SSE$$
 \therefore $SSE = SS - SST - SSB$

 $\mathsf{SSE:}\ \mathsf{Sum}\ \mathsf{of}\ \mathsf{Squares}\ \mathsf{for}\ \mathbf{Errors}$

$$F = \frac{MST}{MSE}$$

where

Test statistic is

$$MST = \frac{SST}{k-1}$$
 $MSE = \frac{SSE}{n-b-k+1}$

ANOVA F-test to compare treatment means in a randomized block design

Theorem (ANOVA F-test for a RBD)

 H_0 : $\mu_1 = \cdots = \mu_k$ H_a : At least one pair of treatment means different.

use the test statistic

 $F = \frac{MST}{MSE}$

If H_0 is **TRUE**

$$F \sim Fisher$$
- $F(k-1, n-b-k+1)$

- this defines the rejection region for significance level $\alpha,$ and the p-value, in the usual way.

After the ANOVA test is complete, and the hypothesis

is rejected, we can proceed with the "post-hoc" tests of

hypotheses $\mu_i = \mu_j$ for $i \neq j$.

terms of response.

Notes:

 $H_0: \mu_1 = \cdots = \mu_k$

 In a RBD, it is not (in general) possible to estimate individual treatment means, that is, x
_i does not estimate μ_i as it is an average across blocks, which are believed to be different in Assumptions:

- 1. Experimental units (between blocks) are independent, and treatments are allocated at random (within blocks).
- 2. Normality
- 3. *bk* block/treatment combinations correspond to populations with equal variances.

ANOVA Table

SOURCE	DF	SS	MS	F
TMTS	k-1	SST	MST	F = MST/MSE
BLOCKS	b-1	SSB	MSB	
ERROR	n-k-b+1	SSE	MSE	
TOTAL	n-1	SS		

2. Testing the Block Means However, we can test whether the block means $\mu_1^{(B)},\ldots,\mu_b^{(B)}$ are significantly different. For

$$H_0: \mu_1^{(B)} = \cdots = \mu_b^{(B)}$$

we use the F statistic

$$F = \frac{MSB}{MSE}$$

where

$$MSB = \frac{SSB}{b-1}$$

If H_0 is **TRUE**

$$F \sim \text{Fisher-F}(b-1, n-k-b+1)$$

100

That is, we treat the blocks as levels of another factor, and test to see whether this factor affects response.

Comment: The "sum of squares" decompositions Example (Soil Analysis (see handout)) CRD: SS = [SST] + SSEResults of two ANOVA F-tests: Test of Conclusion F р RBD: SS = [SST + SSB] + SSESOLVENT No Difference 0.673 0.585 SOIL 10.568 0.001 Difference are both of the form Here SOLVENT is the treatment variable, SOIL is the blocking variable. TOTAL = SYSTEMATIC + RANDOM VARIATION VARIATION VARIATION 3. Remember to check the assumptions (independence, normality, equal variances in each treatment/block combination) For the CRD: SST "SYSTEMATIC" For the RBD: SST + SSBEqual variances may be hard to check as we only have one observation per treatment/block comparison. "RANDOM" For both: SSE

	1.5 Factorial Experiments
We have studied the Randomized Complete Block Design where each block/treatment combination has one experimental unit. An incomplete design could also be considered, where some block/treatment combinations are omitted. However, this design does not lead to straightforward ANOVA analysis. 103	 Designs studied so far: CRD - one factor RBD - one factor, plus one blocking variable, so two factors in total, where one (the blocking variable) is a known source of systematic variation. However, in the RBD, we must assume that the treatments behave in a similar way across blocks.

Let *i* index treatments $(1 \le i \le k)$ and consider block *j*, and two treatment (factor levels) i_1 and i_2 .

In an RBD, we assume that

$$E[X_{i_1j} - X_{i_2j}] = \mu_{i_1} - \mu_{i_2}$$

which does **NOT** depend on j.

That is, the expected difference in response due to the two treatments does not depend on the block.

But perhaps the difference **does** depend on block; perhaps we have **INTERACTION**.

In the current RBD, we do not have enough data to look for this. We now seek to extend the RBD to allow for tests for interaction; we do this by using **replication**.

RBD with Balanced Replication

Suppose we have r observations per block/treatment combination (termed *replicates*), so that we have n = bkr experimental units in total.

Balanced designs have equal numbers of replicates in each block/treatment combination.

In this design, all the quantities

SST, SSB, SSE, SS MST, MSB, MSE

can be defined, and an ANOVA F-test can be carried out - the only difference is that n = bkr.

Sum of Squares for Treatments (SST)

$$SST = \sum_{i=1}^{k} br(\overline{x}_i - \overline{x})^2$$

Sum of Squares for Blocks (SSB)

$$SSB = \sum_{j=1}^{b} kr(\overline{x_{j}^{(B)}} - \overline{x})^{2}$$

Overall Sum of Squares (SS)

$$SS = \sum_{i=1}^{k} \sum_{j=1}^{b} \sum_{t=1}^{r} (x_{ijt} - \overline{x})^2$$

and SSE = SS - SST - SSB

Third index t indexes the replicates.

The RBD with replication does allow the investigation of interaction. The new test is based on the decomposition

$$SS = SST + SSB + SSI + SSE$$

where SSI is the sum of squares for Interaction.

We have SST, SSB and SS as before, and

$$SSI = \sum_{i=1}^{k} \sum_{j=1}^{b} r(\overline{x}_{ij} - \overline{x}_i - \overline{x_j^{(B)}} + \overline{x})^2$$

where

$$\overline{x}_{ij} = \frac{1}{r} \sum_{t=1}^{r} x_{ijt}$$
 $i = 1, \dots, k, \ j = 1, \dots, b$

is the sample mean for replicates in (i, j)th treatment/block combination.

Testing in the RBD with Replication

The three F statistics

$$F = \frac{MST}{MSE}$$
 $F = \frac{MSB}{MSE}$ $F = \frac{MSI}{MSE}$

can be used to test for significant Treatment, Block and Interaction effects respectively.

Now

$$MSE = \frac{SSE}{\text{Error d.f.}}$$

But what is "Error d.f." ? It is a constant that dictates how large SSE should be on average.

The general rule for computing the error d.f. for any model is

Error d.f. =
$$n - p$$

where n is the total sample size and p is the total number of parameters fitted.

How many parameters do we fit ?

► No Interaction

$$p = 1 + (b - 1) + (k - 1)$$

that is, the overall mean $\mu,$ plus the b-1 differences from μ due to the blocks, plus the k-1 differences from μ due to the treatments.

Interaction

$$p = bk$$

that is, one parameter in each cell of the two-way table of blocks by treatments.

Thus

No Interaction

$$p = 1 + (b - 1) + (k - 1) = b + k - 1$$

parameters, so

Error d.f.
$$= n - p = n - b - k + 1$$

• Interaction: we fit p = bk parameters, so

Ha

k parameters \longrightarrow 1 parameter

so there are (k-1) extra parameters, and SST varies on (k-1)

FULL MODEL -

 H_0

→ NULL MODEL

Error d.f.
$$= n - p = n - bk$$

It transpires that if

$$MSI = \frac{SSI}{(b-1)(k-1)}$$

is the Mean Square for Interaction, then

$$F = \frac{MSI}{MSE}$$

yields a test statistic suitable for testing interaction. If there is $\ensuremath{\mathbf{no}}$ interaction, then

$$F \sim \text{Fisher-F}((b-1)(k-1), n-bk)$$

where n = bkr.

Why (b-1)(k-1)? This is the number of **extra** parameters we fit to include the interaction.

For the CRD:

degrees of freedom.

For the RBD: the (i,j)th treatment/block combination has mean $\mu_i + \mu_i^{\mathsf{B}}$

so for testing for a TREATMENT effect

 H_a H_0 FULL MODEL \longrightarrow NULL MODEL k parameters \longrightarrow 1 parameter

so there are (k - 1) extra parameters, and *SST* varies on (k - 1) degrees of freedom.

 $\mu_1,\ldots,\mu_k\longrightarrow \mu$

For testing for a BLOCK effect

Ha		H_0
FULL MODEL	\longrightarrow	NULL MODEL
b parameters	\longrightarrow	1 parameter

so there are (b-1) extra parameters, and SSB varies on (b-1) degrees of freedom.

$$\mu_1^{(B)}, \ldots, \mu_k^{(B)} \longrightarrow \mu^{(B)}$$

These models and tests can be fitted and carried out even if we do not have replication.

With replication, we can investigate the interaction, that is the model where the (i, j)th treatment/block combination has mean

$$\mu_i + \mu_i^B + \mu_i$$

rather than the model where

$$\mu_i + \mu_j^B$$

that is, we wish to test

$$H_0$$
 : $\mu_{ij} = 0$ for all *i* and *j*
 H_a : $\mu_{ij} \neq 0$

ANOVA Table SOURCE DF SS MS In the **full interaction** model: we fit *bk* parameters TMTS k-1SST MST F_T BLOCKS b-1SSB MSB F_B In the restricted, no interaction model: we fit INTERACTION (b-1)(k-1) SSI MSI F_{I} ERROR SSE MSE 1 + (b - 1) + (k - 1) = b + k - 1(n - bk)TOTAL n-1SS parameters. Therefore the differences is where $MST = \frac{SST}{k-1}$ $MSB = \frac{SSB}{b-1}$ bk - (b + k - 1) = bk - b - k + 1 = (b - 1)(k - 1)and SSI varies on (b-1)(k-1) degrees of freedom. $MSI = \frac{SSI}{(b-1)(k-1)}$ $MSE = \frac{SSE}{n-bk}$ and $F_T = \frac{MST}{MSE}$ $F_B = \frac{MSB}{MSE}$ $F_I = \frac{MSI}{MSE}$

Exa	mple: Batteri	es Data (see	han	dout)			
1	Dependent Variab	le: Battery Life					
	Source	Sum of Squares	df	Mean Square	F	Sig.	
	Corrected Model	59,154.000	8	7,394.250	11.103	0.000	
	Intercept 398,792.250 1 398,792.250 598.829 0.00						
	material 10,633.167 2 5,316.583 7.983 0.002						
	temp	temp 39,083.167 2 19,541.583 24					
	material * temp	9,437.667	4	2,359.417	3.543	0.019	
	Error	17,980.750	27	665.954			
	Total	475,927.000	36				
	Corrected Total	77,134.750	35				
	R Squared = $.767$	(Adjusted R Squa	red =	.698)			

For $\alpha = 0.05$, there is a significant **temp** effect (p < 0.001), and a significant **material** effect (p = 0.002), and a significant interaction (p = 0.019)

 $\ensuremath{\mathsf{NB}}$: If we do not have replication, we CANNOT fit the interaction. Recall that

$$\mathsf{Error} \ \mathsf{d}.\mathsf{f}. = \mathit{n} - \mathit{bk}$$

but if r = 1, n = rbk = bk, so the error d.f. is zero.

In fact, SSE = 0 also, so the MSE is not defined.

We now study multifactor designs, to assess the effects and interactions of several factors simultaneously.

We consider all possible combinations of

FACTOR A with *a* levels FACTOR B with *b* levels FACTOR C with *c* levels

to define the treatments in a factorial design.

Factorial Experiments

A complete factorial experiment is one in which every combination of a number of factors is utilized.

i.e. the number of treatments is equal to the total number of factor-level combinations.

