

**STATISTICS****CONFIDENCE INTERVAL**

- Then the  $100(1 - \alpha)\%$  CI for  $\mu_1$  is:

**KNOWN POP VAR**  
 Population variance  $\sigma^2$ :  $\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

**UNKNOWN POP VAR**  
 Sample variance  $s^2$ :  $\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$

where  $t$  has  $v = (n - 1)$  degrees of freedom.

Then the  $(100 - \alpha)\%$  CI for  $(\mu_1 - \mu_2)$  is: MEAN  $\bar{w}$  1 SAMPLE

**KNOWN POP VAR**  
 Known population variances  $\sigma_1^2$  &  $\sigma_2^2$ :  $(\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

**UNKNOWN VAR**  
 Unknown population variances:  $(\bar{x}_1 - \bar{x}_2) \pm t_{\alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$

**Sample Variances**  
 $S_1^2, S_2^2$ :  $\left(\frac{\sum(x_i - \bar{x}_1)^2}{n_1 - 1}\right) + \left(\frac{\sum(x_i - \bar{x}_2)^2}{n_2 - 1}\right)$

Where  $t$  has  $v = n - 1$  degrees of freedom.  
 (round to nearest int)

Before & After MEAN  $\bar{w}$  2 PAIRED SAMPLES

- Then the  $(100 - \alpha)\%$  CI for  $\mu_d = (\mu_1 - \mu_2)$  is: Have?  
 $d \pm t_{\alpha/2} \frac{s_d}{\sqrt{n}}$

Where  $t$  has  $v = (n - 1)$  degrees of freedom.

\*Large sample:  $np \geq 5$  &  $n(1-p) \geq 5$

**ONE-SAMPLE PROPORTION**

Then the  $(100 - \alpha)\%$  CI for  $p$  is:

Have?  
 Size  $n$   
 Success  $\hat{p}$   
 Failure  $1 - \hat{p}$

$$\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

**TWO-SAMPLE DIFFERENCE IN PROPORTIONS**

Then the  $(100 - \alpha)\%$  CI for  $(p_1 - p_2)$  is:

Have?  
 Size  $n_1$   
 Success  $\hat{p}_1$   
 Size  $n_2$   
 Success  $\hat{p}_2$

$$(\hat{p}_1 - \hat{p}_2) \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}}$$

$$= n_1 \hat{p}_1 (1 - \hat{p}_1), n_2 \hat{p}_2 (1 - \hat{p}_2) \geq 5$$

**VARIANCE  $\bar{w}$  1 SAMPLE**

• Then the  $(100 - \alpha)\%$  CI for  $\sigma^2$  is (note the asymmetry):

$$\frac{(n-1)s^2}{X^2_{\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{X^2_{1-\alpha/2}}$$

Where both  $X^2$  values have  $v = n - 1$  degrees of freedom.

\*normally distributed populations

• Then the  $(100 - \alpha)\%$  CI for  $\sigma_1^2/\sigma_2^2$  is:

$$\frac{s_1^2}{s_2^2 f_{\alpha/2}(v_1, v_2)} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{s_1^2 f_{\alpha/2}(v_2, v_1)}{s_2^2}$$

With degrees of freedom  $v_1 = (n_1 - 1)$  and  $v_2 = (n_2 - 1)$ .

**PREDICTION INTERVAL**

\*Normally distributed data

- Then the  $(100 - \alpha)\%$  PI for a future observation  $x_0$  is:

**KNOWN POP VAR**  
 Population variance  $\sigma^2$ :  $\bar{x} \pm z_{\alpha/2} \sqrt{1 + 1/n}$

**UNKNOWN POP VAR**  
 Sample variance  $s^2$ :  $\bar{x} \pm t_{\alpha/2} s \sqrt{1 + 1/n}$

With  $t$  having  $v = (n - 1)$  degrees of freedom.

**ESTIMATES**

To estimate some population parameter  $\theta$ , we want to find a point estimate  $\hat{\theta}$ , which is a single value of a statistic  $\hat{\theta}$ .

A statistic  $\hat{\theta}$  is an unbiased estimator of the parameter  $\theta$  if:

$$E[\hat{\theta}] = \theta$$

**Binomial Experiment**

P - population success proportion

$\hat{p} = \frac{x}{n}$  - Sample success proportion

$\hat{P}$  - statistic =  $\{\hat{p}_1, \hat{p}_2, \dots, \hat{p}_n\}$

An interval estimate of parameter  $\theta$  is an interval

$$\hat{\theta}_L \leq \theta \leq \hat{\theta}_U$$

where  $\hat{\theta}_L, \hat{\theta}_U$  depend on the statistic  $\hat{\theta}$  being sampled as well as its sampling distribution.

