

EVERYTHING QUAL

(just like everything bagel)

Anqi Zhao
08/20/2025

Reviews (from previous students who used this note and took the Qual)

You're invited to leave a review as well if you find this note helpful or if you hate this _(`)_/

I will also add an acknowledgment section if you find any typos/ errors (there will be a lot, but still first come first serve ofc)

Please contact Anqi Zhao (the one who wrote the note and uploaded) at azhao6@ncsu.edu or find me on discord or at LAU

I may add a **table of contents** later here

PhD Qualifying EXAM Syllabus

(Updates to this document must be updated on the Student Copy in the PhD Qualifying Exam **Student Resources** shared drive. The Graduate Services Coordinator can edit that shared drive.)

ST 703-704: Statistical Methods

Representative Texts

- Rao, P.V. Statistical Research Methods in the Life Sciences, Brooks/Cole.
- Ott, R. Lyman and Longnecker, Michael T. Introduction to Statistical Methods and Data Analysis.
- Damon. Jr., Richard A. and Harvey, Walter R. Experimental Design, ANOVA, and Regression, Harper and Row, Publishers, 9187. 508 pp.
- Neeter, John, Wasserman, William, and Kutner, Michael H. Applied Linear Statistical Models 3rd Ed., Richard D. Irwin, Inc., 1990, 1182 pp.
- Ostle, Bernard and Mensing, Richard W. Statistics in Research, 3rd Ed., Iowa State University Press.
- Snedecor, George W. and Cochran, William G. Statistical Methods, 7th Ed., Iowa State University.
- Steel, Robert G.D. and Torrie, James H. Principles and Procedures of Statistics: A Biometrical Approach, 2nd Ed., McGraw – Hill.
- SAS Institute Inc., SAS/STAT User's Guide, Release 6.03 ED., Cary, NC: SAS Institute Inc., 1988. 1028pp.
- SAS Institute Inc., SAS Language: Reference, Version 6, First Ed., Cary, NC: SAS Institute Inc., 1990. 1042 pp.

Topics to Review

- Definition and computation of elementary descriptive statistics
- Populations and samples
- Sampling distributions and the Central Limit Theorem
- Use of the Z, t, Chi Square, and F tables
- Logical basis of confidence and tests of hypothesis
- Inference on mean, variance and proportion of one population
- Inference on means and proportions from two populations – independent and paired samples
- Inference on variances from two populations
- Power and sample size calculations
- Basic concepts of experimental design including the concept of experimental unit, experimental error, replication, relative efficiency, blocking, covariance, and randomization.
- For each of the following experimental designs:
 - Completely Randomized Design (CRD)

- Randomized Complete Block Design (RCBD)
- Split Plot and Repeated Measures designs

Students Should:

- Know the models (including alternative parameterizations) and related distributional assumptions.
- Be able to recognize which design is appropriate for a given experiment and understand the advantages and disadvantages of each design.
- Know how to construct the ANOVA table with degrees of freedom, sums of squares and mean squares.
- Know what hypotheses can be tested and how to test them, including use of expected mean squares to find appropriate denominators.
- Know how to estimate and place confidence intervals on meaningful linear combinations of the fixed effects such as treatment contrasts, treatment means and other linear combinations.
- Know how to estimate the variance components in the model (using method of moments) and how to use them to obtain variances for linear combinations of (estimated) fixed effects (and understand the correlation structure as a function of these variance components.)
- Be able to recognize replication and subsampling, and account for them in the model, ANOVA table and analysis.
- Know how to make multiple comparisons using a number of procedures that adjust for multiple testing, including the Tukey-Kramer, Scheffe, Bonferroni, and Benjamini-Hochberg procedures.
- Know how to account for one or more covariates.
- Basic concepts of treatment designs including:
 - Treatments and treatment combinations
 - Control versus experimental treatments
 - Factorial treatment designs
 - Factors and their levels
 - Main effects and interactions (1st order, 2nd order, etc.)
 - Nested designs
 - Nested-Factorial designs
- For balanced data, partitioning of the treatment sum of squares in the ANOVA table for each of the treatment designs for fixed, random and mixed models and interpret expected mean squares.
- Multiple regression using matrix notation, including

SS

df

- Model and assumptions

SSE $\sum_i (\hat{y}_i - \bar{y}_i)^2$

$n-p$

- Normal equations and parameter estimators
- Properties of the estimators

Lack of fit $\sum_i (\hat{y}_i - \bar{y}_{\cdot i})^2$

$n^* - p$

- Inference in multiple regression, comparing subsetted models
- Lack of fit and pure error

Pure error $\sum_i (y_i - \bar{y}_{\cdot i})^2$

$n - n^*$

- Residual diagnostics, including

- Misspecified mean model
- Misspecified covariance model
- Outliers and influential points
- Multicollinearity
- Model assessment and exploration, including a variety of criteria, k-fold cross validation, and stepwise or subset exploration
- Correlation
- Analysis of Covariance
- Biased Regression
 - Penalized regression (e.g., ridge, lasso, etc.)
 - Reduced dimensions (e.g., principal components, partial least squares, etc.)
- Generalized Linear Models
 - Model and assumptions
 - Inference
 - Logistic regression
 - Poisson modeling
- Linear Mixed Models
 - Models, assumptions, implied covariance structure
 - Subject-specific versus marginal model
 - ANOVA-type mixed models
 - Inference
 - Clustered/repeated measures/longitudinal data
- Sampling Strategies: simple random, stratified random, cluster, and systematic sampling
- Be able to read and interpret code and output from SAS and R for conducting the above analyses

ST 701-702: Statistical Theory

Representative Texts

- Casella, G. and Berger, R.L. *Statistical Inference*, 2nd Ed., Wadsworth/Brooks Cole, Pacific Grove, CA, 2001.
- Hogg, R.V., and Craig, A.T. *Introduction to Mathematical Statistics*, 4th Ed., MacMillan.
- Rohatgi, V.K. *An Introduction to Probability Theory and Mathematical Statistics*, John Wiley & Sons, New York, 1976.

Topics

- Basic probability calculus
- Random variables, probability distributions, density functions and distribution functions
- Discrete probability models: e.g. binomial, Poisson, geometric, negative binomial, hypergeometric, etc.
- Continuous probability models: e.g. uniform, exponential, beta, gamma, normal Weibull, Cauchy, extreme value, log-normal, etc.
- Multivariate probability models: multinomial, bivariate normal

- Expected value, variance, covariance, correlation, moments (about zero and about the mean) and moment-generating functions
- Moments of functions of random variables
- Joint distributions, conditional distributions, marginal distributions and expectations
- Distributions of functions of random variables, order statistics
- Chebyshev's, Markov's and Jensen's inequalities
- Normal theory: joint distribution of the sample mean and variance, central and noncentral distributions for Student t, Chi square and F
- Convergence in probability and the weak law of large numbers
- Convergence in distribution, the central limit theorem, asymptotic normality, Slutsky's theorem and the "delta method"
- Sufficient and minimal sufficient statistics
- Ancillary statistics
- Complete statistics
- Basu's theorem
- Method of moments, maximum likelihood estimation

Jan 2015 Part 1 Q2

Aug 2017 Part 2 Q3

Aug 2019 Part 2 Q3

- Bayesian inference: prior and posterior probability distributions, conjugate priors, Bayes estimators based on squared error loss; hierarchical models
- Properties of estimators -- unbiasedness, mean squared error, Cramer-Rao lower bound, Rao-Blackwell and Lehmann-Scheffe theorems, UMVUE, consistency, asymptotic efficiency
- Logical basis for and properties of hypothesis tests
- Type I and II error, level of significance and power
- Simple and composite hypotheses
- Unbiased tests
- Likelihood ratio tests
- Neyman-Pearson lemma for MP tests; Karlin-Rubin theorem for UMP tests, UMPU tests
- Asymptotic tests: chi-square, Wald, score
- Confidence interval construction by inversion of hypothesis tests
- Confidence interval construction using pivots
- Properties of confidence intervals: shortest length, UMA and UMA unbiased

ST 705: Linear Models

Representative Texts

- Graybill, F.A. Theory and Applications of the Linear Model, Duxbury, N. Scituate, Mass, 1976.
- Monahan, J.F. A Primer on Linear Models, CRC Press, 2008.
- Searle, S.R. Linear Models, John Wiley, New York, 1971.

- Seber, G.A.F. Linear Regression Analysis, Wiley, 2003.

Topics

- Review of linear systems of equations, generalized inverses and projection matrices, vector spaces and subspaces
- Linear statistical models and reparameterization
- Least squares theory and computation, including normal equations and partition of sums of squares
- Estimability and estimable linear functions, restricted models
- Gauss-Markov theorem, BLUE, and generalized least squares
- Theory and application of multivariate normal distribution and related distributions of quadratic forms: central and noncentral chi-squared, central and noncentral F. Cochran's Theorem
- Testing the general linear hypothesis and Likelihood Ratio Tests
- Joint distribution of several BLUES under normality
- Confidence intervals and sets for parameters and predictions
- Random effects, mixed models, and variance component estimation

Confidence interval

General Rules:

① Interval estimator (random):

$$P(\theta \in (\hat{\theta}_L, \hat{\theta}_U)) \geq 1 - \alpha$$

② Interval estimator (fixed):

$(\hat{\theta}_L, \hat{\theta}_U)$ has actual numbers

Confidence interval for μ :

(one sample - known σ^2)

Assumptions: ① y_1, \dots, y_n iid

② $E(y_i) = \mu$; $\text{Var}(y_i) = \sigma^2$

$y_i \stackrel{iid}{\sim} (\mu, \sigma^2)$ \leftarrow may not normally distributed

③ $E(\bar{y}) = \mu$

$$\text{SE}(\bar{y}) = \sqrt{\frac{\sigma^2}{n}} = \frac{\sigma}{\sqrt{n}}$$

	y_i normal	y_i not normal
n large	$\bar{y} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$	$\bar{y} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$
n small	$\bar{y} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$?

(one sample - unknown σ^2)

	y_i normal	y_i not normal
n large	$\bar{y} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$	$\bar{y} \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$
n small	$\bar{y} \pm t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}$?

(two sample (μ_1, μ_2) - independent)

Assumptions: ① Samples from 2 different populations

Group 1: Y_{11}, \dots, Y_{1n_1} } mutually
Group 2: Y_{21}, \dots, Y_{2n_2} } independent

② $Y_{ij} \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$ $Y_{2j} \stackrel{iid}{\sim} N(\mu_2, \sigma^2)$

③ $E(\bar{Y}_1 - \bar{Y}_2) = \mu_1 - \mu_2$

$$\begin{cases} \hat{SE}(\bar{Y}_1 - \bar{Y}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \\ \hat{SE}(\bar{Y}_1 - \bar{Y}_2) = Sp \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \end{cases}$$

← for this case, only consider unknown σ^2

If $\sigma_1^2 \neq \sigma_2^2$:

$$(\bar{Y}_1 - \bar{Y}_2) \pm t_{\alpha/2, v} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

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

↓
Satterthwaite approximate (same type)

If $\sigma_1^2 = \sigma_2^2$:

$$(\bar{Y}_1 - \bar{Y}_2) \pm t_{\alpha/2, n_1+n_2-2} \cdot Sp \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

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

pooled SD

(two-sample (μ_1, μ_2) - paired data)

Assumptions: ① Samples from 2 measurements of same unit

Group 1: Y_{11}, \dots, Y_{1n} } NOT
Group 2: Y_{21}, \dots, Y_{2n} } independent

② $Y_{ij} \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$ $Y_{2j} \stackrel{iid}{\sim} N(\mu_2, \sigma^2)$

③ $E(\bar{Y}_1 - \bar{Y}_2) = \mu_1 - \mu_2$

$$\text{Var}(\bar{Y}_1 - \bar{Y}_2) = \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n} - 2\text{Cov}(\bar{Y}_1, \bar{Y}_2)$$

Let $w_j = Y_{1j} - Y_{2j}$, thus $w_j \stackrel{iid}{\sim} N(\mu_1 - \mu_2, \sigma_w^2)$

$$\sigma_1^2 + \sigma_2^2$$

$$\bar{w} \pm t_{\alpha/2, n-1} \cdot SE(\bar{Y}_1 - \bar{Y}_2)$$

confidence interval for π :

(one sample)

- Assumptions:
- ① $Y_i = 0$ (failure); $Y_i = 1$ (success)
 - ② $Y_1 \dots Y_n$ are iid
 - ③ $E(Y_i) = \pi$; $\text{Var}(Y_i) = \pi(1-\pi)$
 - ④ $E(\bar{Y}) = \pi$; $\text{SE}(\bar{Y}) = \sqrt{\frac{\pi(1-\pi)}{n}}$

Wald - typed CI (from CLT) ← know how to implement this and know its drawbacks

When n is large, $\frac{P(1-P)}{n} \approx \frac{\pi(1-\pi)}{n}$, and $P = \bar{Y}$ is normal:

$$P \pm Z_{\alpha/2} \pm \sqrt{\frac{P(1-P)}{n}}$$

- Issues:
- ① can include negative values or > 1 values
 - ② doesn't work at $p=0$ or $p=1$
 - ③ no coverage guarantee, especially when n is smaller.

solution

Wilson - Score Interval ← at least remember the name and know that it's better

$$\frac{P + Z_{\alpha/2}^2 / 2n}{1 + Z_{\alpha/2}^2 / n} \pm Z_{\alpha/2} \sqrt{\frac{P(1-P)/n + Z_{\alpha/2}^2 / 4n^2}{1 + Z_{\alpha/2}^2 / n}}$$

(two sample - independent)

Assumptions:

- ① Samples from 2 different populations
 - Group 1: $Y_{11} \dots Y_{1n_1}$
 - Group 2: $Y_{21} \dots Y_{2n_2}$
- mutually
independent
- ② $E(Y_{ij}) = \pi_i$, $E(Y_{2j}) = \pi_2$
 - ③ $E(P_1 - P_2) = \pi_1 - \pi_2$
- $$\text{SE}(P_1 - P_2) = \sqrt{\frac{\pi_1(1-\pi_1)}{n_1} + \frac{\pi_2(1-\pi_2)}{n_2}}$$

$$(P_1 - P_2) \pm Z_{\alpha/2} \sqrt{\frac{P_1(1-P_1)}{n_1} + \frac{P_2(1-P_2)}{n_2}}$$

Confidence Interval v.s. Prediction Interval

* The confidence interval provides information about the mean of \hat{Y} at a given $X = x_0$.

- Expect $(1-\alpha)\%$ of intervals to capture $\beta_0 + \beta_1 x_0$.

$$P(\beta_0 + \beta_1 x_0 \in CI) = 1-\alpha$$

* The prediction interval captures Y at $X = x_0$.

- $Y = \beta_0 + \beta_1 x_0 + E$ is random and independent of the collected data

$$P(Y \in PI) = 1-\alpha$$

- Both Y and PI are random quantities.

$$\text{Var}(Y - \hat{Y}) = \text{Var}(Y) + \text{Var}(\hat{Y}) = \sigma^2 \left(1 + \left(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum(x_i - \bar{x})^2} \right) \right)$$

↑ assume Y and \hat{Y} independent, because Y_{n+1} is the future data, \hat{Y} from observed

$$Y - \hat{Y} \sim N(0, \text{Var}(Y - \hat{Y})) \Rightarrow \frac{Y - \hat{Y}}{\text{SE}(Y - \hat{Y})} \sim t_{n-2} \quad \begin{matrix} \text{df of MSE:} \\ \text{df total: } n-1 \\ \text{df Reg: 1} \end{matrix}$$

$$\Rightarrow P(-t_{n-2, \alpha/2} \leq \frac{Y - \hat{Y}}{\sqrt{\text{MSE}(1 + \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum(x_i - \bar{x})^2})}} \leq t_{n-2, \alpha/2}) = 1-\alpha$$

$$\Rightarrow \hat{Y} \pm t_{n-2, \alpha/2} \cdot \sqrt{\text{MSE}(1 + \frac{1}{n} + \frac{(\bar{x} - x_0)^2}{\sum(x_i - \bar{x})^2})}$$

→ if already have $\text{Var}(\hat{Y})$, simply "+ \hat{Y}^2 "

This example is in SLR case only, but the idea remains the same, that is: with $\text{Var}(\hat{Y})$ known, $\text{Var}(Y_{n+1}) = \text{MSE} + \text{Var}(\hat{Y})$.

Shortest CI: Minimum width for fixed coverage

UMA CI: Highest accurate coverage (dual to UMP test)

UMAU CI: Best unbiased CI (dual to UMPU test)

Example:

Normal mean with known variance

Normal mean with unknown variance

Exponential mean

Binomial proportion

CI-Type

Notes

Equal-tailed interval

UMA; shortest

t-interval

UMA unbiased

Shortest interval not symmetric

UMA, not equal tail

Clopper-Pearson interval

Unbiased, not shortest

Hypothesis Testing (I)

General Rules

		Size ↑	
		H ₀ true	H ₀ false
Reject H ₀	Type I error (α)	Power (1 - β)	
	Correct	Type II error (β)	
Accept H ₀			

P-value

Probability $T(\hat{\theta})$ takes on a value as or more extreme than $T(\hat{\theta})$ assuming H_0 .

Pivotal quantity

A pivotal quantity is a function of the data and parameters that has a probability distribution that does not depend on any unknown parameters.

Suppose $X \sim f(x; \theta)$, then

↪ $Q(X, \theta)$ is a pivotal quantity \Leftrightarrow the distribution of Q does not depend on θ
 ↪ Then can construct CI using the distribution:

$$\begin{aligned} P(Dist_{0.975} \leq Q(X, \theta) \leq Dist_{0.05}) &= 0.95 \\ \Rightarrow P((\cdot) \leq \theta \leq (\cdot)) &= 0.95 \end{aligned}$$

Critical value

The critical value is a quantile (cutoff point) of the pivotal quantity's known distribution.

e.g. $Z_{0.025}$, $t_{df_E, 0.025}$, $\chi^2_{df, 0.05}$, $F_{df_U, df_L, 0.05}$

Definition 8.3.5 For $0 \leq \alpha \leq 1$, a test with power function $\beta(\theta)$ is a *size α test* if $\sup_{\theta \in \Theta_0} \beta(\theta) = \alpha$.

Definition 8.3.6 For $0 \leq \alpha \leq 1$, a test with power function $\beta(\theta)$ is a *level α test* if $\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$.

Risk function of a decision rule:

$$d(\hat{\theta}) : \bar{x}, X_{\text{median}}, S_n^2$$

$R_d(\theta) = E_\theta [L(d(x), \theta)]$ the expected loss of decision rule $d(x)$ when θ is true.

Estimation: $R_d(\theta) = E_\theta [(d(x) - \theta)^2]$; $R_d(\theta) = E_\theta [|d(x) - \theta|]$

Testing:
$$\begin{aligned} R_\phi(\theta) &= \begin{cases} \text{loss(false reject)} P_\theta(\phi \text{ rejects } H_0) & \text{if } \theta \in \Sigma_0 \\ \text{loss(false accept)} P_\theta(\phi \text{ accepts } H_0) & \text{if } \theta \in \Sigma_L, \end{cases} \\ &= \begin{cases} \text{loss(false reject)} \Pi_{\phi(\theta)} & \text{if } \theta \in \Sigma_0 \\ \text{loss(false accept)} (1 - \Pi_{\phi(\theta)}) & \text{if } \theta \in \Sigma_L, \end{cases} \end{aligned}$$

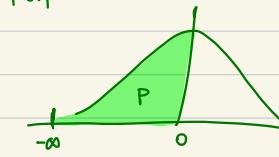
$$\Pi_{\phi(\theta)} = P_\theta(\phi \text{ rejects } H_0) = E_\theta(\phi(x))$$

χ^2 Goodness of fit test

H_0 : Observation follows * distribution

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{\text{#group} - \text{\#parameters} - 1}$$

pdf



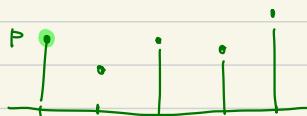
$(-\infty, P)$

$$O_i =$$

$$E_i = n \cdot p$$

$(-\infty, P)$	(P, ∞)
----------------	---------------

pmp



χ^2 independent test

e.g. testing proportion the same across periods:

$$H_0: p_1 = p_2 = p_3$$

$$\chi^2 = \sum_i \sum_j \frac{(E_{ij} - O_{ij})^2}{E_{ij}} \stackrel{H_0}{\sim} \chi^2_{(I-1)(J-1)} \quad \text{with } E_{ij} = \frac{n_{i..} n_{..j}}{n_{..}}$$

T-test / F-test

$$Y_{ij} = \mu + \beta_i + \varepsilon_{ij} \quad i = 1, \dots, 4, \quad \sum \beta_i = 0$$

$$H_0: \beta_1 = \beta_4 \quad (\varepsilon \sim N(0, \sigma^2))$$

① Find $\hat{\beta}_i$ (using restricted normal model or Lagrangian method).

② Must know $\hat{\beta}_1 - \hat{\beta}_4 \sim N(0, \sigma^2)$ under null

$$T = \frac{\hat{\beta}_1 - \hat{\beta}_4}{\sqrt{MSE / k(n)}} \sim t_{n-\text{rank}(X)}$$

$$F = \frac{\hat{\beta}_1^2}{MSE} \sim F_{1, n-\text{rank}(X)}$$

since only one β

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

under $H_0: \beta_1 = \beta_4$

$$SSE = \sum (y_i - \hat{y}_i)^2 \quad | \quad n-p$$

$$SSR = \sum (\hat{y}_i - \bar{y}_i)^2 \quad | \quad p-1$$

$$SST = \sum (y_i - \bar{y})^2 \quad | \quad n-1$$

Inference on variance of two populations

$$H_0: \sigma_1^2 = \sigma_2^2 \quad H_1: \sigma_1^2 \neq \sigma_2^2$$

$$\text{Recall } S_1^2(n_1-1)/\sigma_1^2 \sim \chi^2_{n_1-1}$$

$$S_2^2(n_2-1)/\sigma_2^2 \sim \chi^2_{n_2-1}$$

$$\text{Under } H_0, \quad \frac{S_1^2(n_1-1)/(n_1-1)}{S_2^2(n_2-1)/(n_2-1)} = \frac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1}$$

larger variance

smaller variance

Lack of fit test

n^* : unique x value

p: number of parameters

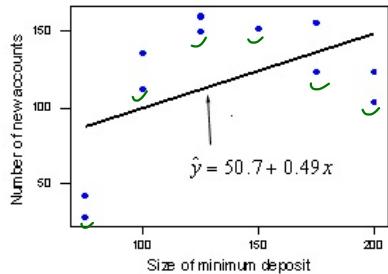
n: total # of observations

H_0 : There is no lack of fit $M_i = \beta_0 + \beta_1 X_i$

$$n^* = 6$$

$$p = 2$$

$$r_i = 2, 2, 2, 1, 2, 2$$



Source	DF	Adj SS	Adj MS	F-value
Regression	p-1	$\sum_{ij} (\bar{Y}_{..} - \hat{Y}_{ij})^2$	SSR / (p-1)	MSR / MSE
Residual Error	n-p	$\sum_{ij} (Y_{ij} - \hat{Y}_{ij})^2$	SSE / (n-p)	
Lack of Fit	n^*-p	$\left\{ \sum_{ij} (\bar{Y}_i - \hat{Y}_{ij})^2 \right\}$	SSLF / (n^*-p)	MSLF / MSPE
Pure Error	$n-n^*$	$\left\{ \sum_{ij} (Y_{ij} - \bar{Y}_{..})^2 \right\}$	SSPE / (n-n^*)	
Total	n-1	$\sum_{ij} (Y_{ij} - \bar{Y}_{..})^2$		

related to fit

doesn't relate to fit

$\sum_{i=1}^{n^*} \sum_{j=1}^{r_i} (Y_{ij} - \bar{Y}_{..})^2$

unique x ↑replicates

← if any point is unique
cannot measure pure error.
Then MSLF = MSE

Try it!

The lack of fit test

Fill in the missing numbers (?) in the following analysis of variance table resulting from a simple linear regression analysis.

Click on the light bulb in each cell to reveal the correct answer.

Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	??	12.597	??	??	0.000
Residual Error	??	??	??		
Lack of Fit	3	??	??	??	??
Pure Error	??	0.157	??		
Total	14	15.522			

Table A.2 Coefficients c_i for orthogonal polynomial trend contrasts

$v = 3$			$v = 4$						
Trend	c_1	c_2	c_3	Trend	c_1	c_2	c_3	c_4	
Linear	-1	0	1	Linear	-3	-1	1	3	
Quadratic	1	-2	1	Quadratic	1	-1	-1	1	
				Cubic	-1	3	-3	1	
$v = 5$									
Trend	c_1	c_2	c_3	c_4	c_5				
Linear	-2	-1	0	1	2				
Quadratic	2	-1	-2	-1	2				
Cubic	-1	2	0	-2	1				
Quartic	1	-4	6	-4	1				
$v = 6$									
Trend	c_1	c_2	c_3	c_4	c_5	c_6			
Linear	-5	-3	-1	1	3	5			
Quadratic	5	-1	-4	-4	-1	5			
Cubic	-5	7	4	-4	-7	5			
Quartic	1	-3	2	2	-3	1			
Quintic	-1	5	-10	10	-5	1			
$v = 7$									
Trend	c_1	c_2	c_3	c_4	c_5	c_6	c_7		
Linear	-3	-2	-1	0	1	2	3		
Quadratic	5	0	-3	-4	-3	0	5		
Cubic	-1	1	1	0	-1	-1	1		
Quartic	3	-7	1	6	1	-7	3		
Quintic	-1	4	-5	0	5	-4	1		
Sextic	1	-6	15	-20	15	-6	1		

Hypothesis testing in ST404

Lack of fit test using deviance

← this is particularly used in GLM.

If Y_i has a distribution "close to normal" with link "close to identity", then $\frac{D_M}{2p} \approx \chi^2_{n-p}$.

(Approximation will not improve as n increases)

e.g. consider n is the # of groups, where n_i is each group's size. want $n_i \uparrow$ instead of n .

$$\text{Bin}(n_i, p_i), i=1, \dots, n$$

H_0 : model M fits the data

H_1 : H_0 is not true

When n_i is reasonably large, can use this

Reject if $\frac{D_M}{2p} > \chi^2_{n-p,\alpha}$

However, this should only be applied when having normal / identity. If used this test under overdispersion (Bin/Poi) Type I error ↑. May perform poorly when expected counts are small.

Solve 1: Fit a quasi-binomial / quasi-poisson model to estimate dispersion

Solve 2: Use a bootstrap or simulation-based test

H_0 : Model M_0 with q reg. parameters

H_1 : Model M (with $M_0 \subseteq M$) with p reg. parameters

$$T = \frac{D_{M_0} - D_M}{2} \xrightarrow{d} \chi^2_{p-q}$$

under H_0

need to distinguish t distribution or χ^2 distribution, and distinguish related df.

$$E(\frac{D_M}{2}) = n-p \Rightarrow \hat{\phi}_{M_0} = \frac{D_M}{n-p}, \text{ when } \hat{\phi} \text{ large, we should concern:}$$

Reason 1: Inadequate linear predictor (underspecification)

Reason 2: Over-dispersion $\text{Var}(Y_i) > \text{Var}(h(y_i))$ where Var is given by family

• Correlated Y_1, \dots, Y_n may lead to over-dispersion

think this

Also see

next page

Asymptotic Wald Test

$$H_0: A\beta = m$$

$$H_1: A\beta \neq m$$

$$Y \sim N(X\beta, \sigma^2 I)$$

$$H_0: A\hat{\beta} = m \quad H_1: A\hat{\beta} \neq m \quad \leftarrow \text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1} \Rightarrow \text{Var}(\hat{\beta}) = \hat{\sigma}^2 (X'X)^{-1}$$

$$T_W = (A\hat{\beta} - m)' (A(\hat{F}' \hat{V}^{-1} \hat{F})^{-1} A')^{-1} (A\hat{\beta} - m) \sim \chi^2_{\text{rank}(A)}$$

$$H_0: h(\theta) = 0$$

$$H_1: h(\theta) \neq 0$$

$$T_W = h(\hat{\theta})' (\hat{H}(\hat{\theta}))^{-1} (\hat{H}(\hat{\theta}))' h(\hat{\theta}) \sim \chi^2_{\text{rank}(A)}$$

Test for correlation r_{xy}

$$r_{xy} = \frac{E((x-\bar{x})(y-\bar{y}))}{\sigma_x \sigma_y}$$

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{s_x s_y}$$

$$(x_i, y_i) \sim_{\text{iid}} N\left(\begin{pmatrix} \bar{x} \\ \bar{y} \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & r_{xy} \sigma_x \sigma_y \\ r_{xy} \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}\right)$$

$$H_0: r_{xy} = 0$$

$$H_1: r_{xy} \neq 0$$

$$Z = \frac{1}{2} \log\left(\frac{1+r_{xy}}{1-r_{xy}}\right) \sim N\left(\frac{1}{2} \log\left(\frac{1+r_{xy}}{1-r_{xy}}\right), \frac{1}{n-3}\right)$$

$$\text{RR: } |Z| \sqrt{n-3} > Z_{\alpha/2}$$

$$\stackrel{H_0}{\sim} N(0, \frac{1}{n-3})$$

Test on $h(\theta)$ with $\hat{\theta}_{MLE}$

This part was discussed under Lmm setting, but may apply on other occasions

$$H_0: h(\theta) = 0 \quad H_1: h(\theta) \neq 0 \quad \text{e.g. } h(\theta) = \underset{rx_1}{AB} - m$$

- $H(\theta) = \frac{\partial h(\theta)}{\partial \theta^T}$ $r \times d$ first partial derivatives of $h(\theta)$
- $\hat{\theta}$ is MLE, $\hat{\theta}_0$ is MLE under H_0
- When H_0 is true, all of $T_W, T_{LR}, T_S \xrightarrow{d} \chi_d^2$ See here, unlike the t-statistics, the df is not dfError, but $\text{rank}(A)$.

Wald statistics:

$$T_W = \frac{h(\hat{\theta})^T \{ H(\hat{\theta}) I(\hat{\theta})^{-1} H(\hat{\theta})^T \}^{-1} h(\hat{\theta})}{\text{Var}(h(\hat{\theta}))} \quad \text{only involves } \hat{\theta}$$

Likelihood ratio statistics:

$$T_{LR} = -2 [l(\hat{\theta}_0; y) - l(\hat{\theta}; y)]$$

Score statistics:

$$T_S = S(\hat{\theta}_0)^T I(\hat{\theta}_0)^{-1} S(\hat{\theta}_0) \quad \text{only involves } \hat{\theta}$$

HT terminologies

- A test ϕ is consistent if $\text{power}(\phi) = \Pi_\phi(\phi) \rightarrow 1$ \leftarrow as data sample size \uparrow
- A test ϕ is unbiased if $\sup_{\theta \in \Theta_0} \Pi_\phi(\theta) \leq \inf_{\theta \in \Theta_1} \Pi_\phi(\theta) \leftarrow p(\text{reject } H_0) \leq p(\text{reject } H_1)$
 $\equiv \text{Power} \geq \alpha$ always
- When $H_0: \theta = a$; $H_1: \theta = b$, try using Neymann-Pearson Lemma show MP(LR)
(when $H_0: \theta = a$; $H_1: \theta > a$ or $\theta < a$, using Karlin-Rubin Lemma show UMP(LR))
Both show LR monotonic with test statistics $T(x)$.

When $H_0: \theta = a$; $H_1: \theta \neq a$ cannot use Neymann Pearson/ Karlin-Rubin Lemma.

Cannot be two-sided H_1

One-sample HT

(HT for μ - known σ^2)

$$H_0: \mu = \mu_0$$

$$\bar{Y} \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

$$Z = \frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} \sim N\left(\frac{\mu - \mu_0}{\sigma/\sqrt{n}}, 1\right) \stackrel{H_0}{\sim} N(0, 1)$$

$$H_1: \mu > \mu_0 \quad RR = Z > C = Z_\alpha$$

$$H_1: \mu < \mu_0 \quad RR = Z < C = -Z_\alpha$$

$$H_1: \mu \neq \mu_0 \quad RR = Z < -Z_{\alpha/2} \text{ or } Z > Z_{\alpha/2}$$

} remember when in favor of H_1 , reject H_0 .

} Thus rejection region depends on H_1 .

(HT for μ - unknown σ^2)

$$H_0: \mu = \mu_0$$

$$\bar{Y} \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim t_{n-1}$$

$$T = \frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} \sim t_{n-1}, \lambda \stackrel{H_0}{\sim} t_{n-1}$$

non-central parameter $\lambda = \frac{\mu - \mu_0}{\sigma/\sqrt{n}}$

$$H_1: \mu > \mu_0 \quad RR = Z > C = t_{n-1, \alpha}$$

$$H_1: \mu < \mu_0 \quad RR = Z < C = -t_{n-1, \alpha}$$

$$H_1: \mu \neq \mu_0 \quad RR = Z < -t_{n-1, \alpha/2} \text{ or } Z > t_{n-1, \alpha/2}$$

Fact: $T(\hat{\theta}) \in RR$ iff p-value < α .

p-value

H_1	known σ^2	unknown σ^2
$\mu > \mu_0$	$1 - \Phi(Z) = \Phi(-Z)$	$1 - F_{t_{n-1,0}}(t) = F_{t_{n-1,0}}(-t)$
$\mu < \mu_0$	$\Phi(Z)$	$F_{t_{n-1,0}}(t)$
$\mu \neq \mu_0$	$2(1 - \Phi(Z))$	$2(1 - F_{t_{n-1,0}}(t))$

when Z small, $\Phi(Z)$

is small, then reject

when Z large, $\Phi(-Z)$ is small, thus reject H_0 .

when $|Z|$ large, $1 - \Phi(|Z|)$ small, thus reject

(HT for π)

Exact test: ← uses binomial dist.

$$Y_i \stackrel{iid}{\sim} \text{Bin}(n, \pi)$$

$$H_0: \pi = \pi_0$$

Test statistics: $Y = \sum_i Y_i \sim \text{Bin}(n, \pi)$

$$P(Y=y) = \binom{n}{y} \pi^y (1-\pi)^{n-y} \quad y=0, 1, \dots, n$$

$$H_1: \pi > \pi_0$$

$$P = P(Y \geq y | \pi_0) \\ = \sum_{j=y}^n \binom{n}{j} \pi_0^j (1-\pi_0)^{n-j}$$

P-values: maybe no closed form for RR.

$$H_1: \pi < \pi_0$$

$$P = P(Y \leq y | \pi_0) \\ = \sum_{j=0}^y \binom{n}{j} \pi_0^j (1-\pi_0)^{n-j}$$

→ observed y

$$H_1: \pi \neq \pi_0$$

$K = \{k : P(Y=k | \pi_0) \leq P(Y=y | \pi_0)\} \rightarrow$ means more extreme

$$p = P(Y \in K | \pi_0) = \sum_{k \in K} P(Y=k | \pi_0)$$

Approximation:

$$p = 2 \cdot \min \{P(Y \geq y | \pi_0), P(Y \leq y | \pi_0)\}$$

RAO-Score test: ← use CLT

$$H_0: \pi = \pi_0$$

$$Z = \frac{P - \pi_0}{\sqrt{\frac{\pi_0(1-\pi_0)}{n}}} \stackrel{H_0}{\sim} N(0, 1) \quad p = \bar{Y}$$

→ be careful for this part

$$H_1: \pi > \pi_0$$

$$\text{RR: } Z > Z_{\alpha}$$

↙ same as one sample procedure with known σ^2

$$H_1: \pi < \pi_0$$

$$\text{RR: } Z < -Z_{\alpha/2}$$

$$H_1: \pi \neq \pi_0$$

$$\text{RR: } Z < -Z_{\alpha/2} \text{ or } Z > Z_{\alpha/2}$$

Two-sample HT

$(\mu_1, \mu_2 \text{ independent})$

$$Y_{ij} \stackrel{iid}{\sim} N(\mu_1, \sigma_1^2) \quad Y_{2j} \stackrel{iid}{\sim} N(\mu_2, \sigma_2^2) \quad \text{indep.}$$

$$\bar{Y}_1 - \bar{Y}_2 \sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$$

$$H_0: \mu_1 - \mu_2 = \Delta_0$$

	Test Statistics	Null Distribution	Alternative Distribution	
$\sigma_1^2 = \sigma_2^2$	$\frac{(\bar{Y}_1 - \bar{Y}_2 - \Delta_0)}{SP\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$	$t_{n_1+n_2-2}$	$t_{n_1+n_2-2, \lambda}$	$: SP = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}}$
$\sigma_1^2 \neq \sigma_2^2$	$\frac{(\bar{Y}_1 - \bar{Y}_2 - \Delta_0)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$	t_v	$t_{v, \tilde{\lambda}}$	$: v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$

$$\lambda = \frac{(\mu_1 - \mu_2 - \Delta_0)}{\sigma^2 \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\tilde{\lambda} = \frac{(\mu_1 - \mu_2 - \Delta_0)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$



when Y_{ij}, Y_{2j} paired, $w_j = Y_{ij} - Y_{2j}$, $H_0: \mu_1 - \mu_2 = \Delta_0$

	Test Statistics	Null Distribution	Alternative Distribution	
	$\frac{(\bar{Y}_1 - \bar{Y}_2 - \Delta_0)}{\sqrt{\frac{Var(w_j)}{n}}}$	t_{n-1}	$t_{n-1, \lambda}$	$\lambda = \frac{(\mu_1 - \mu_2 - \Delta_0)}{\sqrt{w}/\sqrt{n}}$

$(\pi_1, \pi_2 \text{ independent})$

$$Y_{ij} \stackrel{iid}{\sim} \text{Bin}(\pi_1) \quad Y_{2j} \stackrel{iid}{\sim} \text{Bin}(\pi_2) \quad \text{indep.}$$

$$H_0: \pi_1 - \pi_2 = \Delta_0$$

H_0	Test statistics	Null Distribution
$\pi_1 - \pi_2 = 0$	$Z = \frac{p_1 - p_2}{\sqrt{p(1-p)} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$	$N(0, 1)$
$\pi_1 - \pi_2 = \Delta_0$	$Z = \frac{p_1 - p_2 - \Delta_0}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$	$N(0, 1)$

$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2} \text{ estimates } \pi \quad (\text{if } H_0: \pi_1 = \pi_2 = \pi)$$

Power & Sample Size

Power = $P(\text{Reject } H_0 \mid H_1)$

(Power for μ , known σ^2)

$$\begin{aligned} Z &\sim N\left(\frac{\mu - \mu_0}{\sigma/\sqrt{n}}, 1\right) \\ \Rightarrow Z - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} &\sim N(0, 1) \end{aligned}$$

H_1	RR	Power	$P(Z > Z_\alpha \mid \mu > \mu_0)$
$\mu > \mu_0$	$Z > Z_\alpha$	$1 - \Phi\left(Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) = \Phi\left(\frac{\mu - \mu_0}{\sigma/\sqrt{n}} - Z_\alpha\right)$	$= P\left(\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} > Z_\alpha \mid \mu > \mu_0\right)$
$\mu < \mu_0$	$Z < -Z_\alpha$	$\Phi\left(-Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right)$	$= P\left(\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} > -Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} \mid \mu < \mu_0\right)$
$\mu \neq \mu_0$	$ Z > Z_{\alpha/2}$	$\Phi\left(\frac{\mu - \mu_0}{\sigma/\sqrt{n}} - Z_{\alpha/2}\right) + \Phi\left(-Z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right)$	$= 1 - \Phi\left(Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right)$

effect size $\mu - \mu_0$

power depends on
the alternative distribution

Q: what is the smallest sample size to achieve $1 - \beta$ sample size.

$$\begin{aligned} \text{e.g. } 1 - \Phi\left(Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) &\geq 1 - \beta & \Phi\left(Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) &\geq 1 - \beta \\ \Phi\left(Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) &\leq \beta & \Phi\left(Z_\alpha + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) &\leq \beta \\ Z_\alpha - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} &\leq Z_\beta & Z_\alpha + \frac{\mu - \mu_0}{\sigma/\sqrt{n}} &\leq Z_\beta \\ \left(\frac{Z_\alpha - Z_\beta}{\mu - \mu_0/\sigma}\right)^2 &\leq n & \left(\frac{Z_\alpha - Z_\beta}{\mu - \mu_0/\sigma}\right)^2 &\leq n \end{aligned}$$

$$\mu \neq \mu_0 \Rightarrow \left(\frac{Z_{\alpha/2} - Z_\beta}{\mu - \mu_0/\sigma}\right)^2 \leq n$$

(power for μ , unknown σ^2)

Recall when σ^2 unknown $\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} \sim t_{n-1, \lambda} = \frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}}$

H_1	RR	Power
$\mu > \mu_0$	$t > t_{n-1, \lambda}$	$1 - F_{t_{n-1, \lambda}}(t_{n-1, \lambda})$
$\mu < \mu_0$	$t < -t_{n-1, \lambda}$	$F_{t_{n-1, \lambda}}(-t_{n-1, \lambda})$
$\mu \neq \mu_0$	$ t > t_{n-1, \lambda/2}$	$1 - F_{t_{n-1, \lambda}}(t_{n-1, \lambda/2}) + F_{t_{n-1, \lambda}}(-t_{n-1, \lambda/2})$

$$\begin{aligned} P(t > t_{n-1, \lambda} \mid \mu > \mu_0) &= P\left(\frac{\bar{Y} - \mu_0}{\sigma/\sqrt{n}} > t_{n-1, \lambda} \mid \mu > \mu_0\right) \\ &= 1 - F_{t_{n-1, \lambda}}(t_{n-1, \lambda}) \end{aligned}$$

(Power for $\mu_1 - \mu_2$, independent and paired)

$$\text{indep: } \lambda = \frac{\mu_1 - \mu_2 - \Delta_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$\text{pairs: } \lambda = \frac{\mu_1 - \mu_2 - \Delta_0}{\sqrt{n_1 + n_2}}$$

(power for π)

① Power for exact test hard to derive (no alternative)

② Alternative for Rao Score test is Normal.

$$\text{Recall } Z = \frac{\hat{\pi} - \pi_0}{\sqrt{\pi_0(1-\pi_0)/n}} \sim N\left(\frac{\pi - \pi_0}{\sqrt{\pi_0(1-\pi_0)/n}}, \frac{\pi(1-\pi)}{\pi_0(1-\pi_0)}\right)$$

thus,

$$Z \cdot \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} = \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} \sim N(0, 1)$$

H_1	RR	Power
$\pi > \pi_0$	$Z > Z_{\alpha}$	$\Phi\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} - Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}}\right)$
$\pi < \pi_0$	$Z < -Z_{\alpha}$	$\Phi(-Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n})$
$\pi \neq \pi_0$	$ Z > Z_{\alpha}$	$\Phi\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} - Z_{\alpha/2} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}}\right) + \Phi\left(-Z_{\alpha/2} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n}\right)$

$$P(Z > Z_{\alpha} \mid \pi > \pi_0)$$

$$= P\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} > Z_{\alpha} \mid \pi > \pi_0\right)$$

$$= P\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} \cdot \frac{\sqrt{\pi(1-\pi)}}{\sqrt{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} > Z_{\alpha} \cdot \frac{\sqrt{\pi(1-\pi)}}{\sqrt{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} \mid \pi > \pi_0\right)$$

$$= 1 - \Phi\left(Z_{\alpha} \cdot \frac{\sqrt{\pi(1-\pi)}}{\sqrt{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n}\right)$$

$$= \Phi\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} - Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}}\right)$$

H_1	RR	Sample Size
$\pi > \pi_0$	$Z > Z_{\alpha}$	$\left(\frac{Z_{\alpha} \sqrt{\pi(1-\pi)} - Z_{\beta} \sqrt{\pi(1-\pi)}}{\pi - \pi_0}\right)^2$
$\pi < \pi_0$	$Z < -Z_{\alpha}$	$\left(\frac{Z_{\alpha} \sqrt{\pi(1-\pi)} - Z_{\beta} \sqrt{\pi(1-\pi)}}{\pi - \pi_0}\right)^2$
$\pi \neq \pi_0$	$ Z > Z_{\alpha}$	$\left(\frac{Z_{\alpha/2} \sqrt{\pi(1-\pi)} - Z_{\beta} \sqrt{\pi(1-\pi)}}{\pi - \pi_0}\right)^2$

$$\Phi\left(\frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} - Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}}\right) \geq 1 - \beta$$

$$\Phi\left(Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n}\right) \leq \beta$$

$$Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n} \leq Z_{\beta}$$

$$Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - Z_{\beta} \leq \frac{\pi - \pi_0}{\sqrt{\pi(1-\pi)}} \sqrt{n}$$

$$Z_{\alpha} \sqrt{\frac{\pi(1-\pi)}{\pi(1-\pi)}} - Z_{\beta} \frac{\pi - \pi_0}{\pi - \pi_0} \leq \sqrt{n}$$

$$\left(\frac{Z_{\alpha} \sqrt{\pi(1-\pi)} - Z_{\beta} \sqrt{\pi(1-\pi)}}{\pi - \pi_0}\right)^2 \leq n$$

Multiple Testing

Family-wise Type I error rate:

$$1. P(\text{Reject } H_0^{ij} | M_i = M_j) = \alpha_i \leftarrow \text{individual test error rate}$$

2. $H_0: M_1 = \dots = M_k$ is then

$$\alpha = P\left(\bigcup_{i,j} \text{Reject } H_0^{ij} | M_1 = \dots = M_k\right) \leftarrow \text{family-wise error rate (FWER)}$$

$$\leq \sum_{i \neq j} P(\text{Reject } H_0^{ij} | M_1 = \dots = M_k)$$

$$= \binom{k}{2} \alpha_i$$

I'm not sure if this $\binom{k}{2}$ is firm, because $(k-1)$ tests are enough,
but that may be considered as Dunnett method (?)