We focus on two factor experiments

FACTOR A with *a* levels FACTOR B with *b* levels

so there are *ab* treatments in total.

21



 individuals from the same base population are assigned at random to one of the *ab* treatments.

123

► Sum of Squares for Interaction (SSI_{AB})

$$SSI_{AB} = \sum_{i=1}^{a} \sum_{j=1}^{b} r(\overline{x}_{ij} - \overline{x}_{i.} - \overline{x}_{.j} + \overline{x}_{..})^2$$

New notation:

► sample mean for Factor A level i

$$\overline{x}_{i.} = \frac{1}{br} \sum_{j=1}^{b} \sum_{t=1}^{r} x_{ijt} \qquad i = 1, \dots, a$$

• sample mean for Factor B level j

$$\overline{x}_{,j} = \frac{1}{ar} \sum_{i=1}^{a} \sum_{t=1}^{r} x_{ijt} \qquad j = 1, \dots, b$$

• sample mean for replicates in (i, j)th factor combination

h r

$$\overline{x}_{ij} = \frac{1}{r} \sum_{t=1}^{r} x_{ijt} \qquad i = 1, \dots, a, \ j = 1, \dots, b$$

overall sample mean

$$\overline{x}_{..} = \frac{1}{n} \sum_{i=1}^{a} \sum_{j=1}^{b} \sum_{t=1}^{r} x_{ijt}$$

These allow computation of

 $SST_A, SST_B, SSI_{AB}, SS, SSE$

$MST_A, MST_B, MSI_{AB}, MSE$

using the degrees of freedom identical to those in the RBD with replication.

Tests for

- significant effect for Factor A
- ► significant effect for Factor B
- ► significant interaction

will be carried out as before using an ANOVA table.

ANOVA Table

SOURCE	DF	SS	MS	F
FACTOR A	a-1	SST_A	MST_A	F_A
FACTOR B	b-1	SST_B	MST_B	F_B
INTERACTION	(a - 1)(b - 1)	SSI _{AB}	MSI _{AB}	F_{AB}
ERROR	(n — ab)	SSE	MSE	
TOTAL	n-1	SS		

If Factor A is not influential (H_0 specifying no difference between responses at different levels of factor A), then

$$F_A \sim \text{Fisher-F}(a-1, n-ab)$$

Similarly,

No effect of Factor B : $F_B \sim \text{Fisher-F}(b-1, n-ab)$ No Interaction : $F_{AB} \sim \text{Fisher-F}((a-1)(b-1), n-ab)$

SEE EXAMPLES HANDOUT

Note: For two factors A and B, the $\ensuremath{\mbox{main}}$ effects plus interaction model can be written

A + B + A.B

whereas the main effects only can be written

A + A.B

A + B

B + A.B

The models

do not make sense.

28

A+B A+B+A.B interaction	A,B A,B is significan	NONE YES t, the only model you s	hould	► FD: two treatment factors "Blocking" factors are known or strongly believed to have a significant effect on the response.
	A + B + A	<i>4</i> .R		
4	A+B+A.B	A+B+A.B $A,Binteraction is significantA+B+i$	A+B+A.B A,B YES interaction is significant, the only model you sl A+B+A.B	A+B+A.B A,B YES interaction is significant, the only model you should A+B+A.B

Estimating Effect Size

►

In multifactor designs, parameter estimation can be carried out in different parameterizations

For the CRD (one-way layout):

• Natural parameters: μ_1, \ldots, μ_k

Contrast parameters:
$$\beta$$
, β_0 , ..., β_{k-1} where

$$\beta = \mu_k$$
 $\beta_i = \mu_i - \mu_k$, $i = 1, \dots, k - 1$

that is, differences from baseline.

For the two-factor designs (RBD/FD): In the two-way layout, with cells (i, j), i = 1, ..., a, j = 1, ..., b. The cell means are m_{ij} , where

$$m_{ij} = \mu_{i.} + \mu_{.j} + \mu_{ij}$$

where $\mu_{i.}$ gives the Factor A contribution, $\mu_{.j}$ gives the Factor B contribution, and μ_{ij} gives the interaction.

The parameterization used by SPSS is the contrast parameterization is $m_{ij} = \beta_0 \qquad i = a, j = b$ $= \beta_0 + \beta_i^{(A)} \qquad i = 1, \dots, a - 1, j = b$ $= \beta_0 + \beta_j^{(B)} \qquad i = a, j = 1, \dots, b - 1$

$$= \beta_0 + \beta_i^{(A)} + \beta_j^{(B)} + \gamma_{ij}^{(AB)}$$

where

$$egin{array}{rll} eta_i^{(A)} & : & ext{contrasts for factor A} \ eta_j^{(B)} & : & ext{contrasts for factor B} \ \gamma_{ij}^{(AB)} & : & ext{interaction} \end{array}$$

 $i = 1, \ldots, a - 1$

 $i = 1, \ldots, b - 1$

SPSS takes the *a*th level of factor A and the *b*th level of factor B as the baseline, and looks at differences compared to this baseline.

The *ab* parameters are

 $\begin{array}{ll} \beta_{0} & 1 \\ \beta_{1}^{(A)}, \dots, \beta_{a-1}^{(A)} & (a-1) \\ \beta_{1}^{(B)}, \dots, \beta_{b-1}^{(B)} & (b-1) \\ \gamma_{ij}^{(AB)}, i = 1, \dots, a-1, j = 1, \dots, b-1 \quad (a-1)(b-1) \\ \end{array}$ Total

For example: a = 3, b = 4.

			Fact	or B	
		1	2	3	4
A	1	Ð	2	3	$\beta_0 + \beta_1^{(A)}$
ctor	2	4	5	6	$\beta_0 + \beta_2^{(A)}$
Fa	3	$\beta_0 + \beta_1^{(B)}$	$\beta_0 + \beta_2^{(B)}$	$\beta_0 + \beta_2^{(B)}$	β_0

where

134

and so on.

Final Note on ANOVA Estimation is still straightforward: PARAMETER **ESTIMATE** We have studied the simplest design scenarios: extension to ▶ incomplete β_0 \overline{X}_{ab} unbalanced $\overline{x}_{i.} - \overline{x}_{ab}$ $\beta_i^{(A)}$ ► nested ▶ random effect $\overline{x}_{.j} - \overline{x}_{ab}$ $\beta_i^{(B)}$ designs are possible. $\gamma_{ij}^{(AB)} \qquad \overline{x}_{ij} - \overline{x}_{i.} - \overline{x}_{.j} + \overline{x}_{ab}$ Furthermore SPSS has greater functionality: for example, it has the capability to carry out ANOVA-like analyses even for the case of non-equal variances (when Levene's test rejects the hypothesis for i = 1, ..., a, j = 1, ..., b. of equal variances). Other parameterizations can be used. 136

 Part II
 In the previous section, we attempted to explain the variation in an observed response variable by fitting models with one or more factors.

 Factors are discrete variables taking different levels; in this section we will now utilize continuous variables that can similarly explain variation in an observed response.

2.1 Simple Linear Regression

We will investigate models relating two quantities x and y through equations of the form

y = ax + b

where a and b are constants (that is, a straight-line).

Variables x and y will not be treated exchangeably - we will regard y as being a function of x.

Such models are deterministic, that is, if we know x (and the values of the constants), we can compute y exactly without error.

A more useful model allows for the possibility that the system is not observed perfectly, that is, we do not observe (x, y) pairs that are always consistent with a simple functional relationship.

Example (Pharmacokinetic Model)

If a dose of drug is taken at time x = 0, the amount (concentration) of drug still in the bloodstream at time x is often well-modelled by a simple equation. Let

- *D* denote the amount of drug taken at x = 0
- ► x time
- y^* is the amount (concentration per unit volume) in the bloodstream.

Then

$$y^{\star} = \frac{D}{V} \exp\{-\lambda x\}$$

where

- λ is the elimination rate
- ► V is the volume of bloodstream.

2.1.1 Probabilistic Models Example (Pharmacokinetic Model (continued)) In a **probabilistic** model, we allow for the possibility that y is observed with random error, that is, Taking logs of both sides, setting $y = \log y^*$, then y = ax + b + ERROR $y = -\lambda x + \log(D/V) = -\lambda x + (\log D - \log V)$ where *ERROR* is a random term that is present due to imperfect that is, y = ax + b where observation of the system due to (i) measurement error or (ii) ► $a = -\lambda$ missing information. ▶ $b = (\log D - \log V)$ Note that we do not treat x and y exchangeably; x is a fixed observed variable that is measured without error, whereas y is an However, in practice, when we measure concentration, we do so observed variable that is measured with random error. with random error. We model the variation in y as a function of x. We observe pairs $(x_i, y_i), i = 1, \ldots, n.$

A Basic Probabilistic Model • $\beta_1 > 0$ - increasing y with increasing x • $\beta_1 < 0$ - decreasing y with increasing x Terminology:

- ▶ y Dependent variable or response variable
- ▶ x Independent variable, or predictor, or covariate

The model we study takes the form

$$y = \beta_0 + \beta_1 x + \epsilon$$

where ϵ is a random error term, a random variable with mean zero and finite variance $(E[\epsilon] = 0, Var[\epsilon] = \sigma^2)$; it represents the error present in the measurement of y.

- ▶ β_0 *Intercept* parameter
- β_1 *Slop*e parameter

- $\beta_1 = 0$ no relationship between x and y

Note:

$$E[Y|x] = \beta_0 + \beta_1 x$$

where E[Y|x] is the expected value of Y for fixed value of x.

Recall the notation

- ► Y a random variable with a probability distribution
- ▶ y a fixed value that the variable Y can take.

Fundamental Problem: If we believe the straight-line model with error is correct, how do we find the values of parameters β_0 and β_1 . We only have the observed data $\{(x_i, y_i), i = 1, \dots, n\}$.

2.1.2 Least Squares Fitting

We select the best values of β_0 and β_1 by minimizing the *error in fit*. For two data points (x_1, y_1) and (x_2, y_2) , the errors in fit are

$$e_1 = y_1 - (\beta_0 + \beta_1 x_1)$$

 $e_2 = y_2 - (\beta_0 + \beta_1 x_2)$

respectively. But note that, potentially, $e_1 > 0$ and $e_2 < 0$ so there is a possibility that these fitting errors cancel each other out. Therefore we look at **squared** errors (as a large negative error is as bad as a large positive error)

$$e_1^2 = (y_1 - (\beta_0 + \beta_1 x_1))^2$$

$$e_2^2 = (y_2 - (\beta_0 + \beta_1 x_2))^2$$

For n data, we obtain n misfit squared errors

$$e_1^2, \ldots, e_n^2$$

We select β_0 and β_1 as the values of the parameters that minimize SSE, where

$$SSE = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - (\beta_0 + \beta_1 x_i))^2$$

We wish to make the total misfit squared error as small as possible.

 SSE - sum of squared errors - is similar to the SSE for ANOVA. We could write

$$SSE = SSE(\beta_0, \beta_1)$$

to show the dependence of $\ensuremath{\textit{SSE}}$ on the parameters.

Minimization of $SSE(\beta_0, \beta_1)$ is achieved **analytically**.