- The interval  $[\hat{\theta}_L, \hat{\theta}_U]$  is referred to as the  $100(1 - \alpha)\%$  confidence interval (CI) = central point  $\pm$  margin of estimate

**ERROR TYPES**

- Type I error: reject  $H_0$  when it is true  
 to reduce  $\alpha$   
 ↑ Sample size  
 ↓ critical region size
- Type II error: fail to reject (FTR)  $H_0$  when it is false  
 to reduce  $\beta$   
 ↑ Sample size  
 ↑ critical region size
- Statistical Power: the probability of correctly rejecting  $H_0$  =  $(1 - \beta)$

In reality, Ho is:

$H_0$	True	False
Reject	Type I error ( $\alpha$ )	Power $(1 - \beta)$
Do Not Reject	Confidence level $(1 - \alpha)$	Type II error ( $\beta$ )

A hypothesis is a conjecture of 1+ populations

**TESTING PROCEDURE + 1.V 2 sided**

Hypothesis testing procedure

- Come up with a null hypothesis ( $H_0$ ): one that captures the belief we are prepared to accept in absence of data
- Come up with an alternate hypothesis ( $H_A$ ): one that we would accept if the null hypothesis is false
- Based on test statistic, reject or FTR the null hypothesis

A hypothesis is 2-sided if it specifies that the parameter value is exactly equal to a given value

- The mean height of MIE237 students is 175 cm
- The variance of MIE237 student heights is 10 cm

A hypothesis is 1-sided if it specifies that the parameter value is at least OR at most a given value

- The mean height of MIE237 students is at most 177 cm
- The variance of MIE237 student heights is at least 10 cm

if  $p$ -value is greater  $\alpha$  - FTR  $H_0$   
 (not in rejection region)

if  $p$ -value is less than  $\alpha$  - Reject  $H_0$   
 (in rejection region)

To Enter SD Mode MODE MODE 1

Clear input data SHIFT MODE 1 = (clear)(sc)

TYPE num then Mt to Input values

**HYPOTHESIS TESTS****One-Sample**

one-sample z-test

The test statistic  $z$  is given as:

Known Population variance  $\sigma^2$ :  $z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$

Have?  
 Size  $n$   
 mean  $\bar{x}$

We reject  $H_0$  if:

- $H_0$  is  $\mu = \mu_0$  AND either  $z > z_{1-\alpha/2}$  or  $z < z_{\alpha/2}$
- $H_0$  is  $\mu \leq \mu_0$  AND  $z > z_{1-\alpha}$
- $H_0$  is  $\mu \geq \mu_0$  AND  $z < z_\alpha$

**one-sample t-test**

The test statistic  $t$  is given as:

Sample variance  $s^2$ :  $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$

Have?  
 Size  $n$   
 mean  $\bar{x}$

With  $v = n - 1$  degrees of freedom.

We reject  $H_0$  if:

- $H_0$  is  $\mu = \mu_0$  AND either  $t > t_{\alpha/2}$  or  $t < -t_{\alpha/2}$
- $H_0$  is  $\mu \leq \mu_0$  AND  $t > t_\alpha$
- $H_0$  is  $\mu \geq \mu_0$  AND  $t < -t_\alpha$

**Two-Sample****Two-Sample z-test**

Then the  $z$  statistic is:

Known Population variance  $\sigma^2$ :  $z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{(\sigma_1^2/n_1) + (\sigma_2^2/n_2)}}$

Have?  
 Size  $n_1$   
 mean  $\bar{x}_1$   
 Size  $n_2$   
 mean  $\bar{x}_2$

We reject  $H_0$  if:

- $H_0$  is  $(\mu_1 - \mu_2) = d_0$  AND either  $z > z_{1-\alpha/2}$  or  $z < z_{\alpha/2}$
- $H_0$  is  $(\mu_1 - \mu_2) \leq d_0$  AND  $z > z_{1-\alpha}$
- $H_0$  is  $(\mu_1 - \mu_2) \geq d_0$  AND  $z < z_\alpha$

**two-sample t-test**

Then the  $t$  statistic is given as  $t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{(\frac{s_1^2}{n_1}) + (\frac{s_2^2}{n_2})}}$

Have?  
 Size  $n_1$   
 mean  $\bar{x}_1$   
 Size  $n_2$   
 mean  $\bar{x}_2$

We reject  $H_0$  if one of these apply, and FTR otherwise

- $H_0$  is  $(\mu_1 - \mu_2) = d_0$  AND either  $t > t_{\alpha/2}$  or  $t < -t_{\alpha/2}$
- $H_0$  is  $(\mu_1 - \mu_2) \leq d_0$  AND  $t > t_\alpha$
- $H_0$  is  $(\mu_1 - \mu_2) \geq d_0$  AND  $t < -t_\alpha$

**Paired t-test**

Then the  $t$  statistic with  $v = (n - 1)$  degrees of freedom is:

