

THE AITKEN MODEL

- $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim (\mathbf{0}, \sigma^2 \mathbf{V})$
- Identical to the Gauss-Markov Linear Model except that $\text{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}$ instead of $\sigma^2 \mathbf{I}$.
- \mathbf{V} is assumed to be a **known** nonsingular Variance matrix.
- The Normal Theory Aitken Model adds an assumption of normality: $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{V})$
- Observations are now correlated, or have unequal variances.
- but the correlations, or unequal variances, follow a known pattern

Examples - 1

- Analysis of averages

The data to be analyzed are averages of unequal numbers of observations.

Y_i is an average of n_i observations. $\text{Var } Y_i = \sigma^2/n_i$

first 4 rows and columns of $\text{Var } \epsilon$ are:

$$\begin{bmatrix} \sigma^2/n_1 & 0 & 0 & 0 \\ 0 & \sigma^2/n_2 & 0 & 0 \\ 0 & 0 & \sigma^2/n_3 & 0 \\ 0 & 0 & 0 & \sigma^2/n_4 \end{bmatrix} = \sigma^2 \begin{bmatrix} 1/n_1 & 0 & 0 & 0 \\ 0 & 1/n_2 & 0 & 0 \\ 0 & 0 & 1/n_3 & 0 \\ 0 & 0 & 0 & 1/n_4 \end{bmatrix}$$

Examples - 2

- Analysis of data on a pedigree (genetic relationships among parents, children, grandchildren, ...)
- Genetic correlations between parents and children, among children, ..., all known.

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

First four rows of $\text{Var } \epsilon$:

$$\sigma^2 \begin{bmatrix} 1 & \rho_{12} & \rho_{13} & \rho_{14} \\ \rho_{12} & 1 & \rho_{23} & \rho_{24} \\ \rho_{13} & \rho_{23} & 1 & \rho_{34} \\ \rho_{14} & \rho_{24} & \rho_{34} & 1 \end{bmatrix}$$

where ρ_{ij} is the known genetic correlation among individuals i, j

Examples - 3

- Regression on data collected over time: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$
 $X_i = \text{year (1990, 1991, ...)}$
- Assume errors follow an autoregressive process (more later),
first 4 rows and columns of $\text{Var } \epsilon$ are:

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

- Aitken model if ρ known.

Statistical analysis of Aitken model data

- Spectral Decomposition Theorem:
 - any positive definite symm. matrix \mathbf{V} can be written as $\mathbf{V} = \mathbf{U}\mathbf{D}\mathbf{U}'$
 - \mathbf{D} is a diagonal matrix of eigenvalues
 - \mathbf{U} is a matrix of orthonormal eigenvectors with the property that $\mathbf{U}'\mathbf{U} = \mathbf{I}$.
- For any positive definite symmetric \mathbf{V} , there exists a nonsingular symmetric matrix $\mathbf{V}^{1/2}$ such that $\mathbf{V}^{1/2} \mathbf{V}^{1/2} = \mathbf{V}$.
- Given \mathbf{U} and \mathbf{D} for which $\mathbf{U}\mathbf{D}\mathbf{U}' = \mathbf{V}$, $\mathbf{V}^{1/2} = \mathbf{U}\sqrt{\mathbf{D}}\mathbf{U}'$
- $\mathbf{V}^{1/2}$ can be viewed as the “square root” of a matrix
- $\mathbf{V}^{1/2} \mathbf{V}^{1/2} = \mathbf{U}\sqrt{\mathbf{D}}\mathbf{U}'\mathbf{U}\sqrt{\mathbf{D}}\mathbf{U}' = \mathbf{U}\sqrt{\mathbf{D}}\mathbf{I}\sqrt{\mathbf{D}}\mathbf{U}' = \mathbf{U}\mathbf{D}\mathbf{U}' = \mathbf{V}$
- Define $\mathbf{V}^{-1/2}$ as $(\mathbf{V}^{-1})^{1/2}$
- Compute $\mathbf{V}^{-1/2}$ by $\mathbf{U}(1/\sqrt{\mathbf{D}})\mathbf{U}'$,
where $1/\sqrt{\mathbf{D}}$ is the diagonal matrix with elements $1/\sqrt{D_{ii}}$

Converting an Aitken model to GM model

- Our data model is $\mathbf{y} = \mathbf{X}\beta + \epsilon$, $\epsilon \sim (0, \sigma^2 \mathbf{V})$
- Pre-multiply all terms by $\mathbf{V}^{-1/2}$
 $\mathbf{V}^{-1/2}\mathbf{y} = \mathbf{V}^{-1/2}\mathbf{X}\beta + \mathbf{V}^{-1/2}\epsilon.$
- This is a regression model, $\mathbf{Z} = \mathbf{W}\beta + \delta$ with
 $\mathbf{Z} = \mathbf{V}^{-1/2}\mathbf{y}, \quad \mathbf{W} = \mathbf{V}^{-1/2}\mathbf{X}, \quad \delta = \mathbf{V}^{-1/2}\epsilon.$
- Why do this? What is $\text{Var } \delta$? $\text{Var}(\delta) = \text{Var}(\mathbf{V}^{-1/2}\epsilon)$
 $= \mathbf{V}^{-1/2}\sigma^2 \mathbf{V}\mathbf{V}^{-1/2}$
 $= \sigma^2 \mathbf{V}^{-1/2} \mathbf{V}^{1/2} \mathbf{V}^{1/2} \mathbf{V}^{-1/2} = \sigma^2 \mathbf{I}.$
- After transformation, we have a Gauss-Markov Model!