145

146

Two routes: (i) calculus and (ii) geometric methods. It follows that the best parameters $\hat{\beta}_0$ and $\hat{\beta}_1$ are given by

$$\widehat{\beta}_1 = \frac{SS_{xy}}{SS_{xx}} \qquad \qquad \widehat{\beta}_0 = \overline{y} - \widehat{\beta}_1 \overline{x}$$

where

► Sum of Squares *SS_{xx}*:

$$SS_{xx} = \sum_{i=1}^{n} (x_i - \overline{x})^2$$

► Sum of Squares *SS*_{xy}:

$$SS_{xy} = \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})$$

 $\widehat{\beta}_0$ and $\widehat{\beta}_1$ are the least-squares estimates

$$y = \widehat{\beta}_0 + \widehat{\beta}_1 x$$

is the least-squares line of best fit. The fitted-values are

$$\hat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_i \qquad i = 1, \dots, n$$

and the residuals or residual errors are

$$\hat{\mathbf{e}}_i = y_i - \hat{y}_i = y_i - \widehat{eta}_0 - \widehat{eta}_1 x_i \qquad i = 1, \dots, n$$

2.1.3 Model Assumptions for Least-Squares

To utilize least-squares for the probabilistic model

$$Y = \beta_0 + \beta_1 x + \epsilon$$

we make the following assumptions

1. The expected error $E[\epsilon]$ is zero so that

$$E[Y] = \beta_0 + \beta_1 x$$

- The variance of the error, Var[ε], is constant and does not depend on x.
- The probability distribution of
 ϵ is a symmetric distribution about zero (a stronger assumption is that
 ϵ is Normally distributed).
- The errors for two different measured responses are independent, i.e. the error ε₁ in measuring y₁ at x₁ is independent of the error ε₂ in measuring y₂ at x₂.

2.1.4 Parameter Estimation: Estimating σ^2

Using the LS procedure, we can construct an estimate of the *error* or *residual error* variance

Recall that

 $Var[\epsilon] = \sigma^2$

An estimate of σ^2 is

$$\widehat{\sigma}^2 = \frac{SSE(\widehat{\beta}_0, \widehat{\beta}_1)}{n-2} = s^2$$

say.

Note that the denominator n-2 is again a *degrees of freedom* parameter of the form

or n-p, where in the simple linear regression, p=2 ($\widehat{\beta}_0$ and $\widehat{\beta}_1$). Note also that

$$SSE(\widehat{eta}_0,\widehat{eta}_1) = \sum_{i=1}^n (y_i - \widehat{y}_i)^2 = SS_{yy} - \widehat{eta}_1 SS_{xy}$$

where

$$SS_{yy} = \sum_{i=1}^{n} (y_i - \overline{y})^2$$

Estimation and Testing for Slope

In the model where

$$E[Y] = \beta_0 + \beta_1 x$$

it is of interest to test the hypothesis

$$H_0 : \beta_1 = 0$$
$$H_a : \beta_1 \neq 0$$

i.e. H_0 implies that there is no systematic contribution of x to the variation of y.

51

To test H_0 vs H_a we us the test statistic

$$t = \frac{\widehat{\beta}_1}{\text{e.s.e}(\widehat{\beta}_1)} = \frac{\widehat{\beta}_1}{s_{\widehat{\beta}_1}}$$

where e.s.e($\hat{\beta}_1$) is the *Estimated Standard Error* of $\hat{\beta}_1$, computed as

$$\mathsf{e.s.e}(\widehat{\beta}_1) = \frac{\mathsf{s}}{\sqrt{\mathsf{SS}_{xx}}}$$

where ${\it s}$ is the estimate of σ defined previously.

If H_0 is true, and $\beta_1 = 0$, then

$$t = \frac{\hat{\beta}_1}{s/\sqrt{SS_{xx}}} \sim \text{Student}(n-2)$$

so we can carry out a significance test at level α in the usual way (use a *p*-value, or construct the rejection region).

Note: we might also consider a one-sided test, where $H_a: \beta_1 > 0$, say.

• If $H_a: \beta_1 \neq 0$, we use the *two-sided* rejection region, with critical values

$$C_R = \pm t_{n-2}(\alpha/2)$$

► If H_a: β₁ > 0, we use the *one-sided* rejection region, with critical value

$$C_R = +t_{n-2}(\alpha)$$

► If H_a : β₁ < 0, we use the one-sided rejection region, with critical value</p>

$$C_R = -t_{n-2}(\alpha)$$

154

Note: To test

$$H_0 : \beta_1 = b$$
$$H_a : \beta_1 \neq b$$

for any *b*, the test statistic is

$$t = \frac{\widehat{\beta}_1 - b}{s/\sqrt{SS_{xx}}}$$

(for example, b = 1 may be of interest. If H_0 is true

$$t \sim \mathsf{Student}(n-2)$$

Confidence Interval

A 100(1 - $\alpha)\%$ confidence interval for β_1 is

$$\beta_1 \pm t_{n-2}(\alpha/2) \times s_{\widehat{\beta}_1}$$

where

$$t_{n-2}(\alpha/2)$$
 : $\alpha/2$ prob. point of Student $(n-2)$ distn.
 $s_{\widehat{\beta}_1}$: Estimated standard error of $\widehat{\beta}_1$

Note: we could perform a similar analysis for $\beta_0,$ but this is generally of less interest.

The only quantity that needs attention is the estimated standard error of $\widehat{\beta}_{0}.$ It can be shown that

e.s.e.
$$(\widehat{eta}_0) = s_{\widehat{eta}_0} = \sqrt{\frac{1}{n} \left(1 + \frac{n\overline{x}^2}{SS_{xx}}\right)}$$

2.1.5 The Coefficient of Correlation

To measure the *strength of association* between the two variables x and y we can use the

Pearson Product Moment Coefficient Of Correlation

or correlation coefficient which measures the strength of the **linear** relationship between x and y.

The coefficient, r, is defined by

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}}$$

where

$$SS_{xx} = \sum_{i=1}^{n} (x_i - \overline{x})^2 \quad SS_{yy} = \sum_{i=1}^{n} (y_i - \overline{y})^2$$
$$SS_{xy} = \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})$$

Note: $-1 \leq r \leq 1$.

- ► If *r* is close to 1, there is a strong linear relationship between *x* and *y* where *y* **increases** with *x*.
- ► If *r* is close to -1, there is a strong linear relationship between *x* and *y* where *y* **decreases** with *x*.

Note: In the model

$$y = \beta_0 + \beta_1 x$$

 $\beta_1 = 0 \implies r \approx 0$, so tests for $\beta_1 = 0$ can also be used to deduce a lack of correlation between the variables.



Testing Correlation

We use ρ to denote the **true** correlation between X and Y.

We can test the hypothesis that $\rho = 0$ (that is, that X and Y are uncorrelated using r. For testing

$$\begin{array}{rcl} H_0 & : & \rho = 0 \\ H_a & : & \rho \neq 0 \end{array}$$

we can use the test statistic

$$t = \frac{r}{\sqrt{(1-r^2)/(n-2)}}$$

If H_0 is true, then approximately

$$t \sim \text{Student}(n-2)$$

Alternately, we could use

$$z = \frac{1}{2} \log \left(\frac{1+r}{1-r} \right)$$

and then, if H_0 is true, as (approximately)

$$Z \sim N\left(\frac{1}{2}\log\left(\frac{1+
ho}{1-
ho}
ight), \frac{1}{n-3}
ight)$$

when $\rho = 0$, so that (approximately)

$$\sqrt{n-3} Z \sim N(0,1)$$

16

A related quantity is the

Coefficient of Determination

or R^2 Statistic

$$r^2 = \frac{SS_{yy} - SSE}{SS_{yy}} = 1 - \frac{SSE}{SS_{yy}}$$

Note that the *total variation* in y is recorded via

$$SS_{yy} = \sum_{i=1}^{n} (y_i - \overline{y})^2$$

and the random variation is recorded via

Example (Blood Viscosity vs PCV)

• $R^2 = r^2 = (0.879)^2 = 0.772$

We have ► *n* = 32

▶ *r* = 0.879

Test of $\rho = 0$:

$$SSE = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

Therefore the variation explained by the linear regression is

$$SSR = SS_{yy} - SSE$$
 as $SS_{yy} = SSR + SSE$

Thus

$$r^{2} = \frac{SSR}{SS_{VV}} = \frac{\text{Variation explained by Regression}}{\text{Total Variation}}$$

 r^2 is a measure of model adequacy, that is, if $r^2\approx 1,$ then the linear model is a good fit.

64

2.1.6 Prediction

After the linear model is fitted, it can be used for **forecasting** or **prediction**. That is, given a new x value we can predict the corresponding y.

As before, we see that at any value of x_p , the prediction \hat{y}_p is

$$\hat{y}_{p} = \hat{\beta}_{0} + \hat{\beta}_{1} x_{\mu}$$

This is the best predictor of y at this x value.

We can also compute the standard error of this prediction; if the value of the random error variance σ^2 is known, then

 $t = \frac{r}{\sqrt{(1 - r^2)/(n - 2)}} = 10.087$

We compare with a Student $(n-2) \equiv$ Student(30) distribution; the *p*-value is 3.73×10^{-11} , so there is strong evidence that $\rho \neq 0$.

s.e.
$$(\hat{y}_p) = \sigma \sqrt{\frac{1}{n} + \frac{(x_p - \overline{x})^2}{SS_{xx}}}$$

If σ is unknown, we estimate σ by $\hat{\sigma} = s$ as defined previously

$$s^2 = \frac{SSE(\widehat{\beta}_0, \widehat{\beta}_1)}{n-2}$$

so that

e.s.e.
$$(\hat{y}_p) = s \sqrt{\frac{1}{n} + \frac{(x_p - \overline{x})^2}{SS_{xx}}}$$

Note: This prediction is the expected value of y at $x = x_p$. That is, we have worked out

$$Var[\widehat{Y}_p] = Var[\widehat{\beta}_0 + \widehat{\beta}_1 x_p]$$

to compute the s.e. for \widehat{Y}_p .

But we can actually predict an **error corrupted** version of \widehat{Y}_p , \widehat{Y}_p^* say, where

$$\widehat{Y}_p^{\star} = \widehat{Y}_p + \epsilon_p$$

where ϵ_p is a new random error.

Prediction Intervals

But

$$Var[\widehat{Y}_{p}^{\star}] = Var[\widehat{Y}_{p} + \epsilon_{p}] = Var[\widehat{Y}_{p}] + Var[\epsilon_{p}] = Var[\widehat{Y}_{p}] + \sigma^{2}$$

that is, there is an **extra** piece of variation due to ϵ_p .