Have:  
 Mean  
 Paired Sample of size  $n$   
 mean difference  $\bar{d}$   
 SD  $s_d$

$$t = \frac{\bar{d} - d_0}{s_d/\sqrt{n}}$$

We reject  $H_0$  if one of these three apply; otherwise we FTR

- $H_0$  is  $\bar{d} = d_0$  AND either  $t > t_{\alpha/2}$  or  $t < -t_{\alpha/2}$
- $H_0$  is  $\bar{d} \leq d_0$  AND  $t > t_\alpha$
- $H_0$  is  $\bar{d} \geq d_0$  AND  $t < -t_\alpha$

**Chi-squared goodness of fit test**

Then the  $X^2$  statistic is:  $X^2 = \sum_i \frac{(o_i - e_i)^2}{e_i}$  Have: expected & observed dist'n in K bins frequency per bin  $e_i$  &  $o_i$

- With  $v = (k - 1)$  degrees of freedom

We reject  $H_0$  if  $X^2 > X^2_\alpha$  for a pre-selected  $\alpha$

**Chi-squared of independence/homogeneity**

The expected value of each cell is given as:

$$e_{ij} = \frac{(\text{row } i \text{ total})(\text{column } j \text{ total})}{\text{grand total}}$$

\*all  $e_{ij} \geq 5$  for all  $i, j$

The  $X^2$  statistic for independence is given as:

$$X^2 = \sum_{i=1}^{g-1} \sum_{j=1}^{h-1} \frac{(o_{ij} - e_{ij})^2}{e_{ij}} \text{ frequency of cell } e_{ij} \& o_{ij}$$

With  $X^2$  having  $v = (|I| - 1)(|J| - 1)$  degrees of freedom.

We reject  $H_0$  if  $X^2 > X^2_\alpha$

\*Special case: if  $|I| = |J| = 2$  (i.e., a 2x2 table), then: **YATES CORRECTION**

$$X^2 = \sum_{i=1}^{g-1} \sum_{j=1}^{h-1} \frac{(|o_{ij} - e_{ij}| - 0.5)^2}{e_{ij}}$$

# STATISTICS

Connection to experiment design!

$$SSR = \text{sum of squares (regression)} = \sum(\hat{y}_i - \bar{y})^2 = b_1 S_{XY}$$

- Amount of error incorporated by the linear regression model as opposed to simply fitting the mean value of  $y$

$$SSE = \text{sum of squares (error)} = \sum(y_i - \hat{y}_i)^2$$

- Amount of error between the regression model and actual data

$$SST = \text{sum of squares (total)} = \sum(y_i - \bar{y})^2 = SSR + SSE = S_{YY}$$

- Amount of total error in the data if fitting to the mean value of  $y$

## LR MODEL

The general form of a simple linear regression is:

Note:  $X$  &  $t$  are RV  
 $\epsilon$  is normally distributed w/ mean 0, variance  $\sigma^2$   
 $Y = \beta_0 + \beta_1 X + \epsilon$

### ORDINARY LEAST SQUARES

Unconstrained, non-linear program

$$\begin{aligned} b_1 &= \frac{n \sum_i (x_i y_i) - \sum_i (x_i) \sum_i (y_i)}{n \sum_i (x_i^2) - (\sum_i (x_i))^2} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{S_{XY}}{S_{XX}} \\ b_0 &= \frac{\sum_i y_i - b_1 \sum_i x_i}{n} = \bar{y} - b_1 \bar{x} \end{aligned}$$

$$S_{XX} = \sum_{i=1}^n (x_i - \bar{x})^2 = (S_{xx})^2 / (n-1)$$

$$S_{YY} = \sum_{i=1}^n (y_i - \bar{y})^2 = (S_y)^2 / (n-1)$$

$$S_{XY} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n}$$

$$S^2 = \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{n-2} = \frac{S_{YY} - b_1 S_{XY}}{n-2}$$

mean square error

## CALCULATOR

MODE MODE 2 1 enter linear regression mode  
 SHIFT CLR 1 (Sc1) clear statistical memory  
 <x-data> <y> <y> <T> enter statistical data

To recall this type of value: Perform this key operation:

$\Sigma x^2$	SHIFT S-SUM 1
$\Sigma x$	SHIFT S-SUM 2
$n$	SHIFT S-SUM 3
$\Sigma y^2$	SHIFT S-SUM ▶ 1
$\Sigma y$	SHIFT S-SUM ▶ 2
$\Sigma xy$	SHIFT S-SUM ▶ 3
$\bar{x}$	SHIFT S-VAR 1
$x\sigma_n$	SHIFT S-VAR 2
$x\sigma_{n-1}$	SHIFT S-VAR 3
$\bar{y}$	SHIFT S-VAR ▶ 1
$y\sigma_n$	SHIFT S-VAR ▶ 2
$y\sigma_{n-1}$	SHIFT S-VAR ▶ 3
intercept $\beta_0$	SHIFT S-VAR ▶ 1, 1
Regression coefficient A $\beta_0$	SHIFT S-VAR ▶ 1, 2
Regression coefficient B $\beta_1$	SHIFT S-VAR ▶ 2, 2