Global F-test

$$H_0: L\theta = 0$$

① $L\theta$ is estimable

② $\text{rank}(L) = \text{rank}(X) - 1$

③ H_0 true $\Rightarrow \theta$ in the null space of L

$$F = \frac{\text{SSR}/\text{rank}(L)}{\text{SSE}/(n-1-\text{rank}(L))} \leftarrow \text{whatever } L \text{ is used, the test statistic is the same (for global F)}$$

$$\stackrel{H_0}{\sim} F_{\text{rank}(X)-1, n-\text{rank}(X)}$$

Global F-test : Full vs. Reduced model

$$F = \frac{\frac{\text{SSE}(H_0) - \text{SSE}(H_1)}{\text{df}(E(H_0)) - \text{df}(E(H_1))}}{\frac{\text{SSE}(H_1)/\text{df}(E(H_1))}{\text{rank}(X)-1}} = \frac{\frac{\text{SSE}(H_0) - \text{SSE}(H_1)}{\text{rank}(X)-1}}{\frac{\text{SSE}(H_1)/(n-\text{rank}(X))}{\text{rank}(X)-1}} \leftarrow \begin{array}{l} \text{since global} \\ \text{full model} \end{array}$$

General Rules:

1. If all pairwise, use Tukey
2. If pairwise with control, use Dunnett
3. If general $L\theta$ with small $\text{rank}(L)$, use Bonferroni \leftarrow when # of test is small
4. When in doubt, use scheffe
5. When improve Bonferroni, avoid massive rejection, use Benjamini
6. Fisher does nothing

Fisher

$$\hat{\ell}'\theta \pm t_{df_E, \alpha/2} * SE(\hat{\ell}'\theta) \quad \leftarrow \text{obviously incorrect}$$

Bonferroni

$$\hat{\ell}'\theta \pm t_{df_E, \alpha/2p} * SE(\hat{\ell}'\theta) \quad p = \binom{k}{2}$$

k : # of parameters involved in H_0 .

Benjamini

Benjamini-Hochberg Step-Up Procedure (For large k to avoid mass rejection)

1. Compute p -values $p_j = 2 \cdot \Pr(t_{df_{\text{error}}} > |\hat{\theta}_j| / SE(\hat{\theta}_j))$. Goal: Control FDR $< \frac{k\alpha}{k} \alpha$
2. Reject H_j^0 if $p_j < T_{BH}$ where $T_{BH} := \max\{p_{(j)} : p_{(j)} < \alpha \frac{j}{k}, 1 \leq j \leq k\}$

Note: FDR is the expected ratio of the number of *falsely* rejected null hypotheses.

two-sided p -value

rule

Tukey

$$\hat{\ell}'\theta \pm q t_{k, df_E, \alpha} * SE(\hat{\ell}'\theta)$$

$$\text{Tukey HSD: } q = \frac{|\hat{\ell}'\hat{\theta}|}{SE(\hat{\ell}'\hat{\theta})} \Rightarrow \text{e.g. } \frac{|\bar{Y}_i - \bar{Y}_j|}{\sqrt{MSE(\frac{1}{n_i} + \frac{1}{n_j})}}$$

This Tukey HSD follows a studentized range dist

Scheffe:

No interpretation:

$$\hat{\ell}'\theta \pm \sqrt{(\text{rank}(x)-1) F_{\text{rank}(x)-1, n-\text{rank}(x), df_{\text{regression}}}} * \frac{SE(\hat{\ell}'\theta)}{\sqrt{df_{\text{regression}}}}$$

\therefore # of estimable functions
(independent) \curvearrowright max-rank(L)

e.g. one-way ANOVA

Intercept:

$$\hat{\ell}'\theta \pm \sqrt{\text{rank}(x) F_{\text{rank}(x), n-\text{rank}(x), df_{\text{Error}}}} * SE(\hat{\ell}'\theta)$$

for full rank linear model: $(1+p)$

e.g. full rank linear model

or maybe ANOVA with μ ?

Sum of square in testing

$$R(\beta_1 | \beta_0) = SSR(\beta_0, \beta_1) - SSR(\beta_0)$$

$$R(\beta_1, \beta_2 | \beta_0, \beta_1) = SSR(\beta_0, \beta_1, \beta_2, \beta_3) - SSR(\beta_0, \beta_1)$$

$$= R(\beta_1, \beta_2, \beta_3 | \beta_0) + SSR(\beta_0) - (R(\beta_1 | \beta_0) + SSR(\beta_0))$$

$$= R(\beta_1, \beta_2, \beta_3 | \beta_0) - R(\beta_1 | \beta_0)$$

■ Special case: model A sets $p - q$ parameters to 0

■ Arrange model B parameters: $(\beta_0, \dots, \beta_q, \beta_{q+1}, \dots, \beta_p)$

■ H_0 : Model A "true" $\Leftrightarrow \beta_{q+1} = \dots = \beta_p = 0$

■ H_1 : Model B "true" \Leftrightarrow Not all $\beta_{q+1}, \dots, \beta_p$ equal 0

$$\begin{aligned} R(\beta_{q+1}, \dots, \beta_p | \beta_0, \dots, \beta_q) &= SSR(H_1) - SSR(H_0) \\ &= R(\beta_1, \dots, \beta_p | \beta_0) - R(\beta_1, \dots, \beta_p | \beta_0) \end{aligned}$$

Hypothesis Testing (2)

F-test

$$Y \sim N_n(X\beta, \sigma^2 I_n)$$

$$\Rightarrow X'\hat{\beta} \sim N(X'\beta, \sigma^2 X'(X'X)^{-1}X)$$

$$\text{independent of } \hat{\sigma}^2 = \frac{1}{n-r} Y'(I-P_X)Y \text{ where } \frac{(n-r)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-r}$$

$$\frac{Y'(I-P_X)Y}{\sigma^2} \sim \chi^2_{n-r}$$

Theorem: In the model $Y \sim N_n(X\beta, \sigma^2 I_n)$ with unknown β, σ^2 .

$$\hat{\beta}_{MLE} = (X'X)^{-1}X'Y \quad \text{and} \quad \hat{\sigma}_{MLE}^2 = \frac{1}{n-r} Y'(I-P_X)Y$$

Definition:

$$H_0: K'\beta = m$$

$$H_1: K'\beta \neq m \quad \leftarrow \text{estimable if } K \in \text{col}(X')$$

where $K \in \mathbb{R}^{P \times S}$ with $\text{rank}(K) = S$ and every column of K is in $\text{col}(X')$.

The general linear hypothesis is said to be testable iff $K'\beta$ is estimable and has full rank.

Lemma: If $K'\beta$ is estimable, then $H := K'(X'X)^{-1}K \in \mathbb{R}^{S \times S}$ is non-singular.

$$H = K'(X'X)^{-1}K = K'(X'X)^{-1}(X'X)(X'X)^{-1}K = K'(X'X)^{-1}X'W$$

$$\text{rank}(W) \leq \min\{n, s\} \leq S = \text{rank}(K) = \text{rank}(X'X(X'X)^{-1}K) = \text{rank}(X'W) \leq \text{rank}(W)$$

$$S = \text{rank}(W) = \text{rank}(W'W) = \text{rank}(H)$$

Example: $Y \sim N_n(X\beta, \sigma^2 I_n)$ $K'\hat{\beta} \sim N(K'\beta, \sigma^2 K'(X'X)^{-1}K) \equiv N(K'\beta, \sigma^2 H)$

$$E(K'\hat{\beta}) = M \xrightarrow{H_0} K'\hat{\beta} - M \xrightarrow{H_0} N(0, \sigma^2 H)$$

$$(K'\hat{\beta} - M)' H^{-1} (K'\hat{\beta} - M) / \sigma^2 \sim \chi^2_S$$

Lemma: Let $Y_n \sim N_n(\mu, V)$ and A, B symmetric matrices. Then if $BVA = 0$, $Y'Ay$ and $Y'By$ are independent.

$$\rightarrow \sigma^{-2} Y'(I-P_X)Y \sim \chi^2_{n-r}$$

$$\text{Consider } V = \sigma^2 I, A = (I-P_X), H = K'(X'X)^{-1}K = LL', B = L^{-1}K'(X'X)^{-1}K\sigma^{-1}$$

Because $BVA = 0$, BY is independent of $Y'Ay$.

$$F = \frac{(K'\hat{\beta} - M)' H^{-1} (K'\hat{\beta} - M) / \sigma^2 S}{Y'(I-P_X)Y / \sigma^2 (n-r)} \xrightarrow{H_0} F_{S, n-r}$$

$$\xrightarrow{H_1} F_{S, n-r} \left(\frac{1}{\sigma^2} (K'\hat{\beta} - M)' (K'\hat{\beta} - M) \right)$$

Likelihood Ratio Test

$$\Omega_0 := \{(\beta, \sigma^2) : k'\beta = m, \sigma > 0\}$$

$$\Omega := \{(\beta, \sigma^2) : \beta \in \mathbb{R}^p, \sigma > 0\}$$

$$LR := \frac{\max_{\Omega_0} \{L(\beta, \sigma^2)\}}{\max_{\Omega} \{L(\beta, \sigma^2)\}}$$

$$P(LR < c \mid H_0) = \alpha$$

Example: $LR = \left(\frac{\|y - X\hat{\beta}\|^2}{\|y - X\hat{\beta}_H\|^2} \right)^{\frac{n}{2}} = \left(\frac{Q(\hat{\beta})}{Q(\hat{\beta}_H)} \right)^{\frac{n}{2}}$

$\hat{\beta}_H$ is from RNE

$$\begin{pmatrix} X'X & P \\ P' & 0 \end{pmatrix} \begin{pmatrix} \beta_H \\ \theta_H \end{pmatrix} = \begin{pmatrix} X'y \\ s \end{pmatrix}$$

with $P = k$, $s = m$

$$\text{Reject } H_0 \text{ if } \left(\frac{Q(\hat{\beta})}{Q(\hat{\beta}_H)} \right)^{\frac{n}{2}} < c \quad \text{or} \quad \frac{[Q(\hat{\beta}_H) - Q(\hat{\beta})]/s}{Q(\hat{\beta})/(n-r)} > \frac{n-r}{s} (c^{-\frac{2}{n}} - 1)$$

This is actually an F-test.

Theorem: If $k'\beta$ is estimable and $\hat{\beta}_H$ is part of a solution to the RNEs with constraint $k'\beta = m$, then

$$\begin{aligned} Q(\hat{\beta}_H) - Q(\hat{\beta}) &\stackrel{(1)}{=} (\hat{\beta}_H - \hat{\beta})' X'X (\hat{\beta}_H - \hat{\beta}) \\ &\stackrel{(2)}{=} (k'\hat{\beta} - m)' (k'(X'X)^{-1} k) (k'\hat{\beta} - m) \\ &\stackrel{(3)}{=} (k'\hat{\beta} - m)' H^{-1} (k'\hat{\beta} - m) \end{aligned}$$

$$\begin{aligned} \text{Proof: } Q(\hat{\beta}_H) - Q(\hat{\beta}) &= \|y - X\hat{\beta}_H\|^2 - \|y - X\hat{\beta}\|^2 \\ &= \|X(\hat{\beta} - \hat{\beta}_H)\|^2 \end{aligned}$$

$$\begin{aligned} (1) &= (\hat{\beta} - \hat{\beta}_H)' X'X (\hat{\beta} - \hat{\beta}_H) \longrightarrow = (\hat{\beta} - \hat{\beta}_H)' K \hat{\theta}_H \\ &= (\hat{\beta} - \hat{\beta}_H)' X'X (X'X)^{-1} X (\hat{\beta} - \hat{\beta}_H) = (\hat{\beta} - \hat{\beta}_H)' K H^{-1} K (\hat{\beta} - \hat{\beta}_H) \\ &\rightarrow = \hat{\theta}_H' K' (X'X)^{-1} K \hat{\theta}_H \quad (3) = (k'\hat{\beta} - m)' H^{-1} (k'\hat{\beta} - m) \\ &= \hat{\theta}_H' H \hat{\theta}_H \end{aligned}$$

$$X'\hat{\beta}_H + K\hat{\theta}_H = X'y$$

$$X'\hat{\beta}_H = X'y - K\hat{\theta}_H$$

$$\rightarrow X'(\hat{\beta} - \hat{\beta}_H) = K\hat{\theta}_H$$

$$K'(X'X)^{-1} X'(\hat{\beta} - \hat{\beta}_H) = K'(X'X)^{-1} K \hat{\theta}_H$$

$$\begin{aligned} \hat{\theta}_H &= H^{-1} K'(X'X)^{-1} X'(\hat{\beta} - \hat{\beta}_H) \\ &= H^{-1} (A'X)(X'X)^{-1} X'(\hat{\beta} - \hat{\beta}_H) \\ &= H^{-1} K'(\hat{\beta} - \hat{\beta}_H) \\ \rightarrow & \text{Since } K = X'A \end{aligned}$$

Corollary: If $K'\beta$ is estimable and $\hat{\beta}$ solves the normal equations, then the $\hat{\beta}_H$ component of a solution to the RNE solver,

$$X'X\beta = X'y - K(K'(X'X)^{-1}K)\hat{\beta} - m$$

Proof: $\hat{\beta}_H$ solves $X'X\beta + K\hat{\theta}_H = X'y$

$$X'X\beta = X'y - K\hat{\theta}_H$$

$$= X'y - K H^{-1} K (\hat{\beta} - \hat{\beta}_H)$$

$$\rightarrow = X'y - K H^{-1} (K \hat{\beta} - m)$$

Theorem: If $P'\beta$ is a system of linear independent, non-estimable functions and $\hat{\beta}_H$ is part of a solution to the RNEs with constraint $P'\beta = \delta$, then

$$Q(\hat{\beta}) = Q(\hat{\beta}_H) \text{ and } \hat{\theta}_H = 0$$

Proof: $X'X\beta_H + P\hat{\theta}_H = X'y$, and so $P\hat{\theta}_H = X'(y - X\beta_H) \in \text{col}(X')$

Since $P'\beta$ is non-estimable, $\text{col}(P) \cap \text{col}(X') = \{0\}$, which implies $P\hat{\theta}_H = 0$.

P has full column rank, so it must be the case that $\hat{\theta}_H = 0$. Thus,

$$\hat{\beta}_H \in \{X'X\beta = X'y\}$$

So that $Q(\hat{\beta}) = Q(\hat{\beta}_H)$

Confidence Interval & Multiple Comparison

Single estimable function $\lambda'\beta$

$Y \sim N_n(X\beta, \sigma^2 I_n)$, recall if $\lambda'\beta$ is estimable, then for $\hat{\beta} = (X'X)^{-1}X'y$,

$\lambda'\hat{\beta} \sim N_1(\lambda'\beta, \sigma^2 \lambda'(X'X)^{-1}\lambda)$. If we estimate σ^2 with $\hat{\sigma}^2 = \frac{1}{n-r} Y'(I-P)x$,

then

$$t = \frac{\lambda'\hat{\beta} - \lambda'\beta}{\sqrt{\hat{\sigma}^2 \lambda'(X'X)^{-1}\lambda}} \sim T_{n-r}$$

thus

$$P(|t| \leq T_{n-r, \alpha/2}) = 1 - \alpha \quad \leftarrow \text{CI coverage}$$

$$\text{CI : } \lambda'\beta \in [\lambda'\hat{\beta} - T_{n-r, \alpha/2} \cdot \sqrt{\hat{\sigma}^2 \lambda'(X'X)^{-1}\lambda}, \lambda'\hat{\beta} + T_{n-r, \alpha/2} \cdot \sqrt{\hat{\sigma}^2 \lambda'(X'X)^{-1}\lambda}]$$

Multiple estimable function $\Lambda' \beta$

Consider s estimable functions, $\lambda_1' \beta, \lambda_2' \beta, \dots, \lambda_s' \beta$, then let $\Lambda = (\lambda_1 \ \dots \ \lambda_s)$ with linearly independent columns,

$$\Lambda' \beta \sim N_s(\Lambda' \beta, \sigma^2 \Lambda' (\mathbf{X} \mathbf{X})^{-1} \Lambda)$$

$$\Lambda' = \begin{pmatrix} \lambda_1' \\ \vdots \\ \lambda_s' \end{pmatrix}$$

In this case,

$$P(a_j \leq \Lambda_j' \beta \leq b_j) = 1-\alpha, \quad \forall j \in \{1, \dots, s\}, \text{ where}$$

$$a_j := \lambda_j' \hat{\beta} - t_{n-r, \alpha/2} \cdot \sqrt{\hat{\sigma}^2 \lambda_j' (\mathbf{X} \mathbf{X})^{-1} \lambda_j}$$

$$b_j := \lambda_j' \hat{\beta} + t_{n-r, \alpha/2} \cdot \sqrt{\hat{\sigma}^2 \lambda_j' (\mathbf{X} \mathbf{X})^{-1} \lambda_j}$$

However,

$$P\left(\bigcap_{j=1}^s \{a_j \leq \Lambda_j' \beta \leq b_j\}\right) \leq P(a_1 \leq \Lambda_1' \beta \leq b_1) = 1-\alpha$$

thus, need to adjust family-wise Type-I-error

Theorem: Let $\{E_j\}$ be a collection of measurable events, then :

$$(i). P\left(\bigcup_j E_j\right) \leq \sum_j P(E_j)$$

$$(ii). P\left(\bigcap_j E_j\right) \geq 1 - \sum_j P(E_j^c)$$

$$P\left(\bigcap_j E_j\right) = 1 - P\left(\bigcap_j E_j^c\right)$$

$$= 1 - P\left(\bigcup_j E_j^c\right) \quad \text{by DeMorgan's rule}$$

$$\geq 1 - \sum_{j=1}^s P(E_j^c) \quad \text{by (i)}$$

① Bonferroni method: adjust level of each interval as α/s .

Bonferroni method interval :

$$P(a_j \leq \Lambda_j' \beta \leq b_j) = 1-\alpha, \quad \forall j \in \{1, \dots, s\}, \text{ where}$$

$$a_j := \lambda_j' \hat{\beta} - t_{n-r, \alpha/(2s)} \cdot \sqrt{\hat{\sigma}^2 \lambda_j' (\mathbf{X} \mathbf{X})^{-1} \lambda_j}$$

$$b_j := \lambda_j' \hat{\beta} + t_{n-r, \alpha/(2s)} \cdot \sqrt{\hat{\sigma}^2 \lambda_j' (\mathbf{X} \mathbf{X})^{-1} \lambda_j}$$

$$P\left(\bigcap_{j=1}^s \{a_j \leq \Lambda_j' \beta \leq b_j\}\right) \geq 1 - \sum_{j=1}^s P(\{a_j \leq \Lambda_j' \beta \leq b_j\}^c)$$

$$= 1 - \sum_{j=1}^s \frac{\alpha}{s}$$

$$= 1 - \alpha \quad \text{CI coverage} \geq 1-\alpha$$

② Scheffé method Construct a confidence interval $\lambda' \beta$ of the form

$$\Gamma(n, c) = \lambda' \hat{\beta} \pm c \cdot \sqrt{\hat{\sigma}^2 \lambda' (\lambda' \lambda)^{-1} \lambda},$$

where c is such that $P(\lambda' \beta \in \Gamma(n, c)) = 1-\alpha$

CI: $\lambda' \hat{\beta} \pm \sqrt{s \cdot F_{s, n-r, \alpha}} \cdot \sqrt{\hat{\sigma}^2 \lambda' (\lambda' \lambda)^{-1} \lambda}$

\downarrow rank of test \downarrow MSE df

③ Tukey method

Replace $t_{n-r, \alpha/2}$ by $q_{t_{k, n-r, \alpha}}^*$

studentized range distribution quantile
of parameters involved \rightarrow MSE df

$$P(a_j \leq \lambda' \beta \leq b_j) = 1-\alpha, \forall j \in \{1, \dots, s\}, \text{ where}$$

$$a_j := \lambda' \hat{\beta} - q_{t_{k, n-r, \alpha}}^* \cdot \sqrt{\hat{\sigma}^2 \lambda' (\lambda' \lambda)^{-1} \lambda}$$

$$b_j := \lambda' \hat{\beta} + q_{t_{k, n-r, \alpha}}^* \cdot \sqrt{\hat{\sigma}^2 \lambda' (\lambda' \lambda)^{-1} \lambda}$$

Hypothesis Testing (3)

- decision rules d_1 & d_2 .
- d_1 dominates d_2 if $R_{d_1}(\theta) \leq R_{d_2}(\theta) \forall \theta \in \Omega$, with " $<$ " for at least one θ .
- If d is dominated by some rule, then it's inadmissible, otherwise it's admissible.
- d is minimax if $\max_{\theta \in \Omega} R_d(\theta) = \min_{d'} \max_{\theta \in \Omega} R_{d'}(\theta)$

Bayes Decision Rule prior distribution of θ with pdf (pmf) $\Phi(\theta)$.

- Bayes risk is the expected loss:

$$\begin{aligned} r_d(\Phi) &= E[L(d(x), \theta)] \\ &= E_{\Phi}[E[L(d(x), \theta) | \theta]] \\ &= E_{\Phi}[R_d(\theta)] \end{aligned}$$

$$\begin{aligned} r_d(\Phi) &= E[L(d(x), \theta)] \\ &= E[E[L(d(x), \theta) | X]] \\ &= E[\text{expected posterior loss for } d(x)] \end{aligned}$$

- Bayes decision rule:

$$d: r_{d\Phi}(\Phi) = \min_d r_d(\Phi)$$

* may not exist, may not be unique

- Bayes estimator:

$$L(a, \theta) = (a - \theta)^2 \Rightarrow d_{\Phi}(x) = E[\theta | X=x]$$

$$L(a, \theta) = |a - \theta| \Rightarrow d_{\Phi}(x) = \text{median}[\theta | X=x]$$

- Bayes testing:

Expected posterior loss:

$$E[L(a, \theta) | X=x] = \begin{cases} \text{Loss(false accept)} P(\theta = \Omega_1 | X=x) & \text{if accept } H_0 \\ \text{Loss(false reject)} P(\theta = \Omega_0 | X=x) & \text{if reject } H_0 \end{cases}$$

Bayes test:

$$\Phi_{\Phi}(x) = \begin{cases} H_0 & \text{if } ① < ② \\ H_1 & \text{if } ① > ② \end{cases}$$

A unique Bayes rule w.r.t. a general loss function L is admissible w.r.t. L .

Φ admissible $\Rightarrow \Phi$ bayes rule \Rightarrow LRT

Parameter Estimation

Simple Linear Regression

Quantities:

Covariance:

$$\text{Cov}(X_i, Y_i) = E((X_i - \mu_X)(Y_i - \mu_Y))$$

$$\Rightarrow \hat{\text{Cov}}(X_i, Y_i) = \frac{1}{n-1} \sum_i (X_i - \bar{x})(Y_i - \bar{y})$$

Correlation:

$$P_{xy} = \frac{\text{Cov}(X_i, Y_i)}{\sigma_x \sigma_y} = \frac{E((X_i - \mu_X)(Y_i - \mu_Y))}{\sigma_x \sigma_y}$$

$$\uparrow = E\left(\left(\frac{X_i - \mu_X}{\sigma_x}\right)\left(\frac{Y_i - \mu_Y}{\sigma_y}\right)\right)$$

$$\Rightarrow r_{xy} = \frac{1}{n-1} \sum_i \frac{(X_i - \bar{x})}{s_x} \frac{(Y_i - \bar{y})}{s_y}$$

test for this see previous pages

SST & SSR :

$$\underbrace{\sum (Y_i - \bar{y})^2}_{\text{SST}} = \underbrace{\sum (Y_i - \hat{y})^2}_{\text{SSE}} + \underbrace{\sum (\hat{y} - \bar{y})^2}_{\text{SSR}}$$

sum of square regression : computes fitted value under

SLR to fitted value assuming $E(Y|X=x) = \beta_0$

$$\downarrow \frac{\sum (Y_i - \hat{y})^2}{(n-2)}: \text{MSE}$$

$0 \leq R^2 \leq 1$: is the % of total variation explained by the simple linear regression model.

$$R^2 = \frac{\text{SSR}}{\text{SST}} \quad \textcircled{=} (r_{xy})^2$$

in SLR

Residuals & assumption check : $R_i = Y_i - \hat{y}_i$

Assumption : $R_i \stackrel{iid}{\sim} N(0, \sigma^2)$

① independent assumption is impossible to check.

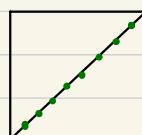
④ Lack of fit: when see a trend in residual

② Residual constant variance.

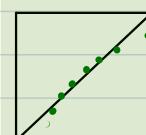
* See previous table for this

③ Normality assumption

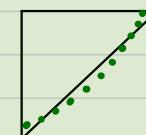
observed
quantile



normal



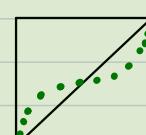
left skewed



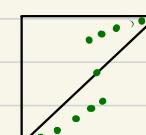
right skewed



light tailed



heavy tailed



bimodal

Theoretical quantile

Estimators

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$> \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

$$SE(\hat{\beta}_0) = \sqrt{MSE\left(\frac{1}{n} + \frac{\bar{X}^2}{\sum(X_i - \bar{X})^2}\right)}$$

$$> \hat{\beta}_1 = \frac{r_{XY}}{S_x} S_y = \frac{S_{XY}}{S_x S_y} \frac{S_y}{S_x} = \frac{S_{XY}}{S_{xx}} = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2}$$

$$SE(\hat{\beta}_1) = \sqrt{MSE\left(\frac{1}{\sum(X_i - \bar{X})^2}\right)}$$

think these as

$$\sigma^2 (\hat{X} \hat{X})^{-1} = \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}^2}{\sum(X_i - \bar{X})^2} \right)^{-1}$$

$$> \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_0$$

$$SE(\hat{Y}) = \sqrt{MSE\left(1 + \frac{(X_0 - \bar{X})^2}{\sum(X_i - \bar{X})^2}\right)} \quad \text{I think here we do NOT assume } \beta_0 \perp \beta_1$$

$$> \tilde{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_0 + \varepsilon_{n+1}$$

$$SE(\tilde{Y}) = \sqrt{MSE\left(1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum(X_i - \bar{X})^2}\right)} \quad \text{see PI in previous pages}$$

Linear Regression v.s. Penalization

(Ordinary) Least Square Regression:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_p X_{ip})^2$$

> $\hat{\beta}$ may or may not be uniquely determined

> Convention: center & scale all $X_{ji} \rightarrow (X_{ji} - \bar{X}_{j.})/s_j$

then $\hat{\beta}_0 = \bar{Y} = \sum_{i=1}^n Y_i/n$ really like centering Y as well

$$> \hat{\beta}_{OLS} = (\hat{X} \hat{X} + \lambda I)^{-1} \hat{X}' Y$$

$\hookrightarrow Y_i \mapsto Y_i - \bar{Y}$ with no intercept column in X

Ridge Regression

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_{i1} - \cdots - \beta_p X_{ip})^2 + \lambda \sum_{j=1}^n \beta_j^2$$

> Balance: minimize SSE v.s. make length of slope vector close to 0. β_1, \dots, β_p

> Highly correlated predictor variables:

- OLS slope is not well determined, have large variance

- shrinking β_1, \dots, β_p to close to 0 helps determination

> Convention: center & scale all $X_{ji} \rightarrow (X_{ji} - \bar{X}_{j.})/s_j$

$Y_i \mapsto Y_i - \bar{Y}$ with no intercept column in X .

$$> \hat{\beta}_{ridge} = (\hat{X} \hat{X} + \lambda I)^{-1} \hat{X}' Y \quad \text{shrinkage estimator}$$

> Bias = $E(\hat{\beta}_{ridge} - \beta) = -\lambda (\hat{X} \hat{X} + \lambda I)^{-1} \beta \sim$ use Sherman-Morrison-Woodbury formula

$$> \text{var}(\hat{\beta}_{ridge}) = \sigma^2 (\hat{X} \hat{X} + \lambda I)^{-1} (\hat{X} \hat{X}) (\hat{X} \hat{X} + \lambda I)^{-1} \quad \text{biased regression}$$

Lasso Regression

$$\hat{\beta} = \underset{\beta}{\operatorname{arg\min}} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \cdots - \beta_p x_{ip})^2 + \lambda \sum_{j=1}^p |\beta_j|$$

> Convention: center & scale all $x_{ji} \rightarrow (x_{ji} - \bar{x}_{j\cdot})/s_j$

$y_i \mapsto y_i - \bar{y}$ with no intercept column in X .

> $\hat{\beta}^{\text{Lasso}}$ has no closed-form as is non-linear.

> Closed form expressions exist in a special case:

① centered & scaled predictors, when X has orthonormal columns : $X^T X = I$

$$\text{LS}(p): \hat{\beta} = (X^T X)^{-1} X^T y \Rightarrow \hat{\beta}_j = x_j^T y$$

LS(k): keep k predictors have the largest impact ← Best Subset

• $(\text{var } \hat{\beta}_j)$ equals

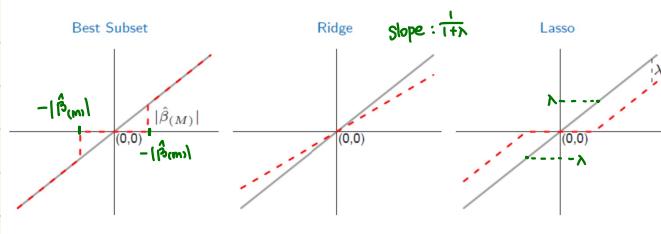
• order $|\hat{\beta}_{(1)}| \geq \cdots \geq |\hat{\beta}_{(k)}| \geq \cdots \geq |\hat{\beta}_{(p)}|$

• if $|\hat{\beta}_{(j)}| \geq |\hat{\beta}_{(k)}|$ report $\hat{\beta}_{(j)}$; otherwise, report 0.

• $\hat{\beta}_j^{\text{LS}(k)} = \hat{\beta}_j \cdot I(|\hat{\beta}_j| \geq |\hat{\beta}_k|) \leftarrow \text{hard threshold}$

$$\text{Ridge}(n): \hat{\beta}_j^{\text{Ridge}} = \frac{x_j^T y}{1+\lambda} = \frac{\hat{\beta}_j}{1+\lambda} \leftarrow \text{proportional shrinkage}$$

$$\text{Lasso}(n): \hat{\beta}_j^{\text{Lasso}} = \text{sign}(\hat{\beta}_j) (\hat{\beta}_j | - \lambda)_{+} = \text{sign}(\hat{\beta}_j) \max(|\hat{\beta}_j| - \lambda, 0) \leftarrow \text{soft thresholding}$$



Properties

Setting: Multicollinearity or p relatively to n

Must center and scale predictors

• shrinkage applied to partial slopes, not intercept

• scaling matters for predictions \hat{y} (?)

As $\lambda \uparrow$, $\hat{\beta} \rightarrow 0$: bias \uparrow variance \downarrow

Ridge v.s. Lasso

Ridge: 1. $\hat{\beta} \rightarrow 0$ not exactly 0

(use when most predictors are important)

2. Shrinks to similar coefficient estimates when correlated

Lasso: 1. $\hat{\beta} \rightarrow 0$ equals to 0

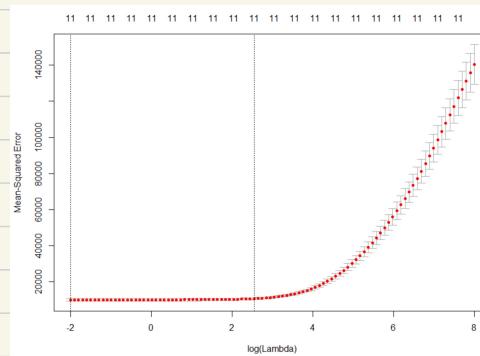
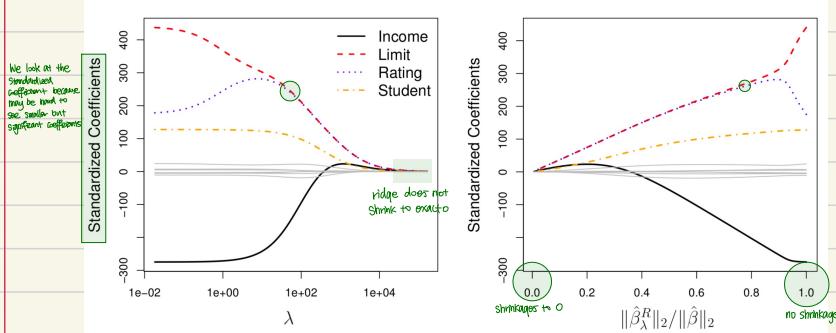
(implicit variable selection)

can be misleading and arbitrary + various adaptations exist

2. Picks one and discard others when correlated

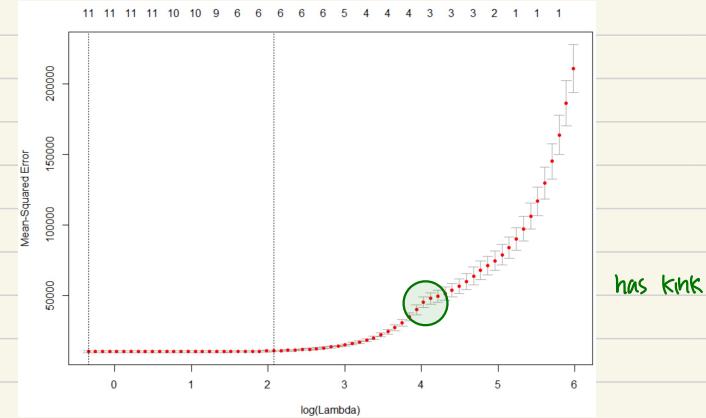
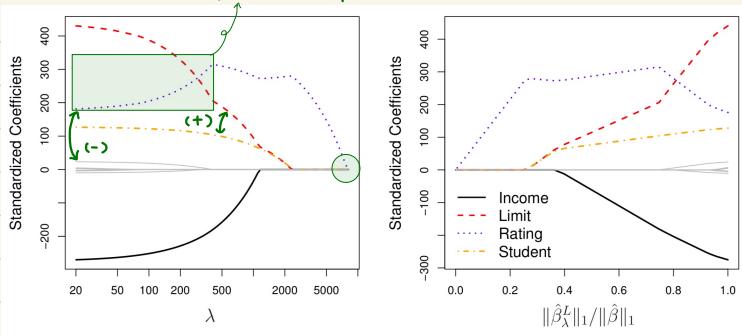
* Both useful for prediction

* Both require care when focus is estimation of β .



Smooth

Reallocation of correlated features



Elastic Net Regression

$$\hat{\beta} = \underset{\beta}{\operatorname{arg\,min}} \text{SSE} + \lambda \left(\frac{(1-\alpha)}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

normally λ is selected α is specified

Dimension Reduction

Convention: center and scale predictors

1. X -space $\rightarrow W$ space (can always be done)
2. Columns of W (scores) = linear combinations of columns of X
3. Columns of W is ordered by relevance
4. Drop irrelevant columns to get W_{reg} -space
5. Perform back conversion $W_{\text{reg}} \rightarrow X$ to do inference.