Thus

$$\text{e.s.e.}(\hat{y}_p^{\star}) = s \sqrt{1 + \frac{1}{n} + \frac{(x_p - \overline{x})^2}{SS_{xx}}} > \text{e.s.e.}(\hat{y}_p)$$

A $100(1 - \alpha)$ % prediction interval for the **mean** value at $x = x_p$ is

$$\hat{y}_{p} \pm t_{n-2} (\alpha/2) s \sqrt{\frac{1}{n} + \frac{(x_{p} - \overline{x})^{2}}{SS_{xx}}}$$

whereas for an individual new value (predicted with error) at $x = x_p$ is

$$\hat{y}_p \pm t_{n-2}(\alpha/2)s\sqrt{1+rac{1}{n}+rac{(x_p-\overline{x})^2}{SS_{xx}}}$$



$$SS = \sum_{i=1}^{n} (y_i - \overline{y})^2$$

$$SSR = \sum_{i=1}^{n} (\hat{y}_i - \overline{y})^2$$

$$SSR = \sum_{i=1}^{n} (\hat{y}_i - \overline{y})^2$$

$$SSE = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
where
$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad i = 1, ..., n$$
Degrees of Freedom
$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
TOTAL $n - 1$ SS
The test of the hypothesis
$$H_0 : E[Y] = \beta_0$$

$$H_0 : E[Y] = \beta_0$$

$$H_0 : E[Y] = \beta_0$$

$$H_0 : E[Y] = \beta_0 + \beta_1 x_i$$
can be completed by using the test statistic
$$F = \frac{MSR}{MSE}$$
If H_0 is true
$$F \sim \text{Fisher-F}(1, n - 2)$$

$$T73$$







R^2 and adjusted R^2

SPSS reports both the R^2 statistic

$$R^2 = 1 - \frac{SSE}{SS}$$

and the **adjusted** R^2 statistic

$$R^2 = 1 - \frac{SSE/EDF}{SS/TDF}$$

where

• EDF = error degrees of freedom = n - 2

• TDF = total degrees of freedom = n - 1

2.1.7 Polynomial Regression

In many practical situations, the simple straight line

$$y = \beta_0 + \beta_1 x$$

is not appropriate. Instead, a model including powers of x

$$x^2, x^3, \ldots, x^k$$

should be considered. For example

$$y = \beta_0 + \sum_{j=1}^k \beta_j x^j = \beta_0 + \beta_1 x + \dots + \beta_k x^k$$

The Polynomial Regression Model

$$Y = \beta_0 + \beta_1 x + \dots + \beta_k x^k + \epsilon$$

where ϵ is a random error term as before can be used to model data.

Two immediate problems:

1. How to choose k

2. How to carry out inference

- estimation
- testingprediction

Jediction

We begin by addressing 2. The estimation of parameters can be again carried out using **Least Squares** provided that the model assumptions listed before are valid. Consider k = 2.

We choose $\beta = (\beta_0, \beta_1, \beta_2)^T$ to minimize the sum of squared errors

$$SSE(\underline{\beta}) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2)^2$$

that is the fitted values for parameters β are

$$\hat{\mathbf{y}}_i = \beta_0 + \beta_1 \mathbf{x}_i + \beta_2 \mathbf{x}_i^2$$

 $\widehat{\beta}$ can be found to minimize *SSE* using calculus techniques (differentiating with respect to the elements of $\underline{\beta}$) to give the minimum SSE

$$SSE(\underline{\beta}) = \sum_{i=1}^{n} (y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i - \widehat{\beta}_2 x_i^2)^2$$

We can also compute the estimated standard errors

$$s_{\widehat{eta}_0}, s_{\widehat{eta}_1}, s_{\widehat{eta}_2}$$

which allow tests of parameters to be carried out, and confidence intervals calculated.

We can also compute prediction intervals.

The best estimate of the residual error variance σ^2 is

$$\widehat{\sigma}^2 = rac{SSE(\widehat{eta})}{n-3}$$

p is the number of parameters estimated equal to three, so we divide by n-3.

υŦ

We can also computeExample (Hooker Pressure Data) \blacktriangleright ResidualsFor the Hooker pressure data, a quadratic polynomial (k = 2)
might be suitable. \blacktriangleright can be used to assess the fit of the model. $Y = \beta_0 + \beta_1 x + \beta_2 x^2$
We need to estimate β_0 , β_1 and β_2 for these data to see if the

- R^2 , Adjusted R^2 statistics
 - used to assess the global fit of the model.
 - used to compare the quality of fit with other models.

We need to estimate β_0 , β_1 and β_2 for these data to see if the model fits better than the straight line model we fitted previously. This can be achieved using SPSS.

It transpires that the quadratic model produces a set of residuals that are patternless, which the straight line model when fitted does not.

See Handout for full details.

Note: It is common to use the Standardized Residuals

$$\widehat{z}_i = \frac{\widehat{e}_i}{\widehat{\sigma}} = \frac{y_i - \widehat{y}_i}{\widehat{\sigma}}$$

where $\hat{\sigma}^2$ is the estimate of σ^2 defined previously, as

$$\operatorname{Var}[\widehat{z}_i] \approx 1$$

if the model fit is good, whereas

 $\operatorname{Var}[\widehat{e}_i] \approx \sigma^2$

which clearly depends on σ . This makes it hard to compare \hat{e}_i across different models when inspecting residuals.

Note: Although the model based on

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

is **not** linear in x, it **is** linear in the parameters. Because of this, we still term this a *linear model*. It is this fact that makes the least-squares solutions easy to find.

This model is no more difficult to fit than the model

$$y = \beta_0 + \beta_1 \frac{x}{1+x} + \beta_2 (1-e^{-x})$$

say - it is still a *linear in the parameters model*. It is in the general class of models

$$y = \beta_0 + \beta_1 g_1(x) + \beta_2 g_2(x)$$

where $g_1(x)$ and $g_2(x)$ are general functions of x.

188

In fact, any model of the form

$$y = \sum_{i=0}^{k} \beta_{j} g_{j}(x) + \epsilon \tag{1}$$

can be fitted routinely using least-squares; if we know x, then we can compute

 $g_0(x), g_1(x), \ldots, g_k(x)$

and plug those values into the formula (1).

Example (Harmonic Regression)

Let

$$\begin{array}{rcl} g_0(x) &=& 1 \\ g_1(x) &=& \left\{ \begin{array}{ll} \cos(\lambda_j x) & j \text{ odd} \\ \sin(\lambda_j x) & j \text{ even} \end{array} \right. \end{array}$$

where k is an even number, k = 2p say.

 $\lambda_j, j = 1, 2, \dots, p$ are constants

$$\lambda_1 < \lambda_2 < \cdots < \lambda_p$$

For fixed x, $cos(\lambda_j x)$ and $sin(\lambda_j x)$ are also fixed, known values.







- because the system of equations based on the derivatives

$$\frac{\partial}{\partial \beta_j} \left\{ SSE(\beta) \right\} = 0 \qquad j = 0, 1, \dots, k$$

can always be solved routinely, so we can always find $\widehat{\beta}$.

In the general model (1), simple formulae for

- ►β
- ► s.e.(β̂)
- ► $\hat{\sigma}^2$

can be found using a matrix formulation.

See handout on website - NOT EXAMINABLE !

Note: One-way ANOVA can be formulated in the form of model (1). Recall

- ► *k* independent groups
- means μ_1, \ldots, μ_k
- y_{ij} *j*th observation in the *i*th group

Let

$$\begin{array}{lll} \beta_0 &=& \mu_k \\ \beta_t &=& \mu_t - \mu_k \qquad t = 1, 2, \dots, k-1. \end{array}$$

Define new data $x_{ij}(t)$ where

$$x_{ij}(t) = \begin{cases} 1 & \text{if } t = i \\ 0 & \text{if } t \neq i \end{cases}$$

10/

2.2 Multiple Linear Regression Then, using the linear regression formulation $y_{ij} = \beta_0 + \sum_{t=1}^{k-1} \beta_t x_{ij}(t) + \epsilon_{ij}.$ Multiple linear regression models model the variation in response yFor any *i*, *j*, $x_{ij}(t)$ is non-zero for only one value of *t*, when t = i. as a function of more than one independent variable. We term this a regression on a factor predictor; it is clear that Suppose we have variables $\beta_0, \beta_1, \ldots, \beta_{k-1}$ can be estimated using least-squares. X_1, X_2, \ldots, X_k This clarifies the link between recording different features of the experimental units. We wish to ANOVA Linear Modelling and model response Y as a function of X_1, X_2, \ldots, X_k . - they are essentially the SAME MODEL formulation. This link extends to ALL ANOVA models; recall that we used the General Linear Model option in SPSS to fit two-way ANOVA.

2.2.1 Multiple Linear Regression Models

Consider the model for datum i

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

where x_{ij} is the measured value of *covariate* j on experimental unit i. That is

$$y_i = \beta_0 + \sum_{j=1}^{\kappa} \beta_j x_{ij} + \epsilon_i$$

where the first two terms on the right hand side are the *systematic* or *deterministic* components, and the final term ϵ_i is the *random* component.

Example (k = 2)

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$$

A three parameter model.

Note: We can also include interaction terms

$$y_{i} = \beta_{0} + \beta_{1}x_{i1} + \beta_{2}x_{i2} + \beta_{12}(x_{i1} \cdot x_{i2}) + \epsilon_{i}$$

where

- The first two terms in x_{i1} and x_{i2} are main effects
- The third term in $(x_{i1} \cdot x_{i2})$ is an interaction

This is a four parameter model.

Multiple Linear Regression Examples	Subgroup analysis , with a factor predictor and a continuous covariate, is a form of interaction modelling; the factor predictor <i>interacts</i> with the covariate to modify the slope across the subgroups, for example.
SEE HANDOUT Multiple regression: Viscosity Example	We can describe the models using the notation previously introduced for ANOVA; consider the single binary factor predictor and single covariate case;
 Factor Regression: Interaction Residuals 	MODEL 0Single horizontal straight line1MODEL 1Two parallel horizontal X_2
 SPSS Instructions 	MODEL 2 Single straight line, X ₁
	MODEL 3 Two parallel straight lines, $X_1 + X_2$ non-zero slope
	MODEL 4 Two non-parallel straight lines $X_1 + X_2 + X_1 \cdot X_2$
1	99 200

For example, for factor predictor X_2 taking two levels and continuous covariate X_1 . When the pooled data are examined, a **positive association** between Y and X_1 is revealed.



Note: Always be on the lookout for *lurking* subgroups (subgroups determined by the levels of an unnoticed factor predictor)

Inferences can change radically when the lurking factor is included in the model

 positive association can be converted into negative association with the continuous covariate.

i.e. increasing X_1 decreases response in subgroup 1, and decreases response in subgroup 2, but appears to increase response overall.

This is known as **Simpson's Paradox in Regression**. It illustrates that pooling data over subgroups must be carried out with care !

 you must fit the factor predictor in the model if you suspect subgroup differences exist.

In the example, the problem arises due to **dependence** between X_1 and X_2 ; all the group with $X_2 = 0$ have **low** values of X_1 , whereas all the group with $X_2 = 1$ have **high** values of X_1

Dependence between covariates and factor predictors makes model fitting and results interpretation complicated.