## LINEAR REGRESSION

### ANOVA

Source of Variation	Sum of Squared Error	Degrees of Freedom	Mean Squared Error	F statistic
Regression	SSR	1	$\frac{SSR}{1} = SSR$	$\frac{SSR}{S^2} = f_{(v_1=1, v_2=n-2)}$
Error	SSE	$n - 2$	$\frac{SSE}{n-2} = s^2$	
Total	SST	$n - 1$		

- We reject  $H_0: \beta_1 = 0$  if  $f > f^*_{(v_1=1, v_2=n-2, \alpha)}$
- Upper-tailed: the rejection region is only if  $f > f^*$  (so  $\alpha$  is not halved)

C1

- The  $100(1 - \alpha)\%$  CI for  $\beta_1$  in the model  $Y = \beta_0 + \beta_1 x$  is:

$$b_1 \pm t_{\alpha/2} \frac{s}{\sqrt{S_{XX}}}$$

where  $t$  has  $(n - 2)$  degrees of freedom.  
 The  $t$  statistic for testing  $\beta_1$  is:

$$t = \frac{b_1 - \beta_{10}}{s/\sqrt{S_{XX}}}$$

HYPOTEST

where  $t$  has  $v = (n - 2)$  degrees of freedom.

We reject  $H_0$  if:

- $H_0: \beta_1 = \beta_{10}$  AND  $t > t_{\alpha/2}$  or  $t < -t_{\alpha/2}$
- $H_0: \beta_1 \leq \beta_{10}$  AND  $t > t_{\alpha}$
- $H_0: \beta_1 > \beta_{10}$  AND  $t < -t_{\alpha}$

FTR  $\beta_1$  suggests there is NO LINEAR relationship b/w x & y

B1

"correlation does not equal causation"  
 Just bc the model fits, does NOT mean that  $x$  results in definite change in  $y$ .  
 LR can be a causal model but isn't always one

C1

Anscombe's quartet is a good counterexample for why we should visualize & NOT rely solely on metrics

All 4 plots on the right have the same:

- Linear regression Line  $\hat{y} = 3 + 0.5x$
- Sample variances  $S_{xx}$  &  $S_{yy}$
- Correlation  $r = 0.82$

## LOGISTIC REG MODEL

Logistic regression function:

- Let's define the odds of success as  $\frac{\hat{p}}{1-\hat{p}}$

$$Y = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k)}}$$

Our fitted model then becomes:

$$\hat{y} = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + \dots + b_k x_k)}}$$

This means that as we change from  $\hat{p}_1$  to  $\hat{p}_2$ , the odds of success changes by a factor of  $\frac{\hat{p}_2/(1-\hat{p}_2)}{\hat{p}_1/(1-\hat{p}_1)}$

### Example

Let's say we are trying to predict whether students pass a course based on how much they study.

Hours 0.50 0.75 1.00 1.25 1.50 1.75 2.00 2.25 2.50 2.75 3.00 3.25 3.50 4.00 4.25 4.50 4.75 5.00 5.25

Pass 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 \text{Hours})}}$$

Our model is:

- Now let's run a logistic regression in R:

- This shows us that:

- $\beta_0 = -4.0777$

- $\beta_1 = 1.5046$

- Note that we don't have a specific F statistic for this output, but our individual coefficients are significant at the 0.05 level

```
> Logistic <- glm(Pass ~ hours, family = binomial)
Call:
glm(formula = Pass ~ hours, family = "binomial")

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -4.0777    0.6387 -6.3877  0.0000 ***
hours        1.5046    0.0325  46.8850  0.0000 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
(Dispersion parameter for binomial family taken to be 1)

Null deviance: 27.729 on 19 degrees of freedom
Residual deviance: 26.860 on 18 degrees of freedom
AIC: 28.06
```

as  $(\text{hours} + 1, \text{Pass}) \sim \text{Logit}$

## OUTLIER ANALYSIS

SD of residuals

$$s_r = \sqrt{\frac{\sum(y_i - \hat{y}_i)^2}{n-2}}$$

Standardized res.

$$StR_i = \frac{e_i}{s_r}$$

Studentized res.

$$StR_i = \frac{e_i}{s_r}$$

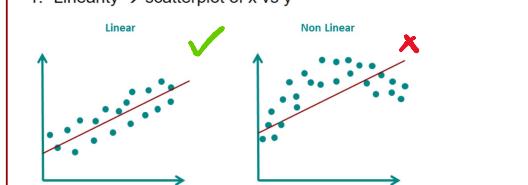
if  $|StR_i| > 3$ , i could be an outlier