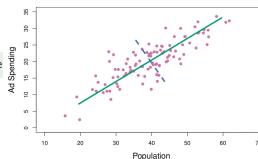
Principal Component Regression (PCR)

- chooses "scores" without info Y

PCR: " W -space" from eigenvalue-eigenvector decomposition of $X'X$

- Eigenvalues $\lambda_1 \geq \dots \geq \lambda_p$. Corresponding orthogonal eigenvectors v_1, \dots, v_p as columns of matrix V
- Xv_1 explains the most variation in the X -space, namely $\left(\frac{\lambda_1}{\lambda_1 + \dots + \lambda_p}\right)$ 100% of the variation in the X -space
- Xv_1 and Xv_2 together explain $\left(\frac{\lambda_1 + \lambda_2}{\lambda_1 + \dots + \lambda_p}\right)$ 100% of the variation in the X -space
- $W = XV$ has p columns, also known as p "principal components"
- The k th principal component is "irrelevant" if its
 - eigenvalue is "small" compared to the largest eigenvalue, i.e., if k th condition index $\delta_k = \sqrt{\frac{\lambda_1}{\lambda_k}}$ is large (say > 10)
 - estimated regression coefficient is statistically zero

Understanding principal components when $p = 2$:



$$pc_1 = .839 * \text{popn} + .544 * \text{ad}$$
$$pc_2 = .544 * \text{popn} - .839 * \text{ad}$$

Fig 6.14 from ISLR

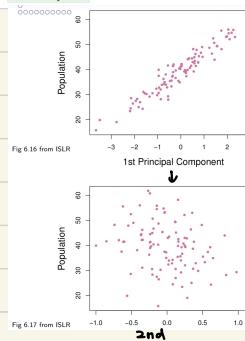


Fig 6.16 from ISLR

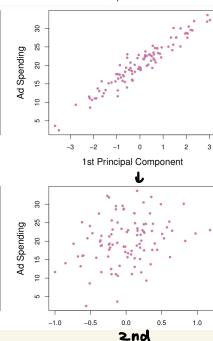


Fig 6.17 from ISLR

Partial Least Squares Regression (PLSR)

- choose "scores" with info Y

PLSR: " W -space" simultaneously seeks high levels of variation in the X -space and strong correlation with Y

- Does not naturally arise from elegant matrix algebra
- More algorithmic, less theoretical
- Can be quite effective

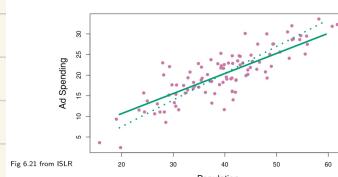


Fig 6.21 from ISLR

Multiple Linear Regression ← weird I don't see much interpretations, but typically

Fitted value: $\hat{y}_i = \hat{\beta}_0 + \sum_j \hat{\beta}_j x_{ij} = x_i^T \hat{\beta}$ are normal χ^2 coefficient means curvature.

$$E(\hat{y}_i | x_i) = x_i^T \hat{\beta}$$

$$\text{Var}(\hat{y}_i | x_i) = \sigma^2 x_i^T (x_i x_i^T)^{-1} x_i < \sigma^2$$

↳ $h_{ii} = x_i^T (x_i x_i^T)^{-1} x_i$ leverage (≤ 1 for models with intercept)

Fitted vector: $\hat{y} = X \hat{\beta} = (X^T X)^{-1} X^T y = Hy$

$$E(\hat{y} | X) = X \hat{\beta}$$

$$\text{Var}(\hat{y} | X) = \sigma^2 H$$

Residual vector: $r = y - \hat{y} = y - Hy = (I - H)y$

Fact 1: $\sum_i r_i = 0$ ✓ models with intercept

Fact 2: Residuals are always correlated

$$\cdot \text{var}(r) = \sigma^2 (I - H)$$

• complicates abilities to check independence assumption

$$SSE = (y - \hat{y})'(y - \hat{y}) = y'(I - H)y$$

Example: Let $x_0^T = (1 \ x_1 \ x_2 \ \dots \ x_p)$ be values of prediction

Prediction value: $\hat{y}_0 = \hat{\beta}_0 + \sum_j \hat{\beta}_j x_{0j} = x_0^T \hat{\beta}$ is normal

$$E(\hat{y}_0 | x_0) = x_0^T \hat{\beta}$$

$$\text{Var}(\hat{y}_0 | x_0) = \sigma^2 x_0^T (x_0 x_0^T)^{-1} x_0$$

$$\underline{\text{CI } x_0^T \hat{\beta}}: \hat{y} \pm t_{\alpha/2, n-p-1} \sqrt{MSE} \sqrt{x_0^T (x_0 x_0^T)^{-1} x_0}$$

$$\underline{\text{PI } x_0^T \hat{\beta}}: \hat{y} \pm t_{\alpha/2, n-p-1} \sqrt{MSE} \sqrt{1 + x_0^T (x_0 x_0^T)^{-1} x_0}$$

Coefficient of Determination: $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$

* However, R^2 always increases as adding predictors.

$$\begin{aligned} \text{adjusted } R^2 &= 1 - \frac{SSE/(n-p-1)}{SST/(n-1)} \\ &= 1 - (1-R^2) \frac{n-1}{n-p-1} \\ &< R^2 \text{ always} \end{aligned}$$

* Adding predictor that does not increase R^2 causing adj R^2 to decrease.

* However, adj R^2 is NOT great for small dataset.

Model Specification

Over-specification

True Model: $E(Y_i | X = x_i) = \sum_{j=0}^{p_i} \theta_j x_{ij} + \sum_{j=p_i+1}^{p-p_i-1} 0 x_{ij}$ OR $E(Y_i | X_i) = x_i \theta_i$

Oracle estimator: $\hat{\theta}_i = (x_i' x_i)^{-1} x_i' y$

$$E(\hat{\theta}_i) = \theta_i$$

$$\text{var}(\hat{\theta}_i) = \sigma^2 (x_i' x_i)^{-1}$$

$$E(\text{MSE}) = \sigma^2 \Rightarrow n - (1+p_i) \text{ DF}$$

Assumed Model:

$$E(Y_i | X = x_i) = \sum_{j=0}^{p_i} \theta_j x_{ij} + \sum_{j=p_i+1}^{p-p_i-1} \theta_j x_{ij} \quad \text{OR} \quad E(Y_i | X_1, X_2) = X \theta = X \theta$$

Estimator: $\hat{\theta} = (X' X)^{-1} X' y$

$$E(\hat{\theta}) = (X' X)^{-1} X' \theta = (\theta)$$

$$\text{var}(\hat{\theta}) = \sigma^2 (X' X)^{-1} > \sigma^2 (x_i' x_i)^{-1}$$

$$E(\text{MSE}) = \sigma^2 \Rightarrow n - (1+p) \text{ DF} \Leftarrow \text{df Error} \downarrow \text{F statistics} \uparrow \text{Power} \downarrow$$

Summary: Both Type I Error & Type II Error are correct.

Under over-specification, correct variables have $\hat{S.E.} \uparrow$ p-value \uparrow $\Leftarrow \text{CI} \uparrow$

Over-specification causes problems as well:

Inflating Estimator Variance:

For assumed model (with $p+1$ parameters):

$$\text{var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_i (x_{ij} - \bar{x}_{ij})^2} \times \frac{1}{1 - R_j^2} \rightarrow R\text{-squared value treating } x_j \text{ as response}$$

↓ variance for SLR with x_j

and other predictors as covariates.

$$R_j^2 \rightarrow 1, \text{ var}(\hat{\beta}_j) \rightarrow \infty$$

Variance Inflation factor (VIF): $\frac{1}{1 - R_j^2}$

$\text{VIF} > 10$ means predictors are highly correlated, complicating ability to identify which variables are important.

(1) model is overly complicated

(2) The predictors are naturally correlated

1. Global F-test p-value \downarrow but individual p-value \uparrow

2. Normally when collinearity happens, VIF \uparrow p-value \uparrow

Under-specification

True Model: $E(Y_i | X = X_i) = \sum_{j=0}^{P_i} \theta_j X_{ij} + \sum_{j=P_i+1}^{P-1} 0 X_{ij}$ OR $E(Y_i | X_1, X_2) = X_1 \theta_1 + X_2 \theta_2 = X \theta$

Oracle estimator: $\hat{\theta} = (X'X)^{-1} X'y$
 $E(\hat{\theta}) = \theta$
 $\text{Var}(\hat{\theta}) = \sigma^2 (X'X)^{-1}$
 $E(\text{MSE}) = \sigma^2 \Rightarrow n - (1+P) \text{ DF}$

Assumed Model:

$E(Y_i | X = X_i) = \sum_{j=0}^{P_i} \theta_j X_{ij} + \sum_{j=P_i+1}^{P-1} 0 X_{ij}$ OR $E(Y_i | X_1) = X_1 \theta_1$

Estimator: $\hat{\theta}_1 = (X_1'X_1)^{-1} X_1'y$
 $E(\hat{\theta}_1) = (X_1'X_1)^{-1} (X_1\theta_1 + X_2\theta_2) = \theta_1 + (X_1'X_1)^{-1} X_2\theta_2$
 $\text{Var}(\hat{\theta}_1) = \sigma^2 (X_1'X_1)^{-1}$
 $E(\text{MSE}) = \sigma^2 + Q(X_2, \theta_2) > \sigma^2 \Rightarrow n - (1+P) \text{ DF} < \text{DF}$

unbiased when
 X_1, X_2 are orthonormal

Summary: Hypothesis Test Type I & Type II errors unpredictable. $\leftarrow \text{CI } x$
 Replication and lack-of-fit can detect underspecified model. pure error may be used to instead of σ^2 . since MSpure always unbiased σ^2 .

Linear Mixed Model

1. In linear models, we consider r.v. ε_i account for the variability in the response not explained by the predictors.

$$2. \hat{\varepsilon}_i = y_i - \hat{y}_i$$

$$\hat{\varepsilon} = (I - X(X^T X)^{-1} X^T) y \sim N(0, (I - X(X^T X)^{-1} X^T) \sigma^2)$$

3. Model formula in R: $y \sim a+b + x:z + I(v^2) - 1$

$$a+b : a+b + a:b$$

$x:z$: interaction of x and z

$I(v^2)$: linear predictor depends on v^2 .

-1: no intercept



$$Model: y = a+b + ab + xz + v^2$$

4. Kronecker product notation:

$$\begin{pmatrix} 1_3 \\ & 1_3 \\ & & 1_3 \end{pmatrix} = I_3 \otimes I_3$$

$$\begin{pmatrix} 1_3, 0 \\ 0, 1_3 \\ 1_3, 0 \\ 0, 1_3 \end{pmatrix} = I_2 \otimes (I_2 \otimes I_3)$$

5. Linear Mixed Model:

When estimating σ_A^2 , Type 3: used to see the tests, but $\hat{\sigma}_A^2 < 0$ with MSR < MSE

ML: Provides $\hat{\sigma}_A^2 \geq 0$, biased

$$\hat{\sigma}_A^2(\text{ML}) = \frac{1}{n-p} (y - \hat{X}\hat{\beta})^T (y - \hat{X}\hat{\beta}) \text{ negatively biased}$$

REML: Provides $\hat{\sigma}_A^2 \geq 0$, $E(\hat{\sigma}_A^2) = \sigma_A^2$

$$\hat{\sigma}_A^2(\text{REML}) = \frac{1}{n-p} (y - \hat{X}\hat{\beta})^T (y - \hat{X}\hat{\beta}) \text{ unbiased}$$

In SAS: proc mixed

In R: lme() for nested

lmer() for nested & crossed

(1|1a): random effect for a

(1|1a:b): random effect for ab interaction (slope)

Prediction

} see following pages

Two varieties ...

- Marginal, aka population-averaged:

$$\hat{Y} = \hat{X}\hat{\beta}$$

- Conditional, aka cluster-specific:

$$\hat{Y} = \hat{X}\hat{\beta} + Z\hat{\alpha}$$

where $\hat{\alpha}$ is best linear unbiased predictor
i.e., $\hat{\alpha}$ is conditional mean of α given observed y

$$\hat{\alpha} = \hat{G}Z^T \hat{V}^{-1} (y - \hat{X}\hat{\beta})$$

BLUP property relies on known G and R

$\alpha \sim N(0, G)$ notation in this case

$$\varepsilon \sim N(0, R)$$

$$\text{cov}(\alpha, \varepsilon) = 0$$

lmor(.) usage and output is crucial

Model specification

The following formula extensions for specifying random-effects structures in R are used by

- lme4
- nlme (nested effects only, although crossed effects can be specified with more work)
- glmmADMB and glmmTMB

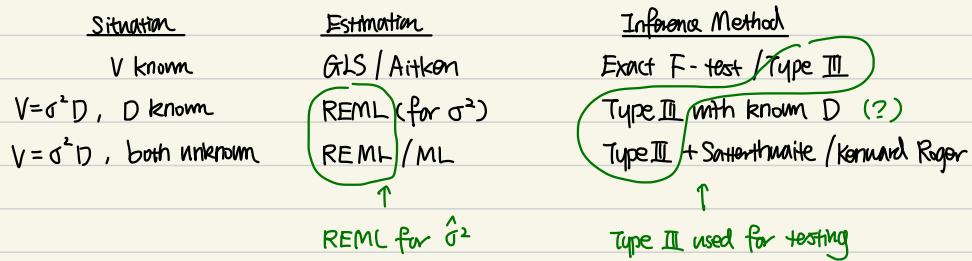
MCMCglmm uses a different specification, inherited from AS-REML.

(Modified from Robin Jeffries, UCLA:)

	formula	meaning
Covered in lecture note	$\{(1 group)$	random group intercept + "random group"
	$\{(x group) = (1+x group)$	random slope of x within group with correlated intercept + "random intercept / subject = group"
	$\{(0+x group) = (-1+x group)$	random slope of x within group: no variation in intercept
	$\{(1 group) + (0+x group)$	uncorrelated random intercept and random slope within group
	$\{(1 site/block) = (1 site)+(1 site:block)$	intercept varying among sites and among blocks within sites (nested random effects)
	$site+(1 site:block)$	fixed effect of sites plus random variation in intercept among blocks within sites
	$(x site/block) = (x site)+(x site:block) = (1 + x site)+(1+x site:block)$	slope and intercept varying among sites and among blocks within sites
	$(x1 site)+(x2 block)$	two different effects, varying at different levels
	$x*site+(x site:block)$	fixed effect variation of slope and intercept varying among sites and random variation of slope and intercept among blocks within sites
	$(1 group1)+(1 group2)$	intercept varying among crossed random effects (e.g. site, year)

Consider under Aitken model, want to test on β but notice $\Sigma \sim N(0, V)$

1. If V is known, $H_0: A\beta = m$ inference is exact.
2. If $V = \sigma^2 D$ with D known, use REML method
3. If $V = \sigma^2 D$ with both σ^2, D unknown, use Type III + Satterthwaite / Kenward Roger find df



Example: Fabric

Two types of yarn are used to weave fabric on looms. In our shop, we have three operators and we randomly select three of our many looms. We will have each operator run each yarn type on each loom twice. The order of these runs will be at random.

The operators are the only three we have who are qualified for this type of weaving. We are interested in comparing them. The yarn types are the only two available for this particular application and are different in the type of fiber used. The looms, on the other hand, are of no particular interest and are thought to be representative of the collection of looms that might be used for this type of application.

The data are average puncture resistance measurements made on five randomly selected spots on the cloth. Here are the data:

Loom	1	1	1	2	2	2	3	3	3
Oper	1	2	3	1	2	3	1	2	3
Yarn 1	59, 52	43, 42	48, 54	88, 86	68, 71	83, 82	50, 43	35, 30	42, 41
Yarn 2	81, 73	54, 56	59, 68	99, 98	82, 88	98, 94	63, 62	49, 40	60, 56

```
proc mixed data=fabric; class loom operator yarn;  
  model resistance=operator|yarn / solution;  
  random loom loom*operator loom*yarn loom*operator*yarn;
```

} same model,
different method

```
proc mixed data=fabric method=ml; class loom operator yarn;  
  model resistance=operator|yarn / solution;  
  random loom loom*operator loom*yarn loom*operator*yarn;
```

```
proc mixed data=fabric method=type3; class loom operator yarn;  
  model resistance=operator|yarn / solution;  
  random loom loom*operator loom*yarn loom*operator*yarn;
```

```
proc mixed data=fabric method=type3; class loom operator yarn;  
  model resistance=operator|yarn / solution ddfm=satterth;  
  random loom loom*operator loom*yarn loom*operator*yarn;
```

```
proc mixed data=fabric method=type3; class loom operator yarn;  
  model resistance=operator|yarn / solution ddfm=kenwardroger;  
  random loom loom*operator loom*yarn loom*operator*yarn;
```

REML: $\hat{\sigma}_L^2 = 406.43$, $\hat{\sigma}_{LO}^2 = 0$, $\hat{\sigma}_{LY}^2 = 0$, $\hat{\sigma}_{LOY}^2 = 0$, $\hat{\sigma}^2 = 12.06$ ML: $\hat{\sigma}_L^2 = 270.84$, $\hat{\sigma}_{LO}^2 = 0$, $\hat{\sigma}_{LY}^2 = 0$, $\hat{\sigma}_{LOY}^2 = 0$, $\hat{\sigma}^2 = 10.23$ type3: $\hat{\sigma}_L^2 = 407.90$, $\hat{\sigma}_{LO}^2 = -1.60$, $\hat{\sigma}_{LY}^2 = -1.97$, $\hat{\sigma}_{LOY}^2 = 0.875$, $\hat{\sigma}^2 = 13.19$

type3 w/ Satterthwaite or Kenward-Roger:

 $\hat{\sigma}_L^2 = 407.90$, $\hat{\sigma}_{LO}^2 = -1.60$, $\hat{\sigma}_{LY}^2 = -1.97$, $\hat{\sigma}_{LOY}^2 = 0.875$, $\hat{\sigma}^2 = 13.19$

Effect	ndf	ddf	F value	Pr>F	Param	Estim	SE	df	Pr> t	
Oper	2	4	68.31	.0008	$O_1 - O_3$	6.83	2.00	4	.0271	
			80.51	.0006			1.85	4	.0208	
			96.59	.0004			1.98	4	.0259	
			96.59	.0004			1.98	7.44	.0096	
Yarn	1	2	159.34	.0062	$O_2 - O_3$	-11.0	2.00	4	.0054	
			187.80	.0053			1.85	4	.0040	
			617.58	.0016			1.98	4	.0051	
			617.58	.0016			1.98	7.44	.0007	
O*Y	2	4	0.60	.5932	effects only has the same F value as type III cas β the same with df the same)					
			0.70	.5474						
			0.48	.6496						
			0.48	.6496						

Computing ... SAS: proc glimmix

Example [Corn]: Consider an experiment to study the effects of 4 planting methods and 3 pesticides on yields of corn. Each pesticide is applied by aerial spraying to 2 large fields. Each large field is subdivided into 4 subplots, and the 4 planting methods are randomly assigned to the 4 subplots. $3 \times 2 \times 6$ fields total

- Completely randomized split-plot design
- Y_{ijk} = yield of j th planting method within k th field receiving i th pesticide

$$Y_{ijk} = \mu + \alpha_i + \gamma_{k(i)} + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}, \quad \gamma_{k(i)} \sim^{iid} N(0, \sigma_{\gamma}^2), \\ \epsilon_{ijk} \sim^{iid} N(0, \sigma^2), \quad \text{Cov}(\gamma_{k(i)}, \epsilon_{ijk}) = 0,$$

$$\Rightarrow \text{Var}(Y_{ijk}) = \sigma_{\gamma}^2 + \sigma^2, \quad \text{Cov}(Y_{ijk}, Y_{ljk}) = \sigma_{\gamma}^2$$

```
proc glimmix data=corn;
  class pesticide field method;
  model yield=pesticide|method / solution;
  random field(pesticide);
```

```
proc glimmix data=fabric;
  class loom operator yarn;
  model resistance=operator|yarn / solution;
  random loom loom*operator loom*yarn loom*operator*yarn;
```

** Results for FABRIC differ (a little) between mixed and glimmix **

Computing ... R

- Function `lme()` in package `nlme` is very popular, but is limited to nested random effects.
- Function `lmer()` in package `lme4` can handle both nested and crossed random effects. It is also newer and better able to deal with large datasets.

But more structure and flexibility in specifying matrices G and R is offered by `lme()` compared to `lmer()`.

This is nested, I think this is because : there is no outside random factor: field.

The interaction cannot stand by itself

```
fit.corn = lmer(yield ~ pesticide * method + (1|field:pesticide), data=corn)
```

```
fit.fabric = lmer(resistance ~ operator * yarn + (1|loom) +  
(1|loom:operator) + (1|loom:yarn)+ (1|loom:operator:yarn), data=fabric)
```

Sleep Study

Interested in **average reaction time** for sleep-deprived subjects.

- Before day 0, subjects had their normal amount of sleep.
- On following nights, subjects were limited to 3 hours of sleep.
- Subjects were given a series of tests on each day.
- This dataset is balanced with no missing observations.

Reaction is average reaction time (ms) on a series of tests given each Days (values 0,1,...,9) to each Subject (there are 18 of them).

Y_{ij} : reaction time for subject i on day j

X_{ij} : day j for subject i on day j

$$E(Y_{ij}) = \beta_0 + \beta_1 X_{ij} \quad \text{expect } \beta_1 > 0$$

But different people may respond differently ...

$$E(Y_{ij}) = \beta_{0i} + \beta_{1i} X_{ij}$$

But different people may respond differently ...

$$E(Y_{ij}) = \beta_{0i} + \beta_{1i}X_{ij}$$

- β_{0i}, β_{1i} as fixed effects limits us to inference on these 18 subjects
- β_{0i}, β_{1i} as random effects allows inference on a broader scale

$$\begin{aligned}\beta_{0i} &\equiv \beta_0 + \alpha_{0i}, \\ \beta_{1i} &\equiv \beta_1 + \alpha_{1i},\end{aligned}$$

$Y_{ij} = \beta_0 + \beta_1 X_{ij} + \alpha_{0i} + \alpha_{1i} X_{ij} + \epsilon_{ij}$

$\alpha_{0i} \sim^{ind} N(0, \sigma_0^2)$
 $\alpha_{1i} \sim^{ind} N(0, \sigma_1^2)$

Called the "random intercepts, random slopes" model.

↓ for different person, effect different

Independently over i , for $i = 1, \dots, N = 18$,

$$\mathbf{Y}_i = \mathbf{X}_i \boldsymbol{\beta} + \mathbf{Z}_i \boldsymbol{\alpha}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\alpha}_i \sim N(\mathbf{0}, \mathbf{D}), \quad \boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \mathbf{R}_i), \quad \boldsymbol{\alpha}_i, \boldsymbol{\epsilon}_i \text{ indep}$$

Clustered Data / Repeated Measures / Longitudinal Data

“Single-level LMM”

Independently over i , for $i = 1, \dots, N$,

$$\mathbf{Y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\boldsymbol{\alpha}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\alpha}_i \sim N(\mathbf{0}, \mathbf{D}), \quad \boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \mathbf{R}_i), \quad \boldsymbol{\alpha}_i, \boldsymbol{\epsilon}_i \text{ indep}$$

- \mathbf{Y}_i is $n_i \times 1$ observed response vector for i th cluster
- \mathbf{X}_i is $n_i \times p$ known covariate matrix for i th cluster
- $\boldsymbol{\beta}$ is $p \times 1$ unknown fixed effects
- \mathbf{Z}_i is $n_i \times q$ known matrix for i th cluster
- $\boldsymbol{\alpha}_i$ is $q \times 1$ random, unobserved, for i th cluster

may not involve all the features
in the random effect.

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_N \end{bmatrix}, \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}, \mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} \\ \vdots & \ddots \\ \mathbf{0} & \mathbf{Z}_N \end{bmatrix}, \boldsymbol{\alpha} = \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_N \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \vdots \\ \boldsymbol{\epsilon}_N \end{bmatrix}, \mathbf{R} = \begin{bmatrix} \mathbf{R}_1 & \mathbf{0} \\ \vdots & \ddots \\ \mathbf{0} & \mathbf{R}_N \end{bmatrix}$$

$$\mathbf{G} = \mathbf{I}_N \otimes \mathbf{D}, \quad n = \sum_{i=1}^N n_i, \quad \mathbf{V} = \text{diag}\{\mathbf{V}_1, \dots, \mathbf{V}_N\}, \quad \mathbf{V}_i = \mathbf{Z}_i \mathbf{D} \mathbf{Z}_i^T + \mathbf{R}_i$$

Single-level LMM: N groups (indexed $i = 1, \dots, N$) each containing n_i observations

$$\mathbf{Y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\boldsymbol{\alpha}_i + \boldsymbol{\epsilon}_i, \quad \boldsymbol{\alpha}_i \sim N(\mathbf{0}, \mathbf{D}), \quad \boldsymbol{\epsilon}_i \sim N(\mathbf{0}, \mathbf{R}_i), \quad \boldsymbol{\alpha}_i, \boldsymbol{\epsilon}_i \text{ indep}$$

Two-level LMM: N first-level groups (indexed $i = 1, \dots, N$) each with n_i second-level groups (indexed $j = 1, \dots, n_i$) containing n_{ij} observations

$$\mathbf{Y}_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + \mathbf{Z}_{1,ij}\boldsymbol{\alpha}_i + \mathbf{Z}_{2,ij}\boldsymbol{\alpha}_{ij} + \boldsymbol{\epsilon}_{ij}, \quad \boldsymbol{\alpha}_i \sim N_{q_1}(\mathbf{0}, \mathbf{D}_1), \quad \boldsymbol{\alpha}_{ij} \sim N_{q_2}(\mathbf{0}, \mathbf{D}_2), \\ \boldsymbol{\epsilon}_{ij} \sim N_{n_{ij}}(\mathbf{0}, \mathbf{R}_{ij}), \quad \boldsymbol{\alpha}_i, \boldsymbol{\alpha}_{ij}, \boldsymbol{\epsilon}_{ij} \text{ indep}$$

- Single-level and two-level LMMs are great for **nested** groupings.
Not for **crossed** random effects.
- Not happy with assuming vectors of random effects are independent?
Use “extended LMM” or “hierarchical specification” :

Replace $\boldsymbol{\alpha}_i, \boldsymbol{\epsilon}_i$ indep with $\boldsymbol{\epsilon}_i | \boldsymbol{\alpha}_i \sim N(\mathbf{0}, \mathbf{R}_i)$

Sleep Study: Software – R

```
> str(sleepstudy)
'data.frame': 180 obs. of 3 variables:
$ Reaction: num 250 259 251 321 357 ...
$ Days     : num 0 1 2 3 4 5 6 7 8 9 ...
$ Subject  : Factor w/ 18 levels "308","309","310",...: 1 1 1 1 1 1 1 1 1 1 ...
> plot(sleepstudy)
> xyplot(Reaction~Days|Subject,sleepstudy)
```

$$Y_i = \beta_0 + \beta_1 Day_i + \alpha_{0i} + \alpha_{1i} Day_i + \varepsilon_i ; \alpha_{0i} \perp \varepsilon_i ; \alpha_{1i} \perp \varepsilon_i$$

fit.sleep = lmer(Reaction ~ Days + (Days|Subject) , sleepstudy)
#includes random intercept; allows correlated REs

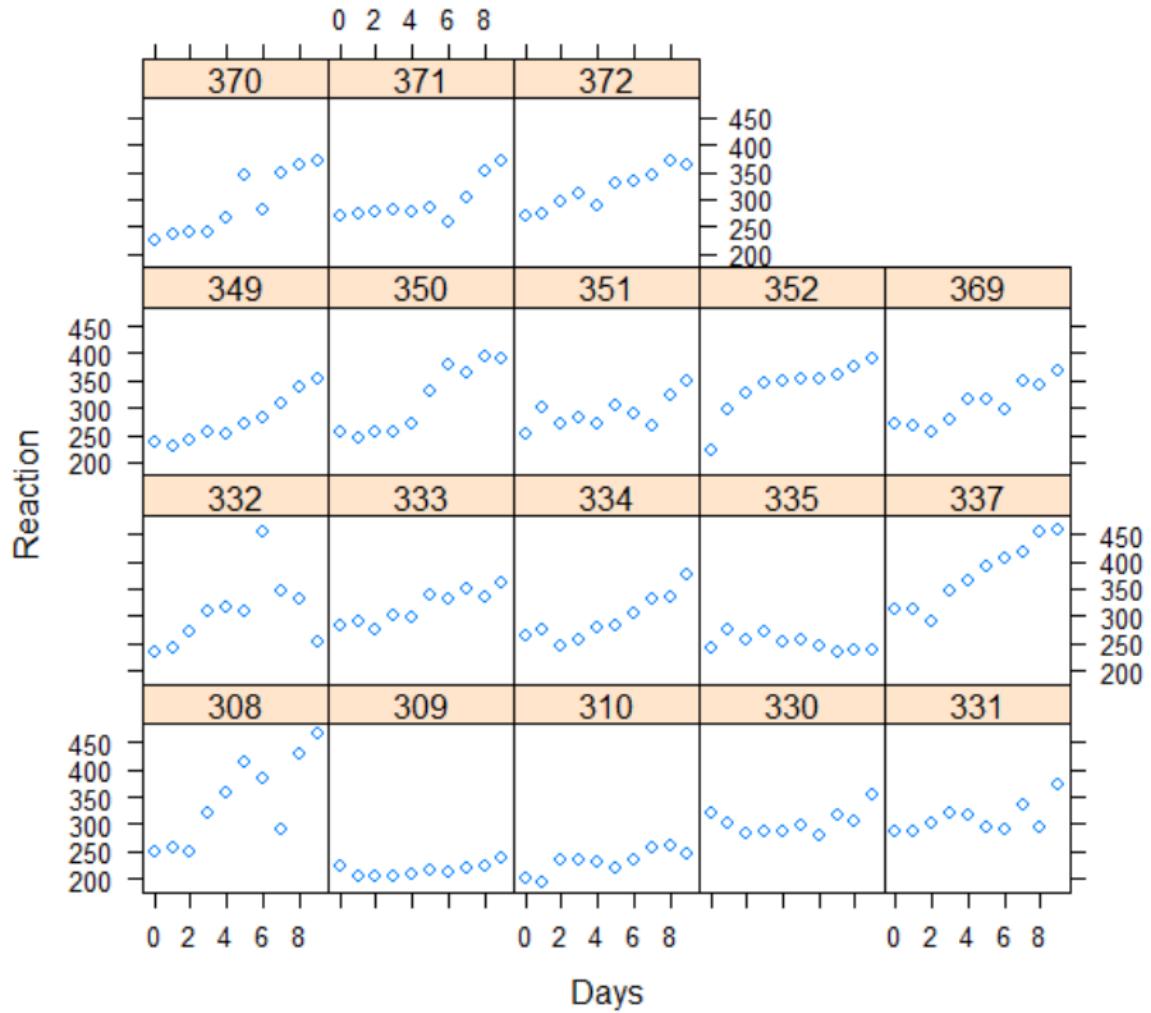
summary(fit.sleep) $Y_i = \beta_0 + \beta_1 Day_i + \alpha_{0i} + \alpha_{1i} Day_i + \varepsilon_i \quad \alpha_{0i} \perp \alpha_{1i} \perp \varepsilon_i$
fit.sleep = lmer(Reaction ~ Days + (Days||Subject) , sleepstudy)
#includes random intercept; REs uncorrelated

summary(fit.sleep) $Y_i = \beta_0 + \beta_1 Day_i + \alpha_{0i} + \alpha_{1i} Day_i + \varepsilon_i \quad \alpha_{0i} \perp \alpha_{1i} \perp \varepsilon_i$
fit.sleep = lmer(Reaction ~ Days + (1|Subject) + (0 + Days|Subject), sleepstudy)
#REs uncorrelated

summary(fit.sleep) $Y_i = \beta_0 + \beta_1 Day_i + \alpha_{0i} + \alpha_{1i} Day_i + \varepsilon_i \quad \alpha_{0i} \perp \alpha_{1i} \perp \varepsilon_i$
fit.sleep = lmer(Reaction ~ Days + (1|Subject) + (-1 + Days|Subject), sleepstudy)
#REs uncorrelated

summary(fit.sleep)

these three were equivalent.



```
> fit.sleep = lmer(Reaction ~ Days + (Days|Subject) , sleepstudy)
> summary(fit.sleep)

Linear mixed model fit by REML ['lmerMod']
Formula: Reaction ~ Days + (Days | Subject)
Data: sleepstudy
```

REML criterion at convergence: 1743.6

Comparison between models fitted by REML, lower is better.

Scaled residuals:

	Min	1Q	Median	3Q	Max
	-3.9536	-0.4634	0.0231	0.4634	5.1793

Standardized residuals: showing the distribution of model errors (should ideally be symmetric & centered at 0)

Random effects:

Groups	Name	Variance	Std.Dev.	Corr
Subject	(Intercept)	612.09	24.740	
	Days	35.07	5.922	0.07
	Residual	654.94	25.592	

Number of obs: 180, groups: Subject, 18

allow both intercept and slope varies by group: $+ A_{0i} + A_{1i}X_{ij}$

subject differs in their baseline Realiim time.

subject differ in how much Days affect their Realiim

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	251.405	6.825	36.838
Days	10.467	1.546	6.771

between intercept and slope

expected reaction time on Day 0 (baseline)

Average Y change by each additional day.

variance not explained by either fixed or random effect.

Correlation of Fixed Effects:

(Intr) for model diagnostic
Days -0.138

$$Y_i = \beta_0 + \beta_1 X_i + \alpha_{0i} + \alpha_{1i} X_i$$

```
> fit.sleep = lmer(Reaction ~ Days + (Days || Subject) , sleepstudy)
```

```
> summary(fit.sleep)
```

Linear mixed model fit by REML [‘lmerMod’]

Formula: Reaction ~ Days + ((1 | Subject) + (0 + Days | Subject))

Data: sleepstudy

REML criterion at convergence: 1743.7

Scaled residuals:

Min	1Q	Median	3Q	Max
-----	----	--------	----	-----

-3.9626	-0.4625	0.0204	0.4653	5.1860
---------	---------	--------	--------	--------

Random effects:

Groups	Name	Variance	Std.Dev.
Subject	(Intercept)	627.57	25.051
Subject.1	Days	35.86	5.988
Residual		653.58	25.565

Number of obs: 180, groups: Subject, 18

Correlation now not estimated

naming convention only: now two separate random effects

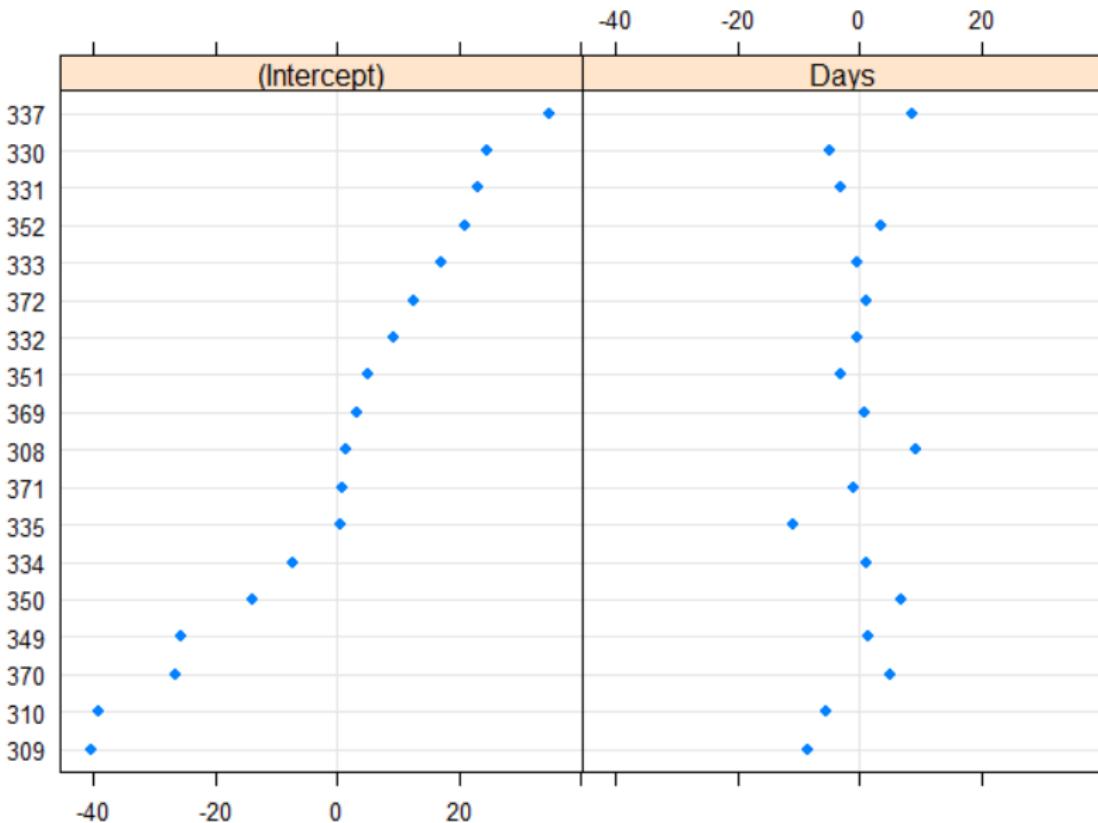
Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	251.405	6.885	36.513
Days	10.467	1.560	6.712

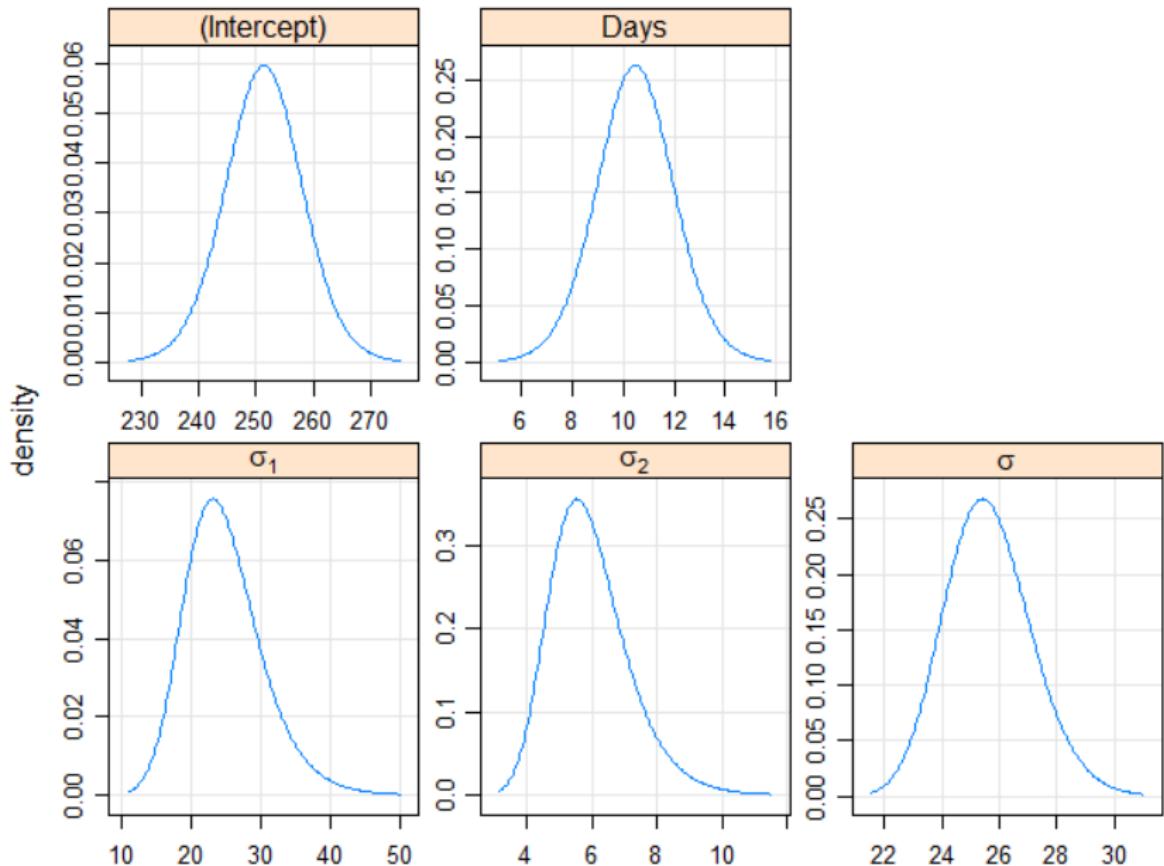
Correlation of Fixed Effects:

(Intr)

Days -0.184

Subject

```
> dotplot(ranef(fit.sleep))
```



```
> fitsleep.profile = profile(fit.sleep); densityplot(fitsleep.profile)
```

Software – SAS

Recall R command

```
> fit.sleep = lmer(Reaction ~ Days + (Days || Subject) , sleepstudy)
```

The following yield equivalent results ...

```
proc mixed data=sleep covtest cl;
  class subject;
  model reaction=days / solution cl;
  random subject days*subject / solution;
proc mixed data=sleep covtest cl;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution;
proc glimmix data=sleep;
  class subject;
  model reaction=days / solution cl;
  random subject days*subject / solution;
proc glimmix data=sleep;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution;
```

$$M + d_i + B_j + \alpha B_{ij} + \varepsilon_{ij}$$

fix slope
 ↑
 $M + d_i + B_j + \alpha B_{ij} + \varepsilon_{ij}$
 ↓
 fix intercept random slope
 ↓ ↓
 random intercept random residual

Correlated **R** & **G** Matrices

`lme4::lmer`: $\mathbf{R} = \sigma^2 \mathbf{I}$ only; \mathbf{G} diagonal (use `||` instead of `|`) or completely unstructured (use `|` instead of `||`)

`proc mixed`: \mathbf{R} controlled by repeated stmt; \mathbf{G} controlled by random stmt

`proc glimmix`: \mathbf{R} controlled by random stmt, with `_residual_` before `/`;
 \mathbf{G} controlled by random stmt

- Autoregressive of order one, i.e., AR(1)

`type=ar(1)`

$$\mathbf{R} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

- Compound symmetry

type=cs

$$\mathbf{R} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{bmatrix}$$

- Toeplitz, two bands

type=toep(2)

$$\mathbf{R} = \sigma^2 \begin{bmatrix} 1 & \rho_1 & 0 & 0 \\ \rho_1 & 1 & \rho_1 & 0 \\ 0 & \rho_1 & 1 & \rho_1 \\ 0 & 0 & \rho_1 & 1 \end{bmatrix}$$

- Toeplitz, three bands

type=toep(3)

$$\mathbf{R} = \sigma^2 \begin{bmatrix} 1 & \rho_1 & \rho_2 & 0 \\ \rho_1 & 1 & \rho_1 & \rho_2 \\ \rho_2 & \rho_1 & 1 & \rho_1 \\ 0 & \rho_2 & \rho_1 & 1 \end{bmatrix}$$

- Unstructured

type=un

$$\mathbf{R} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \sigma_{14} \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} & \sigma_{24} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 & \sigma_{34} \\ \sigma_{14} & \sigma_{24} & \sigma_{34} & \sigma_4^2 \end{bmatrix}$$

- Unstructured, one band

type=un(1)

$$\mathbf{R} = \begin{bmatrix} \sigma_1^2 & 0 & 0 & 0 \\ 0 & \sigma_2^2 & 0 & 0 \\ 0 & 0 & \sigma_3^2 & 0 \\ 0 & 0 & 0 & \sigma_4^2 \end{bmatrix}$$

- Unstructured, two bands

type=un(2)

$$\mathbf{R} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & 0 & 0 \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} & 0 \\ 0 & \sigma_{23} & \sigma_3^2 & \sigma_{34} \\ 0 & 0 & \sigma_{34} & \sigma_4^2 \end{bmatrix}$$

Software – SAS

Recall R command

```
> fit.sleep = lmer(Reaction ~ Days + (Days|Subject) , sleepstudy)
```

The following yield equivalent results . . .

```
proc mixed data=sleep covtest cl;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution type=un;
proc glimmix data=sleep;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution type=un;
covtest diagg / wald cl;
```

Longitudinal modeling that allows correlation across time for each subject . . .

```
proc mixed data=sleep covtest cl;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution;
  repeated / subject=subject type=ar(1);
proc glimmix data=sleep;
  class subject;
  model reaction=days / solution cl;
  random intercept days / subject=subject solution;
  random _residual_ / subject=subject type=ar(1);
  covtest diagr / wald cl;
```

ST704, Sujit K. Ghosh

Generalized Linear Models, Part I

Introduction

Sample questions :

Jan 2024 Part 2 Q₂

Aug 2024 Part 2 Q₃

Aug 2022 Part 1 Q₄

Aug 2021 Part 1 Q₅

Exponential Family

Inference on β

Deviance

Inference on ψ

A Basic Checklist

Some Details

Classical Linear Model (LM):

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- \mathbf{Y} is $n \times 1$, observed
 - \mathbf{X} is $n \times p$, known covariate matrix
 - $\boldsymbol{\beta}$ is $p \times 1$, unknown fixed effects
 - Parameters to be estimated: $\boldsymbol{\beta}, \sigma^2$
-

In other words ...