Recap: we can build general models

$$y_{i} = \beta_{0} + \sum_{j=1}^{k} x_{ij} + \epsilon_{i}$$
We can fit each of these models easily using least-squares to obtain
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variation of y in terms of covariates and factor
predictors x_{1}, \dots, x_{k} .
$$explain the variate x_{1}, \dots, x_{k} .
$$explai$$

Interpreting
$$\hat{\beta}_j$$
 $\hat{\beta}_j$ can be interpreted as the amount of increase in response y
when x_j increases by one unit when the other predictors
 $x_1, x_2, \dots, x_{j-1}, x_{j+1}, \dots, x_k$ Test statistic: $t_j = \frac{\hat{\beta}_j}{s_{\beta_j}} = \frac{\text{ESTIMATE}}{\text{STANDARD ERROR}}$ $i_j, x_2, \dots, x_{j-1}, x_{j+1}, \dots, x_k$ are held fixed.We can test the hypothesis $H_0 : \beta_j = 0$
 $H_0 : \beta_j \neq 0$ using the usual hypothesis testing approach.207

2.2.2 Model Checking	For the multiple regression model, the ANOVA table takes the form
	SOURCE DF SS MS F
	REGRESSION k SSR MSR $F = \frac{MSR}{MSE}$
Using the General Linear Model approach to regression, we can fit models with different numbers of predictors, and	ERROR $n-k-1$ SSE MSE
► assess whether any individual covariate is influential in the model (look at	TOTAL $n-1$ SS
 assess whether there is any explanatory power in the variables combined (look at ANOVA statistics) 	where $MSR = \frac{SSR}{k}$ $MSE = \frac{SSE}{n-k-1}$
	the <i>F</i> statistic is $F = \frac{MSR}{MSE}$
	and if H_0 is true
	$F \sim Fisher-F(k, n-k-1)$
200	

Here

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_a : \text{ At least one } \beta_j \neq 0$$

The model for H_0 has one parameter β_0 . The model for H_a has k + 1 parameters

 $\beta_0, \beta_1, \beta_2, \ldots, \beta_k$

Therefore the number of extra parameters for model H_a is

(k+1)-1=k

i.e. to obtain model H_0 from model H_a we constrain k parameters to be zero.

Because we can constrain model H_a by setting some parameters equal to zero to obtain model H_0 , we say that

Model H_0 is nested inside Model H_a

The number, k, of constraints $\beta_1 = \beta_2 = \cdots = \beta_k = 0$ explains why the ANOVA table Regression degrees of freedom is k - the multiple regression brings in k extra parameters.

In addition, we can use the R^2 or Adjusted R^2 statistic to check overall model adequacy

$$R^2 = 1 - \frac{SSE}{SS_{yy}} = \frac{SS_{yy} - SSE}{SS_{yy}} = \frac{SSR}{SS}$$

which is equal to

VARIATION EXPLAINED BY THE REGRESSION TOTAL VARIATION

Also

Adj.
$$R^2 = 1 - \frac{SSE/(n-k-1)}{SS/(n-1)}$$

 $R^2 > 0.7$ implies that the model is a good fit, that is, most of the variation observed is explained by the regression model.



▶ Three-way Interactions: $x_{j_1} \cdot x_{j_2} \cdot x_{j_3}$

etc.

Dummy Variables

In SPSS, we can use the

General Linear Model \rightarrow Univariate

- pulldown menus to build and fit the model.
 - ▶ To fit factor predictors, we used the Fixed Factor option
 - ► To build models, we use the

 $Model \rightarrow Custom$

selections on the Univariate dialog

Recall that we can fit the factor predictor using the Linear Regression pulldown if we create **dummy variables**.

For example, if factor predictor X has L levels, we create L **new** binary predictors X_1, \ldots, X_L , where, for $I = 1, \ldots, L$

$$X_{l} = \begin{cases} 1 & \text{whenever } X = l \\ 0 & \text{otherwise} \end{cases}$$

We can then include X_1, \ldots, X_L in the regression model.

Example (L = 4) $X_2 \quad X_3$ $X \mid X_1$ X_4 3 0 0 0 1 1 0 0 0 1 3 0 0 1 0 4 0 0 0 1 2 0 1 0 0 2 0 1 0 0 See McClave and Sincich, Section 12.7.

2.2.3 Stepwise Model Selection

We seek a method that allows us to compare nested models.

Suppose we want to compare

MODEL 1 : $y = \beta_0 + \beta_1 x + \beta_2 x^2$ MODEL 2 : $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$

Model 1 is nested inside Model 2 as if we set $\beta_3 = 0$ in Model 2, we get Model 1.

ANOVA tests for Comparing Nested Models Terminology lf ► Complete Model MODEL 1 : $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ $E[Y] = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ MODEL 2 : $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12}(x_1.x_2)$ Reduced Model we can set $\beta_{12} = 0$ in Model 2 to obtain Model 1, so again the models are nested. $E[Y] = \beta_0 + \beta_1 x_1 + \dots + \beta_g x_g$ We can set up a hypothesis test to assess whether the where g < k. The reduced model is obtained from the complete simplification of Model 2 to Model 1 (by setting one or more model by setting parameters equal to zero) is justified by the data. $\beta_{g+1} = \beta_{g+2} = \cdots = \beta_k = 0$

Method

- 1. Fit the **complete model** and obtain the sum of squared errors, *SSE_C*, available from the ANOVA table.
- 2. Fit the **reduced model** and obtain the sum of squared errors, SSE_R , available from the ANOVA table.
- 3. Form the test statistic

$$F = \frac{(SSE_R - SSE_C)/(k - g)}{SSE_C/(n - k - 1)}$$

If H_0 is **true**, then $F \sim \text{Fisher-F}(k - g, n - k - 1)$

Note: k - g is the number of parameters we set equal to zero when moving from complete to reduced model.

Using this F statistic, we can assess whether there is evidence to support the reduced model over the complete model.

The reduced model is nested inside the complete model.

We wish to test the hypothesis

 $\begin{array}{rcl} H_0 & : & \beta_{g+1} = \beta_{g+2} = \cdots = \beta_k = 0 \\ H_a & : & \text{At least one of these } \beta_i \neq 0 \end{array}$

We can test this hypothesis by fitting both models, and combining the results; we focus on the sums of squares quantities.

SOURCE	DF	SS	MS	F	The result holds for comparing any two nested models where t
COMPLETE MODEL	k	SSR _C	MSR _C	F _C	standard model assumptions hold:
ERROR _C	n-k-1	SSE _C	MSE _C		• ϵ uncorrelated
					• ϵ independent of x_1, \ldots, x_k
	<u>n – 1</u>	55			• $\epsilon \sim N(0, \sigma^2)$ Note: It does not allow us to compare non-nested models: for
ced Model ANOVA	n – 1 table:	SS	MS	F	$\bullet \ \epsilon \sim N(0, \sigma^2)$ Note: It does not allow us to compare non-nested models; for example $MODEL \ 1 : \chi = \beta_0 + \beta_1 \chi_1 + \epsilon$
ced Model ANOVA SOURCE REDUCED MODEL	n – 1 table: DF	SS SS SSR _R	MS MSR _R	F F _R	$\bullet \ \epsilon \sim N(0, \sigma^2)$ Note: It does not allow us to compare non-nested models; for example $MODEL \ 1 : y = \beta_0 + \beta_1 x_1 + \epsilon$ $MODEL \ 2 : y = \beta_0 + \beta_2 x_2 + \epsilon$
ced Model ANOVA SOURCE REDUCED MODEL ERROR _R	n-1 table: DF g $n-g-1$	SS SSR _R SSE _R	MS MSR _R MSE _R	F F _R	$\bullet \ \epsilon \sim N(0, \sigma^2)$ Note: It does not allow us to compare non-nested models; for example $MODEL \ 1 : y = \beta_0 + \beta_1 x_1 + \epsilon$ $MODEL \ 2 : y = \beta_0 + \beta_2 x_2 + \epsilon$ $- \text{ NOT NESTED } !$

$$\begin{split} & \mathcal{F} = \frac{(SSE_R - SSE_C)/(k - g)}{SSE_C/(n - k - 1)} = \frac{(1)/2}{(3)/4} \\ & (1) - SSE_R - SSE_C: \text{ this is the improvement in fit when the reduced model is extended to the complete model} \\ & (2) - k - g: \text{ this is the number of extra parameters needed to extend the reduced model to the complete model} \\ & (3) - SSE_C \\ & (4) - n - k - 1 \\ & (3)/(4) - \text{ this is the best guess we have at the true value of σ^2 , that is, the estimate of σ^2 constructed using as much information as possible, once the effects of $\chi_1, \ldots, \chi_k \\ \text{ have been accounted for.} \end{split}$$$

Example (Hooker's Da	ta)		
COMPLETE MODEL	SSR _C SSE _C	2286.933 4.382	
REDUCED MODEL	SSR _R SSE _R	2272.474 18.840	
with $n = 31, k = 2, g =$	1		
$\Longrightarrow k$	-g = 1	, n - k - 1	= 28
So			
$F = \frac{(SSE_R - SSE_C)/}{SSE_C/(n-k)}$	$\frac{(k-g)}{-1)}$	$=\frac{(18.840)}{4.3}$	$\frac{-4.382)/1}{82/28} = 92.383$

Example (Hooker's Data)
We compare F with the
$Fisher-F(k-g,n-k-1) \equiv Fisher-F(1,28)$
distribution. $F_{0.05}(1, 28) = 4.20$
Thus $92.383 = F > F_{0.05}(1,28) = 4.20$
and $H_0: E[Y] = \beta_0 + \beta_1 x$ is REJECTED in favour of $H_a: E[Y] = \beta_0 + \beta_1 x + \beta_2 x^2$.
i.e. the quadratic model fits better than the straight-line model.

NOTE: From the original ANOVA tables, we already know that Model 1 and Model 2 both fit better than the null model

$$\begin{array}{rcl} \text{MODEL 0} & & E[Y] &=& \beta_0 \\ & & y &=& \beta_0 + \epsilon \end{array}$$

where there is no dependence on x.

We have now confirmed that Model 2 fits better than Model 1.

Example (Diabetes Data)

Factor Predictor: group (X_2) Continuous Covariate: loggluf (X_1) Response: logglut (Y)

We have five models to confirm:

Example (Diabetes Data) Example (Diabetes Data) MODEL 4 us the most complex model with 6 parameters In the SPSS parameterization: $\beta_{10}, \beta_{11}, \beta_{20}, \beta_{21}, \beta_{30}, \beta_{31}$ Group 3 Intercept and Slope β_{30}, β_{31} $\beta_{10}=\beta_{30}+\delta_{10}$ Changes in the Intercepts in MODEL 4: $eta_{20}=eta_{30}+\delta_{20}\quad \mbox{Groups 1 and 2 are }\delta_{10}\ \mbox{and }\delta_{20}$ $E[Y] = \begin{cases} \beta_{10} + \beta_{11} x_1 & \text{GROUP 1} \\ \beta_{20} + \beta_{21} x_1 & \text{GROUP 2} \end{cases}$ $\begin{array}{ll} \beta_{11}=\beta_{31}+\delta_{11} & \mbox{Changes in the Slopes in} \\ \beta_{21}=\beta_{31}+\delta_{21} & \mbox{Groups 1 and 2 are } \delta_{11} \mbox{ and } \delta_{21} \end{array}$ Thus the six new parameters are GROUP 3 $\beta_{30}, \beta_{31}, \delta_{10}, \delta_{20}, \delta_{11}, \delta_{21}$ All of the other models are nested inside Model 4; we can obtain them all by setting parameters equal to zero.

MODEL 0 $\beta_{31} = 0$ $\delta_{10} = \delta_{20} = \delta_{11} = \delta_{21} = 0$ MODEL 1 $\beta_{31} = \delta_{11} = \delta_{21} = 0$ MODEL 2 $\delta_{10} = \delta_{20} = \delta_{11} = \delta_{21} = 0$ MODEL 3 $\delta_{11} = \delta_{21} = 0$ Note: $\beta_{31} = 0 \Longrightarrow \delta_{11} = \delta_{21} = 0$, as X_1 is not included in the model. Counting Parameters Whenever we remove a continuous covariate, from a model, we set one parameter to zero. Whenever we remove a factor predictor with L levels from a model, we set L - 1 parameters to zero. Whenever we remove a two-way interaction between these variables from a model, we set 1.(L - 1) = L - 1 parameters to zero.