Generalized Least Squares

- $\hat{\beta}_G = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{V}^{-1}\mathbf{y}$ is a solution to the Aitken Equations:

$$\mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\mathbf{b} = \mathbf{X}'\mathbf{V}^{-1}\mathbf{y}$$

which follow from the Normal Equations

$$\begin{aligned}\mathbf{W}'\mathbf{W}\mathbf{b} &= \mathbf{W}'\mathbf{Z} \\ \Rightarrow \mathbf{X}'\mathbf{V}^{-1/2}\mathbf{V}^{-1/2}\mathbf{X}\mathbf{b} &= \mathbf{X}'\mathbf{V}^{-1/2}\mathbf{V}^{-1/2}\mathbf{y} \\ \Rightarrow \mathbf{X}'\mathbf{V}^{-1}\mathbf{X}\mathbf{b} &= \mathbf{X}'\mathbf{V}^{-1}\mathbf{y}.\end{aligned}$$

- Solving the Normal Equations is equivalent to minimizing

$$(\mathbf{Z} - \mathbf{W}\mathbf{b})'(\mathbf{Z} - \mathbf{W}\mathbf{b}) \text{ over } \mathbf{b} \in \mathbb{R}^p$$

- Now

$$\begin{aligned}(\mathbf{Z} - \mathbf{W}\mathbf{b})'(\mathbf{Z} - \mathbf{W}\mathbf{b}) &= (\mathbf{V}^{-1/2}\mathbf{y} - \mathbf{V}^{-1/2}\mathbf{X}\mathbf{b})'(\mathbf{V}^{-1/2}\mathbf{y} - \mathbf{V}^{-1/2}\mathbf{X}\mathbf{b}) \\ &= (\mathbf{y} - \mathbf{X}\mathbf{b})'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\mathbf{b})\end{aligned}$$

- Thus, $\hat{\beta}_G = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}$ is a solution to this Generalized Least Squares problem.

Estimating $\text{Var } \hat{\beta}_G$ for full rank \mathbf{X}

- If \mathbf{X} full rank, $\hat{\beta}_G = (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y}$

$$\begin{aligned}\text{Var } \hat{\beta}_G &= E (\hat{\beta}_G - \beta_G)(\hat{\beta}_G - \beta_G)' \\&= E ((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y} - (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} E \mathbf{y}) \times \\&\quad ((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{y} - (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} E \mathbf{y})' \\&= ((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} (\mathbf{y} - E \mathbf{y}))((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} (\mathbf{y} - \\&= ((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} (\mathbf{y} - E \mathbf{y}))(\mathbf{y} - E \mathbf{y})' (\mathbf{V}^{-1} \mathbf{X} (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \\&= \sigma^2 ((\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1} \\&= \sigma^2 (\mathbf{X}' \mathbf{V}^{-1} \mathbf{X})^{-1}\end{aligned}$$

- Similar to OLS $\text{Var } \hat{\beta}$, except for extra \mathbf{V}^{-1} in middle
- $\text{Var } \mathbf{C} \hat{\beta}_G = \mathbf{C} \text{Var } \hat{\beta} \mathbf{C}'$
- $E \mathbf{y} = \mathbf{X} \hat{\beta}_G$
- $\text{Var } \hat{\mathbf{y}} = \sigma^2 \mathbf{X} \text{Var } \hat{\beta}_G \mathbf{X}'$

When \mathbf{X} is not full rank:

- estimate $E \mathbf{y}$:
- Under the GM model, the best estimate of $E \mathbf{Z}$ is

$$\begin{aligned}\hat{\mathbf{Z}} &= P_W \mathbf{Z} = \mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{Z} \\ &= \mathbf{V}^{-1/2}\mathbf{X}((\mathbf{V}^{-1/2}\mathbf{X})'\mathbf{V}^{-1/2}\mathbf{X})^{-1}(\mathbf{V}^{-1/2}\mathbf{X})'\mathbf{V}^{1/2}\mathbf{y} \\ &= \mathbf{V}^{-1/2}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1/2}\mathbf{V}^{-1/2}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1/2}\mathbf{V}^{-1/2}\mathbf{y} \\ &= \mathbf{V}^{-1/2}\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}.\end{aligned}$$

- To estimate $E(\mathbf{y}) = \mathbf{V}^{1/2}E(\mathbf{Z})$, use

$$\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}}_G.$$

Estimating $C\beta$ when X not full rank

- $C\beta$ is estimable if C can be written as a linear combination of rows of W , i.e. $C = AW$
- if $C\beta$ estimable, the BLUE is the ordinary least squares (OLS) estimator using Z and W

$$\begin{aligned}C\hat{\beta} &= C(W'W)^{-1}W'Z \\&= C(X'V^{-1/2}V^{-1/2}X)^{-1}X'V^{-1/2}V^{-1/2}y \\&= C(X'V^{-1}X)^{-1}X'V^{-1}y.\end{aligned}$$

- $C(X'V^{-1}X)^{-1}X'V^{-1}y = C\hat{\beta}_G$ is called a Generalized Least Squares (GLS) estimator.