- Linear Predictor, aka Systematic Component:

$$\mathbf{X}\boldsymbol{\beta} \quad \dots \quad \mathbf{x}_i^T \boldsymbol{\beta} \text{ is } i\text{th entry, often } \eta_i$$

- Identity Link:

$$E(Y_i) = \mathbf{x}_i^T \boldsymbol{\beta} \quad \dots \text{ "link" transforms } E(Y_i) \text{ to linear predictor scale}$$

- Random Component:

$$Y_1, \dots, Y_n \text{ indep normal with } \text{var}(Y_i) = \sigma^2$$

Generalized Linear Model (GLM)

- Linear Predictor, aka Systematic Component:

$$\mathbf{X}\boldsymbol{\beta} \quad \dots \quad \mathbf{x}_i^T \boldsymbol{\beta} \text{ is } i\text{th entry, often } \eta_i$$

- Link:

$$g(E(Y_i)) = \mathbf{x}_i^T \boldsymbol{\beta} \quad \dots \quad E(Y_i) = g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta})$$

Note that $g(\cdot)$ is link, $g^{-1}(\cdot)$ is inverse link

- Random Component:

Y_1, \dots, Y_n indep from the **exponential family** with *dispersion parameter* ψ

Parameters to be estimated: $\boldsymbol{\beta}$ and ψ

Properties of link function $g(\cdot)$:

- monotonic and invertible
- maps mean response to a scale where covariate effects are additive
- ensures range restriction for mean response
- distns in exponential family have “canonical” or “natural” link functions
- any suitable link function may be paired with any distribution in the exponential family

Goals for GLM

- Fit the model, i.e., estimate β and ψ
- CI and HT for β : need distribution, standard error
- Inference on functions of β :
 - $H_0 : \mathbf{A}\beta = \mathbf{m}$?
 - Estimate mean response $g^{-1}(\mathbf{x}_i^T \beta)$... CI & HT?
- How well does the model fit?

Exponential family with natural parameter $\theta = \theta(\mu)$, dispersion (scale) parameter ψ , and known weight w : pdf has form ...

$$f(y_i) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{\psi w_i} + c(y_i, \psi w_i) \right\}$$

for some functions $b(\cdot)$ and $c(\cdot)$.

$$N(\mu_i, \sigma^2)$$

$$\begin{aligned} f(y_i) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right\} = \exp \left\{ -\frac{(y_i - \mu_i)^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) \right\} \\ &= \exp \left\{ \frac{-y_i^2 - \mu_i^2 + 2y_i\mu_i}{2\sigma^2} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) \right\} \\ &= \exp \left\{ \underbrace{\frac{y_i\mu_i - \frac{1}{2}\mu_i^2}{\sigma^2}}_{\theta_i = \mu_i, \psi = \sigma^2} - \underbrace{\frac{y_i^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2)}_{c(y_i, \psi w_i)} \right\} \end{aligned}$$

Exponential family with natural parameter $\theta = \theta(\mu)$, dispersion (scale) parameter ψ , and known weight w : pdf has form ...

$$f(y_i) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{\psi w_i} + c(y_i, \psi w_i) \right\}$$

for some functions $b(\cdot)$ and $c(\cdot)$.

Y_i is proportion of successes: $\frac{1}{n_i} Bin(n_i, p_i)$

$$\begin{aligned} f(y_i) &= \binom{n_i}{n_i y_i} p_i^{n_i y_i} (1-p_i)^{n_i - n_i y_i} = \binom{n_i}{n_i y_i} \left(\frac{p_i}{1-p_i} \right)^{n_i y_i} (1-p_i)^{n_i} \\ &= \exp \left\{ n_i y_i \ln \left(\frac{p_i}{1-p_i} \right) + n_i \ln(1-p_i) + \ln \left(\binom{n_i}{n_i y_i} \right) \right\} \\ &= \exp \left\{ \underbrace{\frac{y_i \ln \left(\frac{p_i}{1-p_i} \right) - \ln \left(\frac{1}{1-p_i} \right)}{1/n_i}}_{\theta_i = \ln \left(\frac{p_i}{1-p_i} \right), \psi = 1, w_i = 1/n_i} + \underbrace{\ln \left(\binom{n_i}{n_i y_i} \right)}_{c(y_i, \psi w_i)} \right\} \\ &\quad b(\theta_i) = \ln \left(\frac{1}{1-p_i} \right) = \ln(1 + e^{\theta_i}) \end{aligned}$$

Exponential family with natural parameter $\theta = \theta(\mu)$, dispersion (scale) parameter ψ , and known weight w : pdf has form ...

$$f(y_i) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{\psi w_i} + c(y_i, \psi w_i) \right\}$$

for some functions $b(\cdot)$ and $c(\cdot)$.

Poisson(λ_i)

$$f(y_i) = \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} = \exp \left\{ \begin{array}{l} \underbrace{y_i \ln \lambda_i - \lambda_i}_{\theta_i = \ln \lambda_i, \psi = 1, w_i = 1} \quad \underbrace{-\ln(y_i!)}_{c(y_i, \psi w_i)} \\ b(\theta_i) = \lambda_i = e^{\theta_i} \end{array} \right\}$$

Properties:

- $\mu_i = b'(\theta_i)$ is mean function
- $b''(\theta_i)$ is 'variance' function, with $\text{var}(Y_i) = \psi w_i b''(\theta_i)$

Distribution for Y_i	ψ	w_i	$\theta_i(\mu_i)$ Canonical link $g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}$	Inverse link $g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}) = \mu_i$	$b'(\theta_i)$ Mean function $E(Y_i)$	$b''(\theta_i)$ 'Variance' fnc $\frac{1}{\psi w_i} \text{var}(Y_i)$
$N(\mu_i, \sigma^2)$	σ^2	1	μ_i ; identity	$\mathbf{x}_i^T \boldsymbol{\beta}$; identity	μ_i	1
$\frac{1}{n_i} \text{Bin}(n_i, p_i)$	1	$\frac{1}{n_i}$	$\ln\left(\frac{p_i}{1-p_i}\right)$; logit	$\frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1+e^{\mathbf{x}_i^T \boldsymbol{\beta}}}$; expit	p_i	$p_i(1-p_i)$
$\text{Poisson}(\lambda_i)$	1	1	$\ln \lambda_i$; log	$e^{\mathbf{x}_i^T \boldsymbol{\beta}}$; exp	λ_i	λ_i

Note: The 'variance' function is sometimes denoted as $h(\mu_i)$

Inference on β

- Estimate β by maximum likelihood

- $\sqrt{n} (\hat{\beta} - \beta) \xrightarrow{d} N(\mathbf{0}, [\text{I}_1(\beta)]^{-1})$ where $n\text{I}_1(\beta) = \mathbf{F}^T \mathbf{V}^{-1} \mathbf{F}$

$$\mathbf{V} = \begin{bmatrix} \text{var}(Y_1) & & 0 \\ & \ddots & \\ 0 & & \text{var}(Y_n) \end{bmatrix},$$

$$\mathbf{F} = \frac{\partial \mu}{\partial \beta^T}$$

In other words, $\hat{\beta} \approx N(\beta, [\mathbf{F}^T \mathbf{V}^{-1} \mathbf{F}]^{-1})$

...use this for CIs

- Hypothesis testing for β

- Wald test: $T_W \stackrel{H_0}{\approx} \chi_{\nu}^2$
- Score test: $T_S \stackrel{H_0}{\approx} \chi_{\nu}^2$
- Asymptotic likelihood ratio test: $T_{LR} \stackrel{H_0}{\approx} \chi_{\nu}^2$
- Exact likelihood ratio test

$$T_{LR} = 2\{\ell_{H_a} - \ell_{H_0}\}$$

Deviance

$\ell(\mu, \psi; \mathbf{y})$ denotes loglikelihood

- A *saturated* model has $\hat{\mu} = \mathbf{y}$, and “fits the data perfectly”
- Consider a model M, less complex than *saturated*, with estimated mean response $\hat{\mu}$ based on p regression parameters.

How well does Model M fit, relative to the saturated model?

- Deviance for model M is

$$D_M = \psi \cdot T_{LR}: \text{saturated vs. model M} = \psi \cdot 2 \{ \ell(\mathbf{y}, \psi; \mathbf{y}) - \ell(\hat{\mu}, \psi; \mathbf{y}) \}$$

- Scaled deviance for model M is simply D_M/ψ

$$N(\mu_i, \sigma^2)$$

$$\ell(\mu, \psi; \mathbf{y}) = \sum_{i=1}^n \left\{ -\frac{(y_i - \mu_i)^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) \right\}$$

for normal

$$\Rightarrow \ell(\mathbf{y}, \psi; \mathbf{y}) = \sum_{i=1}^n \left\{ -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) \right\}$$

$$\ell(\hat{\mu}, \psi; \mathbf{y}) = \sum_{i=1}^n \left\{ -\frac{(y_i - \hat{\mu}_i)^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma^2) \right\}$$

* exact !

$$D_M = \psi \cdot 2 \sum_{i=1}^n \left\{ \frac{(y_i - \hat{\mu}_i)^2}{2\sigma^2} \right\} = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 = SSE \dots \dots D_M/\psi \sim \chi^2_{n-p} \text{ is exact}$$



- If Y_i has a distribution “close to normal” with link “close to identity,” then $D_M/\psi \approx \chi^2_{n-p}$

- Approximation will often NOT improve as n increases!
- Suppose data are grouped, n is the # of groups and is fixed. We want the size of each group to be large

$Bin(n_i, p_i), i = 1, \dots, n$: want n_i large

- Lack of fit testing ... H_0 : model M fits the data vs. H_a : not H_0

Reject if $D_M/\psi > \chi^2_{n-p,\alpha}$ Global test

- Consider testing H_0 : model M_0 vs. H_a : model M, where $M_0 \subset M$ is a submodel (i.e., nested) with $q < p$ regression parameters. Note that

$$T_{LR} = \frac{D_{M_0} - D_M}{\psi} \xrightarrow{d} \chi^2_{p-q} \text{ when } H_0 \text{ is true}$$

Distribution D_M

$$N(\mu_i, \sigma^2) \quad \sum_{i=1}^n (y_i - \hat{\mu}_i)^2$$

$$\frac{1}{n_i} Bin(n_i, p_i) \quad 2 \sum_{i=1}^n \left\{ n_i y_i \ln \left(\frac{y_i}{\hat{\mu}_i} \right) + (n_i - n_i y_i) \ln \left(\frac{1-y_i}{1-\hat{\mu}_i} \right) \right\}$$

$$Bin(n_i, p_i) \quad 2 \sum_{i=1}^n \left\{ y_i \ln \left(\frac{y_i}{\hat{\mu}_i} \right) + (n_i - y_i) \ln \left(\frac{n_i - y_i}{n_i - \hat{\mu}_i} \right) \right\}$$

$$Poisson(\lambda_i) \quad 2 \sum_{i=1}^n \left\{ y_i \ln \left(\frac{y_i}{\hat{\mu}_i} \right) - (y_i - \hat{\mu}_i) \right\}$$

Estimating ψ

- Some families set a value for ψ , but it is good practice to estimate it
- Assuming $D_M/\psi \approx \chi^2_{n-p}$, a method of moments estimator of ψ is



$$D_M/\psi \stackrel{\text{set}}{=} E(\chi^2_{n-p}) = n - p \quad \Rightarrow \quad \widehat{\psi}_{\text{MoM}} = \frac{D_M}{n - p}$$

- $\widehat{\psi}$ "large" may be due to
 - inadequate linear predictor, i.e., missing predictors
 - overdispersion, i.e., $\text{var}(Y_i) > \psi w_i h(\mu_i)$ where ψ is given by family
 - Correlation between Y_1, \dots, Y_n can lead to overdispersion
- Recall ???

$$\text{var}(\widehat{\beta}) \doteq \psi (\mathbf{F}^T [\text{diag}(w_1 h(\mu_1), \dots, w_n h(\mu_n))]^{-1} \mathbf{F})^{-1}$$

so over(under)-reporting the value of ψ will lead to SEs that are too large(small).

A Basic Checklist

- Choose family (distribution + link)
- Select covariates and estimate β
 - Lack of fit testing using D_M/ψ if appropriate
 - Compare scaled deviances of nested models
- Estimate ψ
- Report $SE(\hat{\beta}_j)$, $SE(\hat{\mu}_i)$
- Inference on nonlinear function of β will need Taylor's approx for SE
- Model interpretation – odds, odds ratio, log odds, log odds ratio
- Problems with convergence – complete and quasi-complete separation

Inference, treating ψ as known

Maximum likelihood estimation for β : need $\ell(\beta, \psi; \mathbf{y})$

- $\mu_i = g^{-1}(\mathbf{x}_i^T \beta)$ is a function of β
- θ_i is a function of μ_i , and hence a function of β
- $\ell(\mu, \psi; \mathbf{y})$ is loglikelihood function

→ natural parameter

$$\ell(\mu, \psi; \mathbf{y}) = \sum_{i=1}^n \left\{ \frac{y_i \theta_i - b(\theta_i)}{\psi w_i} + c(y_i, \psi w_i) \right\}$$

- Score wrt θ :

$$\frac{\partial \ell(\mu, \psi; \mathbf{y})}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{\psi w_i} = \frac{y_i - \mu_i}{\psi w_i}$$

- Score wrt β :

$$\frac{\partial \ell(\mu, \psi; \mathbf{y})}{\partial \beta} = \mathbf{F}^T \mathbf{V}^{-1} (\mathbf{y} - \mu)$$

is $b''(\theta_i) = h(\mu_i)$?

?

where $\mathbf{V} = \text{diag}(\psi w_1 h(\mu_1), \dots, \psi w_n h(\mu_n))$ has $\text{var}(Y_i)$ on its diagonal, and $\mathbf{F} = \frac{\partial \mu}{\partial \beta}$ has n rows and same number of columns as \mathbf{X}

- Example: $N(\mu_i, \sigma^2)$:

$$\mu_i = \mathbf{x}_i^T \boldsymbol{\beta}, \quad \boldsymbol{\mu} = \mathbf{X} \boldsymbol{\beta}, \quad h(\mu_i) = 1, \quad \psi = \sigma^2, \quad w_i = 1$$

$$\hookrightarrow \mathbf{V} = \sigma^2 \mathbf{I}, \quad \mathbf{F} = \frac{\partial \mathbf{X} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}^T} = \mathbf{X}$$

$$\hookrightarrow \mathbf{X}^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) / \sigma^2 = \mathbf{0} \quad \Rightarrow \quad \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}, \text{ usual normal eqns}$$

- An estimate of the large-sample variance of $\hat{\boldsymbol{\beta}}$ is

$$\widehat{\text{var}}(\hat{\boldsymbol{\beta}}) = \left(\hat{\mathbf{F}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{F}} \right)^{-1} \stackrel{\text{normal, exact}}{=} \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

- Estimate mean response ...

- $\hat{\mu} = g^{-1}(\mathbf{x}^T \hat{\boldsymbol{\beta}})$ usu. biased, even if $\hat{\boldsymbol{\beta}}$ is unbiased
- $SE(\hat{\mu})$: Use Taylor series of $g^{-1}(\mathbf{x}^T \hat{\boldsymbol{\beta}})$ around $\boldsymbol{\beta}^*$ to get linear approximation, then use variance of linearization:

$$\widehat{\text{var}}(\hat{\mu}) \doteq \left[\frac{\partial g^{-1}(\eta)}{\partial \eta} \Big|_{\eta=\hat{\eta}} \right]^2 \mathbf{x}^T \left(\hat{\mathbf{F}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{F}} \right)^{-1} \mathbf{x} \quad \text{delta theorem}$$

- If the canonical link (i.e., $\eta = \theta(\mu)$) is used: $\frac{\partial g^{-1}(\eta)}{\partial \eta} = \frac{\partial \mu}{\partial \theta(\mu)} = h(\mu)$

$$\widehat{\text{var}}(\hat{\mu}) \doteq h(\hat{\mu})^2 \mathbf{x}^T \left(\hat{\mathbf{F}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{F}} \right)^{-1} \mathbf{x}$$

- Tests for functions of β may be conducted using asymptotic Wald, likelihood ratio, and score tests, comparing each of T_W , T_{LR} , T_S to a chi-squared distribution with degrees of freedom matching the number of constraints placed on β .

E.g., for testing $H_0 : \mathbf{A}\beta = \mathbf{m}$, then

$$T_W = (\mathbf{A}\hat{\beta} - \mathbf{m})^T \left(\mathbf{A} \left(\hat{\mathbf{F}}^T \hat{\mathbf{V}}^{-1} \hat{\mathbf{F}} \right)^{-1} \mathbf{A}^T \right)^{-1} (\mathbf{A}\hat{\beta} - \mathbf{m})$$

compares to a $\chi^2_{rank(\mathbf{A})}$ distribution

ST704, Sujit K. Ghosh
GLM, Part II

Logistic Regression

Probit Regression

Interpreting the Logistic Model

Example: Dose Response Modeling, using Logistic Regression

Data and plots

Logistic model

Software

SAS code: genmod, glimmix, logistic

R code: `glm`

Residuals and Diagnostics

Complete and Quasi-complete Separation

A Designed Binomial Study

Poisson Modeling

Contingency Tables via Poisson Modeling

Logistic Regression

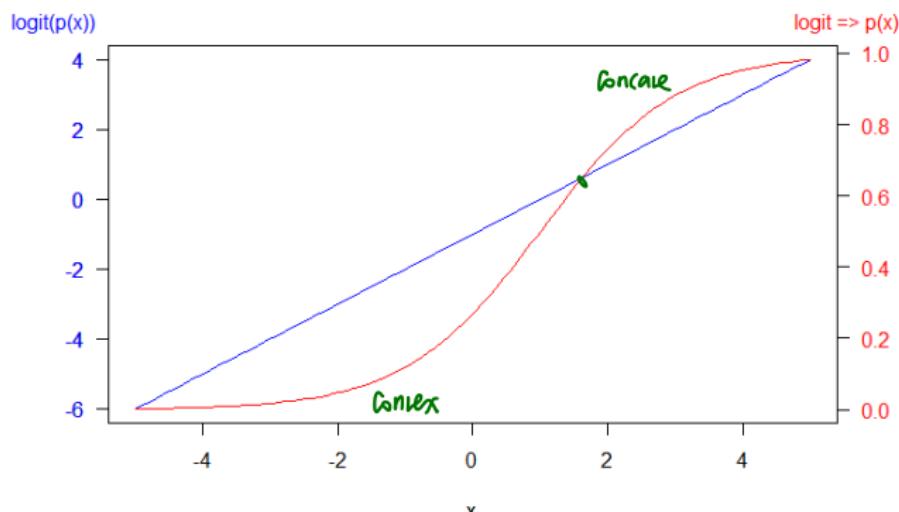
Linear predictor: $x_i^T \beta$

Link function: logit, i.e., $\ln\left(\frac{p_i}{1-p_i}\right)$ also $\ln\left(\frac{p(x)}{1-p(x)}\right)$

Random component: Y_i independent $\frac{1}{n_i} \text{Bin}(n_i, p_i)$, $i = 1, 2, \dots, n$

$$\boxed{\text{logit}(p(x)) = \ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x = -1 + x}$$

$g^{-1}(x)$



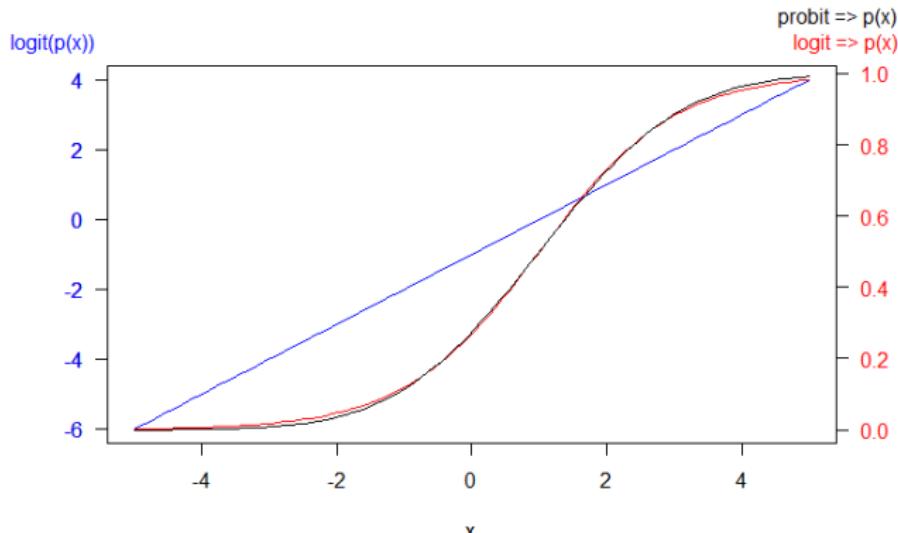
Probit Regression

Linear predictor: $\mathbf{x}_i^T \boldsymbol{\beta}$

Link function: probit, i.e., $\Phi^{-1}(p_i)$, where $\Phi(\cdot)$ is $N(0,1)$ cdf also $\Phi^{-1}(p(x))$

Random component: Y_i independent $\frac{1}{n_i} \text{Bin}(n_i, p_i)$, $i = 1, 2, \dots, n$

$$\Phi^{-1}(p(x)) = \beta_0 + \beta_1 x = -0.6 + 0.6x$$



Interpreting the Logistic Model

$$\ln \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x$$

$$\text{odds}(x) = \frac{p(x)}{1-p(x)}$$

$\equiv 1$ \rightsquigarrow S & F equally likely, aka 'odds are even'

$< 1 \rightsquigarrow$ S less likely

$> 1 \rightsquigarrow$ S more likely

$$\ln \left(\frac{p(x)}{1-p(x)} \right) \stackrel{\beta_0 = \beta_1 = 0}{=} 0$$

\rightsquigarrow odds ($= e^0 = 1$) are even and don't change with x

$$\ln \left(\frac{p(x)}{1-p(x)} \right) \stackrel{\beta_1=0}{=} \beta_0$$

\rightsquigarrow odds ($= e^{\beta_0}$) are not even but don't change with x

* $\beta_0 > 0$: S is more likely odd > 1

* $\beta_0 < 0$: S is less likely add x

$$\left[\frac{p(x+1)}{1-p(x+1)} \right] \Bigg/ \left[\frac{p(x)}{1-p(x)} \right]$$

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

is 'odds ratio for 1 unit increase in x'

$$\ln \left(\frac{p(x)}{1-p(x)} \right) = \beta_0 + \beta_1 x \quad \Rightarrow \quad \frac{p(x)}{1-p(x)} = e^{\beta_0} (e^{\beta_1})^x$$

- The odds at $x = 0$ is e^{β_0}
 - The odds increase multiplicatively by e^{β_1} for 1 unit increase in x
 - $\beta_1 > 0$: $p(x) \uparrow$ as $x \uparrow$
 - $\beta_1 < 0$: $p(x) \downarrow$ as $x \uparrow$

Dose Response Modeling

Dose: % concentration of insecticide

Response: proportion of larvae killed, $\hat{p}_i: y_i = k_i/n_i, n_i = 20, i = 1, \dots, 7$

Assume:

- (1) larvae react independently
- (2) larvae exposed to same conc have equal probability of survival

concentration (%)	.375	.75	1.5	3	6	12	24
# of larvae killed	0	1	8	11	16	18	20

Grouped data

...

$n = 7$ groups

NOT

$n = 7 \times 20 = 140$

- Choose family : $Bin(n_i, p_i)$ + logit link
- Select covariates ... look at plots ... but $\hat{p}_i = 0, \hat{p}_i = 1$ cause problems

empirical logit:

$$\ln\left(\frac{k_i + 0.5}{n_i - k_i + 0.5}\right)$$

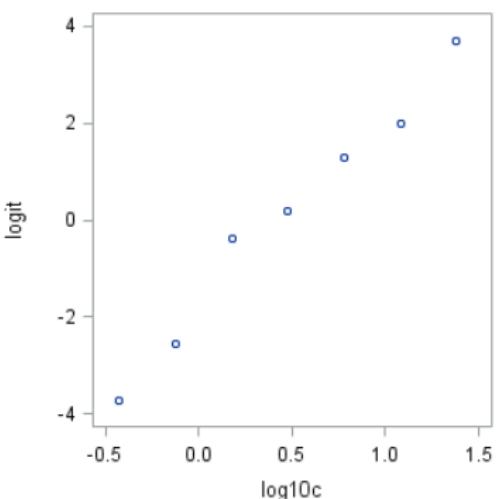
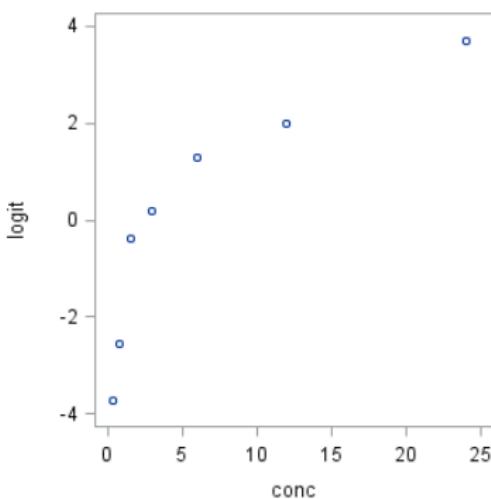
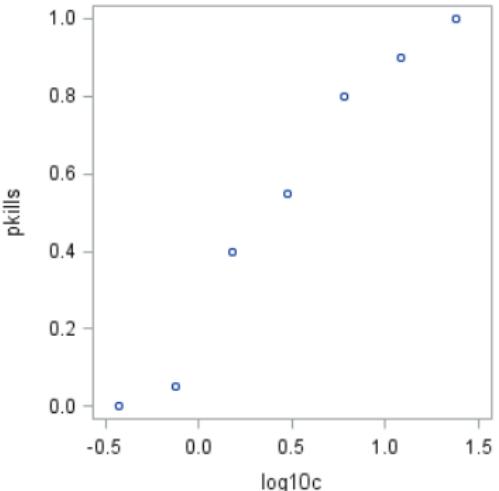
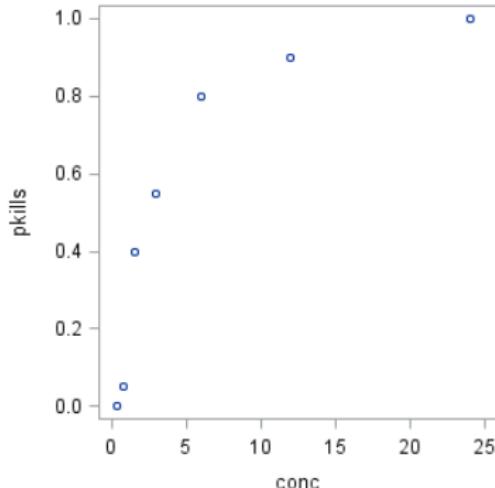
$$\ln\left(\frac{P}{1-P}\right) = \ln\left(\frac{k_i/n_i}{1-k_i/n_i}\right) \\ = \ln\left(\frac{k_i}{n_i - k_i}\right)$$

what is this 0.5



Logistic
O

Probit
O



via Poisson

Logistic Model

Linear predictor: $\beta_0 + \beta_1 x$, $x = \log_{10}(\text{conc})$

Link function: logit, i.e., $\ln\left(\frac{p_i}{1-p_i}\right)$

Random component: Y_i independent $\text{Bin}(n_i = 20, p_i)$, $i = 1, 2, \dots, 7$

$$\ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x$$

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

Fit to data:

$$\ln\left(\widehat{\frac{p(x)}{1-p(x)}}\right) = -1.7305 + 4.1651x = \widehat{\beta}_0 + \widehat{\beta}_1 x$$

$$\ln \left(\widehat{\frac{p(x)}{1-p(x)}} \right) = -1.7305 + 4.1651x = \hat{\beta}_0 + \hat{\beta}_1 x$$

- $e^{\hat{\beta}_0} = 0.177$ estimates “odds at $x = 0$ ” (i.e., odds at $c = 10^0 = 1$)
 ... survival is more likely (more than 5×) than death at % conc of 1

$$\text{CI: LR} \rightarrow (e^{-2.5351}, e^{-1.0557}) = (0.079, 0.348)$$

$$\text{CI: Wald} \rightarrow (e^{-2.4637}, e^{-0.9973}) = (0.085, 0.369)$$

- $e^{\hat{\beta}_1} = 64.399$ estimates “odds ratio for 1 unit increase in x ”
 ... odds of death increase by a factor of 64.399 for each increase of 1% conc

$$\text{CI: LR} \rightarrow (e^{3.0174}, e^{5.6003}) = (20.438, 270.508)$$

$$\text{CI: Wald} \rightarrow (e^{2.8872}, e^{5.4430}) = (17.943, 231.135)$$

Lack of fit?

- $D_M \equiv 4.6206$

Is this “large,” suggesting a bad model?

- Ok to do a formal test because $n_i = 20$ is reasonable

H_0 : model fits data well vs. H_a : not H_0

$$p\text{-value} = \Pr(\chi^2_{n-p} > D_M/\psi) = \Pr(\chi^2_{7-2} > 4.6206) = 0.4639$$

Estimate ψ : $\widehat{\psi} = \frac{D_M}{n-p} = \frac{4.6206}{7-2} = 0.9241$

“close to” theoretical value 1

Inverse regn: What concentration kills 50% of larvae? ... 80%? ... 100δ%?

$$\ln\left(\frac{\delta}{1-\delta}\right) = \beta_0 + \beta_1 x_\delta \quad \Rightarrow \quad x_\delta = \frac{1}{\beta_1} \left\{ \ln\left(\frac{\delta}{1-\delta}\right) - \beta_0 \right\}$$

x_δ is a nonlinear function of β , so need Taylor expansion for standard error!

	Estimate	SE
LD50	2.6030	0.3647
LD80	5.6016	1.0246

SAS code

```

data kills;
  input conc kills;
  trials=20;
  pkills=kills/trials;
  logit=log((kills+.5)/(trials - kills+.5));
  log10c=log10(conc);
datalines;
0.375 0
0.75  1
1.5   8
3.0   11
6.0   16
12.0  18
24.0  20
;
proc sgscatter data=kills;
  plot (pkills logit)*(conc log10c) / columns=2 rows=2;
proc genmod data=kills;
  model kills/trials=log10c / dist=binomial link=logit type1 type3 lrci;
                                         *scale=deviance;
proc glimmix data=kills;
  model kills/trials=log10c / dist=binomial link=logit solution cl;

```

SAS code

```
proc logistic data=kills plots=all;  
  model kills/trials=log10c / link=logit clparm=both clodds=both;
```

Options offered by proc logistic include:

- lots of diagnostics plots, checks
 - forward, backward, stepwise, best subset selection
 - produce a receiver operating characteristic (ROC) curve for fitted model
 - BUT only for independent binomial or multinomial data

SAS code

```
proc nlmixed data=kills df=5;
parameters b0=-1.7 b1=4.0;
p = 1/(1+exp(- b0 - b1*log10c));
model kills ~ binomial(trials,p);
estimate 'LD50' -b0/b1;
estimate 'LD50 original' 10**(-b0/b1);
estimate 'LD80' ( log(0.8/0.2) - b0 ) /b1;
estimate 'LD80 original' 10**(( log(0.8/0.2) - b0 )/b1);
estimate 'OR' exp(b1);
```

R code

```
library(MASS)
library(car)

killsDF = data.frame(conc=c(.375,.75,1.5,3,6,12,24), dead=c(0,1,8,11,16,18,20))
killsDF$log10c = log(killsDF$conc,10)
killsDF$alive = 20 - killsDF$dead
fit = glm(cbind(dead,alive) ~ log10c, family=binomial(link=logit), data=killsDF )
summary( fit, correlation=T )

deviance(fit)
anova(fit)

confint(fit) #profile likelihood CI by default
confint.default(fit) #Wald CI, using z

( dose.p( fit, cf=1:2, p=c(.5,.8)) ) #Does inverse regression

predict(fit, type="response") #estimated mu
predict(fit, type="link") #estimated linear predictor (default)
predict(fit, type="response", se.fit=TRUE)
fitted(fit) #estimated mu
```

```

resid(fit, type="deviance")    #default
resid(fit, type="pearson")

residualPlots(fit)  #loess smooth replaces quadratic
outlierTest(fit)
influencePlot(fit)
crPlots(fit)

library(visreg)
visreg(fit,scale="linear",ylab="log odds (death)",points=list(cex=1))
visreg(fit,scale="response",ylab="Pr(death)",partial=TRUE,points=list(cex=1))

library(glmnet)
fit2 = glmnet( cbind(killsDF$log10c, killsDF$conc),
cbind(killsDF$alive, killsDF$dead),
family="binomial" )

```

Residuals and Diagnostics

- Several types of residuals are commonly used:

• Raw or response: $y_i - \hat{\mu}_i$

not very useful

• Pearson: $r_i^P \equiv \frac{y_i - \hat{\mu}_i}{\sqrt{w_i h(\hat{\mu}_i)}}$ $\rightsquigarrow \mathbb{X}^2 = \sum_{i=1}^n (r_i^P)^2$

very useful for diagnostics

also $\frac{r_i^P}{\sqrt{\psi} \sqrt{1-h_i}}$ where h_i is leverage

• Deviance: $r_i^D \equiv \text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i}$ $\rightsquigarrow D_M = \sum_{i=1}^n d_i = \sum_{i=1}^n (r_i^D)^2$

very useful for diagnostics

also $\frac{r_i^D}{\sqrt{\psi} \sqrt{1-h_i}}$

Note: The chi-squared statistic \mathbb{X}^2 and the scaled chi-squared statistic \mathbb{X}^2/ψ are often used interchangeably with the deviance D_M and scaled deviance D_M/ψ .

- New diagnostic measure: likelihood displacement:

$$LD_i \equiv 2 \left\{ \ell_M(\hat{\theta}; \mathbf{y}) - \ell_M(\hat{\theta}_{(-i)}; \mathbf{y}) \right\}$$

where $\hat{\theta}_{(-i)}$ is MLE from excluding i th observation. Likelihood evaluated with *all* observations.

- Predicted values: $\hat{\mu}_i$, response or $\mathbf{x}_i^T \hat{\beta}$, linear predictor

Complete and Quasi-complete Separation¹

What is it? A linear predictor (almost) completely separates response values

Complete:

y	x1	x2
0	1	3
0	2	2
0	3	-1
0	3	-1
1	5	2
1	6	4
1	10	1
1	11	0

$\Pr(Y = 1|x_1 > 3)$ best estimated as 1

$\Pr(Y = 1|x_1 \leq 3)$ best estimated as 0



$$p(x_1) = \frac{e^{\beta x_1}}{1 + e^{\beta x_1}} \quad x_1 > 3: \beta \rightarrow \infty \quad 1$$

$$p(x_1) = \frac{e^{\beta x_1}}{1 + e^{\beta x_1}} \quad 0 < x_1 \leq 3: \beta \rightarrow -\infty \quad 0$$

Quasi-complete:

y	x1	x2
0	1	3
0	2	2
0	3	-1
1	3	-1
1	5	2
1	6	4
1	10	1
1	11	0

Consequence: Difficulty getting convergence

Fix: Biased regression ... Firth's penalized procedure. Other penalties

Common with: Rare events; very large predictor space; many binary predictors; small sample size

¹<https://stats.idre.ucla.edu/other/mult-pkg/faq/general/>

faqwhat-is-complete-or-quasi-complete-separation-in-logisticprobit-regression-and-how-do-we-deal-with-them/

SAS Code

```
data one;  
input y x1 x2;  
datalines;  
0 1 3  
0 2 2  
0 3 -1  
0 3 -1  
1 5 2  
1 6 4  
1 10 1  
1 11 0  
;  
proc genmod data=one descending; } ①  
model y = x1 x2 / lrci lrcl ;  
proc logistic data=one; } ②  
model y(event='1') = x1 x2 / cl plcl;  
proc logistic data=one; } ③  
model y(event='1') = x1 x2 / cl plcl firth;
```

to model the prob. of 1

10 proc genmod data=one **descending**
 11 model y = x1 x2 / dist=binomial lrci lrcl ; run;

} ①

NOTE: PROC GENMOD is modeling the probability that $y='1'$.

NOTE: The Pearson chi-square and deviance are not computed since the AGGREGATE option is not specified.

NOTE: Algorithm converged.

WARNING: Convergence not attained for at least one side of profile likelihood confidence interval for Prm1. Number of iterations = 50.

WARNING: Convergence not attained for at least one side of profile likelihood confidence interval for Prm2. Number of iterations = 50.

WARNING: Convergence not attained for at least one side of profile likelihood confidence interval for Prm3. Number of iterations = 50.

NOTE: The scale parameter was held fixed. appropriate for binomial models.

NOTE: PROCEDURE GENMOD used (Total process time):

real time	0.05 seconds
cpu time	0.04 seconds

The GENMOD Procedure

①

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Log Likelihood		0.0000	
Full Log Likelihood		0.0000	
AIC (smaller is better)		6.0000	
AICC (smaller is better)		12.0000	
BIC (smaller is better)		6.2383	

model perfectly separates the data,
can find parameters that perfectly
classify $y=0$ vs. $y=1$

Algorithm converged.

Analysis Of Maximum Likelihood Parameter Estimates

Parameter	DF	Estimate	Standard	Likelihood Ratio	95% Confidence Limits	Wald Chi-Square	Pr > ChiSq
			Error	Chi-Square			
Intercept	1	-107.266	4.912E8	-107.266	-107.266	0.00	1.0000
x1	1	25.1805	75202009	25.1805	25.1805	0.00	1.0000
x2	1	9.5189	2.1684E8	9.5189	9.5189	0.00	1.0000
Scale	0	1.0000	0.0000	1.0000	1.0000		

NOTE: The scale parameter was held fixed.

huge : non-meaningfull/unstable

As cannot find proper CI

```

12 proc logistic data=one; * descending;
13   model y(event='1') = x1 x2 / cl plcl; } ②
  
```

Wald ↗ Profile-likelihood ↘

NOTE: PROC LOGISTIC is modeling the probability that y=1.

WARNING: There is a complete separation of data points. The maximum likelihood estimate does not exist.

WARNING: The LOGISTIC procedure continues in spite of the above warning. Results shown are based on the last maximum likelihood iteration. Validity of the model fit is questionable.

NOTE: There were 8 observations read from the data set WORK.ONE.

NOTE: PROCEDURE LOGISTIC used (Total process time):

real time	0.08 seconds
cpu time	0.03 seconds

```

14 proc logistic data=one; * descending;
15   model y(event='1') = x1 x2 / cl plcl firth; } ③
16 run;
  
```

Firth penalized procedure

NOTE: PROC LOGISTIC is modeling the probability that y=1.

NOTE: Convergence criterion (GCONV=1E-8) satisfied for the intercept-only model.

NOTE: Convergence criterion (GCONV=1E-8) satisfied.

NOTE: There were 8 observations read from the data set WORK.ONE.

NOTE: PROCEDURE LOGISTIC used (Total process time):

real time	0.08 seconds
cpu time	0.04 seconds

The LOGISTIC Procedure

2

Complete separation of data points detected.

WARNING: The maximum likelihood estimate does not exist.

WARNING: The LOGISTIC procedure continues in spite of the above warning. Results shown are based on the last maximum likelihood iteration. Validity of the model fit is questionable.

Model Fit Statistics

	Intercept Only	Intercept and Covariates
Criterion		
AIC	13.090	6.005
SC	13.170	6.244
-2 Log L	11.090	0.005

dramatic improvement,
but the "perfect fit" is misleading

5

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	11.0850	2	0.0039
Score	6.8932	2	0.0319
Wald	0.1302	2	0.9370

wald directly depend on SE,
only wald insignificant indicates
large SE.

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	-20.7083	73.7757	0.0788	0.7789
x1	1	4.4921	12.7425	0.1243	0.7244
x2	1	2.3960	27.9875	0.0073	0.9318

very large

unstable

Odds Ratio Estimates

	Point Estimate	95% Wald Confidence Limits
Effect x1	89.311	<0.001 >999.999
x2	10.980	<0.001 >999.999

massive

meaning less

Parameter Estimates and Profile-Likelihood Confidence Intervals

Parameter	Estimate	95% Confidence Limits
Intercept	-20.7083	. -2.2738
x1	4.4921	0.4161
x2	2.3960	.

- Does not compute limits, because likelihood is flat or infinite.