Models 0,1,2,3 are nested inside Model 4.

Two approaches to finding the best model are used:

- 1. Start with Model 0 and try to add terms that improve the model fit (Forward Selection)
- 2. Start with Model 4 and try to remove terms that improve the model fit (Backward Selection)

Note:

- Models 0,1 and 2 are nested inside Model 3.
- Model 0 is nested inside Models 1 and 2.

Therefore we can begin with Model 4, or Model 3 or Model 1 or 2, and simplify to a nested model.

Example (Diabetes Data)

Here n = 144. From SPSS output handouts:

Model	Description	SSE	р
0	1	28.504	1
1	X_2	4.160	3
2	X_1	3.738	2
3	$X_1 + X_2$	1.472	4
4	$X_1 + X_2 + X_1 \cdot X_2$	1.318	6

p is the number of non-zero parameters; k or g is always p - 1 in the following analysis.

236

Backward Selection: Complete Model : Model 4 i.e. Model 4 Reduced Model : Model 3 $X_1 + X_2 + X_1 \cdot X_2$ Here k = 5, g = 3 so k - g = 2, and fits significantly better than Model 3 n-k-1 = 144 - 5 - 1 = 138. $X_1 + X_2$. We have $F = \frac{(SSE_R - SSE_C)/(k - g)}{SSE_C/(n - k - 1)} = \frac{(1.472 - 1.318)/2}{1.318/138} = 8.062$ - we cannot simplify the complete model to the reduced model without the loss of significant explanatory power. The Interaction is Necessary in the Model We compare this with the Backward selection stops here; we cannot simplify further from the Fisher-F(k - g, n - k - 1) = Fisher-F(2, 138)complete model. distribution; we have $F_{\alpha}(2, 138) = 3.061$, so we Reject H_0 at $\alpha = 0.05$

Forward Selection: we start with Model 0 and build up.

Model 1 vs Model 0 F = 412.568

Model 2 vs Model 0 F = 940.846

It seems that Model 2 is the better improvement, so we try the selection path

 $\mathsf{Model}\ 0 \longrightarrow \mathsf{Model}\ 2 \longrightarrow \mathsf{Model}\ 3 \longrightarrow \mathsf{Model}\ 4$

Model	SSE	$SSE_R - SSE_C$
0	28.504	-
2	3.738	24.766
3	1.472	2.266
4	1.318	0.154

ie we work down the table, 28.504 - 3.738 = 24.766 etc.

Comparison $k g SSE_C$ $SSE_R - SSE_C$ F 24.766 940.82 2 vs 0 1 0 3.738 3 vs 2 3 1 1.472 2.266 107.76 4 vs 3 5 3 1.318 0.154 8.06

Recall that n = 144, and

 $F = \frac{(SSE_R - SSE_C)/(k - g)}{SSE_C/(n - k - 1)}$

Under each H_0 ,

$$F \sim \text{Fisher-F}(k-g, n-k-1)$$

- ► F_{0.05}(1, 142) ≃ 3.92 < 940.82 Therefore Model 0 is **NOT** an adequate simplification of Model 2
- ► F_{0.05}(2, 140) ≃ 3.07 < 107.76 Therefore Model 2 is **NOT** an adequate simplification of Model 3
- ► F_{0.05}(2,138) ≃ 3.07 < 8.06 Therefore Model 3 is **NOT** an adequate simplification of Model 4

All of the null hypotheses are **rejected**.

Therefore by both forward and backward selection, we select Model $\ensuremath{4}$

$$X_1 + X_2 + X_1 \cdot X_2$$

as the most appropriate model.

Note: In this sequence of hypothesis tests, the convention is **not** to correct for multiple testing (we don't know how many tests we are going to do), although a correction could be used.

F-tests for Unbalanced Designs	F-tests for Unbalanced Designs
Example (Potato Damage Data)	
The damage to potato plants caused by cold temperatures is to be studied. In this experimental study, three binary factor predictors were used: we label them A, B and C rather than X_1, X_2, X_3 to recall the link with Factorial Designs in ANOVA. Each factor takes two levels: $\frac{Factor}{A} = \frac{Levels}{Potato Variety} = 0$ -Variety 1, 1- Variety 2 B Acclimatization Routine 0- Room Temp, 1- Cold Room C Preparation Treatment 04C, 18C Thus we have a 2 × 2 × 2 three-way factorial design.	However, the design is unbalanced ; we have different numbers of replicates in each of the 8 factor-level combinations. This means we cannot use conventional 3-way ANOVA; the lack of balance means that the stated <i>p</i> -values may be misleading if we perform a standard ANOVA . Thus we are forced to use the General Linear Model F-test approach.
24.	244

Counting the numbers of parameters We begin with the most complex model and do backward selection. Term Parameters Here the most complex model can be written Α (a - 1)1 В (b - 1)1 A + B + C + A.B + A.C + B.C + A.B.CС (c - 1)1 A.B (a-1)(b-1)1 that is, A.C (a-1)(c-1)1 ▶ all main effects (terms 1,2 and 3) (b-1)(c-1)B.C 1 ▶ all two-way interactions (terms 4,5 and 6) A.B.C (a-1)(b-1)(c-1)1 ▶ all three-way interactions (term 7) Total We may write this model where a = b = c = 2. A * B * CWe have 7 parameters in total (excluding the baseline mean) when all terms are considered, so

which is termed the full factorial model.

k = 7

In the following tables columns are: Complete Model Reduced Model SSE_C SSE_R k g F (test statistic) $F_{0.05}(k - g, n - k - 1)$ We denote the critical value by F_{α} and check whether $F > F_{\alpha}$.

Potato Damage Data: ANOVA-F Tests

We compare four models: M_{R_1}, M_{R_2} and M_{R_3} are nested within the complete model $M_C.$

 $\begin{array}{rcl} M_{C} & : & A + B + C + A.B + A.C + B.C + A.B.C \\ M_{R_{1}} & : & A + B + C + A.B \\ M_{R_{2}} & : & A + B + C \\ M_{R_{3}} & : & A + B + A.B \end{array}$

COMP.	RED.	SSE_C	SSE_R	k	g	F	F_{α}
M _C	M_{R_1}	4968.876	5093.746	7	4	0.561	2.76
M_{R_1}	M_{R_2}	5093.746	7183.674	4	3	28.721	3.92
M_{R_1}	M_{R_3}	5093.746	6319.640	4	3	16.846	3.92

Note: The quoted F_α values are approximate as the textbook does not tabulate all Fisher-F distributions. We take $\alpha=0.05$

248

252

Conclusions	
 Taking the comparisons in order: 1. M_C vs M_{R1} : F < F_α. Therefore the result is not significant: Model M_{R1} is an adequate simplification of Model M_C, and we choose M_{R1} over M_C. The model M_{R1} now becomes the complete model. 2. M_{R1} vs M_{R2} : F > F_α. Therefore the result is significant: Model M_{R2} is not an adequate simplification of Model M_{R1} 3. M_{R1} vs M_{R3} : F > F_α. Therefore the result is significant: Model M_{R3} is not an adequate simplification of Model M_{R1} 	Thus the final model is A + B + C + A.B i.e. all main effects, plus the interaction between potato variety and acclimatization routine. We cannot simplify this model further without significant loss in terms of goodness of fit. Note: $R^2 = 0.631$ and Adjusted $R^2 = 0.610$, so we have a reasonable fit. 250

ask Distraction Data			
Example (Task Distraction Data)	Exampl	e (Task Dist	raction Data)
In an experimental study, the number of errors made in performing a specified task was recorded. The experiment investigated the		Group	2 : Delayed smoker 3 : Active smoker
influence of various predictors on the numbers of errors made. There are two factor predictors (A, B) and one continuous covariate (X) .	В	Task	1 : Pattern Recognition 2 : Cognitive Task 3 : Driving Simulation
We have a balanced design with 15 people (replicates) in each factor-level subgroup.	X Di	istraction Leve	el

We compare four models with the complete model

Complete Model : A * B * X

$$A + B + X + A.B + A.X + B.X + A.B.X$$

Number of parameters



For illustration we consider the following sequence of models:

▶ Reduced Model 1: M_{R_1}

$$A + B + X + A.X + B.X$$

► Reduced Model 2: M_{R2}

$$A + B + X + B.X$$

► Reduced Model 3: *M*_{*R*₃}

$$B + X + B.X$$

► Reduced Model 4: *M*_{*R*₄}

B + X

Task Distraction Data: ANOVA-F Tests Conclusions $: A + B + X + A \cdot B + A \cdot X + B \cdot X + A \cdot B \cdot X$ M_C Taking the comparisons in order: : A + B + X + A.X + B.X M_{R_1} 1. $M_C \text{ vs } M_{R_1}$: $F > F_{\alpha}$. Therefore the result is **significant**: M_{R_2} : A+B+X+B.XModel M_{R_1} is not an adequate simplification of Model M_C : B + X + B.X M_{R_3} M_{R_4} : B + X2. $M_{R_1} vs M_{R_2}$: $F < F_{\alpha}$. Therefore the result is not significant: Model M_{R_2} is an adequate simplification of Model M_{R_1} SSE_C SSE_R COMP. RED. k F F_{α} 3. M_{R_2} vs M_{R_3} : $F > F_{\alpha}$. Therefore the result is significant: g \overline{M}_{R_1} 5660.010 7627.479 17 5.084 2.02 M_C 9 Model M_{R_3} is not an adequate simplification of Model M_{R_2} \overline{M}_{R_2} 7627.479 7971.274 9 7 2.817 3 07 M_{R_1} 4. M_{R_3} vs M_{R_4} : $F > F_{\alpha}$. Therefore the result is significant: M_{R_2} M_{R_3} 7971.274 8404.654 7 5 3.452 3.07 Model M_{R_4} is not an adequate simplification of Model M_{R_3} \overline{M}_{R_3} M_{R_4} 8404.654 11154.715 5 3 21.105 3.07

Stepwise Selection in SPSS: Options Follow-up Analysis It is possible to carry out stepwise selection in SPSS using the Linear Regression pulldown menu, and the Method pulldown list. In a follow up analysis (see Handout), it transpires that the model ▶ Enter : All variables in a *block* are entered in a single step. A + B + X + A.B + A.X + B.X▶ Stepwise : At each step, the independent variable not in the equation that has the smallest p-value in the F-test is ie selected. entered, if that probability is sufficiently small. Variables **Note:** $R^2 = 0.863$ and Adjusted $R^2 = 0.831$, so we have a good already in the regression equation are removed if their fit. p-value becomes sufficiently large. The method terminates when no more variables are eligible for inclusion or removal. Note: we must take great care with the sequence of models. ▶ Remove : All variables in a block are removed in a single step.

Stepwise Selection in SPSS: Options

- ► **Backward :** Variables are entered into the equation and then sequentially removed. The variable with the smallest *partial correlation* with the dependent variable is considered first for removal. After the first variable is considered, the variable remaining in the equation with the smallest partial correlation is considered next. The procedure stops when there are no variables in the equation that satisfy the removal criteria.
- ► Forward : Variables are sequentially entered into the model starting from the null model. The first variable considered for entry into the equation is the one with the largest positive or negative correlation with the dependent variable. This variable is entered into the equation only if it satisfies the criterion for entry. If the first variable is entered, the independent variable not in the equation that has the largest partial correlation is considered next. The procedure stops when there are no variables that meet the entry criterion.

