

Experimental Design - controlled, can identify the effect of an action (gets causation)

observational studies can't account for confounding effects (can only get association)

i.e. confounder is smoking, observational study - does drinking lead to lung cancer

response variable = outcome being measured (**no treatment? => observational study**)

factor = explanatory var under the control of the experimenter (level = values of factor)

treatment = combination of experimental conditions applied to experimental unit (incl. control)

experimental unit = the unit on which the treatment is applied (i.e. subject)

blocking variable = nuisance, variations/variable not interested in this study

randomization, placebo, blinding, control in experiments can eliminate bias & confounders

blocking and match-paired designs can reduce variation in the data

- the response variable: presence or absence of a birth defect

the experimental units: the women

the factor: the dose level of folic acid is the factor

the levels of the factor: there are two levels, 0mg and 0.8mg

number of treatments: two

whether blocking was applied: no

1-way ANOVA

- used to compare ≥3 groups (quantitative response, categorical factor/explanatory var)

If sample size is large or data is normal, better to use ANOVA than Kruskal-Wallis (also one factor only)

Assume: residuals (mean-point/within gr) are normal, obs within grp are normal, all grp have same variance

obs within each grp & in diff grp are ind., groups are formed by randomly sampled from population

H0: all population means are the same

Ha: at least one population mean is different from the others

1. find overall sample mean: $(\mu_1 + \mu_2 + \mu_3 + \mu_4) / 4$

- taking group mean (\bar{y}) only works if u have same # of replicates per group (\bar{y}_i is mean of group i)

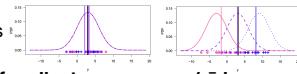
- one-way ANOVAs are for single factor experiments

- factor has g levels/groups (i)

- each group has t observations (j)

		Application method		
		Dipping	Spraying	
Primer type	1	4.0, 4.5, 4.3	5.4, 5.9, 5.6	
	2	5.6, 4.9, 5.4	5.8, 6.1, 6.3	
	3	3.8, 3.7, 4.0	5.5, 5.0, 5.0	

the response variable: adhesive force
the experimental units: the aircraft wings
the factors: primer type and application method
the levels of the factors: two for application method, three for primer type
number of treatments: six



2. find sum of squares: total SS (var) around overall mean = within group SS + between group SS

- if mean rubber strength b/w materials same => diff (H0 → Ha), total SS + between SS would increase

- F statistic Ha > F statistic H0

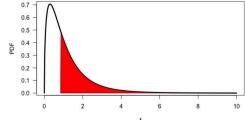
$$X \sim \chi^2_m \cdot X/m \sim F_{m,n} \quad m,n = \text{dof}$$

$$Y \sim \chi^2_n \cdot Y/n \sim F_{n,m}$$

Within group SS $\frac{\sum (y_i - \bar{y})^2}{\sigma^2} \sim \chi^2_{g(t-1)}$ Between group SS $\frac{\sum (\bar{y}_i - \bar{y})^2}{\sigma^2} \sim \chi^2_{g-1}$

3. test statistic = mean square ratio = SS/dof = $\frac{\text{Between group SS}/(g-1)}{\text{Within group SS}/g(t-1)} \sim F_{g-1,g(t-1)}$

- test stat should follow $F_{3,16}$ dist. if H0 is true



Source	Sum of Squares	d.o.f.	MS	F
Between	6.82	3	2.273	
Within	43.02	16	2.689	Between group MS = 0.845
Total	49.84	19		

We reject H0 if: p-value = area under the curve to the right of 0.845

- between SS > within SS

- $F > F_{a,g-1,n-g}$

Multiple comparisons - posthoc method performing set of pairwise comparisons to see which grp differ

- test stat to compare two groups (where variance is MSE) => pairwise comparison

$$t = \left| \frac{\bar{y}_i - \bar{y}_j}{\sqrt{\frac{MSE}{n_i} + \frac{MSE}{n_j}}} \right| > t_{g(n-1)}$$

reject H0 if bigger than t-table val (see bonferroni for t-table val)
there are $(n-1)(n-2)/2 = 6$ possible pairwise comparisons to make within groups