Var $\mathbf{C}\hat{\beta}_G$

- Need to be careful because Var $\hat{\beta}_G$ doesn't exist

$$\begin{aligned}\text{Var } \mathbf{C}\hat{\beta}_G &= \text{Var } \mathbf{C}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y} \\&= \text{Var } \mathbf{C}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{Z} \\&= \text{Var } \mathbf{AW}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{Z} \\&= \text{Var } \mathbf{AP}_W\mathbf{Z} \\&= \mathbf{AP}_W\text{Var } \mathbf{ZP}_W\mathbf{A}' \\&= \mathbf{AP}_W(\sigma^2\mathbf{I})\mathbf{P}_W\mathbf{A}' \\&= \sigma^2\mathbf{AP}_W\mathbf{P}_W\mathbf{A}' = \sigma^2\mathbf{AP}_W\mathbf{A}' \\&= \sigma^2\mathbf{AW}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'\mathbf{A}' \\&= \sigma^2\mathbf{C}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{C}' \\&= \sigma^2\mathbf{C}(\mathbf{X}'\mathbf{V}^{-1/2}\mathbf{V}^{-1/2}\mathbf{X})^{-1}\mathbf{C}' \\&= \sigma^2\mathbf{C}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{C}'\end{aligned}$$

Weighted Least Squares

- When \mathbf{V} is diagonal, the term "Weighted Least Squares" (WLS) is commonly used instead of GLS.
- Define \mathbf{D} = diagonal matrix of inverse weights, $D_{ii} = 1/w_i$
- $(\mathbf{y} - \mathbf{Xb})' \mathbf{D}^{-1} (\mathbf{y} - \mathbf{Xb}) = \sum_{i=1}^n w_i (y_i - \mathbf{X}_i \beta)^2$
- When obs. have unequal variances, i.e. $\text{Var } \mathbf{y} = \text{diag}(\sigma_i^2)$, w_i is proportional to $1/\sigma_i^2$
- WLS assumes weights are known.
- If weights are estimated, e.g. using s_i^2 , WLS analysis is in trouble (when small d.f. for each s_i^2) or approximate (when large d.f. for each s_i^2).

Summary of Aitken model results

- Aitken model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim (0, \sigma^2 \mathbf{V})$
- equivalent to GM model: $\mathbf{Z} = \mathbf{W}\boldsymbol{\beta} + \boldsymbol{\delta}$, by premultiplying by $\mathbf{V}^{-1/2}$
- estimate $E(\mathbf{y}) = \mathbf{V}^{-1/2}E(\mathbf{Z})$ by $\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}$.
- estimate $\mathbf{C}\boldsymbol{\beta}$ by $\mathbf{C}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\mathbf{y}$
- estimate $\text{Var } \mathbf{C}\boldsymbol{\beta}$ by $\hat{\sigma}^2 \mathbf{C}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{C}'$
- estimate $\text{Var } \boldsymbol{\delta}$ by $\hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_G)' \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_G)}{n-k}$
- which means estimating $\text{Var } \mathbf{y}$ by $\hat{\sigma}^2 \mathbf{V}$
- If add normality, all inferential results from Stat 500 follow
- In particular, df. for $\hat{\sigma}^2 = N - \text{rank } \mathbf{X}$

What if V misspecified?

- Data follow Aitken model, $\mathbf{y} = \mathbf{X}\beta + \epsilon$, $\epsilon \sim (0, \sigma_G^2 V)$, but analyzed using $\mathbf{y} = \mathbf{X}\beta + \epsilon$, $\epsilon \sim (0, \sigma^2 I)$?
- Results stated; derivations can be found in Kutner et al. (big white) or most Econometrics textbooks.
- $E(y)$ unbiased, so $\mathbf{C}\hat{\beta}$ unbiased
- OLS estimates not as efficient as GLS, $\text{Var } \mathbf{C}\hat{\beta} > \text{Var } \mathbf{C}\hat{\beta}_G$
- $\text{Var } \mathbf{C}\hat{\beta}$ is **not** $\sigma^2 \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}$. Bias can be considerable.
- Assume \mathbf{X} full rank, so $(\mathbf{X}'\mathbf{X})^{-1}$ exists

$$\begin{aligned}\text{Var } \mathbf{C}\hat{\beta} &= \text{Var } \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' (\text{Var } \mathbf{y}) \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}' \\ &= \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\sigma_G^2 V)\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}' \\ &= \sigma_G^2 \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}' V \mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'\end{aligned}\tag{1}$$