Parameter Estimates and Wald Confidence Intervals

Parameter Estimates and Wald Confidence Intervals			
Parameter	Estimate	95% Confidence Limits	
Intercept	-20.7083	-165.3	123.9
x1	4.4921	-20.4827	29.4669
x2	2.3960	-52.4584	57.2505

- includes 0, insignificant
 - unreliable due to inflated standard errors

The LOGISTIC Procedure

3

Intercept-Only Model Convergence Status
Convergence criterion (GCONV=1E-8) satisfied.

Model Convergence Status

Convergence criterion (GCNV=1E-8) satisfied.

Model Fit Statistics

	Intercept Only	Intercept and Covariates
Criterion		
AIC	7.478	5.995
SC	7.557	6.233
-2 Log L	5.478	-0.005

³
⁵ rare but okay with penalized likelihood

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	5.4830	2	0.0645
Score	6.8932	2	0.0319
Wald	2.6789	2	0.2620

wald tends to be less reliable in small or noisy data

Analysis of Penalized Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald	
				Chi-Square	Pr > ChiSq
Intercept	1	-3.2433	2.1696	2.2348	0.1349
x1	1	0.5163 B ₁	0.3318	2.4212	0.1197
x2	1	0.4777 B ₂	0.6322	0.5711	0.4498

Odds Ratio Estimates

	Point Estimate	95% Wald Confidence Limits		CI includes 1 so neither significant
Effect x1	1.676	0.875	3.211	
x2	1.612	0.467	5.567	

Parameter Estimates and Profile-Likelihood Confidence Intervals

Parameter	Estimate	95% Confidence Limits	
Intercept	-3.2433	-15.6408	-0.1134
x1	0.5163	0.0516	2.1353
x2	0.4777	-0.6628	4.5243

Generally more accurate
than the Wald test.
(when small sample size)

Parameter Estimates and Wald Confidence Intervals

Parameter	Estimate	95% Confidence Limits
Intercept	-3.2433	-7.4956 1.0089
x1	0.5163	-0.1340 1.1667
x2	0.4777	-0.7613 1.7168

R Code

```

> one = data.frame( y=c(0,0,0,0,1,1,1,1),x1=c(1,2,3,3,5,6,10,11),
+                     x2=c(3,2,-1,-1,2,4,1,0) )
> fit1 = glm(y ~ x1+x2, data=one, family=binomial) ←①
Warning message: glm.fit: fitted probabilities numerically 0 or 1 occurred
> summary(fit1)
Call: glm(formula = y ~ x1 + x2, family = binomial, data = one)

Deviance Residuals:
  1      2      3      4      5      6      7      8 
-2.110e-08 -1.404e-05 -2.522e-06 -2.522e-06  1.564e-05  2.110e-08  2.110e-08  2.110e-08

```

separation detected

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-66.098	183471.722	0.000	1
x1	15.288	27362.843	0.001	1
x2	6.241	81543.720	0.000	1

very large

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1.1090e+01 on 7 degrees of freedom

Residual deviance: 4.5454e-10 on 5 degrees of freedom

AIC: 6

fake perfect fit : due to overfitting

Number of Fisher Scoring iterations: 24

$T = \frac{D_0 - D_M}{4p} \sim \chi^2_{df_p - df_M}$
 under completely separation
 thus deviance / LR tests
 are not reliable.
 (wald also highly affected)

birth

```

> library(logistf)
> fit2 = logistf(y ~ ., data=one) ← ②
> summary(fit2)
logistf(formula = y ~ ., data = one)
  
```

Model fitted by Penalized ML

Confidence intervals and p-values by Profile Likelihood Profile Likelihood Profile Likelihood

	coef	se(coef)	lower	0.95	upper	0.95	Chisq	p
(Intercept)	-2.9748886	2.0332566	-15.47721061	-0.1208941	4.2179522	0.03999841		
x1	0.4908484	0.3241088	0.05268297	2.1275832	5.0225056	0.02501994	significant	
x2	0.4313730	0.5941957	-0.65793072	4.4758930	0.7807099	0.37692411		

Likelihood ratio test=5.505687 on 2 df, $p=0.06374636$, n=8

Wald test = 2.569861 on 2 df, p = 0.2766698

Covariance-Matrix:

	[,1]	[,2]	[,3]
[1,]	4.1341324	-0.4970381	-0.6764776
[2,]	-0.4970381	0.1050465	0.0260937
[3,]	-0.6764776	0.0260937	0.3530685

penalized estimates : shrink extreme values

under small sample and separation, use LRT.

A Designed Binomial Study

- An experiment involves a 4×4 factorial design with factors

temperature: T_1, T_2, T_3, T_4 and concentration: 0, 0.1, 1.0, 10

Completely randomized design

- What is the effect on germination probability of seeds?
- For each treatment combination, there are 4 dishes each with 50 seeds. Count number that germinate, Y_{ijk} .
- Assume
 - seeds germinate independently
 - seeds treated similarly have the same probability of germinating
 - Then $Y_{ijk} \sim \text{Bin}(50, \pi_{ij})$
- Some questions:
 - $p_{1j} = p_{2j} = p_{3j} = p_{4j}$
 - $p_{i1} = p_{i2} = p_{i3} = p_{i4}$
 - $p_{1\cdot} = p_{2\cdot} = p_{3\cdot} = p_{4\cdot}$
 - $p_{\cdot 1} = p_{\cdot 2} = p_{\cdot 3} = p_{\cdot 4}$

Code

```
proc genmod data=germrate;
  class temp conc;
  model germ/trials = temp conc temp*conc /
    link=logit dist=binomial type1 type3;

fit = glm( cbind(germ,trials-germ) ~ temp*conc,
           family=binomial(link=logit), data=germrate );
summary( fit, correlation=F )
deviance( fit )
anova( fit )
```

Code

```

fit = glm( cbind(germ,trials-germ) ~ temp*conc,
            family=binomial(link=logit), data=germrates )
summary( fit, correlation=F )
deviance( fit )
anova( fit )
  
```

how much deviance each term
reduces when added sequentially
to the model.

> anova(fit)

Analysis of Deviance Table

Model: binomial, link: logit

Response: cbind(germ, trials - germ)

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev
NULL			63	1193.80
temp	3	763.69	60	430.11
conc	3	282.01	57	148.11
temp:conc	9	92.46	48	55.64

very significant

Good model : most deviance explained

From this, can test null v.s. temp ; temp v.s. temp + conc ; temp + conc v.s. temp + conc + temp*conc

Poisson Modeling

GLMs with log link are often called log-linear models.

AT&T 1988 soldering experiment : R data solder in package rpart, $n = 720$

skips: number of defects (solder skips) on a circuit board [response]

Opening: amount of clearance around the mounting pad (3 levels)

Solder: amount of solder (Thick or Thin)

Mask: type and thickness of the material used for the solder mask (A1.5, A3, A6, B3, B6)

PadType: geometry and size of the mounting pad (10 levels)

Panel: each board was divided into 3 panels

Code

```
library(rpart)
plot(skips~.,solder)
solder$Panel = as.factor(solder$Panel)
summary(solder)

fit = glm(skips~., family=poisson, data=solder); summary(fit)
anova(fit, test="Chisq")
```

```
proc genmod data=solder;
  class opening solder mask padtype panel;
  model skips = opening solder mask padtype panel /
    link=log dist=poisson type1 type3;
```

Poisson Modeling: $n_{ij} \sim \text{Poisson}(\lambda_{ij})$

Assuming independence:

$n_{i\cdot} \sim \text{Poisson}(\lambda_{i\cdot})$	where $\lambda_{i\cdot} = \sum_j \lambda_{ij}$
$n_{\cdot j} \sim \text{Poisson}(\lambda_{\cdot j})$	where $\lambda_{\cdot j} = \sum_i \lambda_{ij}$
$n_{\cdot\cdot} \sim \text{Poisson}(\lambda_{\cdot\cdot})$	where $\lambda_{\cdot\cdot} = \sum_{i,j} \lambda_{ij}$
$\lambda_{ij} = \frac{\lambda_{i\cdot} \lambda_{\cdot j}}{\lambda_{\cdot\cdot}}$	

$$\rightsquigarrow \ln(\lambda_{ij}) = -\ln(\lambda_{\cdot\cdot}) + \ln(\lambda_{i\cdot}) + \ln(\lambda_{\cdot j}) \equiv \mu + \alpha_i + \beta_j$$

Looks like a two-way ANOVA without interaction!

Multinomial Modeling: Assumes $n_{\cdot\cdot}$ is fixed.

$\left(\{n_{ij}\}_{\forall i,j} \right) \sim \text{Multinomial} \left(n_{\cdot\cdot}, \{p_{ij}\}_{\forall i,j} \right)$

Assuming independence: $p_{ij} = p_{i\cdot} p_{\cdot j}$

$$\rightsquigarrow \lambda_{ij} = E(n_{ij}) = n_{\cdot\cdot} p_{i\cdot} p_{\cdot j}$$

Looks like a two-way ANOVA without interaction!

Poisson Modeling: $n_{ij} \sim \text{Poisson}(\lambda_{ij})$

Independence: $\ln(\lambda_{ij}) = \mu + \alpha_i + \beta_j$

Saturated model: $\ln(\lambda_{ij}) = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij}$

Is there an in-between model?

Linear-by-linear association: $\ln(\lambda_{ij}) = \mu + \alpha_i + \beta_j + \gamma u_i v_j$

u_i, v_j represent "scores" e.g., $u_i = i - 2$, $v_j = j - 2.5$

Agreement: $\ln(\lambda_{ij}) = \mu + \alpha_i + \beta_j + \gamma \mathbb{I}(i=j)$

Example: Results of rating the same 236 units by two different raters on an ordinal scale from 1 to 5.

		Rater 1					Total
		1	2	3	4	5	
Rater 2	1	10	6	4	2	2	24
	2	12	20	16	7	2	57
	3	1	12	30	20	6	69
	4	4	5	10	25	12	56
	5	1	3	3	8	15	30
	Total	28	46	63	62	37	236

ST704, Sujit K. Ghosh
Generalized Linear Mixed Models

Introduction

Example

Generalized Linear Mixed Model (GLMM)

$$\text{LM: } E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$$

$$\text{GLM: } E(\mathbf{Y}) = g^{-1}(\mathbf{X}\boldsymbol{\beta})$$

$$\text{LMM: } E(\mathbf{Y}|\boldsymbol{\alpha}) = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha}$$

$$\text{GLMM: } E(\mathbf{Y}|\boldsymbol{\alpha}) = g^{-1}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\alpha})$$

- \mathbf{Y} is $n \times 1$, observed
- \mathbf{X} is $n \times p$, known covariate matrix
- $\boldsymbol{\beta}$ is $p \times 1$, unknown fixed effects
- \mathbf{Z} is $n \times q$, known covariate matrix
- $\boldsymbol{\alpha}$ is $q \times 1$, random effects, where $\boldsymbol{\alpha} \sim N(\mathbf{0}, \mathbf{G})$
- $g(\cdot)$ is link function and $g^{-1}(\cdot)$ is inverse link function
- GLMM: $\text{Var}(\mathbf{Y}|\boldsymbol{\alpha}) = \mathbf{A}^{1/2} \mathbf{R} \mathbf{A}^{1/2}$
 - $\mathbf{R} = \psi \mathbf{I}$ as default
 - $\mathbf{A} = \text{diag}(w_1 h(\mu_1), \dots, w_n h(\mu_n))$

Slight abuse of notation in that functions $g(\cdot)$ and $g^{-1}(\cdot)$ apply to scalars, not vectors.

Standardized Mortality Ratio

Standardized Mortality Ratio (SMR) is a ratio between the observed number of deaths in an study population and the number of deaths would be expected, based on the age- and sex-specific rates in a standard population and the population size of the study population by the same age/sex groups. If the ratio of observed:expected deaths is greater than 1.0, there is said to be "excess deaths" in the study population.

... The SMR is used to compare the mortality risk of an study population to that of a standard population. It is especially applicable where the two populations have dissimilar age distributions, and in cases where direct age standardization may not be appropriate because the study population is small, or when lack of age-specific mortality rates precludes calculation of directly-age-standardized mortality rates. ¹

Y_i : # deaths in region i

E_i : expected # deaths in region i , according to age & sex death rates

SMR_i : $\frac{Y_i}{E_i}$

¹https://ibis.doh.nm.gov/resource/SMR_ISR.html

Source: SAS glimmix "Example 44.3 Smoothing Disease Rates"

Lip cancer in 56 counties of Scotland, 1975-1980

X_i is % of employees in agriculture, fishing, forestry.

May be a surrogate for exposure to sunlight.

Does X_i help explain variability in SMR_i across counties?

Belief: $Y_i \sim Poisson(\lambda_i)$, so $SMR_i \equiv \frac{Y_i}{E_i}$ has mean $\frac{\lambda_i}{E_i}$

Model 1: $\ln\left(\frac{\lambda_i}{E_i}\right) = \beta_0 + \beta_1 x_i \rightsquigarrow \ln(\lambda_i) = \ln(E_i) + \beta_0 + \beta_1 x_i$

Model 2: $\ln\left(\frac{\lambda_i}{E_i}\right) = \beta_0 + \beta_1 x_i + \alpha_i$ random $\rightsquigarrow \ln(\lambda_i) = \ln(E_i) + \beta_0 + \beta_1 x_i + \alpha_i$

```
proc glimmix data=lipcancer plots=(studentpanel residualpanel);
  class county;
  loge = log(expected);
  model observed = employment / dist=poisson offset=log
    solution cl ddfm=none;
  random county;
  covtest zerog / cl(type=profile);
  covtest indep;
  SMR = observed/expected;
  SMR_pred = exp(_zgamma_ + _xbeta_);
  id county employment SMR SMR_pred;
  output out=glimmixout;
proc sgplot data=glimmixout;
  reg x=smr y=smr_pred / datalabel=county clm curvelabel;
run;
```

R function to use: glmer in package lme4

Description

Fit a generalized linear mixed-effects model (GLMM). Both fixed effects and random effects are specified via the model formula.

Usage

```
glmer(formula, data = NULL, family = gaussian, control = glmerControl(),
       start = NULL, verbose = 0L, nAGQ = 1L, subset, weights, na.action,
       offset, contrasts = NULL, mustart, etastart,
       devFunOnly = FALSE, ...)
```

Arguments

formula: a two-sided linear formula object describing both the fixed-effect and random effects part of the model, with the response on the left of a ~ operator and the terms, separated by + operators, on the right.

Random-effects terms are distinguished by vertical bars ("|") separating expressions for design matrices from grouping factors.

Model Design

Simple Linear Model

Model: $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad i=1, \dots, n \quad \varepsilon_i \sim N(0, \sigma^2)$

Estimators: $\hat{\beta}_0 = \bar{Y} - \bar{X}\hat{\beta}_1$

$$\hat{\beta}_1 = r_{xy} \frac{s_y}{s_x} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{s_x s_y}$$

$$s_x = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Source	df	SS	MS	E(MS)	F
Regression	1	$\sum (Y_i - \hat{Y}_i)^2$	SSR	$\sigma^2 + \beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2$	$MSR/MSE \sim F_{1, n-2} \left(\frac{\beta_1^2 \sum_{i=1}^n (x_i - \bar{x})^2}{2 \sigma^2} \right)$
Error	$n-2$	$\sum (Y_i - \hat{Y}_i)^2$	$SSE/(n-2)$	σ^2	
Total	$n-1$	$\sum (Y_i - \bar{Y})^2$			

$H_0: \beta_1 = 0$

$$T = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)} \stackrel{H_0}{\sim} t_{n-2}$$

Multiple Linear Regression

Model: $Y_{ij} = \beta_0 + \beta_1 X_{1j} + \beta_2 X_{2j} + \dots + \beta_p X_{pj} + \varepsilon_{ij} \quad i=1, \dots, p \quad j=1, \dots, n \quad \varepsilon_{ij} \sim N(0, \sigma^2)$

Estimators: $\hat{\beta} = (X'X)^{-1}X'Y$

Source	df	SS	E(MS)	F
Regression	p	$\sum \hat{\varepsilon}_{ij} (\bar{Y}_{..} - \hat{Y}_{ij})^2$	$\sigma^2 + \frac{1}{p} \beta^T X^T X \beta$	$MSR/MSE \sim F_{p, n-p-1} \left(\frac{\beta^T X^T X \beta}{2 \sigma^2} \right)$
Error	$n-(p+1)$	$\sum \hat{\varepsilon}_{ij} (\hat{Y}_{ij} - \bar{Y}_{..})^2$	σ^2	
Total	$n-1$	$\sum \hat{\varepsilon}_{ij} (\bar{Y}_{..} - \bar{Y}_{..})^2$		or $\frac{1}{p} \beta^T X^T (I - P_{\beta}) X \beta$ if non-centered

Cell-Means

Models: $Y_{ij} = M_i + \varepsilon_{ij} \quad i=1, \dots, I \quad j=1, \dots, n_i \quad \varepsilon_{ij} \sim N(0, \sigma^2)$

Estimators: $\hat{M}_i = \bar{Y}_{..}$

Source	df	SS	E(MS)	F
Treat	$I-1$	$\sum \hat{\varepsilon}_{ij} (\bar{Y}_{..} - \bar{Y}_{..})^2 = Y^T (P_X - P_{\beta}) Y$	$\sigma^2 + \frac{1}{I-1} \sum_{i=1}^I n_i (\bar{Y}_{ii} - \bar{M})^2$	$MS_{Treat}/MSE \sim F_{I-1, I(I-1)} \left(\frac{\sum_{i=1}^I n_i (\bar{Y}_{ii} - \bar{M})^2}{2 \sigma^2} \right)$
Error	$I(I-1)$	$\sum \hat{\varepsilon}_{ij} (\hat{Y}_{ij} - \bar{Y}_{..})^2 = Y^T (I - P_{\beta}) Y$	σ^2	"
Total	$IJ-1$	$\sum \hat{\varepsilon}_{ij} (\bar{Y}_{..} - \bar{Y}_{..})^2 = Y^T (I - P_{\beta}) Y$		$\frac{1}{I-1} M^T X^T (I - P_{\beta}) X M$

$H_0: M_1 = \dots = M_I$

$$F = \frac{MS_{Treat}}{MSE}$$

$H_0: M_i = 0$

$$T = \frac{\hat{M}_i}{SE(\hat{M}_i)} = \sqrt{\frac{Y_{..}}{MSE/n_i}}$$

Factorial effect Model (ANOVA)

① Fixed effect one-way (unbalanced)

Models: $Y_{ij} = \mu + d_i + \varepsilon_{ij}$ $i=1, \dots, I$ $j=1, \dots, n_i$ $\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$

notice only pure error

in this case

$\sum Y_{ij} = \bar{Y}_i$.

Models: $Y_{ij} = \mu + d_i + \varepsilon_{ij}$ $i=1, \dots, I$ $j=1, \dots, \frac{N}{I}$ $\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$ Completely randomized design (CRD)

Estimates: $\hat{\mu} + \hat{d}_i = \bar{Y}_i$. $\sum c_i \hat{d}_i = \sum c_i \bar{Y}_i$. s.t. $\sum c_i = 0$

$$\hat{\mu} + \bar{d}_i = \bar{Y}_i$$

Source	df	SS	E(MS)	F
Mean	1	$N\bar{Y}_i^2$	$\gamma^T P_1 Y$	$\frac{MSM}{MSE} \sim F_{1, N-I} \left(\frac{N(\mu + \bar{d})^2}{2\sigma^2} \right)$
Trt	$I-1$	$\sum_j (\bar{Y}_i - \bar{Y}_.)^2$	$\gamma^T (P_X - P_2) Y$	$\frac{MST}{MSE} \sim F_{I-1, N-I} \left(\frac{1}{2\sigma^2} Q^T X^T (I - P_2) X Q \right)$
Error	$N-I$	$\sum_{ij} (Y_{ij} - \bar{Y}_i)^2$	$\gamma^T (I - P_X) Y$	\downarrow
Total	$N-1$	$\sum_j (Y_{ij} - \bar{Y}_.)^2$		$\frac{1}{I-1} Q^T X^T (I - P_2) X Q$

$$H_0: d_1 = \dots = d_I$$

$$F = \frac{MST_{Trt}}{MSE}$$

$$H_0: d_i - d_j = 0$$

$$T = \frac{e' \hat{\theta}}{SE(e' \hat{\theta})} = \frac{e' \hat{\theta}}{\sqrt{MSE} \sqrt{h} h} \sim t_{n-\text{rank}(X)} (N-I)$$

with $e' \hat{\theta} = h Y$

$$F = \frac{e' \hat{\theta}^2 / h h}{MSE} \sim F_{1, n-\text{rank}(X)}$$

$$H_0: \sum_{i=1}^I c_i d_i = 0$$

$$T = \frac{e' \hat{\theta}}{\sqrt{MSE} \sqrt{\sum c_i^2}} = \frac{\sum c_i \bar{Y}_i}{\sqrt{MSE} \sqrt{\sum c_i^2}} \sim t_{n-\text{rank}(X)} = df \text{ error}$$

② Random effect one-way (balanced)

Models: $Y_{ij} = \mu + A_i + \varepsilon_{ij}$ $i=1, \dots, a$ $j=1, \dots, \tilde{n}$ $\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$ } mutually independent
 Estimates: $\hat{\mu} = \bar{Y}_..$ $A_i \stackrel{iid}{\sim} N(0, \sigma_A^2)$ } independent

Source	df	SS	E(MS)	F
Mean	1	$N\bar{Y}_{..}^2 = Y^T P_1 Y$		
A	$a-1$	$\sum_i \bar{Y}_j (\bar{Y}_{..} - \bar{Y}_i)^2 = Y^T (P_X - P_2) Y$	$\sigma^2 + \tilde{n} \sigma_A^2$	$\frac{MSA}{MSE} \left(\frac{\sigma^2}{\sigma^2 + \tilde{n} \sigma_A^2} \right) \sim F_{a-1, a(n-1)}$
Error	$a(\tilde{n}-1)$	$\sum_i \bar{Y}_j (Y_{ij} - \bar{Y}_i)^2 = Y^T (I - P_X) Y$	σ^2	
Total	$n-1$	$\sum_{ij} (Y_{ij} - \bar{Y}_{..})^2 = Y^T (I - P_2) Y$		

$$H_0: \sigma_A^2 = 0$$

$$F = \frac{MSA}{MSE} \stackrel{H_0}{\sim} F_{a-1, a(n-1)}$$

③ Fixed effect (two way full factorial: crossed) ← simple contrasts are NOT estimable

Models: $Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \varepsilon_{ijk}$

$$i=1, \dots, a \quad j=1, \dots, b \quad k=1, \dots, n \quad \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$$

Estimates:

$$\begin{cases} \mu + \bar{\alpha}_i + \bar{\beta}_j + (\bar{\alpha}\bar{\beta})_{ij} = \bar{Y}_{..} \\ \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} = \bar{Y}_{ij} \\ \mu + \alpha_i + \bar{\beta}_j + (\bar{\alpha}\bar{\beta})_{ij} = \bar{Y}_{i..} \\ \mu + \bar{\alpha}_i + \beta_j + (\bar{\alpha}\beta)_{ij} = \bar{Y}_{.j} \end{cases}$$

$$\sum_i C_i (\alpha_i + (\alpha\beta)_{ij}) = \sum_i C_i \bar{Y}_{ij}, \text{ simple effect of } \alpha \text{ at } j$$

$$\sum_i C_i (\alpha_i + (\bar{\alpha}\bar{\beta})_{ij}) = \sum_i C_i \bar{Y}_{i..}, \text{ main effect of } \alpha$$

$$\sum_j d_j (\beta_j + (\alpha\beta)_{ij}) = \sum_j d_j \bar{Y}_{ij}, \text{ simple effect of } \beta \text{ at } i$$

$$\sum_j d_j (\beta_j + (\bar{\alpha}\bar{\beta})_{ij}) = \sum_j d_j \bar{Y}_{.j}, \text{ main effect of } \beta$$

$$\sum_i d_i \sum_j d_j (\alpha\beta)_{ij} = \sum_i d_i \sum_j d_j \bar{Y}_{ij}.$$

when balanced LS Mean = Mean

unbiased and BLUE

Interaction effect of $\alpha\beta$ →
(consistency of simple effect)

Models: $Y_{ijk} = \mu^* + \alpha_i^* + \beta_j^* + (\alpha\beta)_{ij}^* + \varepsilon_{ijk}$

$$i=1, \dots, a \quad j=1, \dots, b \quad k=1, \dots, n \quad \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$\mu^* = \mu - d\bar{\beta}..$$

$$\alpha_i^* = \alpha_i + d\bar{\beta}_i..$$

$$\beta_j^* = \beta_j + d\bar{\beta}..$$

$$\alpha\beta_{ij}^* = d\beta_{ij} - d\bar{\beta}_i.. - d\bar{\beta}_j.. + d\bar{\beta}..$$

Source	df	SS	MS	F
Mean	1	$N \bar{Y}_{..}^2$		
α	$a-1$	$\sum_{ijk} (\bar{Y}_{ij..} - \bar{Y}_{..})^2$	$\sigma^2 + b\tilde{n}Q(\alpha^*)$	$\frac{MS_{\alpha}}{MS_{\text{Error}}} \sim F((a-1) \frac{b\tilde{n}Q(\alpha^*)}{2\sigma^2})$
β	$b-1$	$\sum_{ijk} (\bar{Y}_{.j} - \bar{Y}_{..})^2$	$\sigma^2 + a\tilde{n}Q(\beta^*)$	$\frac{MS_{\beta}}{MS_{\text{Error}}} \sim F((b-1) \frac{a\tilde{n}Q(\beta^*)}{2\sigma^2})$
$\alpha\beta$	$(a-1)(b-1)$	$\sum_{ijk} (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j} + \bar{Y}_{..})^2$	$\sigma^2 + \tilde{n}Q(\alpha\beta^*)$	$\frac{MS_{\alpha\beta}}{MS_{\text{Error}}} \sim F((a-1)(b-1) \frac{\tilde{n}Q(\alpha\beta^*)}{2\sigma^2})$
Error	$n-ab$	$\sum_{ijk} (\bar{Y}_{ij..} - \bar{Y}_{..})^2$	σ^2	
Total	$n-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{..})^2$		

$$Q(\alpha^*) = \frac{1}{a-1} \sum_{i=1}^a (\alpha_i^* - \bar{\alpha}^*)^2$$

$$= \frac{1}{a-1} \sum_{i=1}^a (\alpha_i + d\bar{\beta}_i.. - \bar{\alpha} - d\bar{\beta}..)^2$$

$$H_0: \alpha_1^* = \dots = \alpha_a^*$$

$$\alpha_i^* = \alpha_i + d\bar{\beta}_i.. \leftarrow \bar{Y}_{i..} - \bar{Y}_{..} \Rightarrow \alpha_i^* - \bar{\alpha}.$$

$$H_0: \beta_1^* = \dots = \beta_b^*$$

$$\beta_j^* = \beta_j + d\bar{\beta}.. \leftarrow \bar{Y}_{.j} - \bar{Y}_{..} \Rightarrow \beta_j^* - \bar{\beta}.$$

$$H_0: \alpha\beta_{ij}^* = 0 \quad \forall i, j$$

$$\alpha\beta_{ij}^* = d\beta_{ij} - d\bar{\beta}_i.. - d\bar{\beta}_j.. + d\bar{\beta}.. \leftarrow \bar{Y}_{ij..} - \bar{Y}_{i..} - \bar{Y}_{.j} + \bar{Y}_{..} \Rightarrow \alpha\beta_{ij}^*$$

$$(d_i + \beta_j + d\bar{\beta}_{ij}) - (d_i + \bar{\beta}_i.. + d\bar{\beta}_i..) - (d_j + \beta_j + d\bar{\beta}_j..) + (d.. + \bar{\beta}.. + d\bar{\beta}..)$$

$$= d\beta_{ij} - d\bar{\beta}_i.. - d\bar{\beta}_j.. + d\bar{\beta}..$$

$$= d\beta_{ij}^*$$

④ Mixed two-way effect model (balanced, crossed)

Models: $Y_{ijk} = \mu + \alpha_i + \beta_j + \delta\beta_{ij} + \varepsilon_{ijk}$

$$\begin{array}{cccc} i=1, \dots, a & j=1, \dots, b & k=1, \dots, n \\ \end{array}$$

Estimates: $\hat{\mu} + \hat{\alpha}_i = \bar{Y}_{..}$

$$\hat{\mu} + \hat{\alpha}_i = \bar{Y}_{ii..}$$

$$\left. \begin{array}{l} \delta\beta_{ij} \stackrel{iid}{\sim} N(0, \sigma_{\delta\beta}^2) \\ \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2) \end{array} \right\} \text{mutually independent}$$

Source	df	SS	E(MS)	F
α	$a-1$	$\sum_{ijk} (\bar{Y}_{...} - \bar{Y}_{i..})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + b\tilde{r}\tilde{c}\sigma(\alpha)$	$\frac{MSA}{MSAB} \sim F_{a-1, (a-1)(b-1)} \left(\frac{(a-1)b\tilde{r}\tilde{c}\sigma(\alpha)}{2(\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2)} \right)$
B	$b-1$	$\sum_{ijk} (\bar{Y}_{...} - \bar{Y}_{j..})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + a\tilde{r}\tilde{c}\sigma_\beta^2$	$\frac{MSB}{MSAB} \left(\frac{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + a\tilde{r}\tilde{c}\sigma_\beta^2}{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2} \right) \sim F_{b-1, (a-1)(b-1)}$
$\delta\beta$	$(a-1)(b-1)$	$\sum_{ijk} (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{j..} + \bar{Y}_{...})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2$	$\frac{MSAB}{MSE} \left(\frac{\sigma^2}{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2} \right) \sim F_{(a-1)(b-1), n-ab}$
Error	$n-ab$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{ij.})^2$	σ^2	
Total	$n-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{...})^2$		

$H_0: \alpha_1 = \dots = \alpha_a \quad MSA / MSAB \xrightarrow{H_0} F_{(a-1), (a-1)(b-1)}$
 $H_0: \sigma_B^2 = 0 \quad MSB / MSAB \xrightarrow{H_0} F_{(b-1), (a-1)(b-1)}$
 $H_0: \sigma_{\alpha\beta}^2 = 0 \quad MSAB / MSE \xrightarrow{H_0} F_{(a-1)(b-1), n-ab}$

⑤ Random two-way effect model (balanced, crossed)

Models: $Y_{ijk} = \mu + A_i + B_j + AB_{ij} + \varepsilon_{ijk}$

$$\begin{array}{cccc} i=1, \dots, a & j=1, \dots, b & k=1, \dots, \tilde{n} & \end{array}$$

$$A_i \stackrel{iid}{\sim} N(0, \sigma_A^2) \quad AB_{ij} \stackrel{iid}{\sim} N(0, \sigma_{AB}^2) \quad \left. \begin{array}{l} B_j \stackrel{iid}{\sim} N(0, \sigma_B^2) \\ \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2) \end{array} \right\} \text{independent}$$

Source	df	SS	E(MS)	F
α	$a-1$	$\sum_{ijk} (\bar{Y}_{...} - \bar{Y}_{i..})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + b\tilde{r}\tilde{c}\sigma_A^2$	$\frac{MSA}{MSAB} \left(\frac{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2}{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + b\tilde{r}\tilde{c}\sigma_A^2} \right) \sim F_{a-1, (a-1)(b-1)}$
B	$b-1$	$\sum_{ijk} (\bar{Y}_{...} - \bar{Y}_{j..})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + a\tilde{r}\tilde{c}\sigma_\beta^2$	$\frac{MSB}{MSAB} \left(\frac{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2}{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2 + a\tilde{r}\tilde{c}\sigma_\beta^2} \right) \sim F_{b-1, (a-1)(b-1)}$
$\delta\beta$	$(a-1)(b-1)$	$\sum_{ijk} (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{j..} + \bar{Y}_{...})^2$	$\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2$	$\frac{MSAB}{MSE} \left(\frac{\sigma^2}{\sigma^2 + \tilde{r}\tilde{c}\sigma_{\alpha\beta}^2} \right) \sim F_{(a-1)(b-1), n-ab}$
Error	$n-ab$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{ij.})^2$	σ^2	
Total	$n-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{...})^2$		

$H_0: \sigma_A^2 = 0 \quad MSA / MSAB \xrightarrow{H_0} F_{(a-1), (a-1)(b-1)}$
 $H_0: \sigma_B^2 = 0 \quad MSB / MSAB \xrightarrow{H_0} F_{(b-1), (a-1)(b-1)}$
 $H_0: \sigma_{AB}^2 = 0 \quad MSAB / MSE \xrightarrow{H_0} F_{(a-1)(b-1), n-ab}$

⑥ Fixed effect (two way full factorial: nested)

Models: $Y_{ijk} = \mu + \alpha_i + \beta_{j(i)} + \varepsilon_{ijk}$ $\varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$

$$i=1, \dots, a \quad j=1, \dots, b_i \quad k=1, \dots, n_i$$

Estimates: A main: $\frac{1}{b_i} \sum_{j=1}^{b_i} \sum_{i=1}^a c_i (\alpha_i + \beta_{j(i)}) = \sum_{i=1}^a c_i (\alpha_i + \bar{\beta}_{c(i)})$

B simple for given i: $\sum_{j=1}^{b_i} d_j (\alpha_i + \beta_{j(i)}) = \sum_{j=1}^{b_i} \beta_{j(i)}$

There are $b_i - 1$ linearly independent B simple effects at i (from d_j independent contrasts)

$$B_1', \dots, B_{b_i-1}'$$

Source	df	SS	F(MS)	F
α	$a-1$	$\sum_{jk} (Y_{...} - \bar{Y}_{..})^2$	$\sigma^2 + \frac{1}{a-1} \sum_i \sum_j (\alpha_i^* - \bar{\alpha}^*)^2$	$\frac{MS_{\alpha}}{MSE} \sim F \left(\frac{\sum_i (\alpha_i^* - \bar{\alpha}^*)^2}{2\sigma^2} \right)$
$\beta(\alpha)$	$\sum_{i=1}^a (b_i - 1)$	$\sum_{jk} (\bar{Y}_{..} - \bar{Y}_{ij.})^2$	$\sigma^2 + \frac{1}{\sum(b_i-1)} \sum_i \sum_j (\beta_{j(i)} - \bar{\beta}_{c(i)})^2$	$\frac{MS_{\beta(\alpha)}}{MSE} \sim F \left(\frac{\sum_i (\beta_{j(i)} - \bar{\beta}_{c(i)})^2}{2\sigma^2} \right)$
Error	$n - \sum_i b_i$	$\sum_{jk} (Y_{ijk} - \bar{Y}_{ij.})^2$	σ^2	
Total	$n-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{...})^2$		

$$H_0: A' = \dots = A^{a-1} = 0 \quad \text{main effect } A$$

$$H_0: B_1' = \dots = B_{b_i-1}' = 0 \quad \text{overall test for factor } B.$$

$$B_2' = \dots = B_{b_i-1}' = 0$$

$$B_a' = \dots = B_{b_i-1}' = 0$$

⑦ Mixed effect (nested, balanced)

Models: $Y_{ijk} = \mu + \alpha_i + \beta_{j(i)} + \varepsilon_{ijk}$ $\varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$

$i=1, \dots, a \quad j=1, \dots, b \quad k=1, \dots, n_i \quad \beta_{j(i)} \stackrel{iid}{\sim} N(0, \sigma_{B(\alpha)}^2)$

} mutually

} independent

$$\text{Estimates: } M + \bar{\alpha}_i = \bar{Y}_{...}$$

$$M + \hat{\alpha}_i = \bar{Y}_{i..}$$

Source	df	SS	F(MS)	F
α	$a-1$	$\sum_{jk} (Y_{...} - \bar{Y}_{..})^2$	$\sigma^2 + \tilde{\alpha} \sigma_{B(\alpha)}^2 + b \tilde{\alpha} \sigma_{B(\alpha)}^2$	$\frac{MS_{\alpha}}{MS_{B(\alpha)}} \sim F_{a-1, ab-1} \left(\frac{(a-1)b\tilde{\alpha} \sigma_{B(\alpha)}^2}{2(\sigma^2 + \tilde{\alpha} \sigma_{B(\alpha)}^2)} \right)$
$\beta(\alpha)$	$a(b-1)$	$\sum_{jk} (\bar{Y}_{..} - \bar{Y}_{ij.})^2$	$\sigma^2 + \tilde{\alpha} \sigma_{B(\alpha)}^2$	$\frac{MS_{\beta(\alpha)}}{MSE} \sim F_{ab-1, n-ab} \left(\frac{\sigma^2}{\sigma^2 + \tilde{\alpha} \sigma_{B(\alpha)}^2} \right)$
Error	$n - ab$	$\sum_{jk} (Y_{ijk} - \bar{Y}_{ij.})^2$	σ^2	
Total	$n-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{...})^2$		

$$H_0: \sum_i C_i \alpha_i = 0 \quad \text{main effect of } A \text{ with contrast } C$$

$$H_0: \sigma_{B(\alpha)}^2 = 0$$

$$H_0: \alpha_1 = \dots = \alpha_a$$

⑧ Random effect (nested)

Models: $Y_{ijk} = \mu + A_i + B_{j(i)} + \varepsilon_{ijk}$

$i=1, \dots, a \quad j=1, \dots, b \quad k=1, \dots, n$

$\varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$
 $B_{j(i)} \stackrel{iid}{\sim} N(0, \sigma_{B(A)}^2)$
 $A_i \stackrel{iid}{\sim} N(0, \sigma_A^2)$

mutually independent

Source	df	SS	E(MS)	F
A	a-1	$\sum_k (Y_{..} - \bar{Y}_{..})^2$	$\sigma^2 + \tilde{n}\sigma_{B(A)}^2 + b\tilde{n}\sigma_A^2$	$\frac{MSA}{MSB(A)} \left(\frac{\sigma^2 + \tilde{n}\sigma_{B(A)}^2}{\sigma^2 + \tilde{n}\sigma_{B(A)}^2 + b\tilde{n}\sigma_A^2} \right) \sim F_{a-1, ab-n}$
B(A)	a(b-1)	$\sum_k (Y_{ij..} - \bar{Y}_{ij..})^2$	$\sigma^2 + \tilde{n}\sigma_{B(A)}^2$	$\frac{MSB(A)}{MSE} \left(\frac{\sigma^2}{\sigma^2 + \tilde{n}\sigma_{B(A)}^2} \right) \sim F_{ab-1, n-ab}$
Error	n-ab	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{ijk})^2$	σ^2	
Total	n-1	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{..})^2$		

$$H_0: \sigma_A^2 = 0$$

$$H_0: \sigma_{B(A)}^2 = 0$$

ANCOVA

Models: $Y_{ij} = \mu + \alpha_i + \beta X_{ij} + \varepsilon_{ij} \quad i=1, \dots, k \quad j=1, \dots, n_i \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$

Estimates: $\hat{\beta}$ by SLR

$$\hat{\mu} + \hat{\alpha}_i = \bar{Y}_i - \hat{\beta} \bar{X}_i.$$

$$\Sigma_i c_i \hat{\alpha}_i = \Sigma_i c_i \bar{Y}_i - \hat{\beta} \Sigma_i c_i \bar{X}_i.$$

$$\hat{\mu} + \hat{\alpha}_i + \hat{\beta} \bar{X}_{..} = \bar{Y}_i$$

$$\hat{\mu} + \hat{\alpha}_i + \hat{\beta} (\bar{X}_i - \bar{X}_{..}) = \bar{Y}_i$$

Source	df	SS	E(MS)	F
α	k-1	$\sum_j (\bar{Y}_i - \bar{Y}_{..})^2$	$\sigma^2 + \frac{1}{k-1} \sum_j (\alpha_i - \bar{\alpha}_{..})^2$	
β	1	$\sum_j (\bar{Y}_i - \bar{Y}_{..})^2$?	
Error	N-k-1	$\sum_{ij} (Y_{ij} - \bar{Y}_{ij})^2$	σ^2	
Total	N-1	$\sum_{ij} (Y_{ij} - \bar{Y}_{..})^2$		

$$H_0: \beta = 0 \text{ and } \alpha_1 = \dots = \alpha_k \text{ Global F test}$$

$$H_0: \alpha_1 = \dots = \alpha_k \text{ treatment test}$$

$$H_0: \beta = 0 \text{ covariate test}$$

Two-way effect model

① Two-way effect model (crossed)

Models: $Y_{ijk} = \mu + \alpha_{ij} + \epsilon_{ijk}$

$$i=1, \dots, a \quad j=1, \dots, b \quad k=1, \dots, n_{ij} \quad \epsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$$

if $n_{ij} = \bar{n}$, balanced \Rightarrow complete

Estimates: $\hat{\mu} + \alpha_{ij} = \bar{Y}_{ij}$.

$$\hat{\mu} + \bar{\alpha}_{..} = \bar{Y}_{..}$$

$$\mu + \bar{\alpha}_{..} = \bar{Y}_{..}$$

$$\mu + \bar{\alpha}_{.j} = \bar{Y}_{.j}$$

$$\hat{\alpha}_j = \sum_{i=1}^a C_i \alpha_{ij} \Rightarrow \hat{\alpha}_j = \sum_{i=1}^a C_i \bar{Y}_{ij}$$

$$\hat{\alpha}_{..} = \frac{1}{b} \sum_{j=1}^b \hat{\alpha}_j \Rightarrow \hat{\alpha}_{..} = \frac{1}{b} \sum_{j=1}^b \sum_{i=1}^a C_i \bar{Y}_{ij}$$

$$AB = \sum_{j=1}^b \alpha_j A_j \Rightarrow \hat{AB} = \sum_j C_j \hat{\alpha}_j \bar{Y}_{ij}$$

when $a > 2$ or $b > 2$, more than one contrasts available.