- multiple comparisons problem: each pairwise comparisons are not ind., chance of T1 error increases

- Bonferroni's equality: divide alpha by all possible pairwise comparisons => now the chance of making a T1 error is less than 5% (T1 error = rejecting H0 when true = alpha)

- corrected sig level: $0.05/6 = 0.00833 \Rightarrow$ two-sided alternative: $1 - 0.00833/2 = 0.9958 \Rightarrow$ plug to t-table

Colony (l, m) $\bar{y}_l - \bar{y}_m$, t, Significant?

Colony (l, m)	$\bar{y}_l - \bar{y}_m$	t	Significant?
1, 2	-29.4	6.22	Y
1, 3	-32.2	6.81	Y
1, 4	-30.6	6.47	Y
2, 3	-2.8	0.59	N
2, 4	-1.2	0.25	N
3, 4	1.6	0.34	N

- t-table val (dof 16, 0.9958) = 3
- colony 1 differs from other colonies
- on average, colony 1 depart about 30 days earlier than the others
- note: $\text{Var}(\bar{y}_l - \bar{y}_m) = \frac{\sigma^2}{n_l} + \frac{\sigma^2}{n_m}$ $\text{ESD}(\bar{y}_l - \bar{y}_m) = \sqrt{\frac{MSE}{n_l} + \frac{MSE}{n_m}}$

Contrasts - extending two sample t-test for general multi-level single factor

- i.e. want to compare group bio with level A, B and group non-bio with level C

$$1. \bar{y}_A + \bar{y}_B - 2\bar{y}_C = 0.46 \quad \text{or} \quad \frac{\bar{y}_A + \bar{y}_B - 2\bar{y}_C}{\sqrt{2}} = 0.46$$

$$2. \text{Var}(\bar{y}_A + \bar{y}_B - 2\bar{y}_C) = \text{Var}(\bar{y}_A) + \text{Var}(\bar{y}_B) + 2\text{Var}(\bar{y}_C) \\ = \text{Var}(\bar{y}_A) + \text{Var}(\bar{y}_B) + 2^2\text{Var}(\bar{y}_C) \\ = \sigma^2/n + \sigma^2/n + 2^2(\sigma^2/n) \\ = 6(\sigma^2/5)$$

$$3. t^* = (\bar{y}_A + \bar{y}_B - 2\bar{y}_C) / \sqrt{6(MSE/5)}$$

- why use contrast? if we pooled data from 2 grp into 1, we can't account for small diff these grp have

- $C_A + C_B - 2C_C = 0$, this constraint on coeff/weights with grp mean (y) of resulting stat => contrast

- under H0, mean of ALL groups should be same, which means coeff/weights should sum to 0

- if n is not equal amongst groups, use this formula:
this is step 2, so use c properly => $\text{MSE} = \text{variance } \sigma^2$

2-way ANOVA - used to look at the effect of 2 factors on a continuous response variable

$\hat{\sigma}^2$ = Within group mean square

► Also known as: Mean square error (MSE)

$$\text{MSE} = \frac{\sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{ij})^2}{g(t-1)}$$

- **main effect:** constant effect of one factor across all values the other factors

$$\begin{array}{ll} A - A + & \text{Main effect of A: } \bar{y}_{A+} - \bar{y}_{A-} = \frac{50+40}{2} - \frac{30+20}{2} = 20 \\ B - 20 & \\ B + 30 & \end{array}$$

$$\begin{array}{ll} & \text{Main effect of B: } \bar{y}_{B+} - \bar{y}_{B-} = \frac{50+30}{2} - \frac{40+20}{2} = 10 \end{array}$$

- **no interaction = slope is equal** (lines will never touch)