- Can't simplify!
- All inference is suspect, unless $(\mathbf{X}' V \mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}$ close to I

Sandwich estimator

- However, equation (1) can be used to advantage!
- Can estimate $\text{Var } \mathbf{y} = \text{Var}$ from the empirical variance-covariance matrix of the residuals, $\text{Var } \boldsymbol{\epsilon} = \boldsymbol{\epsilon} \boldsymbol{\epsilon}'$, or a modeled version of that.
- If you suspect a problem with homogeneous variances, independent errors, estimate $\text{Var } \mathbf{C}\hat{\boldsymbol{\beta}}$ using (1).
- called “White’s heteroscedastic consistent variance estimator” in econometrics
- Recent statistical literature has called this the “sandwich” estimator, Because of the meat ($\mathbf{X}' \mathbf{V} \mathbf{X}$) between the two slices of bread ($\mathbf{X}' \mathbf{X}$).

Example data analysis using Aitken model

- Study comparing metabolite concentrations in two genotypes of plants (A, B)
- Seedlings commonly grown in flats, here 9 seedlings per flat
- 4 flats per genotype. All 9 seedlings same genotype.
- Measure size of each seedling 2 weeks after germination.
- review of 500/401: what is the observational unit (ou)?
what is the experimental unit?

Analysis - 2

- Seedlings too small to measure metabolite concentration on each.
- Combine together all plants in a flat. One measurement per flat.
- This is a physical average of the metabolite concentration in seedling.
- Goal: estimate mean difference between genotypes in size and in metabolite concentration.
- Problem: some seedlings died. Some responses are averages of 9 seedlings; some are averages of 5 seedlings.
- Standard analyses assume that death unrelated to size or metabolite conc.
- Many ways to model (and hence to analyze) data like this.

Analysis - 3

- For now, only consider metabolite data.
- **A** model (not the model) assumes that the variation among responses is due only to variation among seedlings.
- Y_{ij} is metabolite concentration measured in flat j of genotype i
- this is an average of n_{ij} seedlings.
- If only variation among seedlings, $\text{Var } Y_{ij} = \sigma^2 / n_{ij}$
- Aitken model with $\mathbf{V} = \text{diag}(1/n_{ij})$

Analysis - 4

- The data:

| | | | | | | | | |
|-------------|-------|-------|-------|-------|-------|-------|-------|-------|
| Flat | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| Genotype | A | B | B | A | A | B | B | A |
| # seedlings | 9 | 6 | 7 | 9 | 6 | 6 | 8 | 5 |
| Metabolite | .0231 | .0275 | .0521 | .0001 | .0063 | .0138 | .1061 | .0482 |

- The results:

| Parameter | GLS | | OLS | |
|-----------------------|--------|--------|--------|--------|
| | est. | s.e. | est. | s.e. |
| Group A mean | 0.0168 | 0.0159 | 0.0194 | 0.0163 |
| Difference | 0.0373 | 0.0229 | 0.0305 | 0.0230 |
| rMSE = $\hat{\sigma}$ | 0.085 | | 0.032 | |

- Heteroscedasticity not very large ($\max \text{Var} / \min \text{Var} < 2$)
- Estimates similar, $\hat{\sigma}^2$ very different!

Analysis - 5

- Different parameters: Aitken model: $\text{Var } Y_{ij} = \sigma_G^2/n_i$, GM model: $\text{Var } Y_{ij} = \sigma^2$
- σ_G^2 is the variance among measurements of individual seedlings
- σ^2 from OLS is variance among flat means
- In fact, average $\hat{\sigma}_G^2/n_i = 0.085^2 \cdot 0.1487 \approx 0.032^2$
- Very dependent on assumption of no variation other than seedling-seedling
- Assumes no additional variation among flats
- Common experience is that there is both flat-flat variation and seedling-seedling variation.
- Accounting for both requires a mixed model (coming up).
- Can estimate both variance components even though did not measure individual seedlings

R code for the Aitken model and weighted LS

```
# the metabolite example
# enter the data (genotype, # seedlings, and ave. conc)
genotype <- c('A', 'B', 'B', 'A', 'A', 'B', 'B', 'A')
gen <- as.factor(genotype)

nsdl <- c(9, 6, 7, 9, 6, 6, 8, 5)
metab <- c(0.0231, 0.0275, 0.0521, 0.0001, 0.0063, 0.0138,
           0.1061, 0.0482)

# the ols analysis
ols.lm <- lm(metab ~ gen)
summary(ols.lm)
```

```
# an Aitken model analysis, by hand
# v is an 8 x 8 matrix with 1/nsdl on the diagonal
v <- diag(1/nsdl)
# diag(vector) creates a matrix with vector on diag
# diag(matrix) extracts the diagonal

# need to get the inverse square root matrix of v
# use eigen function to do that

temp <- eigen(v)
# returns a list with two components: values, a vector
# and vectors, a matrix
```

```

u <- temp$vector
d <- diag(temp$values)
# convert e-vals to a diagonal matrix

round(v - u %*% d %*% t(u), 5) # check  $u d t' = v$ 
all(v == u %*% d %*% t(u))
# another possible check, more sensitive to num. error

svi <- u %*% diag(1/sqrt(temp$values)) %*% t(u)
# inverse square root matrix

# a check that svi does what we want it to do
v %*% svi %*% t(svi)
#  $svi svi' = v^{-1}$ , and  $v * v^{-1} = I$ 