Causes of outliers:

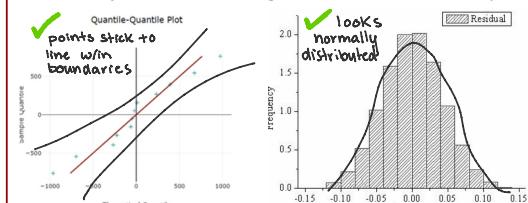
- Measurement error REMOVE + RERUN LR
- Extreme but part of normal process KEEP
- Extreme and not part of normal process REMOVE + RERUN LR

## ASSUMPTIONS OF LR

- Linearity  $\rightarrow$  scatterplot of  $x$  vs  $y$



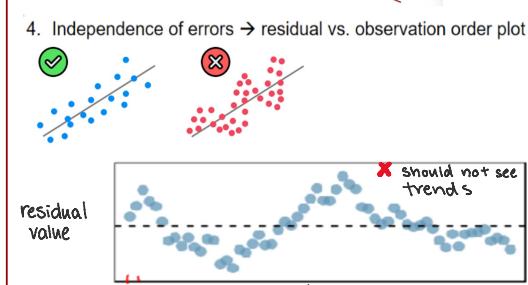
- Normality of residuals  $\rightarrow$  histogram of residuals, Normal QQ plot



- Constant variance  $\rightarrow$  residual vs. fitted value plot



- Independence of errors  $\rightarrow$  residual vs. observation order plot



# STATISTICS

## MULTIPLE LINEAR REGRESSION

### MULTIPLE LINEAR REG

The multiple linear regression model is given as:

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon$$

Where we now have  $k$  many explanatory variables and therefore  $k$ -many  $\beta$  coefficients

Explanatory variables have power of 1 & no interaction is each other

### ORDINARY LEAST SQUARES

Let the general MLR equation be:

$$y = b^T X + e$$

$$\text{Where } y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, b = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_k \end{bmatrix}, X = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{k,1} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1,n} & \dots & x_{k,n} \end{bmatrix}, e = \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}$$

Then the solution of the point estimates vector  $b$  is:

$$b = (X^T X)^{-1} X^T y$$

$$S^2 = \frac{SSE}{n-k-1} = \frac{\sum_i^n (y_i - \hat{y}_i)^2}{n-k-1} \quad \text{Variance of residuals}$$

- The  $100(1-\alpha)\%$  CI for  $\beta_i$  is:

$$b_i \pm t_{v,\alpha/2} \sqrt{S^2 c_{ii}}$$

Where  $t$  has  $v = n - k - 1$  degrees of freedom.

The  $t$  statistic with  $v = n - k - 1$  degrees of freedom is:

$$t = \frac{b_i - \beta_{i0}}{\sqrt{c_{ii}}}$$

We reject  $H_0$  if:

- $H_0: \beta_i = \beta_{i0}$  AND  $t > t_{v,\alpha/2}$  or  $t < -t_{v,\alpha/2}$
- $H_0: \beta_i \leq \beta_{i0}$  AND  $t > t_{v,\alpha}$
- $H_0: \beta_i > \beta_{i0}$  AND  $t < -t_{v,\alpha}$

### ASSUMPTIONS OF MLR

- Linearity
- Normality of residuals
- Constant variance
- Independence of errors

However, now we have an additional assumption of:

- No multicollinearity → matrix plot of correlations (correlogram)

avoid correlations  $> |0.7|$

Correlation between  $x_i$  and  $x_j$ :

$$r_{x_i x_j} = \frac{S_{x_i x_j}}{\sqrt{S_{x_i x_i} S_{x_j x_j}}}$$

### ANALYSIS OF VARIANCE APPROACH

$$SSR = \text{sum of squares (regression)} = \sum(\hat{y}_i - \bar{y})^2$$

$$SSE = \text{sum of squares (error)} = \sum(y_i - \hat{y}_i)^2$$

$$SST = \text{sum of squares (total)} = \sum(y_i - \bar{y})^2 = SSR + SSE$$

Source of Variation	Sum of Squared Error	Degrees of Freedom	Mean Squared Error	F statistic
Regression	SSR	$k$	$\frac{SSR}{k} = MSR$	$\frac{MSR}{MSE} = f_{(v_1=k, v_2=n-k-1)}$
Error	SSE	$n - k - 1$	$\frac{SSE}{n - k - 1} = MSE$	
Total	SST	$n - 1$		

$H_0: \beta_1 = \dots = \beta_k = 0$   
 $H_A: \text{At least one } \beta_i \neq 0$   
 We reject  $H_0$  if  $f > f_{v_1, v_2, \alpha}^*$

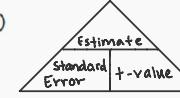
> summary(Full\_Model)