- when there is interaction, it makes no sense to talk about individual factor(s) or main effects (effect of A may not be 20 whatever level of B is)

$$y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2)$$

* overall mean

$$\hat{\mu} = \bar{y}$$

* α_i : group (or treatment) effect for group i

$$\hat{\alpha}_i = \bar{y}_i - \bar{y}$$

* ϵ_{ij} : independent error variables

$$\hat{\epsilon}_{ij} = y_{ij} - \bar{y}_{ij}$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{ij})^2}{g(t-1)}$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{ij})^2}{g(t$$

If Y didn't depend on X , plot scatters about mean $y_{\bar{y}}$ & only var is due to noise ($b_0 \sim \text{noise}$, $b_1 \sim 0$)

- if relationship exists: model/regression SS > error SS	Source	DoF	SS	MS	F
- if relationship DNE: model/regression SS < error SS	Model	1	$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$	MSM	MSM/MSSE
- assumption $\epsilon \sim N(0, \sigma^2)$ works well if n is big	Error	$n-2$	$\sum_{i=1}^n (y_i - \hat{y}_i)^2$	MSE	
Under H_0 , $F \sim F_{1,n-2}$	Total	$n-1$	$\sum_{i=1}^n (y_i - \bar{y})^2$		

Under H_0 , $F > F_{1,n-2}$ reject H_0 , test stat falls in far upper tail of F -dist

- i.e. $b_1 = 0$, relationship exists, MSM/MSE to be > 1

Regression Estimators - quantify uncertainty around a parameter estimate

$$Y = \beta_0 + \beta_1 X + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad \text{- params are } \beta_0, \beta_1, \sigma$$

- estimators b_0, b_1 and σ^2 are all random variables (fxns of the data which are subject to random error according to the model)

If $\text{Var}(Y_i | X = x_i) = \sigma^2$ for all i , it can be shown that

$$\text{Var}(b_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \begin{aligned} &\text{- Var}(b_1) \text{ is inversely proportional to } \sigma^2 \\ &\text{- more spread out } x \text{ (bigger } x), \text{ lower the variance of the estimate of the slope} \end{aligned}$$

- predictor, explanatory, or independent variables = x

- response, observable, or dependent variables = y

$$\text{Var}(b_0) = \frac{\sigma^2 \sum x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}$$

- to reduce $\text{Var}(b_1)$, we may want a large range of x values (more spread out)

- replace known variance σ^2 with estimated variance $s^2 \Rightarrow$ formula for $\text{sd}(b_1) \Rightarrow b_1 = \text{estimate}$

$$\text{CI: estimate} \pm t_{n-2}^* \times \text{se}(\text{estimate})$$

$$b_1 = 0.0731$$

$$t_{n-2}^* = 2.01$$

$$\text{se}(b_1) = \sqrt{\frac{s^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$s^2 = \frac{\sum e_i^2}{n-2} = \text{MSE} = 61.47$$

ANOVA test for regression

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Test statistic: $F = \text{MSM}/\text{MSE} \sim F_{1,n-2}$

Alternative test (t-test)

$$H_0: \beta_1 = 0$$

$$\text{Test statistic: } t = \frac{b_1}{\text{se}(b_1)} \sim t_{n-2}$$

$$H_a: \beta_1 \neq 0 \rightarrow t_{48,0.975}^* = 2.01$$

$$H_a: \beta_1 > 0 \rightarrow t_{48,0.95}^* = 1.68$$

$$H_a: \beta_1 < 0 \rightarrow t_{48,0.05}^* = -1.68$$

- one-sided: $H_a: \beta_1 < 0$ or $H_a: \beta_1 > 0$

- two-sided: $H_a: \beta_1 \neq 0$

ANOVA test (f-val) OR t-test

- $t^2 = F$ bc a random variable t_m^2 is a random variable $F_{1,m}$

- For $H_0: \beta_1 = 0$ t-test and F-test are equivalent

(t-test is always two-tailed bc it accounts for sq vals only)

(f-test is always positive, so it can't be equivalent to 1-tail)

Parameter estimates $\hat{\beta}_0 = b_0$ and $\hat{\beta}_1 = b_1$ and their standard errors $\text{se}(b_0)$ and $\text{se}(b_1)$ can be used for inference

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

$$\triangleright t = b_1/\text{se}(b_1)$$

Why should we not test for H_0 : intercept/ $\beta_0 = 0$?