Contrast table: (2×3 factorial)

A	1	2
C ₁	1	-1
C ₂	1	0

B	1	2	3
d ₁₁	1	-1	0
d ₁₂	1	0	-1

	d ₁₁₁	d ₁₁₂	d ₁₁₃	d ₁₂₁	d ₁₂₂	d ₁₂₃	fraction
A ¹	1	1	1	-1	-1	-1	1/3
B ¹	1	-1	0	1	-1	0	1/2
B ²	1	0	-1	1	0	-1	1/2
A'B ¹	1	-1	0	-1	1	0	
A'B ²	1	0	-1	-1	0	1	

Source	df	SS	F(MS)	F
Model	ab-1	$\sum_{ijk} (\bar{Y}_{..} - \bar{Y}_{ij})^2$	$\sigma^2 + \bar{n} Q(\alpha)$	$\frac{MSM}{MSE} \sim F_{ab-1, n-ab} ((ab-1) \frac{\bar{n} Q(\alpha)}{2 \sigma^2})$
Error	n-ab	$\sum_{ijk} (\bar{Y}_{..} - Y_{ijk})^2$	σ^2	
Total	n-1	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{..})^2$		
$H_0: A^1 = 0$	No A main effect	$df = a-1 = 1$		
$H_0: B^1 = B^2 = 0$	No B main effect	$df = b-1 = 2$		
$H_0: A'B^1 = A'B^2 = 0$	No AB interaction effect	$df = (a-1)(b-1) = 2$		

② Two-way effect model (nested)

Models: $Y_{ijk} = \mu + \gamma_j(c_{ij}) + \epsilon_{ijk}$ $\epsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma^2)$

$$i=1, \dots, a \quad j(i) = 1, \dots, b_i \quad k=1, \dots, n_{ij}$$

Estimates: $\sum_i C_i \bar{\gamma}_j(c_{ij}) = \sum_i C_i \bar{Y}_{i..}$ A main effect since nested, B main effect

$$\sum_j d_j \bar{\gamma}_j(c_{ij}) = \sum_j d_j \bar{Y}_{ij} \quad B \text{ simple effect} \leftarrow \text{doesn't make sense.}$$

Block Design :

Models: $Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$ $i=1, \dots, I$ $j=1, \dots, \frac{N}{I}$ $\epsilon_{ij} \sim N(0, \sigma^2)$ Completely randomized design (CRD)

Models: $Y_{hi} = \mu + \beta_h + \tau_i + \epsilon_{hi}$ $h=1, \dots, b$ $i=1, \dots, t$ $\epsilon_{hi} \sim N(0, \sigma^2)$ Randomized complete block design (RCBD)

↑ block effect ↑ treatment effect

no replicates per combination
thus no interaction

↓ treatment applied on all blocks

Estimators: $\hat{\Sigma}_i C_i \bar{\tau}_i = \Sigma_i C_i \bar{Y}_i$ (BLUE)

Source	df	SS	MS	F
Block	$b-1$	$\sum_{hi} (\bar{Y}_{..} - \bar{Y}_{hi})^2$	$\sigma^2 + t\bar{Q}(\beta)$	Typically not interested in this
Treatment	$t-1$	$\sum_{hi} (\bar{Y}_{..} - \bar{Y}_{..i})^2$	$\sigma^2 + b\bar{Q}(\tau)$	$MSE \sim F_{t-1, (b-1)(t-1)} \frac{b(\bar{Q}(\tau))}{2\sigma^2}$
Error	$(t-1)(b-1)$	$\sum_{hi} (\bar{Y}_{..} - \bar{Y}_{hi} - \bar{Y}_{..i} + \bar{Y}_{hi})^2$	σ^2	
Total	$tb-1$	$\sum_{hi} (\bar{Y}_{hi} - \bar{Y}_{..})^2$		only in the blocking case, SSE isn't in easy form.

Models: $Y_{hi} = \mu + \beta_h + \tau_i + \epsilon_{hi}$ $\epsilon_{hi} \sim N(0, \sigma^2)$ Incomplete block design (IBD)

$h=1, \dots, b$ $i=1, \dots, t$ when $n_{hi} = 1$

Estimators: Relies on software to construct BLUE for $\Sigma_i C_i \bar{\tau}_i$

- Problems:
1. May not be able to estimate all treatment contrast.
 2. May have pairwise contrasts with different variances.

Models: Balanced incomplete block design (BIBD)

Rules: b : block number

k : experimental units in one block

t : treatment number

Each treatment occurs in exactly $r = \frac{bk}{t}$ blocks

Each pair of treatment appears in $\lambda = \frac{r(c-1)}{t-1}$ blocks

df same as RCBD

Models: $Y_{ijk} = \mu + \beta_i + \gamma_j + \tau_k + \epsilon_{ijk}$ $\epsilon_{ijk} \sim N(0, \sigma^2)$ Latin Square Designs (LSD)

$i=1, \dots, k$ $j=1, \dots, k$ $k=1, \dots, t$ where $n_{ijk} = 1$ $\sum_i n_{ijk} = \sum_j n_{ijk} = 1$; $\sum_{ij} n_{ijk} = k$

Estimators:

$$\hat{\mu} + \hat{\beta}_i + \hat{\gamma}_j + \hat{\tau}_k = \bar{Y}_{..k}$$

$$\sum_k C_k \tau_k = \sum_k C_k \bar{Y}_{..k}$$

Thus all pairwise contrasts

have equal variance

Split-plot design (CRD-RCBD)

Models: Whole plot: $Y_{ijk} = \mu + \alpha_i + W_{j(i)} + \varepsilon_{ijk}$

$\begin{cases} W_{j(i)} \stackrel{iid}{\sim} N(0, \sigma_w^2) \\ \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma_s^2) \end{cases}$ } mutually independent

$i=1, \dots, t \quad j=1, \dots, r \quad k=1, \dots, s$

WP main effect

Sub-plot: $Y_{ijk} = \mu + W_{j(i)} + \beta_k + \varepsilon_{ijk}$

$\begin{cases} W_{j(i)} \stackrel{iid}{\sim} N(0, \sigma_w^2) \\ \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma_s^2) \end{cases}$ } mutually independent

$i=1, \dots, t \quad j=1, \dots, r \quad k=1, \dots, s$

SP main effect

Full model: $Y_{ijk} = \mu + \alpha_i + W_{j(i)} + \beta_k + \alpha\beta_{ik} + \varepsilon_{ijk}$

$i=1, \dots, t \quad j=1, \dots, r \quad k=1, \dots, s$

$\begin{cases} W_{j(i)} \stackrel{iid}{\sim} N(0, \sigma_w^2) \\ \varepsilon_{ijk} \stackrel{iid}{\sim} N(0, \sigma_s^2) \end{cases}$ } mutually independent

Source	df	SS	MS	F
α	$t-1$	$\sum_{jk} (\bar{Y}_{..} - \bar{Y}_{..})^2$	$\sigma_s^2 + S\sigma_w^2 + rsQ(\alpha^*)$	MSA/MSW
Error	$t(r-1)$	$\sum_{ijk} (\bar{Y}_{ij..} - \bar{Y}_{..})^2$	$\sigma_s^2 + S\sigma_w^2$	MSW/MSE
β	$s-1$	$\sum_{jk} (\bar{Y}_{..k} - \bar{Y}_{..})^2$	$\sigma_s^2 + trQ(\beta^*)$	MSB/MSE
$\alpha\beta$	$(s-1)(t-1)$	$\sum_{ijk} (\bar{Y}_{ijk} - \bar{Y}_{ij..} - \bar{Y}_{..k} + \bar{Y}_{..})^2$	$\sigma_s^2 + rQ(\alpha\beta^*)$	MSAB/MSE
Error	$t(r-1)(s-1)$	$\sum_{ijk} (\bar{Y}_{ijk} - \bar{Y}_{ij..} - \bar{Y}_{..k} + \bar{Y}_{..})^2$	σ_s^2	
Total	$trs-1$	$\sum_{ijk} (Y_{ijk} - \bar{Y}_{..})^2$		

Source	A	B	AB	W(AB)		
df	$a-1$	$b-1$	$(a-1)(b-1)$	$ab(r-1)$		
Source	C	AC	BC	A β C	Error	Total
df	$c-1$	$(a-1)(c-1)$	$(b-1)(c-1)$	$(a-1)(b-1)(c-1)$	$abc(r-1)(c-1)$	$abcr-1$

$$H_0: \sigma_w^2 = 0$$

$$H_0: \alpha_i^* = 0 / \beta_j^* = 0 / \alpha\beta_{ij}^* = 0 \leftarrow \text{be careful with this}$$

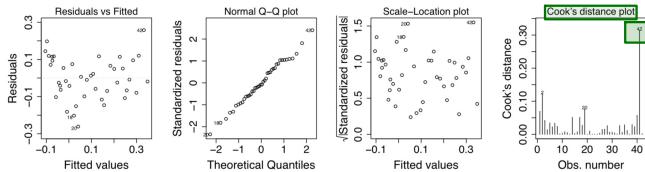
Effect Principles

1. Effect hierarchy: main effect and lower-order interactions more likely to be significant
2. Effect heredity: If interaction significant, keep parent terms
3. Effect sparsity: Few effects are significant
4. Factor sparsity: Few factors are significant

Model Diagnostics

Example Model Diagnostics

- Some default residual plots are produced by `plot(gw.mod1)`.



- There is a trend in the mean of the residuals, violating **independence**.
- The QQ plot is close to a straight line, so **normality** is OK.
- The residual magnitudes seem consistent with **constant variance**.
- The 42nd observation has a very high influence on the results.

Unusual data point

Definition:

$$\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y = PY \leftarrow \text{projection of } Y \text{ onto the column space of } X, \hat{Y} \in \text{col}(X)$$

P: hat matrix

$$\hookrightarrow \hat{Y}_i = P_{ii}Y_1 + \dots + P_{i,i-1}Y_{i-1} + P_{ii}Y_i + P_{i,i+1}Y_{i+1} + \dots + P_{in}Y_n$$

By examining the elements of P, get to know the influence of individual response values on fitted values.

P_{ii}: leverage of ith component: leverage ↑ influence on fit ↑

$$1. \frac{1}{n} \leq P_{ii} \leq 1, P_{ii} = \sum_j P_{ij}^2$$

$$\hookrightarrow \text{in SLR, } P_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum(x_i - \bar{x})^2}$$

centered design

matrix with

intercept in model

$$2. \text{ If } P_{ii} = 1, \text{ then rest of row (& column) equals 0: } \hat{Y}_i = Y_i$$

$$3. \text{ If } P_{ii} = 0, \text{ then rest of row (& column) equals 0: } \hat{Y}_i = 0$$

$$4. -1 \leq P_{ij} \leq 1$$

5. P is determined entirely by the design matrix X, and not at all by Y.

P properties:

1. P is symmetric even if $(X'X)^{-1}$ is not. \hookrightarrow orthogonal projection

2. P is invariant of $(X'X)^{-1}$ choice since P unique orthogonal projection, $(X'X)^{-1}$ varies

$$3. X = PX \Rightarrow X(X'X)^{-1}X'X = X \Rightarrow (X'X)^{-1}X' = X^{-1}$$

$$X'P = X' \Rightarrow X'(X'X)^{-1}X' = X' \Rightarrow X(X'X)^{-1} = (X')'$$

$$4. \text{rank}(P) = \text{trace}(P) = \text{rank}(X), \text{ where } \text{trace}(P) = \sum_i P_{ii}$$

$$5. \text{If } 1 \in \text{col}(X), \text{ then each row (column) of P sums to 1. } \sum_i P_{ij} = \sum_j P_{ij} = 1$$

\hookrightarrow with interaction

cook's distance

Potential issues in regression

1. Assumption violation

Assumptions:

- 1.1: X is known, without error

- 1.2: $E(\varepsilon) = 0$, i.e. each model is correctly specified

- 1.3: $\text{Var}(\varepsilon) = \sigma^2 I$, i.e. equal variance, uncorrelated

- 1.4: $\varepsilon \sim N(\cdot)$

2. Unusual Data

3. Computational instability

1.1. X observed with error (measurement error, errors-in-variables models)

- Estimators biased : usually towards 0
- Affects everything
- instrumental variable may help

1.2. Mean model is misspecified

- underspecification
- Model not additive
- Nonlinear relationship
- Diagnose : Plot e_i vs. \hat{y}_i should show random scatter.
 - (a). $\sum e_i = 0$ if model includes intercept
 - (b). $\sum e_i \hat{y}_i = 0$ (a) & (b) $\Rightarrow \text{Corr}(e, \hat{y}) = 0 \Rightarrow$ only implies no linear relationship
 - (c). $\sum e_i x_{ij} = 0$ (a) & (c) $\Rightarrow \text{Corr}(e, x_j) = 0$

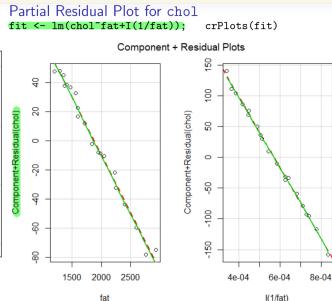
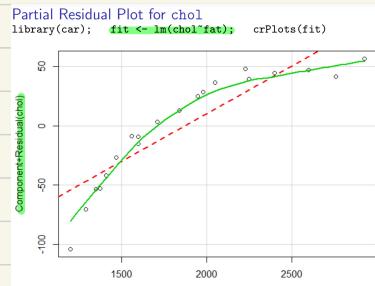
Plot e_i v.s. X_j should show random scatter.

Partial Residual (leverage) plots should show straight lines.

has more information than usual residual plot (when no much multicollinearity)
e.g. straight line says add X_j to model

Nonlinear relationship says to transform X_j before adding to model.

$$\begin{aligned}\text{Partial Residual for } X_j: e^* &= Y - (\beta_0 + \beta_1 x_1 + \dots + \beta_{j-1} x_{j-1} + \cancel{\beta_j x_j} + \beta_{j+1} x_{j+1} + \dots + \beta_p x_p) \\ &= Y - \hat{Y} + \hat{\beta}_j x_j \\ &= e + \hat{\beta}_j x_j\end{aligned}$$

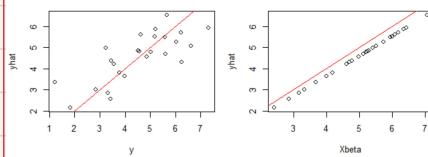
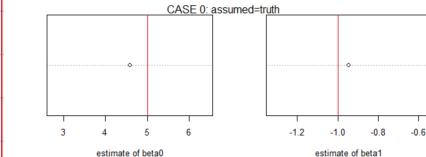


Simulation Study: Effect of Misspecifying the Mean Model Simulation Setup:

Truth: $\mathbf{Y} = \mathbf{X}\beta + \mathbf{Z}\gamma + \epsilon \Rightarrow \mathbf{Y} \sim N(\mathbf{X}\beta + \mathbf{Z}\gamma, \sigma^2 I)$
 Assume: $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ to get $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

Then

	$\gamma = 0$	$Z = XB$	$X^T Z = 0$	other
$E(\hat{\beta})$	β	✓		✓
$E(\hat{Y})$	$\mathbf{X}\beta + \mathbf{Z}\gamma$		✓	
$E(\hat{Y})$	$\mathbf{X}\beta$	✓		✓



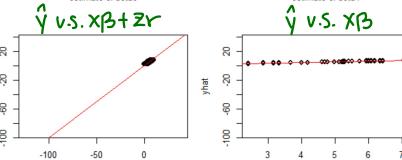
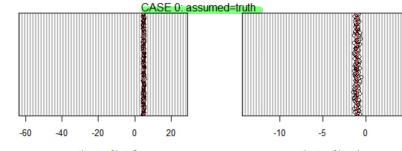
$$n = 25, \quad \mathbf{X} \text{ is } n \times 2, \quad \beta = \begin{bmatrix} 5 \\ -1 \end{bmatrix}, \quad \mathbf{Z} \text{ is } n \times 1, \quad \gamma = 5$$

Case 0: $\gamma = 0$... demonstrates usual properties

Case 1: $\gamma = 5$, \mathbf{Z} in column space of \mathbf{X} ... $Z = X \begin{bmatrix} 4 \\ 1 \end{bmatrix}$

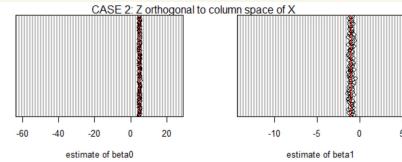
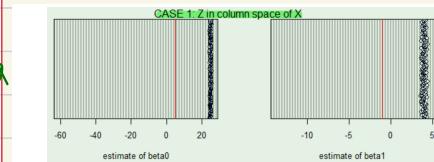
Case 2: $\gamma = 5$, columns of \mathbf{Z} orthogonal to columns of \mathbf{X} ... $Z = (I - P_X)(4 + N(0, 1))$

Case 3: $\gamma = 5$, other ... $Z = [-1, -2, \dots, -n]^T$

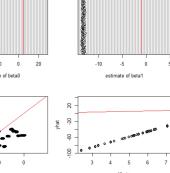
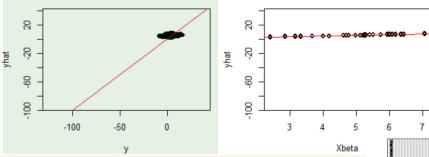
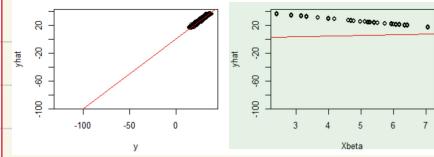


unbiased $\hat{\beta}$

\hat{Y} vs. $E(Y)$



unbiased $\hat{\beta}$



$\hat{\beta}$ biased

Note: 1. If want to have $E(\hat{\beta}) = \beta$, design $x'z = 0$

- $E(\hat{\beta}) = \beta$

- $E(\hat{Y}) = XB$

- e can be used to estimate σ^2 as $e \sim N(zr, \sigma^2(I-P))$

2. If want to have $E(\hat{Y}) = XB + zr$, let $z \in \text{col}(x)$

Thus "estimate β " and "estimate $E(Y)$ " have different requirements

1.3 (i) Heterogeneity i.e. unequal variance

- estimator still unbiased, but not best

- ANOVA sum of squares still ok

- standard errors are wrong, tests, CI wrong

- Transformation may fix *

Truth: $Y = XB + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 v)$

$$Y = XB + \varepsilon \quad \varepsilon \sim N(0, \sigma^2 I)$$

(ii). Correlated errors

- estimator still unbiased, but not best

- ANOVA sum of squares still ok

- standard errors are wrong, tests, CI wrong

- e.g.: time series, spatial, split plot, subsampling

- Other estimation procedure may help

1.4 Checking for normality

Why not use $e = \hat{Y} - Y$ to check for normality?

1. Residuals are not independent

$$\text{cov}(e_i, e_j) = -\sigma^2 P_{ij}$$

2. Residuals are not identically distributed

$$\text{Var}(e_i) = \sigma^2 (1 - P_{ii}) \text{ not constant}$$

3. Residuals are not a random sample from error distribution ε

Standardized residual: $r_i = e_i / \hat{\sigma} \sqrt{1 - P_{ii}}$

Studentized residual: $r_i = e_i / \hat{\sigma}_{(i)} \sqrt{1 - P_{ii}}$

Deleted/Jackknife / LOOCV results:

1. Assumptions violated:

1.3 (i) Heterogeneity, i.e., unequal variances

1.3 (ii) Correlated errors

Truth: $\mathbf{Y} = \mathbf{X}\beta + \epsilon$, $\epsilon \sim N(\mathbf{0}, \sigma^2 V)$ $\Rightarrow \mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 V)$

Assume: $\mathbf{Y} = \mathbf{X}\beta + \epsilon$, $\epsilon \sim N(\mathbf{0}, \sigma^2 I)$ to get $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

Then

- $E(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\beta = \beta$ unbiased
- $E(\hat{Y}) = \mathbf{X}E(\hat{\beta}) = \mathbf{X}\beta$ unbiased
- $E(\mathbf{e}) = (\mathbf{I} - \mathbf{P})E(\mathbf{Y}) = \mathbf{0}$ unbiased for 0
- $\text{var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 V \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$

is not as “small” as possible among unbiased estimators

★ $\text{cov}(\hat{Y}, \mathbf{e}) = \mathbf{P}\sigma^2 V(\mathbf{I} - \mathbf{P}) \neq 0$ correlated

even though $\hat{Y}^T \mathbf{e} = \mathbf{Y}^T \mathbf{P} \cdot (\mathbf{I} - \mathbf{P}) \mathbf{Y} = 0$ orthogonal

Notation
oooo

Inference
oo

Problems
o

Mean model misspecified
oooooooooooo

Covariance misspecified
ooo●oooo

Normality
oooo

Possible remedial actions . . .

- Use Weighted Least Squares (WLS)
if know $\text{var}(Y) = a\sigma^2$, know a , know independent
- Use Iteratively Reweighted Least Squares (IRLS)
if know $\text{var}(Y) = f[\text{E}(Y)]$, know independent
- Use Estimated Generalized Least Squares (EGLS)
- Use Generalized Linear Model (GLM)
- **Transform!**

Notation
ooooInference
ooProblems
oMean model misspecified
ooooooooooooCovariance misspecified
oooo●ooooNormality
oooo

There is the Box-Cox Family of Transformations:

$$Y_i^{(\lambda)} = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda Y} & \lambda \neq 0 \\ Y \ln(Y_i) & \lambda = 0 \end{cases} \quad \text{where } \bar{Y} = (Y_1 Y_2 \cdots Y_n)^{1/n}$$

Denominator and -1 in numerator are just for scaling.

Converts scale back to original units, and thus allows direct comparison of SSE across models with different powers.

- Fit ANOVA for several values of λ , e.g., $\lambda = -1, -0.9, \dots, 1$. Record $SSE^{(\lambda)}$ for each value of λ .
- Plot $\{\lambda, SSE^{(\lambda)}\}$ and determine $SSE_{\min}^{(\lambda)}$
- In the end, choose any value of λ that causes

$$SSE^{(\lambda)} \leq SSE_{\min}^{(\lambda)} \left[1 + \frac{t_{df_e, \alpha/2}^2}{df_e} \right],$$

where df_e is the degrees of freedom associated with any $SSE^{(\lambda)}$

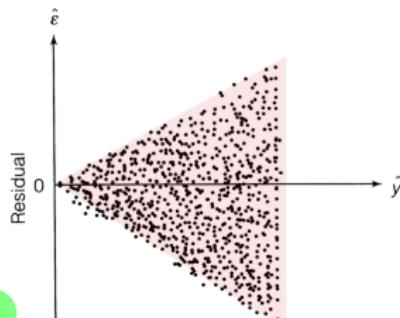
Note: $\hat{\beta}^{(\lambda)}$ can be very different as λ changes! Also sensitive to $\mathbf{X}\beta$.

Other Guidance on Choosing a Transformation

If $V(Y) \propto [E(Y)]^{2k}$ then transform to Y^{1-k} , where " $Y^0 = \ln(Y)$ ".

$$k=1 : V(Y) \propto \{E(Y)\}^2$$

- Use transformation $Y^{1-k} = Y^0 = \ln(Y)$
- Great if $Y \sim \text{Gamma}(\alpha, \beta)$, with $V(Y) = [E(Y)]^2/\alpha$
 - Example: Y = survival time of mice subjected to a treatment

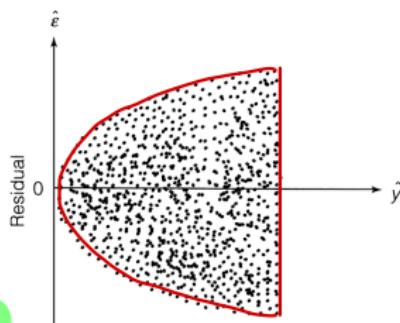


- Plot e_i versus \hat{Y}_i is fan-shaped
- Possible alternative approach: *generalized linear model*

If $V(Y) \propto [E(Y)]^{2k}$ then transform to Y^{1-k} , where " $Y^0 = \ln(Y)$ ".

$k = 0.5 : V(Y) \propto E(Y)$

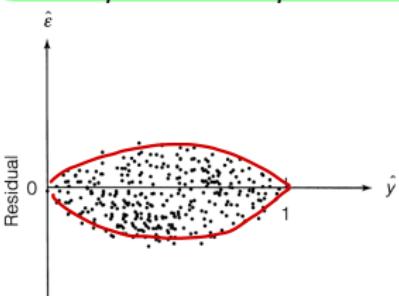
- Use transformation $Y^{1-k} = \sqrt{Y}$
- Great if $Y \sim \text{Poisson}$, with $V(Y) = E(Y)$
 - Example: $Y = \# \text{ trees in 1000 acres of a forested area}$
This is a count having a very large (possibly "infinite") upper bound



- Plot e_i versus \hat{Y}_i is fan-shaped
- Possible alternative approach: *loglinear regression model*

$$V(Y) \propto \{E(Y)\}\{1 - E(Y)\}$$

- Use transformation $\arcsin(\sqrt{Y}) = \sin^{-1}(\sqrt{Y})$
- Great if $nY \sim \text{Binomial}$, with
 $V(Y) = [E(Y)][1 - E(Y)]/n$
 - Example: Y = proportion of 30 trees that are afflicted by a fungus
 - (arcsin good when $E(Y) < 0.3$ or $E(Y) > 0.7$)
- Plot e_i versus \hat{Y}_i has bulge for \hat{Y} close to 0.5



- Possible alternative approach: logistic regression

Unequal variances often coexist with nonnormality!

Notation
oooo

Inference
oo

Problems
o

Mean model misspecified
oooooooo
oooooooo

Covariance misspecified
oooooooo
●oo

Normality
oooo

1. Assumptions violated:

1.3 (ii) Correlated errors:

- Estimators are still unbiased, but not best.
- ANOVA sum-of-squares are still ok
- Standard errors are wrong, and hence tests and CIs are wrong
- Creates the problem: Time series, spatial, split plot, subsampling
- Tests
- Possible fix: Estimation procedures other than OLS

How to diagnose?

- Be guided by the type of data
- Look for patterns among residuals over time/space/etc.

Detecting Correlation

Positively(Negatively) correlated data can lead to *standard errors that are seriously under-(over-)estimated*, thus drastically affecting hypothesis tests and confidence intervals.

Durbin-Watson Test of Autocorrelation

Durbin-Watson test statistic:

$$d = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} \approx 2(1 - \hat{\rho}) \approx \begin{cases} 0 & \text{strong positive autocorrelation} \\ 2 & \text{no autocorrelation} \\ 4 & \text{strong negative autocorrelation} \end{cases},$$

where $\hat{\rho}$ is the sample correlation between e_i and e_{i-1}

- Ordering of the data matters. Looking for correlation with immediate neighbors, following sequence.
- Null distribution complicated

Example: 35-year sales history of a company

Y : annual sales, in thousands of dollars [SALES]

X : year [T]

OLS Regression:

$$\hat{\beta}_0 = 0.402, \text{ with } se(\hat{\beta}_0) = 2.206$$

$$\hat{\beta}_1 = 4.296, \text{ with } se(\hat{\beta}_1) = 0.107$$

AR(1) Regression:

$$\hat{\beta}_0 = 0.422, \text{ with } se(\hat{\beta}_0) = 3.670$$

$$\hat{\beta}_1 = 4.295, \text{ with } se(\hat{\beta}_1) = 0.179$$

```
proc reg data=sales35;
  model sales=t / dwprob; **d=0.821, rho=0.590, pval[Ha:+ve corr] is <.0001;
proc arima data=sales35;
  identify var=sales crosscorr=t;
  estimate p=1 input=t; run;
```

```
fit = lm(SALES ~ T, data=SALES35); summary(fit)
library(car); durbinWatsonTest(fit)
```

```
# AR(1) regression:
fit = arima(SALES35$SALES, order=c(1, 0, 0), xreg = SALES35$T); fit
```

Check for Normality

- Checking for normality ranks low relative to other checks
 - Expectation & variance of estimators and sums of squares are unaffected by nonnormality
 - ANOVA F-test is reasonably robust to nonnormality
 - HTs and CIs are more affected by nonnormality, but robust in large samples
- Testing for normality can be overkill. Instead,
 - use histogram, with normal curve overlaid
 - use normal quantile-quantile (Q-Q) plot and look for pattern¹
 - straight line \rightsquigarrow normal distribution (intercept is mean, slope is std. dev.)
 - "S" \rightsquigarrow symmetric, light-tailed distribution
 - "tangent" \rightsquigarrow symmetric, heavy-tailed dist'n. **Unequal variances?**
 - **Outliers?**
 - "J" \rightsquigarrow positively skewed dist'n. **Log transformation?**
 - "r" \rightsquigarrow negatively skewed dist'n
 - line not through origin \rightsquigarrow missing important predictor variable

¹heavily dependent on sample size!

- Is e a good choice for testing normality?

- Is it true that e_1, \dots, e_n forms a random sample?

No because $\text{var}(e_i) = \sigma^2(1 - P_{ii})$ & $\text{Cov}(e_i, e_j) = -\sigma^2 P_{ij}$

- internally studentized residual

(R=standardized residual, SAS=studentized residual):

$$r_i = \frac{e_i}{s\sqrt{1 - P_{ii}}}, \quad s = \sqrt{MSE}, \quad \text{MSE from regn with all } n \text{ obs}$$

- * $\text{var}(r_i) \approx 1$.
- * r_i, r_j are likely dependent
- * $r_i \approx t_{df_e}$ (dependent numerator & denominator)

- externally studentized residual or studentized deleted residual

(R & SAS=rstudent residual):

$$r_i^* = \frac{Y_i - \hat{Y}_{i(i)}}{s_{(i)} \sqrt{1 + \mathbf{x}_i (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i^T}} = \frac{e_i}{s_{(i)} \sqrt{1 - P_{ii}}}$$

- * $\text{var}(r_i^*) \approx 1$.
- * r_i^*, r_j^* are likely dependent
- * $r_i^* \approx t_{df_e-1}$ (better approx than r_i) (df_e is from regression with all n observations)
- * r_i^* reflects large values more dramatically than r_i (Atkinson 1983)

- **deleted/Jackknife/LOOCV results:** run regression n times, with the i th observation excluded during the i th run: (Belsley, Kuh, Welsch 1980)

- $\mathbf{X}_{(i)}$ is new design matrix from omitting the i th observation
- $s_{(i)}^2 = MSE_{(i)}$ is the MSE from omitting the i th observation
No need to rerun the regn: $(n - p - 1)s_{(i)}^2 = (n - p)s^2 - r_i^2/(1 - P_{ii})$
- $\hat{\beta}_{(i)}$ is estimate of full p -dimensional vector β from omitting the i th obs.
- $\hat{\beta}_{j(i)}$ is estimate of parameter β_j from omitting the i th obs.
- $\hat{\mathbf{Y}}_{(i)}$ is prediction of full n -dimensional vector \mathbf{Y} from omitting the i th obs.
- $\hat{Y}_{i(i)}$ is prediction of i th observation Y_i from omitting the i th obs.
- $Y_i - \hat{Y}_{i(i)}$ is called the i th deleted residual
- r_i^* is the t-statistic for agreement between Y_i and $\hat{Y}_{i(i)}$:

$$r_i^* = \frac{Y_i - \hat{Y}_{i(i)}}{\text{s.e.}(Y_i - \hat{Y}_{i(i)})} = \frac{Y_i - \hat{Y}_{i(i)}}{s_{(i)} \sqrt{1 + \mathbf{x}_i (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i^T}} = \frac{Y_i - \hat{Y}_i}{s_{(i)} \sqrt{1 - P_{ii}}}$$

Notation
ooooInference
ooProblems
oMean model misspecified
oooooooo
ooooooooCovariance misspecified
oooooooo
oooNormality
ooo●

Creating and interpreting a normal Q-Q plot:

- y-axis: residual, e.g. r_i^* , in increasing order
x-axis: quantiles from a normal distribution
For j th ordered residual, plot $\left\{ \Phi^{-1} \left[\frac{j-3/8}{n+1/4} \right], j\text{th ordered residual} \right\}$
- Good when $n \geq 30$, better when $n \geq 50$. (theory says ok when $n \geq 5$)
- Look for pattern

How to correct nonnormality?

TRANSFORMATION, e.g. Box-Cox

Make variances equal

Address outliers

Illustration

Salary as a function of years of experience

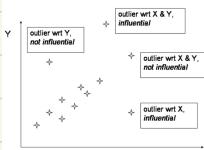
3. Unusual Data

(i) Outlier

(ii) Influential point

Outliers and Influential data points ... *

- Outlier:** data point that is "unusual" when compared to other points in the dataset. May not be influential. Can be outlier wrt X (aka high leverage) or wrt Y .
- Influential data points:** has undue influence in that inference can change drastically depending on whether that single point is included in the analysis.



What to do about outlier wrt Y ?

- improved model is needed
- claim that point has other features driving response that model doesn't account for

What to do about outlier wrt X ?

- try to fill in gap in X space, then refit

What to do about influential points?

Fit model with and without points. Change in inference?

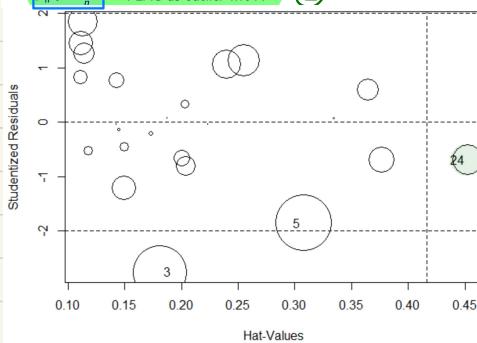
- Yes : report both results
- No : report results from full set

Outlier wrt X (aka High Leverage Point):

Which data points are far from the centroid of the X -space?

Does a data point have "high leverage," i.e., large value of P_{ii} ?

- P_{ii} is weight of i th observation in determining \hat{Y}_i , since $\hat{Y} = \mathbf{P}\mathbf{Y}$
- $\frac{1}{n} \leq P_{ii} \leq \frac{1}{c}$, c is number of rows same as i th row
 - $P_{ii} = \frac{1}{c}$ only when $x_{ij} = \bar{x}_{j..}, j = 1, \dots, p$, i.e., the i th point is at the centroid
 - model should capture this point well, if it is a good model
 - relatively speaking, this point deserves little scrutiny
 - omitting this point will cause little change in the analysis
- $P_{ii} = 1$ only when x_i is far from the centroid
 - regression line fits the point exactly, i.e., $\hat{Y}_i = Y_i$. Indicative of good model? Not really
 - this point may or may not be in line with the other points
 - relatively speaking, this point deserves heavy scrutiny
 - omitting this point may cause great change in the analysis
- $\sum_{i=1}^n P_{ii} = r(\mathbf{X})$ since $\text{tr}(\mathbf{P}) = r(\mathbf{X})$
- $P_{ii} > \frac{2r(\mathbf{X})}{n} \rightsquigarrow \text{FLAG as outlier wrt } X$ (2)



Outlier wrt Y :

Is Y_i well predicted?

- A "no" answer could result from outlier wrt Y or an inadequate model!
- Recall that r_i^* is the t-statistic for agreement between Y_i and $\hat{Y}_{(i)}$:

$$r_i^* = \frac{Y_i - \hat{Y}_{(i)}}{\text{s.e.}(Y_i - \hat{Y}_{(i)})} = \frac{Y_i - \hat{Y}_{(i)}}{s_{(i)}\sqrt{1 + \mathbf{x}_i(\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i^T}} = \frac{Y_i - \hat{Y}_{(i)}}{s_{(i)}\sqrt{1 - P_{ii}}} \quad \boxed{\text{Studentized residual}}$$

- $|r_i^*| > 3 \rightsquigarrow \text{FLAG as outlier wrt } Y$

• proc reg in SAS applies threshold of 2 by default – too low

- Bonferroni: $|r_i^*| > t_{df_e-1, \frac{\alpha}{2n}} \rightsquigarrow \text{FLAG as outlier wrt } Y$

We're actually making inferences about n residuals!

Must make adjustment for multiple testing.

• R has a function for this, namely car::outlierTest

Influential Data Points:

Examine effect on $\hat{\beta}$, \hat{Y} (also possible: \hat{Y}_i , $\hat{\beta}_j$, $\text{Var}(\hat{\beta})$)

- Effect on $\hat{\beta}$, \hat{Y} : Cook's D(istance)

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})' (\mathbf{X}^T \mathbf{X}) (\hat{\beta}_{(i)} - \hat{\beta})}{r(\mathbf{X}) s^2} = \frac{(\hat{Y}_{(i)} - \hat{Y})' (\hat{Y}_{(i)} - \hat{Y})}{r(\mathbf{X}) s^2} \text{ computing } \frac{r_i^2}{r(\mathbf{X})} \left(\frac{P_{ii}}{1 - P_{ii}} \right)$$

$D_i > 1$ (some use > 0.8) $\rightsquigarrow \text{FLAG as influential}$

$D_i > F_{r(\mathbf{X}), n-r(\mathbf{X}), 0.05} \rightsquigarrow \text{FLAG as influential}$ (3)

• proc reg in SAS applies a different threshold by default – too low

Cook's distance

3. Computational Instability, aka Multicollinearity:

- Do different predictor variables provide redundant information?
- Affected by choice of values for predictors X_1, X_2, \dots, X_p

Definition:

Multicollinearity exists when two or more of the predictor variables used in regression are moderately or highly correlated.

$$\mathbf{X} = \begin{bmatrix} 1 & 4 & 4.01 \\ 1 & 6 & 5.98 \\ 1 & 7 & 7.02 \\ 1 & 8 & 7.99 \end{bmatrix}, \quad r(\mathbf{X}) = 3$$

But solving $\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{Y}$ will be highly unstable!

NEEDS:

- Recognize when multicollinearity is a problem
- Take corrective action

Diagnosis ... this IS very important

Indirect:

- Look for symptoms, although they can happen for other reasons:
 - Large changes in $\hat{\beta}$ when predictor variable is added/omitted
 - $\hat{\beta}_j$ has a sign opposite what is expected
 - Model has significant F but many non-significant t tests
 - Sequential SS and partial SS are very different
 - Variances of $\hat{\beta}_j$ are very large
 - High correlation between $\hat{\beta}_j$ and $\hat{\beta}_k$. This may violate simple interpretation of regression coefficients as measuring change in $E(Y)$ when a given predictor variable is increased by 1 while all others are held constant.

Direct:

- Large (simple) correlations between predictor variables
- Large variance inflation factors
- Condition number and indices, i.e., scrutinize structure of $\mathbf{X}^T \mathbf{X}$ (won't cover)

Variance Inflation Factors (VIF) ...

- One for each predictor variable in the model, except intercept
- $VIF_j = \frac{1}{1-R_j^2}$, R_j^2 is coefficient of determination from regressing X_j on all other predictor variables (including intercept)
 - $0 \leq R_j^2 \leq 1$ implies $1 \leq VIF_j$
 - $R_j^2 \approx 0$ implies $VIF_j \approx 1$ and $R_j^2 \approx 1$ implies VIF_j very large

X_j involved in linear dependency with other non-intercept predictor variables
 $\Rightarrow R_j^2 \approx 1 \Rightarrow VIF_j$ very large
- $\text{var}(\hat{\beta}_j) = \sigma^2 \cdot VIF_j^\dagger$: j th variable \perp all others $\Rightarrow VIF_j = 1$
- FLAG: $VIF_j > 10$ indicates problem, $VIF_j > 30$ severe problem
- Unaffected by centering predictor variables
- Detects overall collinearity problems with more than just the intercept
- No direct indication of:
 - number of linear dependencies
 - which other variables are involved in linear dependency with X_j

† assumes centered & scaled predictors

Possible Corrective Actions ...

- Drop one or more of the correlated predictor variables
- If you decide to keep all predictor variables in the model:
 - Avoid making inferences about the individual β values, and don't try to determine "relative importance" of the predictor variables
- Sacrifice unbiasedness to get smaller variance \rightsquigarrow **Biased Regression:**
 - Principal Components Regression (PCR)
 - Partial Least Squares Regression (PLSR)
 - Ridge Regression
 - LASSO, ElasticNet, etc.

Example: Cigarettes [dataset FTCCIGAR.txt]

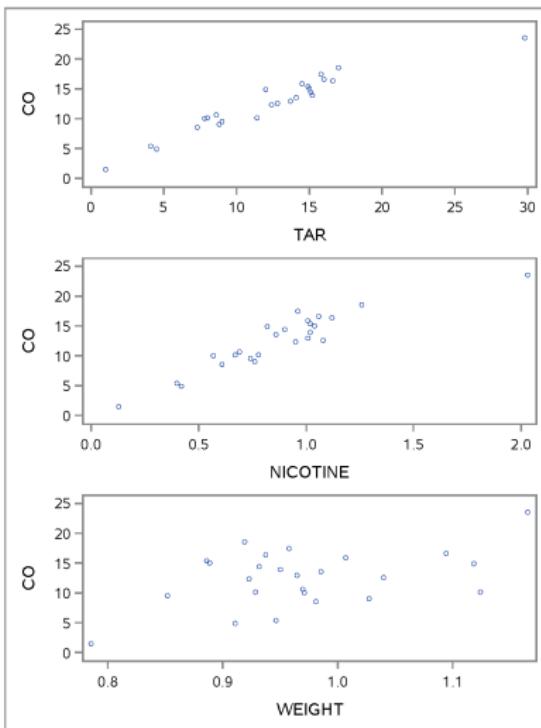
Can we model the carbon monoxide content of cigarettes as a function of their tar content, their nicotine content, and their weight?