Call:

lm(formula = y ~ x1 + x2 + x3, data = d)

Residuals:

Min 1Q Median 3Q Max  
 -1.8532 -1.4495 -0.3219 0.5919 3.2121



Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	$b_0$	5.8871	$t_0$	9.36e-05 ***
x1	$b_1$	0.1909	$t_1$	0.000479 ***
x2	$b_2$	0.2673	$t_2$	6.58e-05 ***
x3	$b_3$	0.6171	$t_3$	0.591572
---				

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1  
 Residual standard error:  $s$  on 9 degrees of freedom  
 Multiple R-squared: 0.9117, Adjusted R-squared: 0.8823  
 F-statistic:  $t\text{-value}$  on K and  $V_2$  DF, p-value: 4.496e-05

B<sub>i</sub>

HYPOTHESIS TEST

### QUALITY OF FIT & MODEL SELECTION

As the number of fitted explanatory variables goes up,  
 $R^2 \rightarrow 1$

This is because the value of  $SSE = \sum(y_i - \hat{y}_i)^2$  will decrease, while  
 $SST = \sum(y_i - \bar{y})^2$  remains constant

The **Adjusted  $R^2$**  value divides the terms in  $R^2$  (i.e., adjusts) by their degrees of freedom:

$$R_{adj}^2 = 1 - \frac{SSE/(n - k - 1)}{SST/(n - 1)}$$

Factors to consider:  
 if ↓ in SSE doesn't match ↑ in DOF,  $R^2_{adj}$  decreases to indicate poorer fit

- Higher  $R_{adj}^2 >$  lower  $R_{adj}^2$
- Simpler models > more complicated models

### INTERACTION TERMS

Models that contain different types of interaction terms

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{12} x_1 x_2 + \beta_{13} x_1 x_3 + \beta_{23} x_2 x_3 + \beta_{123} x_1 x_2 x_3$$

If main effect terms are not significant, then unlikely interaction terms will be either

- Analyze partial models first and then add interaction terms

a)  $\alpha = P(\text{reject } H_0 \mid H_0 \text{ is true}) = P(\bar{x} < 14.9 \mid \mu \geq 15)$   
 $= P(z < \frac{\bar{x}-\mu_0}{\sigma/\sqrt{n}}) = P(z < \frac{14.9-15}{0.5/\sqrt{50}}) = P(z < -1.41)$   
 From Normal table,  $P(z < -1.41) = 0.0793 = \alpha$

b)  $\beta = P(\text{fail to reject } H_0 \mid H_0 \text{ is false}) = P(\bar{x} > 14.9 \mid \mu_A = 14.8)$   
 $= P(z > \frac{\bar{x}-\mu_A}{\sigma/\sqrt{n}}) = P(z > \frac{14.9-14.8}{0.5/\sqrt{50}}) = P(z > 1.41)$   
 From Normal table,  $P(z > 1.41) = 1 - P(z < 1.41) = 1 - 0.9207 = 0.0793 = \beta$

### CATEGORICAL VARIABLES

Add indicator or dummy variables

- Example:

$$z_A = \begin{cases} 1 & \text{if Yeast A} \\ 0 & \text{otherwise} \end{cases}$$

$$z_B = \begin{cases} 1 & \text{if Yeast B} \\ 0 & \text{otherwise} \end{cases}$$

If adding indicator variables:

- Indicator main effect will affect the intercept values
- Indicator interaction effect will affect the slope values

A bakery is interested in modelling the amount of bread baked using the pH of the water and one of three types of yeast.

If we use Yeast A, our model is:

$$\hat{y}_A = -161.897 + 54.294(\text{pH}) + 89.998 + 0$$

$$= -71.899 + 52.294(\text{pH})$$

If we use Yeast B, our model is:

$$\hat{y}_B = -161.897 + 54.294(\text{pH}) + 0 + 27.166$$

$$= -134.731 + 52.294(\text{pH})$$

If we use Yeast C, our model is:

$$\hat{y}_C = -161.897 + 54.294(\text{pH}) + 0 + 0$$

$$= -161.897 + 54.294(\text{pH})$$

> Model12 = lm(y2 ~ ph + z1 + z2)  
> summary(Model12)

Call:  
 lm(formula = y2 ~ ph + z1 + z2)

Residuals:  
 Min 1Q Median 3Q Max  
 -31.038 -13.640 3.601 10.798 26.374

Coefficients:  
 (Intercept) -161.897 37.433 -4.325 0.000699 \*\*\*

ph 54.298 4.755 11.417 1.77e-05 \*\*\*

z1 89.998 11.052 8.143 1.11e-06 \*\*\*

z2 27.166 11.018 2.467 0.027127 \*

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 19.05 on 14 degrees of freedom  
 Multiple R-squared: 0.9404, Adjusted R-squared: 0.9277  
 F-statistic: 73.68 on 3 and 14 DF, p-value: 8.34e-09