```

```
# use svi to transform y, and X
# notation follows that in notes
z <- svi %*% metab
w <- svi %*% model.matrix(ols.lm)
  # could also use any other equivalent X matrix
  # w includes a column for intercept

gls.lm <- lm(z~-1+w) # need to suppress default intercept
summary(gls.lm)
```



```
# could also do this analysis using weighted least
#   squares, since V diagonal
#   weights are 1/variance= nsdl

wls.lm <- lm(metab~gen, weight=nsdl)
summary(wls.lm)

# some useful hints / tricks to construct different
#   sorts of V matrices

# AR(1) structure. V has bands.
# This sort of matrix is a toeplitz matrix

toeplitz(1:5)
```

```
# to generate a 8 x 8 ar(1) matrix with specified rho
rho <- 0.7
v <- toeplitz(rho^(0:7))
# sequence starts at 0 so diag = 1, ends at rho^7

# compound symmetry structure (obs within groups)
# we haven't yet talked about this model,
# it's here because it fits here, but you won't
#   need it for a while.
grp <- c(1,1,2,2,3,3,4,4,4)
# or use rep(1:4,c(2,2,2,3))
v <- outer(grp, grp, '==')
# true (i.e. 1) if grp[i] = grp[j]
k <- 2/1.5
# ratio of sigma^2_p / sigma^2_c
v <- v + diag(rep(k,length(grp)))
# add k to diag., also coerces T/F to 1/0 (good)
```

Maximum Likelihood

- Suppose $f(\mathbf{w}|\theta)$ is the probability density function (*pdf*) or probability mass function (*pmf*) of a random vector \mathbf{w} , where θ is a $k \times 1$ vector of parameters.
- Given a value of the parameter vector θ , $f(\mathbf{w}|\theta)$ is a real-valued function of \mathbf{w} .
- The likelihood function $L(\theta|\mathbf{w}) = f(\mathbf{w}|\theta)$ is a real-valued function of θ for a given value of \mathbf{w} .
- $L(\theta|\mathbf{w})$ is not a pdf. $\int_{\theta} L(\theta|\mathbf{w})d\theta \neq 1$.

- For any potential observed vector of values \mathbf{w} , define $\hat{\theta}(\mathbf{w})$ to be a parameter value at which $L(\theta|\mathbf{w})$ attains its maximum value. $\hat{\theta}(\mathbf{w})$ is a maximum likelihood estimator (MLE) of θ .
- Invariance property of MLES: The MLE of a function of θ , say $\mathbf{g}(\theta)$, is the function evaluated at the MLE of θ : $\widehat{\mathbf{g}(\theta)} = \mathbf{g}(\hat{\theta})$
- Often much more convenient to work with $l(\theta|\mathbf{w}) = \log L(\mathbf{w}|\theta)$.
- If $l(\theta|\mathbf{w})$ is differentiable, candidates for the MLE of θ can be found by equating the score function

$$\frac{\partial l(\theta|\mathbf{w})}{\partial \theta} \equiv \begin{bmatrix} \frac{\partial l(\theta|\mathbf{w})}{\partial \theta} \\ \vdots \\ \frac{\partial l(\theta|\mathbf{w})}{\partial \theta_k} \end{bmatrix} \text{ to } \mathbf{0} \text{ and solving for } \theta$$

- The score equations are $\frac{\partial l(\theta|\mathbf{w})}{\partial \theta} = \mathbf{0} \Rightarrow \frac{\partial l(\theta|\mathbf{w})}{\partial \theta_j} = 0 \quad \forall j = 1, \dots, k$
- One strategy for obtaining an MLE is to find solution(s) of the score equations and verify that at least one such solution maximizes $l(\theta|\mathbf{w})$.
- If the solution(s) to the score equations lie outside the appropriate parameter space, they are not MLE's

- Example: Normal Theory Gauss-Markov Linear Model

$$\underbrace{\mathbf{y}}_{n \times 1} = \underbrace{\mathbf{X}}_{n \times p} \underbrace{\boldsymbol{\beta}}_{p \times 1} + \underbrace{\boldsymbol{\epsilon}}_{n \times 1} \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}) \quad \underbrace{\boldsymbol{\theta}}_{(p+1) \times 1} = \begin{bmatrix} \boldsymbol{\beta} \\ \sigma^2 \end{bmatrix}$$

$$\begin{aligned} f(\mathbf{y}|\boldsymbol{\theta}) &= \frac{\exp\left\{\frac{-1}{2}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\sigma^2\mathbf{I})^{-1}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})\right\}}{(2\pi)^{n/2}|\sigma^2\mathbf{I}|^{1/2}} \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{\frac{-1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})\right\} \\ l(\boldsymbol{\theta}|\mathbf{y}) &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