2.2.5 Pitfalls of Regression Modelling

Five issues to bear in mind in ANOVA, Regression and General Linear Modelling.

- 1. Model assumptions
- 2. Data transformations
- 3. Model selection
- 4. Multicollinearity
- 5. Predicting beyond the range of the covariates

See Handout.



25

3.1 Distribution-free tests for Categorical Data	Doll and Hill Data
Categorical data are data in which experimental units are allocated to one of a number of categories according to their characteristics. The categories are defined by one or more factors Examples: Female/Male - two categories Smoker/Former Smoker/Non Smoker - three categories. 	Table 13.11.Smokers and non-smokers among male cancer patients and controls (Doll and Hill 1950) Smokers Non-smokers TotalLung cancer 6472Controls62227649

Juvenile Delinquency and Spectacle-Wearing

Table 10.14Spectacle wearing among juvenile delinquenand non-delinquents who failed a vision test (Weindlingal1986)

		Juvenile delinquents	Non delinquents	Total
Spectacle wearers	Yes No	1 8	5 2	6 10
	Tota	9	7	16

The data are **counts** of experimental units that fall into each category. Suppose

- 1. There are n experimental units in the study
- 2. There are k categories
- 3. The probabilities of the k outcomes are p_1, \ldots, p_k , where

 $p_1 + \cdots + p_k = 1$

- 4. The experimental units are independent
- 5. The counts in the k categories are n_1, \ldots, n_k , where

 $n_1 + \cdots + n_k = n$

The experimental design is termed a Multinomial Experiment Note: The categories can arise as combinations of factor levels; we What kinds of tests can be carried out for such data ? can have 1. Tests about p_1, \ldots, p_k • one-way classification (categories of a single factor, A) • H_0 : $p_1 = \cdots = p_k = 1/k$ ► two-way classification (categories defined by combinations of ▶ H_0 : p_1, \ldots, p_k determined by some parametric distribution levels of two factors, A and B) (Normal, Poisson etc.) and so on. The counts table is often called a contingency table 2. Tests about the factors A and Band the entries in the table are called **cells**. ▶ are A and B dependent ? ▶ i.e. does classification by A influence classification by B. The idea can be extended to larger numbers of factors (A, B, C, \ldots) to produce a multi-way table. We will focus on at most two-way tables, with r rows and c columns, yielding an $r \times c$ table.

Chi-Squared Test

For one-way tables: suppose that a null hypothesis **completely specifies** p_1, \ldots, p_k , that is, we have H_0 of the form

$$H_0$$
: $p_1 = p_1^{(0)}, \ldots, p_k = p_k^{(0)}$

where $p_1^{(0)}, \ldots, p_k^{(0)}$ are fixed probabilities. For example, for k = 3,

$$H_0$$
: $p_1 = p_2 = p_3 = 1/3$

or

$$H_0$$
: $p_1 = 1/2, p_2 = p_3 = 1/4$

To test this hypothesis against the general alternative hypothesis

$$H_a$$
 : H_0 not true.

we use the test statistic

$$X^{2} = \sum_{i=1}^{k} \frac{\left(n_{i} - np_{i}^{(0)}\right)^{2}}{np_{i}^{(0)}}$$

If H_0 is true,

$$X^2 \sim \mathsf{Chi} ext{-squared}(k-1).$$

that is, X^2 is approximately distributed as Chi-squared(k-1).

In this formula

- ► *n_i* is the **observed** count in cell *i*
- $np_i^{(0)}$ is the **expected** count in cell *i* if H_0 is **true**.

Sometimes the formula is written

$$X^{2} = \sum_{i=1}^{k} \frac{(O_{i} - E_{i})^{2}}{E_{i}}$$

where O_i is the observed count, and E_i is the expected count.

lf

$$X^2 > \text{Chisq}_{\alpha}(k-1)$$

then we reject H_0 at the α significance level, where $\text{Chisq}_{\alpha}(k-1)$ is the $1 - \alpha$ (right-hand) tail critical value of the Chi-squared distribution with k - 1 degrees of freedom.

This method can be extended in the one-way case to test distribution assumptions, that is, for example

 H_0 : Data Normally distributed

or

H_0 : Data Poisson distributed

Unfortunately this facility is not available in SPSS; direct calculation is possible but involved.

For the two-way table, we can test the hypothesis

 H_0 : Factor A and Factor B levels are assigned independently

that is, classification by factor A is independent of classification by factor B. We use the same test statistic that can be rewritten

$$X^{2} = \sum_{i=1}^{r} \sum_{i=1}^{c} \frac{(n_{ij} - \widehat{n}_{ij})^{2}}{\widehat{n}_{ij}}$$

where

$$\widehat{n}_{ij} = rac{n_{i.}n_{.j}}{n}$$
 $n_{i.} = \sum_{j=1}^{c} n_{ij}$ $n_{.j} = \sum_{i=1}^{r} n_{ij}$

The terms $n_{i.}$ and $n_{.j}$ are the row and column totals for row i and column j respectively.

If H_0 is true

$$X^2 \sim \text{Chi-squared}((r-1)(c-1))$$

i.e. the degrees of freedom quantity is (r-1)(c-1). Otherwise the test proceeds as before: we check whether

$$X^2$$
 > Chisq _{α} (($r-1$)($c-1$))

and if so, we reject H_0 .

274

Example (DNA Sequence Data)

Counts of Nucleotides A,C,G,T in a genomic segment related to the breast cancer gene BRCA2.

Example (Eye and Hair Colour Data)

The assignment of hair and eye colour in a sample of humans

See handout.

Note: For the Chi-squared test to be valid, we need the expected cell counts

$$np_i^{(0)}$$
 $i=1,\ldots,k$

or

$$\hat{n}_{ij}$$
 $i = 1, ..., r, j = 1, ..., c$

to be sufficiently large. The convention is to require the expected value to be greater than $\ensuremath{\textit{five}}.$

Note: If r = c = 2 we have a 2×2 table, and another **exact** test can be used which does not rely on the large sample approximation

Fisher's Exact Test

- another test for independence of assignment of the row and column factor levels
- ► test statistic and null distribution are complicated (based on the hypergeometric distribution)
- ► SPSS computes test statistic and *p*-value.

Example (Juvenile Delinquency and Spectacle Wearing)

and a *p*-value of 0.013. Therefore we reject H_0 .

 $X^2 = 6.112$

 $\mathsf{Chi}\operatorname{-squared}_{0.05}(1) = 3.841$

Fisher's Exact Test: p-value is 0.035 (1-sided) or 0.024 (2-sided).

Thus we reject H_0 and we have evidence of association between

Compare with Chi-squared((r-1)(c-1)) =Chi-squared(1); we

Chi-squared Test:

have

the factors.

Example (Juvenile Delinquency and Spectacle Wearing) Is there any association between the two factors ? A : Spectacle Wearing (Yes/No) B : Juvenile Delinquent (Yes/No) Delinguent Yes No $| n_i |$ 5 Yes 1 6 Spectacles No 2 10 8 9 7 16 n_{.i}

A **case-control** study is an observational study where participants are selected for the study with regard to their **disease status**.

► a sample of **cases** (disease sufferers)

Case-Control Studies

► a sample of **controls** (healthy patients)

We investigate the possible association between disease status and a factor that takes two levels. A 2×2 table of counts is formed for all combinations of disease status/factor level.



The Chi-squared test is potentially not valid here because of the design. An alternative test statistic is based on the ${\bf odds\ ratio}$

$$\mathsf{O.R.} = \frac{n_{11}n_{22}}{n_{12}n_{21}} = \hat{\psi}$$

say. The test statistic is

$$Z = \frac{\log \widehat{\psi}}{\text{s.e.}(\log \widehat{\psi})}$$

where

s.e.
$$(\log \hat{\psi}) = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

That is,

$$Z = \frac{\log n_{11} + \log n_{22} - \log n_{12} - \log n_{21}}{\sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}}$$

Under

 H_0 : No association between factor and disease status

it follows that

 $Z \rightsquigarrow N(0,1)$

Here log means In or **natural log**.

Example (BCG Vaccination and Leprosy)

$$n_{11} = 101, n_{12} = 554, n_{21} = 159, n_{22} = 446$$

Therefore

 $\widehat{\psi}$

$$=\frac{n_{11}n_{22}}{n_{12}n_{21}}=0.511\qquad \log \widehat{\psi}=-0.671$$

and

so

s.e.
$$(\log \hat{\psi}) = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}} = 0.142$$

$$Z = \frac{-0.671}{0.142} = -4.717$$

For a text at lpha= 0.05, the two-sided critical value is ± 1.96 , so we

Reject H₀.

284

3.2 Single Population Tests Example (Smoking and Lung Cancer) $n_{11} = 647, n_{12} = 622, n_{21} = 2, n_{22} = 27$ Therefore $\log \widehat{\psi} = \log \frac{647 \times 27}{2 \times 622} = 2.642$ We seek non-parametric or distribution-free tests for hypotheses relating to single samples, the equivalents of one-sample Z- or T-tests, which rely on the **normality** of the samples. and s.e. $(\log \hat{\psi}) = \sqrt{\frac{1}{647} + \frac{1}{2} + \frac{1}{622} + \frac{1}{27}} = 0.735$ Normally these tests are formulated in terms of ranks of the data to give so $Z = \frac{2.642}{0.735} = 3.590$ **Rank Tests** For a text at $\alpha =$ 0.05, the two-sided critical value is ± 1.96 , so we Reject H₀ and report evidence for association.

For example, if the data are

0.24 3.16 1.97 2.10 0.92

we sort them into ascending order, and assign ranks in order

The tests depend on the behaviour of statistics computed in terms of the ranks, and rely on a **large sample** justification.

Rather than test the **mean**, we test the **median**, x_{MED} , where

$$\Pr[Observation \leq x_{MED}] = \frac{1}{2}$$

i.e. the halfway point of the distribution.

The sample median is the halfway point of the sorted sample.

Let $\boldsymbol{\eta}$ denote the population median. We wish to test, for example,

$$H_0$$
 : $\eta = \eta_0$

See Handout



Note: For the MWW test	
► Textbook convention : Label the samples so that n ₁ > n ₂ (i.e. sample 1 is the one with the larger sample size)	Other two sample tests are available: Kolmogorov-Smirnov Test
► SPSS convention : Label the samples such that x _{MED1} < x _{MED2}	 Moses Extreme Reactions Test Wald-Wolfowitz Runs Test None make distributional assumptions, all perform best when the
(i.e. sample 1 is the one with the smaller median) and only test $H_0~:~\eta_1=\eta_2$	sample size is large.
291	297

3.4 Comparing Two Dependent Samples	3.4 Comparing Three or More Populations
Suppose we have repeat measurements on the same experimental units. In this case, the within-subject data are dependent ; we have pairing of observations.	We now seek non-parametric equivalents to ANOVA useful for different designs. We study tests for (a) the Completely Randomized Design (CRD) (b) the Randomized Block Design (RBD) For (a) we use the
We can use the	Kruskal-Wallis Test
Wilcoxon Signed Rank Test	and for (b) we use the
See Handout	Friedman Test.
	See Handout

Summary of the Non-Parametric Tests

- ► Chi-Squared Test : Goodness of Fit/independence in contingency tables
- **Sign Test** : One Sample (equivalent of one sample *t*-test)
- ► Mann-Whitney-Wilcoxon : Two Sample (equivalent of two sample *t*-test)
- Wilcoxon Signed Rank : Paired Data
- ► Kruskal-Wallis : one-way layout, multigroup comparison equivalent of ANOVA for CRD.
- ► Friedman : two-way blocked layout, equivalent of two-way ANOVA for RBD.