$\downarrow$  compare t vs t<sub>df</sub>

54

MIC 247

## ONE - FACTOR EXPERIMENT

We assume that each of the  $k$  populations are:

- Independent
- No outliers
- Normally distributed with individual mean  $\mu_i$  and common variance  $\sigma^2$

ANOVA statistical model:

$$\text{value of level } i \rightarrow y_{ij} = \bar{\mu} + \alpha_i + \epsilon_{ij}$$

grand deviation of level  $i$   
mean of level  $i$   
 $\epsilon_{ij}$  effect of level  $i$

## ASSUMPTIONS

$$\sum_i \alpha_i = 0$$

$\epsilon_{ij}$  is normally distributed w/ mean 0 & variance  $\sigma^2$

If we reject  $H_0$ : find which means don't match 3 ways

1. Aggregate groups using linear contrasts

$$W = \sum_i c_i \mu_i \text{ where } \sum_i c_i = 0$$

indicate the 'side' the mean is on & strength of composition of group

2. Conduct Hypothesis Test  $H_0: \sum_i c_i \mu_i = 0$  Two contrasts  $w_a = \sum_i c_i \mu_i$  and  $w_b = \sum_i d_i \mu_i$  are orthogonal  
 $H_A: \sum_i c_i \mu_i \neq 0$  if and only if  $\sum_i \frac{c_i d_i}{n_i} = 0$

$$SSW = \frac{(\sum_i c_i \bar{y}_i)^2}{\sum_i \left( \frac{c_i^2}{n_i} \right)} \quad \text{Sum of Sq's for contrast } W$$

$$\bar{y}_i = \text{mean}$$

$$\bar{y}_{i..} = \text{total n-i size of } i$$

$$MSE = \frac{\sum_i (y_{ij} - \bar{y}_{i..})^2}{N-k} = \frac{SSE}{N-k}$$

$$f_w = \frac{SSW}{MSE} \quad \text{We reject } H_0 \text{ if } f > f_{(v_1=1, v_2=N-k)}$$

$F$  - test for contrast  $W$

## PLANNED COMPARISONS (contrasts)

If can guess which means are going to differ in advance can use pairwise individual 2-sample t-tests

$$H_{0(ij)}: \mu_i = \mu_j \text{ or alternatively } \mu_i - \mu_j = 0$$

$$H_{A(ij)}: \mu_i \neq \mu_j \text{ or alternatively } \mu_i - \mu_j \neq 0$$

\* EXPERIMENT-WISE TYPE I ERROR RATE  
 $1 - (1 - \alpha)^r$ , for  $r$ -many pairwise comparisons

# of pairwise comparisons ↑  
 prob of committing type I error ↑

if only have posterior knowledge of significant differences

$$H_0: \mu_i = \mu_j$$

$$H_A: \mu_i \neq \mu_j$$

$$SRSE_{i,j} = \sqrt{\frac{MSE}{2} \left( \frac{1}{n_i} + \frac{1}{n_j} \right)}$$

## TUKEY'S TEST

(posterior pairwise comparisons)

- Recall that  $H_0: \mu_i = \mu_j$
- If  $q > q^*$ , then we reject  $H_0$
- If  $q \leq q^*$ , then we fail to reject  $H_0$

$$q_{ij} = \frac{|\bar{y}_i - \bar{y}_j|}{SRSE_{i,j}}$$

Compare  $q$  to  $q^*$  from Tukey's Table

$$H_0: \sigma_1^2 = \dots = \sigma_k^2$$

- $H_A$ : at least one variance is not equal

$$S_{\text{pooled}}^2 = \frac{1}{N-k} \sum_{i=1}^k (n_i - 1) s_i^2$$

$$b = \frac{1}{S_{\text{pooled}}^2} \left[ (s_1^2)^{n_1-1} \times \dots \times (s_k^2)^{n_k-1} \right]^{\frac{1}{N-k}}$$

Compare  $b$  to  $b_k(\alpha, n_1, \dots, n_k)$ , where  $b_k(\alpha, n_1, \dots, n_k) \approx \frac{\sum_{i=1}^k n_i b_k(\alpha, n_i)}{N}$   
 and  $b_k(\alpha, n_i)$  is from the Bartlett table

## BARTLETT'S TESTS

↳ test for homogeneity of variance

Reject  $H_0$  if  $F > F^*$

## RANDOMIZED COMPLETE BLOCK DESIGN

- SSA: sum of squares (effect) =  $b \sum_i (\bar{y}_i - \bar{y}_{..})^2$
- SSB: sum of squares (blocks) =  $k \sum_j (\bar{y}_{.j} - \bar{y}_{..})^2$
- SSE: sum of squares (error) =  $\sum_i \sum_j (y_{ij} - \bar{y}_i - \bar{y}_{.j} + \bar{y}_{..})^2$
- SST: sum of squares (total) =  $\sum_i^k \sum_j (y_{ij} - \bar{y}_{..})^2 = SSA + SSB + SSE$