Predicting the response for a particular value of the explanatory variable

► Here, predicting the number of goals for a given number of shots

For a given x^* value, do we want to predict

► the mean response?

► Here, μ_x the mean number of goals for players taking x^* shots

► The response for an individual?

► Here, \bar{y} , the number of goals for a single player that had x^* shots

For both, the point estimate is: $b_0 + b_1 x^*$

But CIs are different because $\text{se}(\cdot)$ are different

$$\triangleright \text{se}(\bar{y}) < \text{se}(y)$$

- point estimate for mean and point $y_{\bar{y}}$ will be same but confidence interval will have different width bc $\text{sd}(y_{\bar{y}}) > \text{sd}(\text{mean}) \Rightarrow$ extra variability in a single response compared to mean (mean doesn't have σ)

Multiple Linear Regressions

- Visualizing relationship b/w Y and (X_1, X_2) is difficult, especially if X_1 and X_2 has correlations/interaction

- minimize least squares simple linear regression multiple linear regression

$$\text{pink} = y_{\bar{y}}$$

$$\sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_p x_{pi}))^2$$

- what is an interpretation of estimate β_j ? We estimate that for day with a common POL, the average number of EMV increases by β_2 units for every one increase in TMP. $\text{ERV} = \beta_0 + \beta_1 \text{POL} + \beta_2 \text{TMP} + \epsilon$

- β_1 CI (POL) has lower endpoint is positive \Rightarrow likely high POL causes more ERV

β_2 CI (TMP) has lower endpoint is positive but narrow \Rightarrow less likely high TMP causes more ERV

- β CI includes 0 \Rightarrow no evidence that there exists a relationship between Y and X (accept H_0)

=> good chance that the true population value is close to 0

$$\text{Residual SS} \quad \sum_{i=1}^n e_i^2 \quad \text{Residual MS / MSE} \quad \frac{\sum_{i=1}^n e_i^2}{n - (p+1)} \quad \text{Multiple R}^2 \quad R^2 = 1 - \frac{\text{Residual SS}}{\text{Total SS}}$$

- when >1 predictor variable, we cannot use scatterplots because there are >2 variables to plot together

- relationship b/w y and x_1 may change if we add x_2 in the model (i.e. scatterplot of y and x_1 may not show true relationship in a multi regression model because x_1 and x_2 could be correlated)

- interpretation of params: β_1 is the change in Y when X_1 is changed by 1 unit (while other predictors fixed)

Curve Fitting in Regression - using the multiple linear regression machinery to fit curves to data

- $b_1 = \frac{dy}{dx} = \frac{db_0 + b_1 x}{dx} =$ - b_1 is slope = rate of change

- adding an explanatory variables (X) will decrease RSS

- exceptions where RSS will be same (these cannot be fitted in practice):

1. If $y = \beta_0 + \beta_1 x_1$ with no error (i.e. perfect fit, RSS will be 0, so cannot be lower)
2. If x_1 and x_2 are perfectly correlated

- increasing power of regression line \Rightarrow will never increase residual (will always decrease or stay the same)

$$y = b_0 + b_1 x + b_2 x^2 \quad \text{- we can't say, } b_1 \text{ is change in } Y \text{ when } x \text{ is changed by 1 unit while we hold } x^2$$

- instead, change in Y is $b_1 + 2b_2 x$ when x is changed by 1 unit

For any value c , $\hat{\beta}_1 + 2\hat{\beta}_2 c$ is the estimated rate of change in the conditional mean of Y given X as X

The model

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

postulates a linear relationship between X and Y , whereas

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$$

describes a quadratic relationship.

But: statisticians consider both these models as linear models

- Because in both cases the expected value of Y given X is modeled as a linear function of the parameters, the regression coefficients
- Can use linear model software to fit both models

polynomial models with order q $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_q X^q + \epsilon$

Which polynomial order model is the most consistent with the data?