- The score function is

$$\frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \boldsymbol{\beta}} \\ \frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \sigma^2} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2}(\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\boldsymbol{\beta}) \\ \frac{(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})}{2\sigma^4} - \frac{n}{2\sigma^2} \end{bmatrix}$$

- The score equations are

$$\frac{\partial l(\boldsymbol{\theta}|\mathbf{y})}{\partial \boldsymbol{\theta}} = \mathbf{0} \Leftrightarrow \mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{y} \quad \sigma^2 = \frac{(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})'(\mathbf{y}-\mathbf{X}\boldsymbol{\beta})}{n}$$

- Here, any solution lies in the appropriate parameter space:
 $\beta \in \mathbb{R}^p, \sigma^2 \in \mathbb{R}^+.$
- A solution to the score equations is

$$\begin{bmatrix} \hat{\beta} \\ \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n} \end{bmatrix},$$

where $\hat{\beta}$ is a solution to the normal equations

- Need to show this is a maximum of the likelihood function
- We already know that any solution to the normal equations, $\beta = \hat{\beta}$, minimizes $(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$ for $\beta \in \mathbb{R}^p$.
- Thus,

$$\forall \sigma^2 > 0, \quad l\left(\begin{bmatrix} \hat{\beta} \\ \sigma^2 \end{bmatrix} \mid \mathbf{y}\right) \geq l\left(\begin{bmatrix} \beta \\ \sigma^2 \end{bmatrix} \mid \mathbf{y}\right) \quad \forall \beta \in \mathbb{R}^p$$

- The 2nd derivative of $l(\theta|\mathbf{y})$ with respect to σ^2 is < 0 so

$$\begin{bmatrix} \hat{\beta} \\ \frac{(\mathbf{y}-\mathbf{x}\hat{\beta})'(\mathbf{y}-\mathbf{x}\hat{\beta})}{n} \end{bmatrix} \text{ is an MLE of } \theta = \begin{bmatrix} \beta \\ \sigma^2 \end{bmatrix}$$

- Thus, if $C\beta$ is estimable, the MLE of $C\beta$ is $C\hat{\beta}$ (by the Invariance Property of MLEs), which is the BLUE of $C\beta$
- Note that the MLE of σ^2 is not the unbiased estimator we have been using.

$$E\left[\frac{(\mathbf{y}-\mathbf{x}\hat{\mathbf{b}})'(\mathbf{y}-\mathbf{x}\hat{\mathbf{b}})}{n}\right] = E\left(\frac{SSE}{n}\right) = \frac{n-p}{n}\sigma^2 < \sigma^2$$

Thus, the MLE underestimates σ^2 on average.

- Now consider the general linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \boldsymbol{\Sigma}),$$

Where $\boldsymbol{\Sigma}$ is a positive definite covariance matrix whose entries depend on unknown parameters in some vector $\boldsymbol{\gamma}$.

- For example,

$$\boldsymbol{\Sigma} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{bmatrix}, \boldsymbol{\gamma} = \begin{bmatrix} \sigma^2 \\ \rho \end{bmatrix}$$

- Not Aitken model when ρ unknown.

- In general,

$$\theta = \begin{bmatrix} \beta \\ \gamma \end{bmatrix} \text{ and } f(\mathbf{y}|\theta) = \frac{\exp \left\{ -\frac{1}{2}(\mathbf{y} - \mathbf{x}\beta)'(\boldsymbol{\Sigma})^{-1}(\mathbf{y} - \mathbf{x}\beta) \right\}}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}}$$

$$l(\theta|\mathbf{y}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{y} - \mathbf{x}\beta)' \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{x}\beta) - \frac{n}{2} \log(2\pi)$$

- We know that for any positive definite covariance matrix Σ , $(\mathbf{y} - \mathbf{X}\beta)' \Sigma^{-1} (\mathbf{y} - \mathbf{X}\beta)$ is minimized over $\beta \in \mathbb{R}^p$ by the GLS estimator $\hat{\beta}_g = (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1} \mathbf{y}$.
- Thus, for any γ such that Σ is a positive definite covariance matrix,

$$l\left(\begin{bmatrix} \hat{\beta}_g \\ \gamma \end{bmatrix} \mid \mathbf{y}\right) \geq l\left(\begin{bmatrix} \beta \\ \gamma \end{bmatrix} \mid \mathbf{y}\right) \quad \forall \beta \in \mathbb{R}^p$$

- We define the profile log likelihood for γ to be

$$l^*(\gamma \mid \mathbf{y}) = l\left(\begin{bmatrix} \hat{\beta}_g \\ \gamma \end{bmatrix} \mid \mathbf{y}\right)$$