Source of Variation	Sum of Squared Error	Degrees of Freedom	Mean Squared Error	F statistic
Effects	SSA	$k - 1$	$\frac{SSA}{k-1} = MSA$	$\frac{MSA}{MSE} = f_{(v_1=k-1, v_2=(k-1)(b-1))}$
Subjects	SSB	$b - 1$	$\frac{SSB}{b-1} = MSB$	
Error	SSE	$(k-1)(b-1)$	$\frac{SSE}{(k-1)(b-1)} = MSE$	
Total	SST	$kb - 1$		

Completely Randomized design

Randomly assign subjects into any of the  $K$  treatments being analyzed

## Randomized Block Design

used when there is a second factor to stratify by  
 1. stratify by factor B into blocks

2. randomize subjects into treatments along factor A

## Randomized complete block design

assign each treatment once to each block, same as two factor experiment w/o interaction

## Interpolation v. Extrapolation

Interpolation: estimating a value within the domain of data points  
 ex: predicting an intermediate value in LR

Extrapolation: estimating a value beyond the data's domain  
 ex: LR beyond the data

Two contrasts  $w_a = \sum_i c_i \mu_i$  and  $w_b = \sum_i d_i \mu_i$  are orthogonal if and only if  $\sum_i \frac{c_i d_i}{n_i} = 0$

Orthogonality means that both  $SSW_a$  and  $SSW_b$  form independent components of  $SSA$

- This means that we are separating  $SSA$  into their  $w_a$  and  $w_b$  "components"
- If  $w_a$  and  $w_b$  are not orthogonal, then their  $SSW_i$  will overlap

## TWO - FACTOR EXPERIMENT &amp; INTERACTION

\* assume each of the  $a, b$  level combinations has exactly  $n$  observations

Same assumptions with a two-way ANOVA, but with  $ab$  groups:

- Independence of observations
- Equality of variances
- No outliers
- Normality

## ANOVA MODEL

$$y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}$$

let  $\alpha_i$  represent effect of A  
 let  $\beta_j$  represent effect of B

## Constraints

$$\sum_i \alpha_i = 0, \sum_j \beta_j = 0, \text{ and } \sum_i \sum_j (\alpha\beta)_{ij} = 0$$

Interaction effect of AB:

- $H_{0(AB)}: (\alpha\beta)_{1,1} = \dots = (\alpha\beta)_{a,b} = 0$
- $H_{A(AB)}: \text{at least one } (\alpha\beta)_{i,j} \text{ is not equal}$

$$SSA = bn \sum_i (\bar{y}_{i..} - \bar{y}_{...})^2$$

$$SSB = an \sum_j (\bar{y}_{.j} - \bar{y}_{...})^2$$

$$SS(AB) = n \sum_i \sum_j (\bar{y}_{ij} - \bar{y}_{i..} - \bar{y}_{.j} + \bar{y}_{...})^2$$

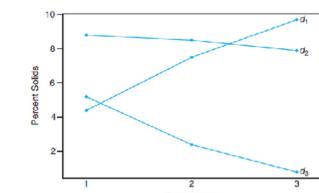
$$SSE = \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{ij.})^2$$

$$SST = \sum_i \sum_j \sum_k (y_{ijk} - \bar{y}_{...})^2 = SSA + SSB + SS(AB) + SSE$$

Source of Variation	Sum of Squared Error	Degrees of Freedom	Mean Squared Error	F statistic
Main Effect A	SSA	$a - 1$	$\frac{SSA}{a-1} = MSA$	$\frac{MSA}{MSE} = f_{A(v_1=(a-1), v_2=ab(n-1))}$
Main Effect B	SSB	$b - 1$	$\frac{SSB}{b-1} = MSB$	$\frac{MSB}{MSE} = f_{B(v_1=(b-1), v_2=ab(n-1))}$
Interaction	SS(AB)	$(a-1)(b-1)$	$\frac{SS(AB)}{(a-1)(b-1)} = MS(AB)$	$\frac{MS(AB)}{MSE} = f_{AB(v_1=(a-1)(b-1), v_2=ab(n-1))}$
Error	SSE	$ab(n-1)$	$\frac{SSE}{ab(n-1)} = MSE$	
Total	SST	$abn - 1$		

Assess each  $H_0$  independently based on its f-statistic REJECT  $H_0$  if  $f > f^*(v_1, v_2)$

## INTERACTION PLOTS



If there are  $k$  treatments, then there are at most  $k - 1$  orthogonal linear contrasts, where:

$$SSA = SSW_1 + SSW_2 + \dots + SSW_{k-1}$$

## PLANNED COMPARISONS