Pros:

- No distributional assumptions
- Applicable for most sorts of data
- ► Large sample approximations make them easy to implement

Cons:

- ► Low power compared to parametric tests (i.e. often do not reject H_0 when they should - prone to Type II Error)
- ► Small sample null distributions difficult to compute.

3.6 Rank Correlation

To measure the association between two variables, we previously used the *correlation coefficient*, r; for data x_1, \ldots, x_n and

 $y_1, ..., y_n$,

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}}$$

where

$$SS_{xy} = \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y}) \quad SS_{xx} = \sum_{i=1}^{n} (x_i - \overline{x})^2 \quad SS_{yy} = \sum_{i=1}^{n} (y_i - \overline{y})^2$$

r is a measure of the linear association between X and Y

Pearson Product Moment Coefficient of Correlation

A more general measure of association is the

Spearman Rank Correlation Coefficient

We compute this as follows:

1. For each sample separately, compute the ranks of the data, denote the ranks for the data x_1, \ldots, x_n and y_1, \ldots, y_n by u_1, \ldots, u_n and v_1, \ldots, v_n respectively.

2. Compute

$$r_{S} = \frac{SS_{uv}}{\sqrt{SS_{uu}SS_{vv}}}$$

ie the Pearson correlation between the ranks.

 r_S is the **Spearman Correlation**.

Notes:

1. If there are no ties in the data

$$c_S = 1 - rac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)}$$

where $d_i = u_i - v_i$.

2. r_{S} is potentially a measure of the **non-linear** association between X and Y.

The calculation can be applied directly to rank data i.e. u_1, \ldots, u_n and v_1, \ldots, v_n can be preference ranks given by two observers.

Tests for r_{S}

To test

 $H_0 : \rho = 0$

vs

(1) H_a : $\rho > 0$ (2) H_a : $\rho < 0$ (3) H_a : $\rho \neq 0$

We may use r_{S} as a test statistic. The distribution of r_{S} under H_{0} is tabulated in McClave and Sincich.

The Role of Randomization and Permutation Tests

If Spearman_ α is the α tail quantile of the null distribution, we have the following rejection regions:

- (1) : Reject H_0 if $r_S > \text{Spearman}_{\alpha}$
- (2) : Reject H_0 if $r_S < -\text{Spearman}_{\alpha}$
- (3) : Reject H_0 if $|r_S| > \text{Spearman}_{\alpha/2}$

Randomization or **Permutation** procedures are useful for computing **exact** null distributions for non-parametric test statistics when sample sizes are small.

We focus first on two sample comparisons; suppose that two data samples $x_1 \ldots, x_{n_1}$ and $y_1 \ldots, y_{n_2}$ (where $n_1 \ge n_2$) have been obtained, and we wish to carry out a comparison of the two populations from which the samples are drawn. The Wilcoxon test statistic, W, is the sum of the ranks for the second sample. The permutation test proceeds as follows:

1. Let $n = n_1 + n_2$. Assuming that there are no ties, the pooled and ranked samples will have ranks

1 2 3 ... *n*

- 2. The test statistic is $W = R_2$, the rank sum for sample two items. For the observed data, W will be the sum of n_2 of the ranks given in the list above.
- 3. If the null hypothesis
 - H_0 : No difference between population 1 and population 2

were **true**, then we would expect **no pattern** in the arrangements of the group labels when sorted into ascending order. That is, the sorted data would give rise a **random** assortment of group 1 and group 2 labels.

30

- 4. To obtain the exact distribution of W under H₀ (which is what we require for the assessment of statistical significance), we could compute W for all possible permutations of the group labels, and then form the probability distribution of the values of W. We call this the **permutation null distribution**.
- But W is a rank sum, so we can compute the permutation null distribution simply by tabulating **all possible subsets** of size n₂ of the set of ranks {1, 2, 3, ..., n}.

6. There are

$$\binom{n}{n_2} = \frac{n!}{n_1! n_2!} = N$$

say possible subsets of size n_2 . For example, for n = 6 and $n_2 = 2$, the number of subsets of size n_2 is

$$\binom{8}{2} = \frac{8!}{6! \ 2!} = 28$$

However, the number of subsets increases dramatically as n increases; for $n_1 = n_2 = 10$, so that n = 20, the number of subsets of size n_2 is

$$\binom{20}{10} = \frac{20!}{10! \ 10!} = 184756$$

7. The exact rejection region and *p*-value are computed from the permutation null distribution. Let W_i , i = 1, ..., N denote the value of the Wilcoxon statistic for the *N* possible subsets of the ranks of size n_2 . The probability that the test statistic, W, is less than or equal to w is

$$\Pr[W \le w] = \frac{\text{Number of } W_i \le w}{N}$$

We seek the values of w that give the appropriate rejection region, \mathcal{R} , so that

$$\Pr[W \in \mathcal{R}] = \frac{\text{Number of } W_i \in \mathcal{R}}{N} = \alpha$$

It may not be possible to find critical values, and define \mathcal{R} , so that this probability is **exactly** α as the distribution of W is **discrete**.

	Ranks W	Ranks	W Ranks W Ranks W
	1 2 3 6	1 7 8	16 2 7 10 19 4 6 7 17
	1 2 4 7	1 7 9	17 2 8 9 19 4 6 8 18
	1 2 5 8	1 7 10	18 2 8 10 20 4 6 9 19
	1 2 6 9		
Simple Example	1 2 10 13	2 3 5	
Suppose $n_1 = 7$ and $n_2 = 3$. There are	1 3 4 8	2 3 6	
Suppose $n_1 = r$ and $n_2 = 5$. There are	1 3 5 9	2 3 7	12 3 4 10 17 4 9 10 23
	1 3 6 10	2 3 8	13 3 5 6 14 5 6 7 18
(10) 101	1 3 7 11	2 3 9	14 3 5 7 15 5 6 8 19
$\binom{10}{10} - \frac{10!}{10} - \frac{120}{10}$	1 3 8 12	2 3 10	15 3 5 8 16 5 6 9 20
$\left(\frac{3}{3}\right) = \frac{7}{7131} = 120$	1 3 9 13	2 4 5	
(3) 1:3:	1 3 10 14		
	1 4 5 10	2 4 7	
subsets of the ranks $\int 1/2$ $\int 1/2$ $\int 1/2$ of size 3. The subsets are	1 4 7 12	2 4 9	
	1 4 8 13	2 4 10	
listed below, together with the rank sums.	1 4 9 14	2 5 6	13 3 7 8 18 5 9 10 24
	1 4 10 15	2 5 7	14 3 7 9 19 6 7 8 21
	1 5 6 12	2 5 8	15 3 7 10 20 6 7 9 22
	1 5 7 13	2 5 9	16 3 8 9 20 6 7 10 23
	1 5 8 14	2 5 10	17 3 8 10 21 6 8 9 23
	1 5 9 15		
	1 5 10 16		
	1 6 9 16	2 6 10	
	1 6 9 16	$\begin{vmatrix} 2 & 0 & 10 \\ 2 & 7 & 8 \end{vmatrix}$	
	1 6 10 17	2 7 9	

There are 22 possible rank sums, $\{6, 7, 8, \ldots, 25, 26, 27\}$; the number of times each is observed is displayed in the table below, with the corresponding probabilities and cumulative probabilities.

W	6	7	8	9	10	11	12	13	14	15	16
Frequency	1	1	2	3	4	5	7	8	9	10	10
Prob.	0.008	0.008	0.017	0.025	0.033	0.042	0.058	0.067	0.075	0.083	0.083
Cumulative Prob.	0.008	0.017	0.033	0.058	0.092	0.133	0.192	0.258	0.333	0.417	0.500
W	17	18	19	20	21	22	23	24	25	26	27
W Frequency	17 10	18 10	19 9	20 8	21 7	22 5	23 4	24 3	25 2	26 1	27 1
W Frequency Prob.	17 10 0.083	18 10 0.083	19 9 0.075	20 8 0.067	21 7 0.058	22 5 0.042	23 4 0.033	24 3 0.025	25 2 0.017	26 1 0.008	27 1 0.008

Thus, for example, the probability that W = 19 is 0.075, with a frequency of 9 out of 120. From this table, we deduce that

 $\Pr[8 \le W \le 25] = 0.983 - 0.033 = 0.950$

implying that the two-sided rejection region for $\alpha=$ 0.05 is the set $\mathcal{R}=\{6,7,26,27\}.$

Placenta Permeability Data

Placenta Permeability Data

Example

Thus the Wilcoxon statistic is

$$W = R_2 = 2 + 5 + 6 + 8 + 9 = 30$$

Now, here $n_1 = 10$ and $n_2 = 5$. There are

$$\binom{15}{5} = \frac{15!}{10! \, 5!} = 3003$$

subsets of the ranks $\{1, 2, 3, \dots, 15\}$ of size 5.

In the permutation null distribution, the possible values of W are $\{15, 16, \ldots, 64, 65\};$ the probabilities are given below.

Placenta Permeability Data

Example

W	15	16	17	18	19	20	21	22	23	24	25	26	27
Frequency	1	1	2	3	5	7	10	13	18	23	30	36	45
Prob.	0.000	0.000	0.001	0.001	0.002	0.002	0.003	0.004	0.006	0.008	0.010	0.012	0.015
Cumulative Prob.	0.000	0.001	0.001	0.002	0.004	0.006	0.010	0.014	0.020	0.028	0.038	0.050	0.065
W	28	29	30	31	32	33	34	35	36	37	38	39	40
Frequency	53	63	72	83	92	103	111	121	127	134	137	141	141
Prob.	0.018	0.021	0.024	0.028	0.031	0.034	0.037	0.040	0.042	0.045	0.046	0.047	0.047
Cumulative Prob.	0.082	0.103	0.127	0.155	0.185	0.220	0.257	0.297	0.339	0.384	0.430	0.477	0.523
W	41	42	43	44	45	46	47	48	49	50	51	52	53
Frequency	141	137	134	127	121	111	103	92	83	72	63	53	45
Prob.	0.047	0.046	0.045	0.042	0.040	0.037	0.034	0.031	0.028	0.024	0.021	0.018	0.015
Cumulative Prob.	0.570	0.616	0.661	0.703	0.743	0.780	0.815	0.845	0.873	0.897	0.918	0.935	0.950
W	54	55	56	57	58	59	60	61	62	63	64	65	
Frequency	36	30	23	18	13	10	7	5	3	2	1	1	
Prob.	0.012	0.010	0.008	0.006	0.004	0.003	0.002	0.002	0.001	0.001	0.000	0.000	
Cumulative Prob.	0.962	0.972	0.980	0.986	0.990	0.994	0.996	0.998	0.999	0.999	1.000	1.000	

Placenta Permeability Data

Example

By inspection of the table, we see that

$$\Pr[25 \le W \le 55] = 0.972 - 0.038 = 0.934$$

and

$$\Pr[24 \le W \le 56] = 0.980 - 0.028 = 0.952$$