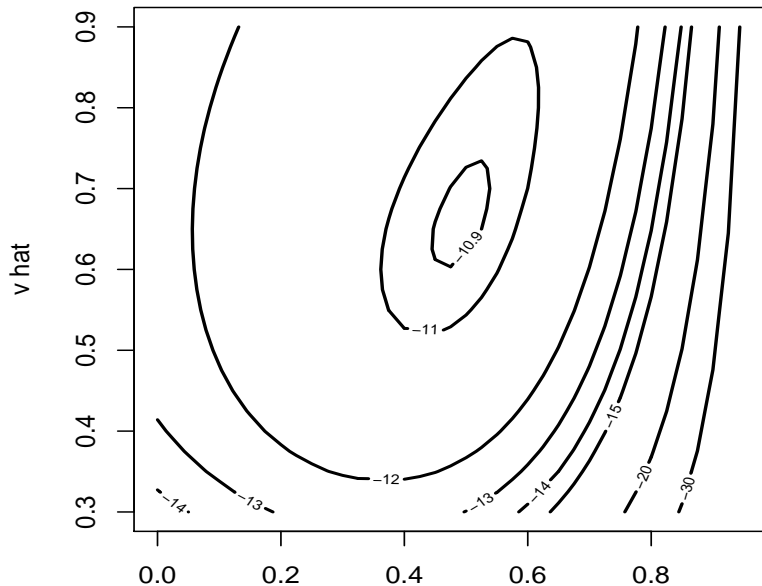
- The MLE of θ is

$$\hat{\theta} = \begin{bmatrix} \hat{\beta}_{\hat{g}} \\ \hat{\gamma} \end{bmatrix}$$

- Where $\hat{\gamma}$ is a maximizer of $l^*(\gamma|\mathbf{y})$ and $\hat{\Sigma}$ is Σ with $\hat{\gamma}$ in place γ .
- In general, numerical methods are required to find $\hat{\gamma}$, a maximizer of $l^*(\gamma|\mathbf{y})$

- Numerical maximization algorithms are iterative.
 - Require a starting value of γ
 - Attempt to find a better value, γ^* , in the sense that $l^*(\gamma^*|\mathbf{y}) > l^*(\gamma|\mathbf{y})$.
- Newton-Raphson algorithm:
 - estimate gradient vector and Hessian matrix at current γ .
 - This gives a quadratic approximation to log-likelihood.
 - Analytic maximum of quadratic gives γ^* .
 - replace γ by γ^*
 - Repeat until no improvement.
- Many details, will cover near end of semester.
- Will mention one now.

• Example: AR(1) model with unknown ρ



- Problem appears to be maximization over two parameters.
- Doesn't have to be. Remember the Aitken model:

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

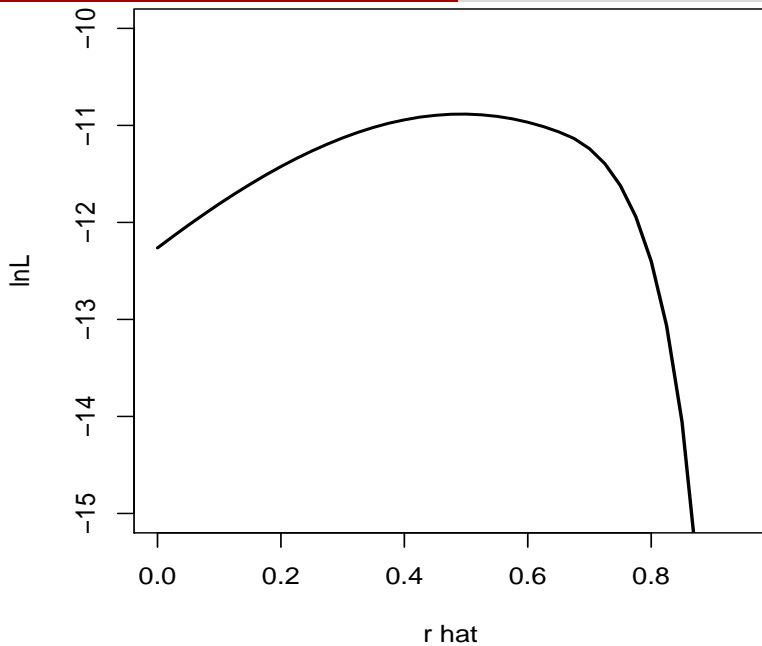
where \mathbf{V} is a function of known constants

- e.g. for AR(1)

-

$$\boldsymbol{\Sigma} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

- $\hat{\boldsymbol{\beta}}_g$ depends on ρ but not σ^2
- Use ρ as the parameter, maximize InL over ρ
- Sometimes called “profiling out” the error variance



- The MLE of the variance component vector γ is often biased.
- For example, for the case of $\epsilon = \sigma^2 I$, where $\gamma = \sigma^2$, the MLE of σ^2 is $\frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n}$ with expectation $\frac{n-p}{n}\sigma^2$.
- The MLE of σ^2 is often criticized for "failing to account for the loss of degrees of freedom needed to estimate β ."