Y : carbon monoxide content from cigarette smoke [CO]

X_1 : tar content of cigarette [TAR]

X_2 : nicotine content of cigarette [NICOTINE]

X_3 : weight of cigarette [WEIGHT]



```
proc reg data=cigar plots=none;
model co=tar nicotine weight /vif;
```

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	3	495.25781	165.08594	78.98	<.0001
Error	21	43.89259	2.09012		
Corrected Total	24	539.15040			

Root MSE	1.44573	R-Square	0.9186
Dependent Mean	12.52800	Adj R-Sq	0.9070
Coeff Var	11.53996		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	3.20219	3.46175	0.93	0.3655	0
TAR	1	0.96257	0.24224	3.97	0.0007	21.63071
NICOTINE	1	-2.63166	3.90056	-0.67	0.5072	21.89992
WEIGHT	1	-0.13048	3.88534	-0.03	0.9735	1.33386

```
proc reg data=cigar plots=none;
model co=nicotine weight /vif;
```

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	2	462.25639	231.12820	66.13	<.0001
Error	22	76.89401	3.49518		
Corrected Total	24	539.15040			

Root MSE	1.86954	R-Square	0.8574
Dependent Mean	12.52800	Adj R-Sq	0.8444
Coeff Var	14.92290		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	1.61398	4.44663	0.36	0.7201	0
NICOTINE	1	12.38812	1.24473	9.95	<.0001	1.33366
WEIGHT	1	0.05883	5.02395	0.01	0.9908	1.33366

```
proc reg data=cigar plots=none;
model co=tar weight /vif;
```

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	2	494.30638	247.15319	121.25	<.0001
Error	22	44.84402	2.03836		
Corrected Total	24	539.15040			

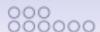
Root MSE	1.42771	R-Square	0.9168
Dependent Mean	12.52800	Adj R-Sq	0.9093
Coeff Var	11.39618		

Parameter Estimates						
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Variance Inflation
Intercept	1	3.11433	3.41620	0.91	0.3718	0
TAR	1	0.80415	0.05904	13.62	<.0001	1.31726
WEIGHT	1	-0.42287	3.81299	-0.11	0.9127	1.31726

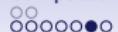
```
proc corr data=cigar
plots=matrix(histogram);
var co tar nicotine weight;
```

Pearson Correlation Coefficients, N = 25 Prob > r under H0: Rho=0				
	CO	TAR	NICOTINE	WEIGHT
CO	1.00000	0.95749 <.0001	0.92595 <.0001	0.46396 0.0195
TAR	0.95749 <.0001	1.00000	0.97661 <.0001	0.49077 0.0127
NICOTINE	0.92595 <.0001	0.97661 <.0001	1.00000	0.50018 0.0109
WEIGHT	0.46396 0.0195	0.49077 0.0127	0.50018 0.0109	1.00000

Unusual Data Points

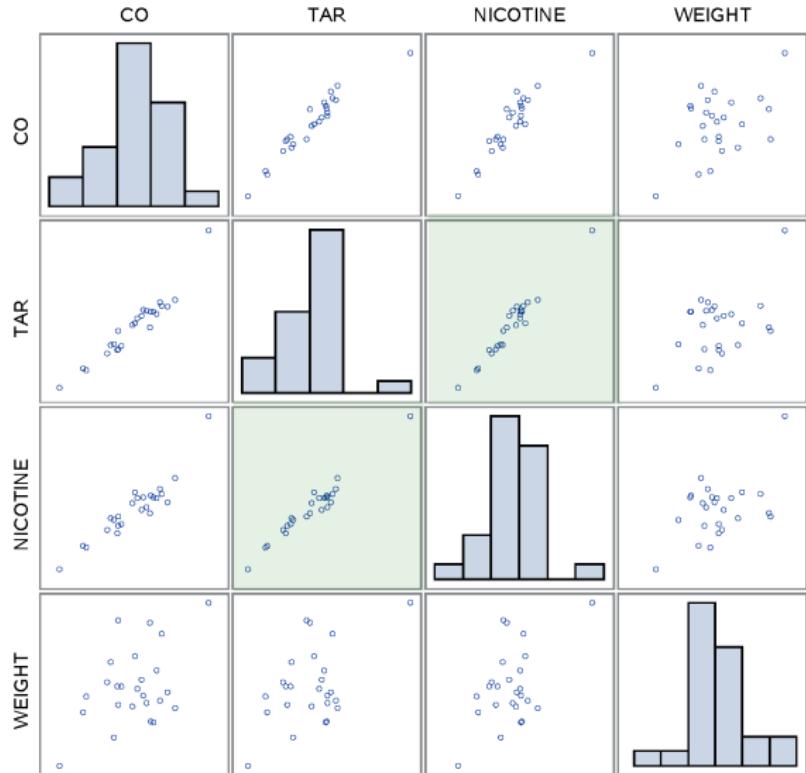


Computational Instability



The CORR Procedure

Scatter Plot Matrix



Using R ...

```
> fit = lm(CO ~ TAR + NICOTINE + WEIGHT, data=FTCCIGAR)
```

```
> install.packages("car")
```

```
> library(car)
```

```
> vif(fit)
```

TAR	NICOTINE	WEIGHT
21.630706	21.899917	1.333859

```
> cor(FTCCIGAR)
```

```
          TAR  NICOTINE    WEIGHT      CO
TAR      1.0000000 0.9766076 0.4907654 0.9574853
NICOTINE 0.9766076 1.0000000 0.5001827 0.9259473
WEIGHT   0.4907654 0.5001827 1.0000000 0.4639592
CO       0.9574853 0.9259473 0.4639592 1.0000000
```

```
>
```

Recommendation: Drop nicotine from the model.
from R^2 & R^2_{adj}

Diagnostic Measures for GLM

Recall: Some residuals

- Is e a good choice for testing normality?

No because $\text{var}(e_i) = \sigma^2(1 - P_{ii})$ & $\text{Cov}(e_i, e_j) = -\sigma^2 P_{ij}$

- internally studentized residual** (2)

(R=standardized residual, SAS=studentized residual):

$$r_i = \frac{e_i}{\sqrt{1 - P_{ii}}} \quad \text{S.E.} = \sqrt{\text{MSE}}, \quad \text{MSE from regn with all } n \text{ obs}$$

* $\text{var}(r_i) \approx 1$

* r_i, r_j are likely dependent

* $r_i \approx t_{df-1}$ (dependent numerator & denominator)

- externally studentized residual or studentized deleted residual** (3)

(R & SAS=rstudent residual):

$$r_i^* = \frac{Y_i - \hat{Y}_{(i)}}{s_{(i)} \sqrt{1 + \mathbf{x}_i (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i^T}} = \frac{e_i}{s_{(i)} \sqrt{1 - P_{ii}}}$$

* $\text{var}(r_i^*) \approx 1$

* r_i^*, r_j^* are likely dependent

* $r_i^* \approx f_{df-1}$ (better approx than r_i) (d_f is from regression with all n observations)

* r_i^* reflects large values more dramatically than r_i (Atkinson 1983)

- deleted/Jackknife/LOOCV results:** run regression n times, with the i th observation excluded during the i th run. (Belsley, Kuh, Welsh 1990)

* $\mathbf{X}_{(i)}$ is new design matrix from omitting the i th observation

* $s_{(i)}^2 = \text{MSE}_{(i)}$ is the MSE from omitting the i th observation

No need to rerun the regn: $(n - p - 1)s_{(i)}^2 = (n - p)^2 - r_i^2 / (1 - P_{ii})$

* $\beta_{(i)}$ is estimate of full p -dimensional vector β from omitting the i th obs.

* $\hat{\beta}_{(i)}$ is estimate of parameter β_j from omitting the i th obs.

* $\hat{Y}_{(i)}$ is prediction of full n -dimensional vector \mathbf{Y} from omitting the i th obs.

* $\hat{Y}_{(i)}$ is prediction of i th observation Y_i from omitting the i th obs.

* $Y_i - \hat{Y}_{(i)}$ is called the i th deleted residual

* r_i^* is the t-statistic for agreement between Y_i and $\hat{Y}_{(i)}$:

$$r_i^* = \frac{Y_i - \hat{Y}_{(i)}}{\text{s.e.}(Y_i - \hat{Y}_{(i)})} = \frac{Y_i - \hat{Y}_{(i)}}{s_{(i)} \sqrt{1 + \mathbf{x}_i (\mathbf{X}_{(i)}^T \mathbf{X}_{(i)})^{-1} \mathbf{x}_i^T}} = \frac{Y_i - \hat{Y}_i}{s_{(i)} \sqrt{1 - P_{ii}}}$$

GLM Diagnostic

1. + Deviance lack of fit test

2. Residual

- Several types of residuals are commonly used:

(1) **Raw or response**: $y_i - \hat{\mu}_i$

$$\sim \chi^2 = \sum_{i=1}^n (r_i^P)^2$$

not very useful for diagnostics

(2) **Pearson**: $r_i^P \equiv \frac{y_i - \hat{\mu}_i}{\sqrt{\text{var}(y_i | \hat{\mu}_i)}}$

$$\sim \chi^2$$

also $\frac{r_i^P}{\sqrt{\psi_i / (1 - h_i)}}$ where h_i is leverage

(3) **Deviance**: $r_i^D \equiv \text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i}$

$$\sim D_M = \sum_{i=1}^n d_i = \sum_{i=1}^n (r_i^D)^2$$

very useful for diagnostics

also $\frac{r_i^D}{\sqrt{\psi_i / (1 - h_i)}}$

Note: The chi-squared statistic χ^2 and the scaled chi-squared statistic χ^2/ψ are often used interchangeably with the deviance D_M and scaled deviance D_M/ψ .

3. Influence (analog Cook's distance)

- New diagnostic measure: likelihood displacement:

$$LD_i \equiv 2 \{ \ell_M(\hat{\theta}; \mathbf{y}) - \ell_M(\hat{\theta}_{(-i)}; \mathbf{y}) \}$$

where $\hat{\theta}_{(-i)}$ is MLE from excluding i th observation. Likelihood evaluated with all observations.

- Predicted values: $\hat{\mu}_i$, response or $\mathbf{x}_i^T \hat{\beta}$, linear predictor

4. Complete / Quasi Complete Separation

- non-convergence, high SE, LR ≈ 0 .

- happens when rare events, large prediction space many binary prediction, small sample size.

- biased regression (Firth's procedure) (to fix)

- quick check: deviance / dfError

≈ 1 model fits well

< 1 possibly overfitting

> 1 underfitting, overdispersion

Exponential Family

General Forms & Rules

Definition 1:

- $$l(y_i, \phi; \theta) = \sum_{i=1}^n \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi w_i} + C(y_i, \phi; w_i) \right\}$$
- ← this is in a very weird form,
it doesn't put $T(x)$ altogether.
- $M_i = g'(x_i; \beta)$ is some function of β
 - θ is some function of M (thus β)

Definition 2:

$$Y = EF(\theta, \phi)$$

$$f(y; \theta, \phi) = \exp \left\{ \frac{T(y)g(\theta) - b(\theta)}{\alpha(\phi)} \right\} h(y, \phi)$$

- ϕ dispersion parameter
- $E(T(y)) = b'(\theta)$
- $\text{Var}(T(y)) = \alpha(\phi) b''(\theta)$

Definition 3: $f(y; \theta) = \exp \{ T(y)g(\theta) - b(\theta) \} h(y)$

$g(\theta)$: natural (canonical) parameter

$T(y)$: sufficient statistics

* Assume $\theta = g(\theta)$ & the canonical parameter :

$$\int \exp \{ T(y)g(\theta) - b(\theta) \} h(y) dy = 1$$

$$\int e^{T(y)g(\theta)} e^{-b(\theta)} h(y) dy = 1$$

$$\int e^{T(y)g(\theta)} h(y) dy = e^{b(\theta)}$$

$$\Rightarrow \int e^{T(y)g(\theta)} h(y) dy = e^{b(\theta)} * g(\theta) = \theta$$

$$\Rightarrow \int T(y)e^{T(y)\theta} h(y) dy = b'(\theta) e^{b(\theta)} * \text{derivative w.r.t. } \theta \text{ on both sides}$$

$$\int T(y)e^{T(y)\theta} h(y) dy = b'(\theta)$$

$$\Rightarrow E(T(y)) = b'(\theta)$$

$$\text{Var}(T(y)) = b''(\theta) \text{ similarly}$$

Univariate Normal

$$\begin{aligned}
 f(x) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \\
 &= \exp\left(-\frac{n}{2}\log(2\pi\sigma^2) - \frac{\sum(x_i-\mu)^2}{2\sigma^2}\right) \\
 &= \exp\left(-\frac{n}{2}\log(2\pi\sigma^2) - \frac{\sum x_i^2 - 2\mu\sum x_i + n\mu^2}{2\sigma^2}\right) \quad \rightarrow \text{Complete \& sufficient statistics } (\sum x_i, \sum x_i^2) \\
 &= \exp\left(-\frac{n}{2}\log(2\pi\sigma^2) + \frac{n\mu\sum x_i - n\mu^2}{\sigma^2} - \frac{\sum x_i^2}{2\sigma^2}\right) \\
 &= \exp\left(\frac{n\sum x_i - n\mu^2}{\sigma^2} - \frac{n}{2}\log(\sigma^2) - \frac{n}{2}\log(2\pi) - \frac{\sum x_i^2}{2\sigma^2}\right) \\
 \phi &= \sigma^2 \quad \theta = \mu \quad b(\theta) = \mu^2/2 \\
 w_i = 1 \quad c(y, \phi, w_i) &= -\frac{n}{2}\log(\sigma^2) - \frac{n}{2}\log(2\pi) - \frac{\sum x_i^2}{2\sigma^2}
 \end{aligned}$$

* MGF derivation

$$\begin{aligned}
 M_X(t) &= E(e^{tX}) \\
 &= \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2 - 2x\mu + \mu^2}{2\sigma^2} + \frac{2tx}{2\sigma^2}} dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2 - 2x(M+t\sigma^2) + M^2}{2\sigma^2}} dx \\
 &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-(M+t\sigma^2))^2}{2\sigma^2} + \frac{M^2 + t^2\sigma^2}{2}} dx \\
 &= e^{t\mu + \frac{t^2\sigma^2}{2}}
 \end{aligned}
 \quad EX = \mu \quad \text{Var } X = \sigma^2$$

Multivariate Normal

$$\begin{aligned}
 f_X(x) &= \frac{1}{\det(2\pi\Sigma)} \exp\left(-\frac{1}{2}(x-\mu)' \Sigma^{-1} (x-\mu)\right) \\
 &= \exp\left(-\log\det|2\pi\Sigma| - \frac{1}{2}(x-\mu)' \Sigma^{-1} (x-\mu)\right) \\
 &= \exp\left(-\log\det|2\pi\Sigma| - \frac{1}{2}(x' \Sigma^{-1} x - 2\mu' \Sigma^{-1} x + \mu' \Sigma^{-1} \mu)\right) \\
 MGF &= e^{\frac{1}{2}t' \Sigma t + t' \mu}
 \end{aligned}$$

* Conditional Normal

$$\begin{aligned}
 \begin{pmatrix} x \\ y \end{pmatrix} &\sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}\right) \\
 x | y &\sim N(\mu_1 + \Sigma_{12}' \Sigma_{22}^{-1} (y - \mu_2), \Sigma_{11} - \Sigma_{12}' \Sigma_{22}^{-1} \Sigma_{21})
 \end{aligned}$$

Gamma

1. Shape α , Scale B ← this is the most common case

$$f_X(x) = \frac{1}{\Gamma(\alpha)B^\alpha} x^{\alpha-1} e^{-x/B} \quad 0 < x$$

$$f(x) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)B^\alpha} x_i^{\alpha-1} e^{-x_i/B} \quad 0 < x_i$$

$$= \exp(-n \log(\Gamma(\alpha)) - n \alpha \log(B) + (\alpha-1) \sum_{i=1}^n \log(x_i) - \frac{1}{B} \sum_{i=1}^n x_i)$$

$$= \exp(-\frac{1}{B} \sum_{i=1}^n x_i - n \log(\Gamma(\alpha)) - n \alpha \log(B) + (\alpha-1) \sum_{i=1}^n \log(x_i))$$

$$\varphi = 1 \quad \theta = -\frac{1}{B} \quad b(\theta) = -n \log(\Gamma(\alpha)) - n \alpha \log(B)$$

$$\omega_i = 1 \quad a(y, \varphi, \omega_i) = (\alpha-1) \sum_{i=1}^n \log x_i$$

$$E(X) = \int_0^\infty x \frac{1}{\Gamma(\alpha)B^\alpha} x^{\alpha-1} e^{-x/B} dx$$

$$= \int_0^\infty \frac{1}{\Gamma(\alpha)B^\alpha} x^{(\alpha+1)-1} e^{-x/B} dx$$

$$= \int_0^\infty \frac{1}{\Gamma(\alpha+1)B^{\alpha+1}} \cdot \alpha B \cdot x^{(\alpha+1)-1} e^{-x/B} dx$$

$$= \alpha B$$

Sufficient & complete statistics
 $T(X) = (\sum_{i=1}^n X_i, \sum_{i=1}^n \log X_i)$

2. Shape α & rate θ * this is not being tested too much.

As long as always remember $\theta = \frac{1}{B}$.

Beta

$$f_X(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \quad \text{with } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \quad 0 \leq x \leq 1$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}$$

$$= \exp(\log(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}) + (\alpha-1)\log(x) + (\beta-1)\log(1-x))$$

$$= \exp((\alpha-1)\log(x) + (\beta-1)\log(1-x) + \log(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}))$$

$$\varphi = 1 \quad \omega = 1 \quad \theta = 0 \quad b(\theta) = \log(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)})$$

$$a(y, \varphi, \omega_i) = (\alpha-1)\log x + (\beta-1)\log(1-x)$$

$$E(X) = \frac{\alpha}{\alpha+\beta}$$

$$\text{Var}(X) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

Sufficient & complete statistics

$$T(X) = (\sum \log x_i, \sum \log(1-x_i))$$

$$\text{Beta}(1, 1) = \text{Unif}(0, 1)$$

Exponential * time between poisson events

1. Scale = mean = β

$$f_X(x) = \frac{1}{\beta} e^{-x/\beta} \sim \text{Gamma}(\text{shape} = 1, \text{scale} = \beta) \quad 0 < x$$

$$= \exp(-\log \beta - \frac{1}{\beta} x)$$

sufficient & complete statistics
 $T(X) = \sum_{i=1}^n X_i$

$$\phi = 1 \quad \omega = 1 \quad b(\theta) = -\log \beta \quad \theta = \frac{1}{\beta}$$

$$\text{acy}(\varphi, \omega) = 0$$

$$E(X) = \int_0^\infty x \cdot \frac{1}{\beta} e^{-x/\beta} dx$$

$$= -xe^{-x/\beta} + \int_0^\infty e^{-x/\beta} dx$$

$$= -xe^{-x/\beta} - \beta e^{-x/\beta} \Big|_0^\infty$$

$$= \beta$$

$$u = x \quad v' = \frac{1}{\beta} e^{-x/\beta}$$

$$u' = 1 \quad v = -e^{-x/\beta}$$

$$E(X) = \beta$$

$$\text{Var}(X) = \beta^2$$

$$M_X(t) = E(e^{tx})$$

$$= \int_0^\infty e^{tx} \frac{1}{\beta} e^{-x/\beta} dx$$

$$= \frac{1}{\beta} \int_0^\infty e^{-x(\frac{1}{\beta} - t)} dx$$

$$= \frac{1}{\beta} \int_0^\infty e^{-x(\frac{1-t}{\beta})} dx \quad \text{lets think } \frac{1}{\beta^*} = \frac{1-t}{\beta}$$

$$= \frac{1}{\beta} \cdot \frac{\beta}{1-t\beta}$$

$$= \frac{1}{1-t\beta}$$

$$\text{then } \frac{1}{\beta^*} \cdot \int \cdot dx = 1$$

$$\Rightarrow \int \cdot dx = \beta^*$$

2. rate = λ ($= 1/\beta$)

$$f_X(x) = \lambda e^{-\lambda x}$$

$$= \exp(\log(\lambda) - \lambda \sum x_i)$$

$$\theta = -\lambda \quad \phi = 1 \quad \omega = 1$$

$$b(\theta) = \log(\lambda)$$

sufficient & complete statistics is $\sum X_i$:

$$E(X) = 1/\lambda$$

$$\text{Var}(X) = 1/\lambda^2$$

$$E(X) = \int_0^\infty x \lambda e^{-\lambda x} dx$$

$$= -xe^{-\lambda x} + \int_0^\infty e^{-\lambda x} dx$$

$$= -xe^{-\lambda x} - \frac{1}{\lambda} e^{-\lambda x} \Big|_0^\infty$$

$$= \frac{1}{\lambda}$$

$$u = x \quad v' = \lambda e^{-\lambda x}$$

$$u' = 1 \quad v = -e^{-\lambda x}$$

$$\begin{aligned}
 M_X(t) &= E(e^{tx}) \\
 &= \int_0^\infty e^{tx} \lambda e^{-x\lambda} dx \\
 &= \int_0^\infty \lambda e^{-x(\lambda-t)} dx \\
 &= \frac{\lambda}{\lambda-t}
 \end{aligned}$$

Weibull* not exp family

* Just know some relationship.

$$\exp(\beta) = \text{Weibull}(1, \beta) = \text{Gamma}(1, \beta)$$

$X \sim \text{Weibull}(a, b) \Leftrightarrow X^a \sim \exp(b^a)$ both scale - liked.

Cauchy* not exp family

* Just know some relationship.

* Undefined mean & variance

$$X \text{ indep } Y \sim N(0, 1) ; \frac{X}{Y} \sim \text{Cauchy}(0, 1)$$

Log-normal* not exp family (?)

$$\begin{aligned}
 f(x) &= \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{\log x - \mu}{2\sigma^2}\right) \\
 &= \exp\left(-\log(x) - \frac{1}{2}\log(2\pi\sigma^2) - \frac{\log x - \mu}{2\sigma^2}\right)
 \end{aligned}$$

sufficient and complete statistics $\sum_i \log(x_i)$ (?)

Uniform* not exponential family

$$X \sim \text{Unif}(0, 1)$$

$$f(x) = 1 \quad I(0 \leq x \leq 1)$$

$$f(x) = \prod_{i=1}^n 1 \quad I(0 \leq x_i \leq 1) \\ = I(0 < x_{(1)}, x_{(n)} < 1)$$

$$P(X_{(1)} \leq x) = 1 - P(X_{(1)} \geq x) \\ = 1 - P(X_1 \geq x)^n$$

$$f_{X_{(1)}}(x) = \frac{\partial}{\partial x} [1 - (1-x)]^n \\ = n \cdot (1-x)^{n-1} \\ = n(1-x)^{n-1} \sim \text{Beta}(1, n)$$

$$P(X_{(n)} \leq x) = P(X_{(1)} \leq x)^n \\ = x^n$$

$$f_{X_{(n)}}(x) = \frac{\partial}{\partial x} P(X_{(n)} \leq x) \\ = nx^{n-1} \sim \text{Beta}(n, 1)$$

$$E(X) = \frac{1}{2}$$

$$\text{Var}(X) = \frac{1}{12}$$

Sufficient statistics $X_{(1)}, X_{(n)}$

Not complete.

$(X_{(1)} + X_{(n)}) - (a+b)$ has mean 0
everywhere, but is non-trivial.

Poisson * counts in time / space

$$\prod_{i=1}^n P(X_i=x_i) = \prod_{i=1}^n \frac{\lambda^x e^{-\lambda}}{x!} \Leftrightarrow \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda}$$

$$= \exp(\sum_{i=1}^n \log(x_i!) + \sum_{i=1}^n x_i \log \lambda - n\lambda)$$

sufficient & complete statistics $\sum X_i$

$$\Phi = 1 \quad w = 1$$

$$\theta = \log \lambda \quad b(\theta) = \lambda = e^{\theta}$$

$$EX = \sum_{x=0}^{\infty} x \cdot \lambda^x e^{-\lambda} / x! \quad EX = \lambda \quad Var X = \lambda$$

$$= \sum_{x=1}^{\infty} \frac{\lambda^x e^{-\lambda}}{(x-1)!}$$

$$= \lambda \sum_{x=1}^{\infty} \frac{\lambda^{x-1} e^{-\lambda}}{(x-1)!}$$

$$= \lambda$$

$$M_X(t) = E(e^{tx})$$

$$= \sum_{x=0}^{\infty} e^{tx} \frac{\lambda^x e^{-\lambda}}{x!}$$

$$= \sum_{x=0}^{\infty} (e^{t\lambda})^x \frac{e^{-\lambda}}{x!}$$

$$= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(e^{t\lambda})^x}{x!}$$

$$= e^{-\lambda} e^{te^{\lambda}}$$

$$= e^{\lambda(e^t - 1)}$$

Bernoulli * occurrence of event in one trial

$$\prod_{i=1}^n p(X_i=x_i) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$= \exp(\sum x_i \log(\frac{p}{1-p}) + n \log(1-p))$$

sufficient & complete statistics $\sum X_i$

$$\theta = \log(\frac{p}{1-p})$$

$$EX = \sum_{x=0}^1 x p^x (1-p)^{1-x} \quad EX = p \quad Var X = p(1-p)$$

$$= p$$

$$M_X(t) = E(e^{tx})$$

$$= \sum_{x=0}^1 e^{tx} p^x (1-p)^{1-x}$$

$$= p e^t + (1-p)$$

Binomial* Occurrence of event in n trials

$$\prod_{i=1}^n p(X_i=x_i) = \prod_{i=1}^n \binom{n}{x_i} p^{x_i} (1-p)^{n-x_i}$$

$$= \exp \left\{ \sum_{i=1}^n \log \binom{n}{x_i} + \sum x_i \log \left(\frac{p}{1-p} \right) + n \log(1-p) \right\}$$

$$\theta = \log \left(\frac{p}{1-p} \right) b(\theta) = \log(1-p) = \ln(1 + e^\theta) \quad \text{if } \sum x_i$$

$$4p = 1 \quad w = n$$

$$\begin{aligned} E[X] &= \sum_{x=0}^n \binom{n}{x} \cdot x \cdot p^x (1-p)^{n-x} \\ &= \sum_{x=1}^n \frac{n!}{x!(n-x)!} p^{x-1} p (1-p)^{n-x} \\ &= n \cdot p \sum_{x=1}^n \frac{(n-1)!}{(x-1)!(n-x)!} p^{x-1} (1-p)^{n-x} \\ &= np \end{aligned}$$

$$E[X] = np$$

$$\text{Var}[X] = np(1-p)$$

$$M_X(t) = E(e^{tX})$$

$$= E(e^{t \sum Y_i}) \quad \text{with } Y_i \stackrel{iid}{\sim} \text{Ber}(p)$$

$$= E(e^{tY_1}) E(e^{tY_2}) \dots$$

$$= (pe^t + (1-p))^n \quad \leftarrow \text{using Bernoulli MGF}$$

Multinomial

$$\prod_{i=1}^n p(X_i=x_i) = \frac{n!}{x_1! \cdots x_n!} p_1^{x_1} \cdots p_n^{x_n}$$

$$= \exp \left\{ \log \left(\frac{n!}{x_1! \cdots x_n!} \right) + x_1 \log(p_1) + \cdots + n \log(p_n) \right\}$$

$$E(X_i) = \sum_{x_i=0}^n x_i \frac{n!}{x_1! \cdots x_n!} p_1^{x_1} \cdots p_n^{x_n}$$

$$= n \cdot p_i$$

$$E[X_i] = np_i$$

$$\text{Var}[X_i] = np_i(1-p_i)$$

$$M_X(t) = \left(\sum_{i=1}^n p_i e^{t_i} \right)^n \quad \leftarrow \text{inspired by binomial dist.}$$

Geometric

1. # of failure until first success at prob of success p (include success)

$$P(X=x) = (1-p)^{x-1} p \quad x=1, \dots$$

$$\prod_{i=1}^n P(X_i=x_i) = (1-p)^{\sum x_i - n} p^n$$

$$= \exp\left\{ \sum x_i \log(1-p) + n \log\left(\frac{p}{1-p}\right) \right\}$$

$$\theta = \log(1-p) \quad b(\theta) = -\log\left(\frac{p}{1-p}\right) \quad \text{is } \sum x_i$$

$$\Phi = I \quad \omega = \frac{1}{n}$$

$$\begin{aligned} EX &= \sum_{x=1}^{\infty} x (1-p)^{x-1} p \\ &= p \sum_{x=1}^{\infty} x (1-p)^{x-1} \\ &= p \left(-\frac{\partial}{\partial p} \sum_{x=1}^{\infty} (1-p)^x \right) \\ &= p \left(-\frac{\partial}{\partial p} \sum_{x=1}^{\infty} (1-p)^{x-1} - 1 \right) \\ &= p \left(-\frac{\partial}{\partial p} \frac{1}{p} - 1 \right) \\ &= p \frac{1}{p^2} \\ &= \frac{1}{p} \end{aligned}$$

Recall Geometric Summation:

$$\sum_{x=0}^{\infty} ar^x = \frac{ar}{1-r} \quad 0 < r < 1$$

$$EX = \frac{1}{p}$$

$$VarX = \frac{1-p}{p^2}$$

$$\begin{aligned} Mx(t) &= \sum_{x=1}^{\infty} e^{tx} (1-p)^{x-1} p \\ &= \sum_{x=1}^{\infty} (e^t(1-p))^x \frac{p}{1-p} \\ &= \frac{p}{1-p} \sum_{x=1}^{\infty} (e^t(1-p))^x \\ &= \frac{p}{1-p} \left(\sum_{x=1}^{\infty} (e^t(1-p))^{x-1} - 1 \right) \quad \text{converges when } 0 < e^t(1-p) < 1 \\ &= \frac{p e^t}{1 - (1-p)e^t} \quad \text{when } e^t(1-p) < 1 \\ &\qquad e^t < \frac{1}{1-p} \\ &\qquad t < \ln\left(\frac{1}{1-p}\right) \end{aligned}$$

2. Number of failures before success (does not include success) ★ may use rule of this case waiting time (discrete) to hit the first success (only count for failures)

$$P(X=x) = (1-p)^x p \quad x=0, \dots$$

$$E(X) = \sum_{x=0}^{\infty} x (1-p)^x p$$

$$= p(1-p) \sum_{x=0}^{\infty} x (1-p)^{x-1}$$

$$= p(1-p) \left(-\frac{1}{sp} \sum_{x=0}^{\infty} (1-p)^x\right)$$

$$= p(1-p) \left(-\frac{1}{sp} \left(\frac{1}{1-p}\right)\right)$$

$$= p(1-p) \frac{1}{p^2}$$

$$= \frac{1-p}{p}$$

$$E(X) = \frac{1-p}{p}$$

$$\text{Var}(X) = \frac{1-p}{p^2}$$

$$M_X(t) = E(e^{tx})$$

$$= \sum_{x=0}^{\infty} e^{tx} p (1-p)^x$$

$$= p \sum_{x=0}^{\infty} (e^t (1-p))^x$$

$$= p \cdot \frac{1}{1 - e^t (1-p)}$$

$$= \frac{p}{1 - e^t (1-p)} \quad \text{with } 0 < e^t (1-p) < 1$$

$$t < \log(\frac{1}{1-p})$$

Negative-Binomial

Number of failures to hit the r^{th} success (does not include r)

Waiting time (discrete) to hit the r^{th} success (only count for failures)

$$P(X=x) = \binom{x+r-1}{r-1} p^r (1-p)^x \quad x=0, 1, \dots$$

$$E(X) = \sum_{x=0}^{\infty} x \binom{x+r-1}{r-1} p^r (1-p)^x$$

$$E(X) = r \frac{(1-p)}{p}$$

$$= \sum_{x=1}^{\infty} \frac{(x+r-1)!}{(r-1)!(x-1)!} p^r (1-p)^x$$

$$\text{Var}(X) = r \frac{(1-p)}{p^2}$$

$$= \sum_{x=1}^{\infty} \frac{(x+r-1)!}{r!(x-1)!} \cdot r p^{r+1} \cdot \frac{1}{p} (1-p)^{x-1} (1-p)$$

$$= r \frac{1-p}{p}$$

$$M_X(t) = \left(\frac{p}{1-(1-p)e^t} \right)^r \quad \leftarrow \text{inspired from geometric MGF}$$

Since Geometric(p) \equiv Neg-binomial(1, p)

Hypergeometric * Like binomial but without replacement

$$P(X=k) = \frac{\binom{k}{k} \binom{N-k}{n-k}}{\binom{N}{n}}$$

= k success, $n-k$ fail
among all possible results

N : population size

k : number of success

n : number of draw

k : number of observed success

$$\begin{aligned} E[X] &= \sum_{x=0}^k x \cdot \frac{\binom{k}{x} \binom{N-k}{n-x}}{\binom{N}{n}} \\ &= \sum_{x=1}^k \frac{1}{\binom{N}{n}} \frac{k!}{(x-1)!(k-x)!} \frac{(N-k)!}{(n-x)!(N-n)!} \\ &= \sum_{x=1}^k \frac{1}{\binom{N}{n}} \frac{(k-1)!}{(x-1)!(k-x)!} \frac{k}{(n-x)!(N-n)!} \\ &= \sum_{x=1}^k \frac{(n-1)!(N-n)!}{(n-1)!} \cdot \frac{n}{N} \frac{(k-1)!}{(x-1)!(k-x)!} \frac{k}{(n-x)!(N-n)!} \\ &= \sum_{x=1}^k \frac{(k-1) \binom{N-k}{n-x}}{\binom{N-1}{n-1}} \cdot \frac{n k}{N} \\ &= \frac{n k}{N} \end{aligned}$$

Distribution Theory

Multivariate normal distribution

$$Y \in \mathbb{R}^p \sim N_p(\mu, \Sigma) \Rightarrow V'Y \sim N(V'\mu, V'\Sigma V) \text{ with } V'\Sigma V > 0$$

$$f_Y(y) = \det(2\pi\Sigma)^{-\frac{1}{2}} \exp(-\frac{1}{2}(y-\mu)' \Sigma (y-\mu)) \Rightarrow M_Y(t) = e^{t'\mu + \frac{1}{2}t'\Sigma t}$$

$$E(Y) = \mu$$

$$\text{Var}(Y) = \Sigma$$

t distribution

* $X \sim N(\mu, 1)$ indep. $U \sim \chi^2_p$
 $X / \sqrt{U/p} \sim T_p(\mu)$

$$X \sim N(0, \sigma^2)$$

$$\frac{X}{\sqrt{\sigma^2}} \sim N(0, 1)$$

$$\frac{S^2(p)}{\sigma^2} \sim \chi^2_{p-1} \Rightarrow \frac{S^2}{\sigma^2} \sim \frac{1}{p} \chi^2_p$$

$$\frac{N(0, \sigma^2)}{\sqrt{\chi^2_p/p}} \sim T_p$$

F distribution

① $U_1 \sim \chi^2_{p_1}$ indep. $U_2 \sim \chi^2_{p_2}$
 $F = \frac{U_1/p_1}{U_2/p_2} \sim F_{p_1, p_2}$

② $U_1 \sim \chi^2_{p_1}(\phi)$ indep. $U_2 \sim \chi^2_{p_2}$
 $F = \frac{U_1/p_1}{U_2/p_2} \sim F_{p_1, p_2}(\phi)$

Chi-squared distribution

① $Z_i \stackrel{iid}{\sim} N(0, 1) \Rightarrow Z \sim N_p(0, I)$
 $V : \sum_{i=1}^p Z_i^2 \sim \chi^2_p \Rightarrow M_V(t) = (1-2t)^{-\frac{p}{2}}$
 $E(V) = p$
 $\text{Var}(V) = 2p$

② $Z_i \stackrel{iid}{\sim} N(\mu_i, 1) \Rightarrow Z \sim N_p(\mu, I)$
 $V : \sum_{i=1}^p Z_i^2 \sim \chi^2_p(\frac{1}{2} \sum_{i=1}^p \mu_i^2) \Rightarrow M_V(t) = (1-2t)^{-\frac{p}{2}} e^{\frac{2\phi t}{1-2t}}$ with $t < \frac{1}{2}$
 $E(V) = p + 2\phi$
 $\text{Var}(V) = 2p + 4\phi$

③ $U_i \stackrel{iid}{\sim} \chi^2_{p_i}(\phi_i)$

$$U := \sum_{i=1}^m U_i \sim \chi^2_{\sum p_i}(\sum \phi_i)$$

* ④ $X \sim N_p(u, v) \Rightarrow X'v^{-1}X \sim \chi^2_p(\frac{1}{2}u'v'u)$

Lemma: $A \in \mathbb{R}^{n \times n}$ is symmetric and idempotent matrix with $\text{rank}(A) = s$ iff $A = G\tilde{G}'$ for some matrix $G \in \mathbb{R}^{n \times s}$ with $G'\tilde{G} = I_s$ and $\text{rank}(G) = s$.

Theorem: Let $X \sim N_n(\mu, I)$. If A is a symmetric and idempotent matrix with $\text{rank}(A) = S$, then $X'AX \sim \chi^2_s(\frac{1}{2}\mu'A\mu)$

Proof: $A = GG'$ $\Rightarrow G'X \sim N_s(G'\mu, G'IG) = N_s(G'\mu, I_s)$
 $\Rightarrow XGG'X \equiv XAX \sim \chi^2_s(\frac{1}{2}\mu'G\mu) \equiv \chi^2_s(\frac{1}{2}\mu'A\mu)$

Usage: $Y \sim N(X\beta, \sigma^2 I)$

$$\frac{Y}{\sigma} \sim N\left(\frac{X\beta}{\sigma}, I\right)$$

$$Y'(P_X - P_{X\beta})Y \sim \chi^2_{\text{rank}(P_X - P_{X\beta})} \left(\frac{\beta'(X(P_X - P_{X\beta}))X\beta}{2\sigma^2} \right)$$

Theorem: Let $X \sim N_n(\mu, V)$. If A is a symmetric and AV idempotent with $\text{rank}(AV) = s$, then $X'AX \sim \chi^2_s(\frac{1}{2}\mu'A\mu)$

Proof: $V = LL'$ for V symmetric & P.O. $\Rightarrow Y := L'X \sim N_p(L'\mu, I_p)$

$$X'AX = X'(L')L'ALL'L'X$$

$$= Y'L'ALY$$

$$= Y'BY \sim \chi^2_s\left(\frac{1}{2}(L'\mu)'B(L'\mu)\right)$$

$$\equiv \chi^2_s\left(\frac{1}{2}\mu'A\mu\right)$$

\Rightarrow since $B' = (L'AL)' = L'AL$ is symmetric

$$BB = L'AL'L'AL = L'AVAL = L'AVAL \cdot L'L' =$$

$$= L'AVAV(L')' = L'AV(L')' = L'AL$$

is idempotent

$$\hookrightarrow \text{rank}(B) = \text{trace}(B) = \text{trace}(L'AL) = \text{trace}(LL'A) = \text{trace}(LA) = \text{rank}(AV) = s$$

Theorem: Let $X \sim N_p(\mu, V)$ and A symmetric, $\text{rank}(A) = s$. If $BVA = 0$ for a given B , then BX and $X'AX$ are independent.

Corollary: Let $X \sim N_p(\mu, V)$, let A be a symmetric matrix with $\text{rank}(A) = r$. let B be a symmetric matrix with $\text{rank}(B) = s$. If $BVA = 0$, then $X'BX$ and $X'AX$ are independent.

Cochrane's theorem: Let $Y \sim N_n(\mu, \sigma^2 I_n)$ and let A_i be symmetric and idempotent with $\text{rank}(A_i) = s_i$ $\forall i = (1, \dots, k)$. If $\sum_{i=1}^k A_i = I_n$, then $\frac{1}{\sigma^2} Y'A_i Y \sim \chi^2_{s_i}(\frac{1}{\sigma^2} \mu'A_i\mu)$, $\sum_{i=1}^k s_i = n$, and $\frac{1}{\sigma^2} Y'A_1 Y, \dots, \frac{1}{\sigma^2} Y'A_k Y$ are independent.