1. Set $q = 1$.

2. Fit model q .

3. If the 95% CI for β_q includes zero then stop, and output model $q-1$ as the answer. Otherwise, increase q by one, and go to Step 2.

- note that this algorithm can give us $q=0$ for the answer ("intercept-only")

- methods for multiple regression can be used in polynomial regression

- **limitation** for a polynomial model: it may be unstable if the order is higher than 2 (i.e. relation may change if diff samples are taken from the same population)

All linear models are empirical models - may be useful to describe the observed data, but may not reflect true relationship b/w predictors & response \Rightarrow predictions outside observed data range can be dangerous (i.e. extrapolation)

Residuals in Regression - check if assumpt. of model is reasonable for data in a linear regression, we assume (residual plots + normal QQ plots)

- mean of the response is a linear function of the predictors

- errors are independent and normal RV with mean 0 and constant variance

- look at residuals (estimated errors) to assess if conditions of a multiple linear regression model seem to hold for our data

- residuals (obs - fit) vs fitted values ($y_{\hat{y}}$)

- we want patternless "band" centred around 0

- pattern (i.e. funnel) suggest inadequate model

- inadeq? may want to change response variable

- residuals vs predictor variables (x)

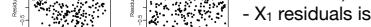
- we want patternless "cloud" centred around 0

- only accounts for 1 var, not generalizable

- normal score plots of residuals

- we want an approximately straight line thru origin

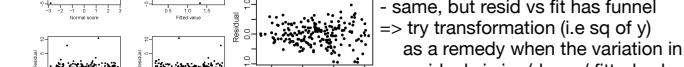
- check departures from normality (skewness, heavy tails)



- data w/ X_1, X_2

- X_1 residuals is nonlinear, add X_2

- same, but there's 3 outliers \Rightarrow remove outliers



- same, but resid vs fit has funnel \Rightarrow try transformation (i.e. sq of y) as a remedy when the variation in residuals is inc/dec w/ fitted vals

- residual = what's left after model is removed \Rightarrow should resemble noise

- H_0 : residual is normal w/ σ^2 that does not change w/ the explanatory variable

$\log Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

Results in

$$\triangleright \text{Implies: 1 unit increase in } X_1 \quad \star e^{\beta_1} = e^{-0.0273} = 0.973$$

$$\triangleright Y = e^{\beta_0} e^{\beta_1 X_1} e^{\beta_2 X_2} e^{\epsilon} \quad \star \text{ (for a constant } X_2 \text{ value)}$$

$$\triangleright \text{So } a (0.973 - 1) \times 100 = -2.7\% \text{ rate of change}$$

Residual plots allow us to check:

- goodness-of-fit: does the model fit the data well?

- constant variance: is the variance of y constant?

- outliers: are there any outliers?

- transformations: any transformations needed for y and x_i 's?

Normal scores plots allow us to check:

- Whether e follow Normal distributions or not

- A near straight line shows that the Normality assumption may be OK

E.g., let's compare 3 groups using a regression

>Create 3 dummy variables: X_1 and X_2

Group 1: $X_1 = 0 \& X_2 = 0 \Rightarrow$ the reference (or baseline)

Group 2: $X_1 = 1 \& X_2 = 0$

Group 3: $X_1 = 0 \& X_2 = 1$

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

Y-axis: Odds ratio

X-axis: X_1, X_2

Group 1: $Y = 1 | X_1 = 0, X_2 = 0$

Group 2: $Y = 1 | X_1 = 1, X_2 = 0$

Group 3: $Y = 1 | X_1 = 0, X_2 = 1$

Odds ratio = $\frac{P(Y=1|X_1=1,X_2=1)}{P(Y=1|X_1=0,X_2=0)}$

$= \frac{P(Y=1|X_1=1) \cdot (1-P(Y=1|X_1=0,X_2=0))}{P(Y=1|X_1=0) \cdot (1-P(Y=1|X_1=0,X_2=0))}$

</div