Distribution relationship

$$N(0, 1)^2 \sim \chi^2_1 ; \quad N(\mu, 1)^2 \sim \chi^2_{1, (\frac{1}{2}\mu^2)} ; \quad N(\mu, \sigma^2)^2 \sim \sigma^2 \chi^2_{1, (\frac{\mu^2}{\sigma^2})}$$

$$\frac{N(0, 1)}{\sqrt{n-p} \chi^2_{n-p}} \sim t_{n-p} ; \quad t_p^2 \sim F_{1,p} ; \quad \chi^2_a \text{ indep. } \chi^2_b , \quad \frac{\chi^2_a / a}{\chi^2_b / b} \sim F_{a,b}$$

Scale

$$\begin{aligned} \chi^2_k &\equiv \text{Gamma}(\frac{k}{2}, 2) & ; \quad \stackrel{iid}{\sum_i} \chi^2_{ai(bi)} &\equiv \chi^2_{\sum ai(bi)} \\ &\equiv \text{Gamma}(\frac{k}{2}, \frac{1}{2}) \end{aligned}$$

rate

$$\begin{aligned} \stackrel{iid}{\sum_i} \exp(\lambda) &\equiv \text{Gamma}(n, \lambda) & ; \quad \frac{1}{n} \stackrel{iid}{\sum_i} \exp(\lambda) &\equiv \text{Gamma}(n, \frac{\lambda}{n}) \\ \text{scale} = 1/\text{rate} & & &\equiv \text{Gamma}(n, \frac{n}{\lambda}) \\ f_x = \frac{1}{x} e^{-x/\lambda} & \quad \lambda: \text{scale} \quad \lambda: \text{rate} \end{aligned}$$

$$\stackrel{iid}{\sum_i} \text{Gamma}(a_i, b) \equiv \text{Gamma}(\sum_i a_i, b) ; \quad \frac{1}{b} \text{Gamma}(a, b) \equiv \text{Gamma}(a, 1)$$

Poi(μ) indep. Poi(λ), $\text{Pois}(\mu) + \text{Pois}(\lambda) = \text{Pois}(\mu + \lambda)$

For poisson, rate = scale = mean = λ

$$\exp(\beta) = \text{Weibull}(1, \beta) \approx \text{Gamma}(1, \beta)$$

$X \sim \text{Weibull}(a, b) \Leftrightarrow X^a \sim \exp(b^a)$ both scale - liked.

$$X \text{ indep. } Y \sim N(0, 1) ; \quad \frac{X}{Y} \sim \text{Cauchy}(0, 1)$$

$$\text{Beta}(1, 1) = \text{Unif}(0, 1)$$

$$\text{neg-binomial}(1, p) = \text{geometric}(p)$$

Inequalities

Boole's Inequality

$$P(\bigcup_{i=1}^{\infty} E_i) \leq \sum_{i=1}^{\infty} P(E_i)$$

Triangular Inequality

$$\|x+y\| \leq \|x\| + \|y\|$$

Product Inequality

matrix norm

$$\|x \times y\| \leq \|x\| \times \|y\|$$

Markov's Inequality

$$P(x > \varepsilon) \leq \frac{E(x)}{\varepsilon}$$

Chebyshov's Inequality

$$P(|X - E(X)| > \varepsilon) \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

Jensen's inequality

$$E(g(x)) \geq g(E(x)) \quad \text{if } g(\cdot) \text{ is convex}$$

$$\sum \theta_i f(x_i) \geq f(\sum \theta_i x_i)$$

$$E(g(x)) \leq g(E(x)) \quad \text{if } g(\cdot) \text{ is concave}$$

* Equality met when :

1. X is constant almost surely :

$$P(X=c) = 1 \Rightarrow E(\phi(x)) = \phi(E(x))$$

2. ϕ is linear (affine) over the range of x :

$$\phi(x) = ax+b$$

a.k.a. if ϕ is differentiable and convex, then

$$\phi(x) \geq \underbrace{\phi(EX) + \phi'(EX)(x-EX)}$$

{ by Taylor expansion

{ first-order condition in convexity

with equality iff \forall tangent line of $g(\cdot)$ at EX ,

$$P(\phi(x) = \phi(EX) + \phi'(EX)(x-EX)) = 1$$

Cauchy-Schwartz Inequality

$$\text{Cov}(U, V) \leq \sqrt{\text{Var}(U) \text{Var}(V)} \Leftrightarrow (E(XY))^2 \leq E(X^2)E(Y^2) \quad \text{equal iff } x = \alpha y$$

Set Operations

Properties of set operations

Commutativity:

$$A \cup B = B \cup A ; A \cap B = B \cap A ; A \Delta B = B \Delta A$$

$$A \Delta B = (A \setminus B) \cup (B \setminus A) = (A \cup B) \setminus (A \cap B) = (A \cap B^c) \cup (B \cap A^c)$$

Associativity:

$$(A \cup B) \cup C = A \cup (B \cup C) ; (A \cap B) \cap C = A \cap (B \cap C) ;$$

$$(A \Delta B) \Delta C = A \Delta (B \Delta C)$$

Distributivity:

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

$$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$$

De Morgan's law:

$$(A \cup B)^c = A^c \cap B^c ; (A \cap B)^c = A^c \cup B^c$$

Variable transformation

General Methods:

1. By CDF, MGF

2. By pdf / pmf :

1. Law of total probability (when in \mathbb{R})

2. Jacobian (when in \mathbb{R}^F)

↑
(when sign changes)

be aware the range of new parameters ← from the range of odds.

3. By distribution relationship recognition

Example

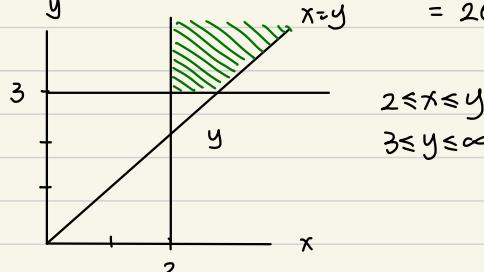
Problem

Find the value of the constant c which make the following function a valid probability density function of (X, Y) . And then calculate $P(X \geq 2, Y \geq 3)$

$$f(x, y) = ce^{-x-y} I_{(0 \leq x \leq y)}$$

$$\begin{aligned} & \iint_{\mathbb{R}^2} f(x, y) dx dy \stackrel{\text{set}}{=} 1 \\ &= \int_0^\infty \int_0^y ce^{-x-y} dx dy \\ &= \int_0^\infty c e^{-y} (-e^{-x}) \Big|_0^y dy \\ &= \int_0^\infty c e^{-y} (1 - e^{-y}) dy \\ &= c \int_0^\infty e^{-y} - e^{-2y} dy \\ &= c (-e^{-y} + \frac{1}{2} e^{-2y}) \Big|_0^\infty \\ &= c (-0 + 0 + 1 - \frac{1}{2}) \\ &= \frac{c}{2} \stackrel{\text{set}}{=} 1 \\ &\Rightarrow c = 2 \end{aligned}$$

$$\begin{aligned} P(X \geq 2, Y \geq 3) &= \int_3^\infty \int_2^y 2e^{-x-y} dx dy \\ &= \int_3^\infty 2e^{-y} (-e^{-x}) \Big|_2^y dy \\ &= \int_3^\infty 2e^{-y} (-e^{-y} + e^{-2}) dy \\ &= -2 \int_3^\infty e^{-2y} dy + 2e^{-2} \int_3^\infty e^{-y} dy \\ &= -2 \cdot (-\frac{1}{2}) e^{-2y} - 2e^{-2} \Big|_3^\infty \\ &= e^{-2y} - 2e^{y-2} \Big|_3^\infty \\ &= -e^{-6} + 2e^{-5} \\ &= 2e^{-5} - e^{-6} \end{aligned}$$



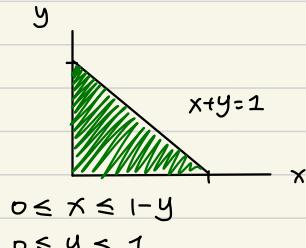
Problem

When a random vector (X, Y) has the following joint pdf

$$f(x, y) = 120xy(1-x-y)I_{(x \geq 0, y \geq 0, x+y \leq 1)}$$

, find the value of $\mathbb{E}(XY)$

$$\begin{aligned}
 E(XY) &= \iint_{\substack{x,y \\ x+y \leq 1}} xy f(x,y) dy dx \\
 &= \int_0^1 \int_0^{1-y} xy \cdot 120xy(1-x-y) dx dy \\
 &= \int_0^1 \int_0^{1-y} 120x^2y^2(1-x-y) dx dy \\
 &= \int_y \int_x^{1-y} 120x^2y^2 dx dy - \int_y^1 \int_x^{1-y} 120x^3y^2 dx dy \\
 &\quad - \int_y \int_x^{1-y} 120x^2y^3 dx dy \\
 &= \int_y 40x^3y^2 \Big|_x^{1-y} dy - \int_y 30x^4y^2 \Big|_x^{1-y} dy - \int_y 40x^3y^3 \Big|_x^{1-y} dy \\
 &= \int_y 40(1-y)^3y^2 - 30(1-y)^4y^2 - 40(1-y)^3y^3 dy \\
 &= \int_y (1-y)^3y^2(40 - 30(1-y) - 40y) dy \\
 &= \int_y (1-y)^3y^2(10 - 10y) dy \\
 &= 10 \int_y (1-y)^4y^2 dy \quad t = 1-y \quad dt = -dy \\
 &= -10 \int_t^1 t^4(1-t)^2 dt \quad y = 1-t \\
 &= -10 \int_t^1 t^4(1-2t+t^2) dt \quad 0 \leq y \leq 1 \\
 &= -10 \int_t^1 t^4 - 2t^5 + t^6 dt \quad t \downarrow \quad t \downarrow \\
 &= -10 \left(\frac{1}{5}t^5 - \frac{1}{3}t^6 + \frac{1}{7}t^7 \right) \Big|_1^0 \\
 &= 10 \left(\frac{1}{5} - \frac{1}{3} + \frac{1}{7} \right) \\
 &= 10 \left(\frac{21}{105} - \frac{35}{105} + \frac{15}{105} \right) \\
 &= 10 \left(\frac{11}{105} \right) \\
 &= \frac{110}{105} \\
 &= \frac{22}{21}
 \end{aligned}$$



Question

If a random vector $(X_1, X_2, X_3)^\top$ has the following joint pdf,

$$f(x_1, x_2, x_3) = 120x_1(1 - x_1 - x_2 - x_3)I_{(x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_1 + x_2 + x_3 \leq 1)}$$

calculate $\mathbb{E}(X_1 X_2)$

$$\begin{aligned} \mathbb{E}(X_1 X_2) &= \int_0^1 \int_0^{1-x} \int_0^{1-x-y} xy \cdot 120x(1-x-y-z) dz dy dx \\ &= \int_0^1 120x^2 \int_0^{1-x} y \int_0^{1-x-y} (1-x-y-z) dz dy dx \\ &= \int_0^1 120x^2 \int_0^{1-x} y (z - xz - yz - \frac{1}{2}z^2) \Big|_0^{1-x-y} dy dx \\ &= \int_0^1 120x^2 \int_0^{1-x} y (1-x-y - x(1-x-y) - y(1-x-y) - \frac{1}{2}(1-x-y)^2) dy dx \\ &= \int_0^1 120x^2 \int_0^{1-x} y (1-x-y - x + x^2 + xy - y + xy + y^2 - \frac{1}{2}(1-2xy - 2x-2y + x^2 + y^2)) dy dx \\ &= \int_0^1 120x^2 \int_0^{1-x} -2xy + x^2 + 2xy^2 - y^2 + y^3 - \frac{1}{2}y + xy^2 + xy + y^2 - \frac{1}{2}x^2y - \frac{1}{2}y^3 dy dx \\ &= -xy + \frac{1}{2}x^2y + 3xy^2 + \frac{1}{2}y^3 - \frac{1}{2}y \end{aligned}$$

easy to start with z
last to integrate x

⋮

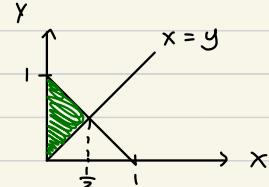
Problem

When a random vector (X, Y) has the following joint pdf

$$f_{1,2}(x, y) = 10xy^2 I_{(0 < x < y < 1)}$$

calculate $P(X + Y \leq 1)$

$$\begin{aligned} P(X + Y \leq 1) &= \int_0^{\frac{1}{2}} \int_X^{1-x} 10xy^2 dy dx \\ &= \int_0^{\frac{1}{2}} \frac{10}{3}xy^3 \Big|_X^{1-x} dx \\ &= \int_0^{\frac{1}{2}} \frac{10}{3}(x(1-x)^3 - x^4) dx \\ &= \int_0^{\frac{1}{2}} \frac{10}{3}(x(1-2x+x^2)(1-x) - x^4) dx \\ &= \int_0^{\frac{1}{2}} \frac{10}{3}(x - 3x^2 + 3x^3 - 2x^4) dx \\ &= \frac{5}{3}x^2 - \frac{10}{3}x^3 + \frac{5}{2}x^4 - \frac{4}{3}x^5 \Big|_0^{\frac{1}{2}} \end{aligned}$$



$$\begin{aligned} &\left\{ \begin{array}{l} x < y < 1-x \\ 0 < x < \frac{1}{2} \end{array} \right. \\ &+ \left\{ \begin{array}{l} 0 < x < y \\ 0 < y < \frac{1}{2} \\ 0 < x < 1-y \\ \frac{1}{2} < y < 1 \end{array} \right. \end{aligned}$$

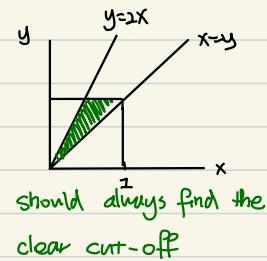
Problem

When a random vector (X, Y) has the following joint pdf

$$f_{1,2}(x, y) = cx^2yI_{(0 < x < y < 1)}$$

find the constant c , and calculate $P(Y \leq 2X)$

$$\begin{aligned} \int_0^y f(x,y) dx dy &\stackrel{\text{set}}{=} 1 & P(Y \leq 2X) \\ &= C \int_0^1 \int_0^y x^2 y dx dy &= \int_0^1 \int_{\frac{y}{2}}^y 15x^2 y dx dy \\ &= C \int_0^1 y \int_0^y x^2 dx dy &= \int_0^1 5x^3 y \Big|_{\frac{y}{2}}^y dy \\ &= C \int_0^1 y \frac{1}{3} x^3 \Big|_0^y dy &= \int_0^1 5y^4 - \frac{5}{8}y^4 dy \\ &= \frac{C}{3} \int_0^1 y^4 dy &= y^5 - \frac{1}{8}y^5 \Big|_0^1 \\ &= \frac{C}{15} y^5 \Big|_0^1 &= \frac{7}{8} \\ &= \frac{C}{15} \stackrel{\text{set}}{=} 1 & \\ C &= 15 \end{aligned}$$



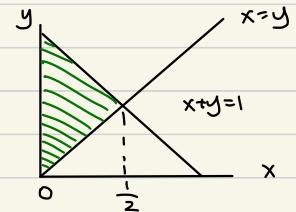
Problem

When a random vector (X, Y) has the following joint pdf

$$f_{1,2}(x, y) = e^{-y} I_{(0 < x < y < \infty)}$$

calculate $P(X + Y \leq 1)$

$$\begin{aligned} P(X + Y \leq 1) &= \int_0^{1/2} \int_x^{1-x} e^{-y} dy dx \\ &= \int_0^{\frac{1}{2}} -e^{-y} \Big|_x^{1-x} dx \\ &= \int_0^{\frac{1}{2}} -e^{-(1-x)} + e^{-x} dx \\ &= -e^{-(1-x)} - e^{-x} \Big|_0^{\frac{1}{2}} \\ &= -e^{-\frac{1}{2}} - e^{-\frac{1}{2}} + e^{-1} + 1 \\ &= 1 + e^{-1} - 2e^{-\frac{1}{2}} \end{aligned}$$



Reparameterization & BLUE

Reparameterization

Definition (Reparameterization): Let $X \in \mathbb{R}^{n \times p}$ and $w \in \mathbb{R}^{n \times q}$

The linear models $Y = X\beta + u$ and $Y = wr + u$ are reparameterizations of each other, if
 $\text{col}(X) = \text{col}(w)$

Theorem: Suppose $\text{col}(X) = \text{col}(w)$, then $P_X = P_w$.

Proposition: If $\text{col}(X) = \text{col}(w)$, then \exists a matrix S s.t. $X = wS$

Similarly, $\exists T$ s.t. $w = XT$

Theorem: Let $\text{col}(X) = \text{col}(w)$.

$$W = XT$$

Suppose $w = XT$ and \hat{r} solves the normal equations in w . is when $\text{col}(X) \subseteq \text{col}(w)$

1). The fitted data & residuals are the same:

$$\begin{aligned}\hat{y} &= P_X y = P_w y \leftarrow \underset{\hat{\beta}_w}{X(X^T X)^{-1} X^T y} = \underset{\hat{\beta}_w}{w(W^T W)^{-1} W^T y} \\ \hat{e} &= (I - P_X)y = (I - P_w)y\end{aligned}$$

2). $\hat{\beta} = T \hat{r}$ solves the normal equations in X

$$Y = X\beta = WR$$

$$= XT\hat{r} \Rightarrow \hat{\beta} = T\hat{r}$$

Theorem: If $\lambda'\beta$ is estimable in the model with design X , and \hat{r} solves the normal equation in design w , then $\lambda' T \hat{r}$ is the least square estimator of $\lambda'\beta$.

Theorem: If $\delta'r$ is estimable in design w ; then $\delta'S\beta$ is estimable in design X , and its least square estimator is $\delta'T\hat{r}$, where $\hat{r} \in \{w^Tw = w'y\}$

$$\begin{aligned}Y &= Wr = X\beta \\ &= WSB \\ \Rightarrow \hat{r} &= S\hat{\beta}\end{aligned}$$

Constraint model

$$\text{col}(x') \perp \text{col}(c')$$

Lemma: Suppose $C \in \mathbb{R}^{(p-r) \times p}$, $r := \text{rank}(x)$, and $\text{col}(x') \cap \text{col}(c') = \{0\}$.

Then the following systems are equivalent:

$$(1) \quad \begin{pmatrix} x'x \\ C'C \end{pmatrix} \beta = \begin{pmatrix} x'y \\ 0 \end{pmatrix}$$

$$(2) \quad \begin{pmatrix} x'x \\ C \end{pmatrix} \beta = \begin{pmatrix} x'y \\ 0 \end{pmatrix}$$

$$(3) \quad (x'x + C'C)\beta = x'y$$

Theorem: Let $C \in \mathbb{R}^{(p-r) \times p}$ with $\text{rank}(C) = p-r$ and $\text{col}(x') \cap \text{col}(c') = \{0\}$. Then:

(1) The matrix $x'x + C'C$ is non-singular

(2) $(x'x + C'C)^{-1}x'y$ uniquely solves $x'x\beta = x'y$, $C'\beta = 0$

(3) $(x'x + C'C)^{-1}$ is a generalized inverse of $x'x$

(4) $C(x'x + C'C)^{-1}x' = 0$

(5) $C(x'x + C'C)^{-1}C' = I$

Definition: The function $\lambda'\beta$ is estimable in the restricted model if and only if there exists a scalar c and vector a

$$\text{s.t. } E(c+a'y) = \lambda'\beta, \forall \beta \in \{\beta : P'\beta = 0\}$$

$$\begin{pmatrix} x'x & P \\ P' & 0 \end{pmatrix} \begin{pmatrix} \beta \\ \theta \end{pmatrix} = \begin{pmatrix} x'y \\ 0 \end{pmatrix}$$

Restricted Model: $y = x\beta + u$, $E(u) = 0$, $\beta \in \{\beta : P'\beta = 0\}$
full-rank

Remember
this

Theorem: In the restricted model, $c+a'y$ is unbiased for $\lambda'\beta$ if and only if $\exists d$

$$\text{s.t. } \lambda = x'a + Pd \text{ and } c = d's.$$

$$(\Rightarrow) \quad E(c+a'y) =$$

$$= d's + a'x\beta$$

$$= d's + a'x\beta + d'P'\beta - d'P'\beta$$

$$= d's + \lambda'\beta - d's$$

$$= \lambda'\beta$$

$$(\Leftarrow) \quad \lambda'\beta = d'P'\beta + a'x\beta$$

$$= d's + a'E(y)$$

Theorem: If $\delta \in \text{col}(P')$, then there exists a solution to the RNEs.

$$\text{restriction: } P'\beta = \delta \Leftrightarrow \delta \in \text{col}(P')$$

Theorem: If $\hat{\beta}_H$ denotes the first component of a solution to true RNE's:

$$\begin{pmatrix} \hat{\beta}_H \\ \theta_H \end{pmatrix} \in \left\{ \begin{pmatrix} X'X & P \\ P' & 0 \end{pmatrix} \begin{pmatrix} \beta \\ \theta \end{pmatrix} = \begin{pmatrix} X'y \\ \delta \end{pmatrix} \right\}$$

then $\hat{\beta}_H$ minimizes $Q(\beta) = \|y - X\beta\|^2$ over $\{P\beta = \delta\}$

Theorem: If $\hat{\beta}_H$ denotes the first component of a solution to true RNE's:

$$\begin{pmatrix} \hat{\beta}_H \\ \theta_H \end{pmatrix} \in \left\{ \begin{pmatrix} X'X & P \\ P' & 0 \end{pmatrix} \begin{pmatrix} \beta \\ \theta \end{pmatrix} = \begin{pmatrix} X'y \\ \delta \end{pmatrix} \right\}$$

if $\beta \in \{P\beta = \delta\}$, then $Q(\beta) = Q(\hat{\beta}_H)$ iff β also solves RNE.

Not only RNE gives a L.S. solution, but a L.S. solution MUST solves RNE.

Best Linear Unbiased Estimator

* a.k.a unbiased estimator that achieves CRLB.

Gauss-Markov Model

$$\begin{aligned} Y &= X\beta + U \\ E(U) &= 0 \\ \text{Cov}(U) &= \sigma^2 I \quad \text{for some } \sigma > 0 \end{aligned} \quad \left. \right\} \text{Gauss-Markov Assumption}$$

Theorem (Gauss-Markov): Under the Gauss-Markov assumptions/model, if $X\beta$ is estimable, then $X'\hat{\beta}$ is the best (minimum variance) linear unbiased estimator of $X\beta$, $\forall \hat{\beta} \in \{X'X\beta = X'y\}$.

(For any unbiased linear estimator, the one with least variance is the LSE).

Theorem: An unbiased estimator of σ^2 is $y'(I - P_X)y / (n-p)$

Aitken Model

$$Y = X\beta + U$$

$$E(U) = 0$$

$$\text{Var}(U) = \sigma^2 V$$

$V > 0$ known

$$W'VW > 0 \quad \forall w \in \mathbb{R}^n$$

Theorem (AITKEN)

So this is a Gauss-Markov Model. Consequently, the BLUE for any estimable $\lambda'\beta$ is $\hat{\beta}_{\text{GLS}}$ where $\hat{\beta}_{\text{GLS}} \in \{X'V^{-1}X\beta = X'V^{-1}Y\}$

Remember
this \Rightarrow

Theorem: The estimator $t'y$ is BLUE for $E(t'y)$ if and only if $t'y$ is unrelated with all unbiased estimators of zero.

Corollary: Under the Aitken Model, $t'y$ is BLUE for $E(t'y)$ iff $\forall t \in \text{col}(X)$

Theorem: Under the Aitken model, $\hat{\beta}_{\text{OLS}}$ is the BLUE for estimable $\lambda'\beta$ iff

$$\exists Q \text{ s.t. } VQ = XQ$$

Bayes prior & posterior

Conjugate Priors

Prior	Data	Posterior
Beta	Bernoulli	Beta
Beta	Binomial	Beta
Beta	Geometric	Beta
Beta	Neg-Binomial	Beta
Gamma	Poisson	Gamma
Gamma	exponential	Gamma
Gamma	Gamma	Gamma
Dirichlet	Multinomial	Dirichlet
Dirichlet	Hypergeometric	Dirichlet
Normal	Normal	Normal

think parameters in these are p.
with $0 < x < 1$ continuous fits well.

may think mirroring the Beta

remember those two separately

Convergence Orders

Notation	Deterministic or Stochastic?	Formal Definition	Intuition	Example
$O(b_n)$	Deterministic	$\exists C > 0 \text{ s.t. } \limsup_{n \rightarrow \infty} a_n \leq C$	a_n	$n = O(n^2)$
$o(b_n)$	Deterministic	$\frac{a_n}{b_n} \rightarrow 0 \text{ as } n \rightarrow \infty$	Negligible compared to b_n	$n = o(n^2)$
$O_p(a_n)$	Stochastic	$\forall \varepsilon > 0, \exists M > 0 \text{ s.t. } \Pr(X_n - a_n > M) < \varepsilon \text{ for large } n$	$X_n - a_n$	$X_n \sim N(0, 1/n^2), \text{ then } X_n = o_p(n^{-1/2})$
$o_p(a_n)$	Stochastic	$\frac{X_n - a_n}{a_n} \xrightarrow{P} 0$	Small compared to a_n with high probability	$X_n \sim N(0, 1/n^2), \text{ then } X_n = o_p(n^{-1/2})$

Ancillary & Completeness

Show ancillary by family

Location family ①

$$f(x; \mu) = f_0(x - \mu)$$

$N(\mu, \sigma^2)$ with known σ^2

Laplace with fixed scale

Cauchy with fixed scale

Uniform $(\mu, \mu + c)$

Logistic with fixed scale

Scale family ②

$$f(x; \sigma) = \frac{1}{\sigma} f_0\left(\frac{x}{\sigma}\right) : \sigma \text{ scale}$$

exponential

chi-squared with ϕ

Gamma with fixed shape

Weibull with fixed shape

location & Scale family ③

$$f(x; \mu, \sigma) = \frac{1}{\sigma} f_0\left(\frac{x-\mu}{\sigma}\right)$$

$N(\mu, \sigma^2)$

Cauchy

Laplace

Logistic

Uniform $(\mu - c\sigma, \mu + c\sigma)$

Location family ancillary statistics ①

• Sample spacing : $(X_{(2)} - X_{(1)}, X_{(3)} - X_{(2)}, \dots, X_{(n)} - X_{(n-1)})$

• Sample range : $X_{(n)} - X_{(1)}$

• Sample variance : $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$

Example 1:

Location family, want to show $X_{(n)} - X_{(1)}$ ancillary.

Show $X - \theta \sim \exp(1)$ invariant of θ , thus $(X_{(n)} - \theta) - (X_{(1)} - \theta)$ is also invariant of θ

Thus $X_{(n)} - X_{(1)}$ is ancillary.

Scale family ancillary statistics ②

• Sample ratio : $(\frac{X_{(2)}}{X_{(1)}}, \dots, \frac{X_{(n-1)}}{X_{(n)}})$

• The t-statistics : $\frac{\bar{X}_n}{S_n}$ OR robust $\tilde{t} = \frac{\bar{X}_n}{X_{(\frac{n}{4})} - X_{(\frac{n}{4})}}$

medium

location-scale family ancillary statistics ③

• normalized sample spacing : $(\frac{X_{(2)} - X_{(1)}}{X_{(n)} - X_{(1)}}, \dots, \frac{X_{(n)} - X_{(n-1)}}{X_{(n)} - X_{(1)}})$

• sample range / sample s.d. ratio : $\frac{X_{(n)} - X_{(1)}}{S_n}$

Uniqueness of power series \star use this to proof completeness.

Example 1

$X_i \stackrel{iid}{\sim} \text{Poisson}(\lambda)$; $\sum X_i \stackrel{iid}{\sim} \text{Poisson}(n\lambda)$

$$E(g(t)) = \sum_{t=0}^{\infty} g(t) e^{-n\lambda} \frac{(n\lambda)^t}{t!} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sum_{t=0}^{\infty} g(t) \frac{(n\lambda)^t}{t!} = 0$$

$$\Rightarrow \sum_{t=0}^{\infty} \frac{g(t)}{t!} (n\lambda)^t = 0$$

$\Rightarrow \frac{g(t)}{t!} = 0 \quad \forall t$ by uniqueness of power series

$$\Rightarrow g(t) = 0 \quad \forall t$$

Example 2

$X_i \stackrel{iid}{\sim} \exp(\lambda)$; $\sum X_i \stackrel{iid}{\sim} \text{Gamma}(n, \lambda)$

$$E(g(x)) = \int_0^{\infty} g(x) \frac{\lambda^x}{\Gamma(n)} x^{n-1} e^{-\lambda x} dx \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \int_0^{\infty} g(x) x^{n-1} e^{-\lambda x} dx = 0$$

$$\frac{\partial}{\partial \lambda} \Rightarrow \int_0^{\infty} g(x) x^{n-1} (-x) e^{-\lambda x} dx = 0 \Rightarrow \int_0^{\infty} g(x) (-x^n) e^{-\lambda x} dx = 0$$

$$\frac{\partial^2}{\partial \lambda^2} \Rightarrow \int_0^{\infty} g(x) x^{n+1} e^{-\lambda x} dx = 0$$

:

thus $g(x) e^{-\lambda x} = 0$ everywhere $\forall x$

thus $g(x) = 0$

Asymptotic Theory

Cramér-Rao Lower bound (RCLB)

← called Cramér condition

Under regularity conditions, the inverse of the Fisher information is attained:

1. Model is correctly specified
2. The log-likelihood is sufficiently smooth (e.g. differentiable up to second order)
3. The Fisher information is positive definite and finite
4. The MLE exists and consistent
5. The score function has mean zero and finite variance

* Violation: Score has to be a linear function of an unbiased estimator.

But for MLE estimators, the CRLB always met asymptotically.

Wilk's Theorem assumptions

$$H_0: \theta \in \Theta_0 \quad H_1: \theta \in \Theta \setminus \Theta_0$$

$$\Lambda = \frac{\sup_{\theta \in \Theta_0} L(\theta)}{\sup_{\theta \in \Theta} L(\theta)} \quad -2 \log \Lambda \xrightarrow{d} \chi_k^2 ; k = \dim(\Theta) - \dim(\Theta_0)$$

by Wilk's theorem

Assumptions:

1. $\theta \in \Theta$ uniquely determine the distribution of data

$$f(x|\theta_1) \neq f(x|\theta_2) \text{ if } \theta_1 \neq \theta_2$$

2. True parameter θ_0 must lie in the interior of parameter space Θ .

If on the boundary, the limiting distribution might not be chi-square.

3. Regularity of Likelihood

Twice differentiable in θ

Well-defined, non-singular Fisher information

Expectation of derivatives exists and can be interchanged with integration

4. Large Sample $n \rightarrow \infty$ to use Wilk

5. MLE behavior Consistent & asymptotic normal

$$\hat{\theta}_{\text{MLE}} \xrightarrow{P} \theta_0 ; \sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, I'(\theta_0))$$

6. Nested Model

$\Theta_0 \subseteq \Theta$ ← ensures the likelihood ratio is well-defined

Violations: zero-inflation Poisson & mixture distribution

Cannot test for $H_0: \pi = 0, 1$; $H_0: p = \text{any value } \in [0, 1]$

Can test for $H_0: \pi = \text{some value } \in (0, 1)$

Law of large number (weak)

Suppose $X \sim D(\cdot)$, $E(X) = \mu$, $\text{Var}(X) = \sigma^2 < \infty$
Thus $\bar{X} \xrightarrow{P} \mu$

Convergence in Probability

$$P(|X_n - X| > \varepsilon) \rightarrow 0 \quad \text{can be some scalar}$$

1. By definition

2. Use WLLN (But must know at least mean & variance)

3. If $\lim \text{Var}(X_n) = 0$, and $\lim E(X_n) = a$, then $X_n \xrightarrow{P} a$

$$P(|X_n - a| > \varepsilon) \leq \frac{E(X_n - a)^2}{\varepsilon^2}$$

when can't get rid of 1.1, $= \frac{E(X_n^2) - 2E(X_n)a + a^2}{\varepsilon^2}$

use $E(1 \cdot 1^2)$ or $\text{Var}(1 \cdot 1)$, otherwise, $= \frac{E(X_n^2) - 2E(X_n)a + a^2}{\varepsilon^2} \rightarrow \text{Var}(X_n)$

use probability method.

$$\rightarrow \frac{\text{Var}(X_n)}{\varepsilon^2}$$

$$\rightarrow 0$$

Convergence in distribution

$X_n \xrightarrow{d} X$ if $F_{X_n} \rightarrow F_X \forall x$

1. Central Limit Theorem (CLT)

2. Delta theorem

3. Slutsky theorem (usually CLT + WLLN)

4. By definition

Central Limit Theorem

1. For sample mean

$$X_i \stackrel{iid}{\sim} D \quad E(X_i) = \mu \quad \text{Var}(X_i) = \sigma^2 < \infty$$

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2)$$

2. For some estimator

$$E(\hat{\theta}_n) = \theta_0 \text{ or bias} = O\left(\frac{1}{\sqrt{n}}\right) \quad \hat{\theta}_n \xrightarrow{P} \theta_0 \quad \sigma^2 < \infty$$

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, \sigma^2)$$

Delta Theorem

$$\sqrt{n}(X - \mu) \xrightarrow{d} N(0, \sigma^2)$$

Consider $g(x)$ differentiable and $g'(\mu) \neq 0$.

$$\sqrt{n}(g(X) - g(\mu)) \xrightarrow{d} N(0, g'(\mu)^2 \sigma^2)$$

Slutsky Theorem

Suppose $X_n \xrightarrow{d} X$, $Y_n \xrightarrow{P} c$, then

$$\cdot aX_n + bY_n \xrightarrow{d} aX + bc$$

$$\cdot X_n Y_n \xrightarrow{d} cX$$

Continuous Mapping Theorem

$\{X_n\}$ be a k -dimensional random vectors. Let $g: \mathbb{R}^k \rightarrow \mathbb{R}^k$ be a continuous function.

$$\text{Then } X_n \xrightarrow{P} X \Rightarrow g(X_n) \xrightarrow{P} g(X)$$

$$X_n \xrightarrow{d} X \Rightarrow g(X_n) \xrightarrow{d} g(X)$$

$$X_n \xrightarrow{as} X \Rightarrow g(X_n) \xrightarrow{as} g(X)$$

Examples:

Question

Let's define $Y_n \triangleq \min_{1 \leq i \leq n} X_i$, where $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f(x; \theta) = e^{-(x-\theta)} I_{(\theta, \infty)}(x)$ then,

$$Y_n \xrightarrow[n \rightarrow \infty]{P} \theta$$

$$\begin{aligned} & P(|Y_n - \theta| > \varepsilon) \\ &= P(|\min X_i - \theta| > \varepsilon) \\ &= P(\min X_i > \varepsilon + \theta) \quad \text{as } \theta < X_i \\ &= P(X_i > \varepsilon + \theta)^n \\ &= \left(\int_{\varepsilon+\theta}^{\infty} e^{-(x-\theta)} dx \right)^n \\ &= \left(e^\theta (-e^{-x}) \Big|_{\varepsilon+\theta}^{\infty} \right)^n \\ &= \left(e^\theta e^{-(\varepsilon+\theta)} \right)^n \\ &= e^{-n\varepsilon} \longrightarrow 0 \quad \text{as } n \rightarrow \infty \end{aligned}$$

Question

If $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Beta}(1, \alpha)$, what is the pdf of the limit distribution of $n^{1/\alpha}(1 - \max_{1 \leq i \leq n} X_i)$?

$$\begin{aligned}
 & P(n^{\frac{1}{\alpha}}(1 - \max X_i) \leq x) \leftarrow \text{CDF} \\
 & = P(1 - \max X_i \leq x \cdot n^{-\frac{1}{\alpha}}) \\
 & = P(\max X_i \geq 1 - x \cdot n^{-\frac{1}{\alpha}}) \\
 & = 1 - P(\min X_i \leq 1 - x \cdot n^{-\frac{1}{\alpha}}) \\
 & = 1 - P(X_i \leq 1 - x \cdot n^{-\frac{1}{\alpha}})^n \\
 & = 1 - (1 - (1 - \underbrace{(1 - \frac{x}{n^{\frac{1}{\alpha}}})^{\alpha}}_{\rightarrow \frac{x^\alpha}{n}}))^n \\
 & \longrightarrow 1 - e^{-x^\alpha}
 \end{aligned}$$

$$\begin{aligned}
 & X \sim \text{Beta}(1, \alpha) \\
 & f_X(x) = \frac{\Gamma(\alpha+1)}{\Gamma(1)\Gamma(\alpha)} (x)^\alpha (1-x)^{\alpha-1} \\
 & = \alpha \cdot (1-x)^{\alpha-1} \\
 & F_X(x) = \int_0^x \alpha (1-x)^{\alpha-1} dx \\
 & = -(1-x)^\alpha \Big|_0^x \\
 & = 1 - (1-x)^\alpha \quad 0 < x < 1
 \end{aligned}$$

$$\begin{aligned}
 f_X(x) &= \frac{\partial}{\partial x} (1 - e^{-x^\alpha}) \\
 &= \frac{\partial g}{\partial h} \frac{\partial h}{\partial x} \\
 &= (-1)(-e^{-x^\alpha}) \alpha x^{\alpha-1} \\
 &= \alpha x^{\alpha-1} e^{-x^\alpha} \quad I(0 < x < \infty)
 \end{aligned}$$

cheat sheet

Law of total expectation (and related definitions)

$$\begin{aligned} E(Y) &= E_x(E(Y|X)) \\ &= P(X>x)E(Y|X>x) + P(X\leq x)E(Y|X\leq x) \\ &= \int_0^\infty P(Y>t)dt \end{aligned}$$

Law of total variance

$$\begin{aligned} \text{Var}(Y) &= E_x(\text{Var}(Y|X)) + \text{Var}_x(E(Y|X)) \\ &= E_y(Y - EY)(Y - EY) \\ \text{Cov}(X, Y) &= \text{Cov}(X, E(Y|X)) \\ &= E_x(XE(Y|X)) - E(X)E(Y) \end{aligned}$$

Taylor Expansion:

$$\begin{aligned} g(x) &= \sum_{k=0}^{\infty} \frac{g^{(k)}(a)(x-a)^k}{k!} \\ \Rightarrow g(x) &\geq g(a) + \nabla g(a)^T(x-a) \quad \text{when } g(\cdot) \text{ convex and differentiable.} \end{aligned}$$

Geometric summation:

$$\begin{aligned} \ln(1+x) &= \sum_{n=0}^{\infty} \frac{(-1)^n}{n+1} x^{n+1} \quad \text{with } x \in (-1, 1] \\ \text{two instances} \quad \curvearrowright \frac{1}{1-x} &= \sum_{n=0}^{\infty} x^n \quad \text{with } x \in (-1, 1) \\ \frac{1}{1+x} &= \sum_{n=0}^{\infty} (-x)^n \quad \text{with } x \in (-1, 1) \\ e^x &= \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n = \sum_{k=0}^{\infty} \frac{x^k}{k!} \\ \sum_{i=1}^n ar^{i-1} &= \frac{a(1-r^n)}{1-r} \Rightarrow \sum_{i=1}^{\infty} ar^{i-1} = \frac{a}{1-r} \\ \sum_{i=1}^n (zi-1) &= n^2 \end{aligned}$$

Useful Property Given $Z \sim E(Z)=\mu, \text{Var}(Z)=V ; E(Z'AZ) = \text{tr}(AV) + \mu'A\mu$

$$\begin{aligned} \uparrow \quad E(Z'AZ) &= E((Z-\mu+\mu)'A(Z-\mu+\mu)) \\ \text{use to find the} \quad &= E((Z-\mu)'A(Z-\mu) + \mu'A(Z-\mu) + (Z-\mu)'A\mu + \mu'\mu A) \\ \text{expectation of SSR} \quad &= E((Z-\mu)'A(Z-\mu)) + 0 + 0 + \mu'A\mu \\ \text{or SSE in terms} \quad &= E(\text{trace}((Z-\mu)'A(Z-\mu))) + \mu'A\mu \\ \text{that } Y'A Y &= E(\text{trace}(A(Z-\mu)(Z-\mu)')) + \mu'A\mu \\ &= \text{tr}(E(A(Z-\mu)(Z-\mu)')) + \mu'A\mu \\ &= \text{tr}(AV) + \mu'A\mu \end{aligned}$$

Gamma function

$$\begin{aligned}\Gamma(x) &= (x-1)\Gamma(x-1) \\ &= (x-1)!\end{aligned}$$

$$\Gamma(\frac{1}{2}) = \sqrt{\pi}$$

$$\Gamma(\frac{3}{2}) = \frac{1}{2} \cdot \Gamma(\frac{1}{2}) = \frac{\sqrt{\pi}}{2}$$

MGF for Normal

$$M_X(t) = e^{\frac{1}{2}t^2\sum t + t' \mu}$$

MGF for χ^2

$$M_X(t) = (1-2t)^{-\frac{p}{2}} e^{\frac{z\phi t}{1-2t}}$$

$$E(X) = p + z\phi$$

$$\text{Var}(X) = 2p + 4\phi$$

Theorem (Gradient)

$$a, b \in \mathbb{R}^n \quad A \in \mathbb{R}^{n \times n}$$

$$(a). \nabla_b a^T b = a$$

$$(b). \nabla_b b^T A b = Ab + A^T b = (A + A^T)b$$

Sequential Sum of Square

If intercept : P_1

TABLE 7.1 ANOVA Table for Sequential SS

Source	df	Projection	SS	noncentrality
b_0	$r(X_0)$	P_{X_0}	$R(b_0)$	$(2\sigma^2)^{-1}(Xb)^T P_{X_0}(Xb)$
b_1 after b_0	$r(X_1^*) - r(X_0)$	$P_{X_1^*} - P_{X_0}$	$R(b_0, b_1) - R(b_0)$	$(2\sigma^2)^{-1}(Xb)^T (P_{X_1^*} - P_{X_0})(Xb)$
...	...			
b_j after b_0, \dots, b_{j-1}	$r(X_j^*) - r(X_{j-1}^*)$	$P_{X_j^*} - P_{X_{j-1}^*}$	$R(b_0, \dots, b_j) - R(b_0, \dots, b_{j-1})$	$(2\sigma^2)^{-1}(Xb)^T (P_{X_j^*} - P_{X_{j-1}^*})(Xb)$
...	...			
b_k after b_0, \dots, b_{k-1}	$r(X_k^*) - r(X_{k-1}^*)$	$P_{X_k^*} - P_{X_{k-1}^*}$	$R(b_0, \dots, b_k) - R(b_0, \dots, b_{k-1})$	$(2\sigma^2)^{-1}(Xb)^T (P_X - P_{X_{k-1}^*})(Xb)$
Error	$N - r(X)$	$I - P_X$	$y^T y - R(b)$	0
Total	N	I	$y^T y$	$(2\sigma^2)^{-1}(Xb)^T (Xb)$