$$E\left[\frac{(\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta})}{n}\right] = \frac{n-p}{n}\sigma^2$$

$$< \sigma^2 = E\left[\frac{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)}{n}\right]$$

- A variation, REML, is unbiased in simple models and less biased in many others.
- We'll see REML in detail when we talk about mixed models

Large N theory for MLE's

- Suppose θ is a $k \times 1$ parameter vector.
- Let $l(\theta)$ denote the log likelihood function.
- Under regularity conditions discussed, e.g., Casella and Berger, we have the following:
 - There is an estimator $\hat{\theta}$ that solves the likelihood equations $\frac{\delta l(\theta)}{\delta \theta} = \mathbf{0}$ and is a consistent estimator of θ , i.e., $\lim_{n \rightarrow \infty} Pr[|\hat{\theta} - \theta| > \epsilon] = 0$ for any $\epsilon > 0$.
 - For sufficiently large n , $\hat{\theta} \sim N(\theta, I^{-1}(\theta))$, where

$$\begin{aligned}
 I(\theta) &= E \left[\left(\frac{\partial l(\theta)}{\partial \theta} \right) \left(\frac{\partial l(\theta)}{\partial \theta} \right)' \right] \\
 &= -E \left[\frac{\partial^2 l(\theta)}{\partial \theta \partial \theta'} \right]
 \end{aligned}$$

- Or, in scalar terms,

$$I(\theta)_{ij} = -E \left[\frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j} \right], \quad i = 1, \dots, k, \quad j = 1, \dots, k.$$

- $I(\theta)$ is known as the Fisher Information matrix.
- $I(\theta)$ can be approximated by $\frac{-\partial^2 l(\theta)}{\partial \theta \partial \theta'} \big|_{\theta=\hat{\theta}}$
- Often called “observed information”

Wald tests and confidence intervals

- Suppose for large n that $\hat{\theta} \sim N(\theta, \mathbf{V})$ and $\hat{\mathbf{V}}$ is a consistent estimator of \mathbf{V} .
- Then, in suitable large samples, $\hat{\mathbf{V}}^{-1/2}(\hat{\theta} - \theta) \sim N(\mathbf{0}, I)$ and $(\hat{\theta} - \theta)' \hat{\mathbf{V}}^{-1}(\hat{\theta} - \theta) \sim \chi_k^2$.
- An approximate $100(1 - \alpha)\%$ confidence interval for θ_i is $\hat{\theta}_i \pm Z_{1-\alpha/2} \sqrt{\hat{\mathbf{V}}_{ii}(\theta)}$, where $Z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of $N(0,1)$ and $\hat{\mathbf{V}}_{ii}(\theta)$ is element (i,i) of $\hat{\mathbf{V}}(\theta)$.
- An approximate p-value for testing $H_0 : \theta = \theta_0$ is $P[\chi_k^2 \geq (\hat{\theta} - \theta_0)' \hat{\mathbf{V}}^{-1}(\hat{\theta} - \theta_0)]$, where χ_k^2 is a χ^2 random variable with k degrees of freedom and k is the number of restrictions in the null hypothesis.
- The above confidence interval and test are based on the asymptotic normality of $\hat{\theta}$

Likelihood ratio based inference

- Suppose we wish to test the null hypothesis that a reduced model provides an adequate fit to a dataset relative to a more general full model that includes the reduced model as a special case.
- Under the regularity conditions mentioned previously,
(REDUCED MODEL DEVIANCE) - (FULL MODEL DEVIANCE)
 $\stackrel{H_0}{\sim} \chi^2_{k_f - k_r}$, where k_f and k_r are the number of free parameters under the full and reduced models, respectively.
- This approximation can be reasonable if n is "sufficiently large".
- Note that the test statistic is equal to $-2\log\Lambda$, where

$$\Lambda = \frac{\text{LIKELIHOOD MAXIMIZED UNDER REDUCED MODEL}}{\text{LIKELIHOOD MAXIMIZED UNDER THE FULL MODEL}}$$

- Λ is known as the likelihood ratio, and tests based on $-2\log\Lambda$ are called likelihood ratio tests.

LRT's and Confidence Regions for a Subvector of θ :

- Suppose θ is a $k \times 1$ vector and is partitioned into vectors θ_1 $k_1 \times 1$ and θ_2 $k_2 \times 1$, where $k = k_1 + k_2$ and $\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$
- Consider a test of $H_0 : \theta_1 = \theta_{10}$
- Suppose $\hat{\theta}$ is the MLE of θ and $\hat{\theta}_2(\theta_1)$ maximizes $l = \left(\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \right)$ over θ_2 for any fixed value of θ_1 .
- Then $2 \left[l(\hat{\theta}) - l \left(\begin{bmatrix} \theta_{10} \\ \hat{\theta}_2(\theta_{10}) \end{bmatrix} \right) \right] \stackrel{H_0}{\sim} \chi_{k_1}^2$ by our previous result when n is "sufficiently large."

- Also, $Pr \left\{ 2[l(\hat{\theta}) - l \left(\begin{bmatrix} \theta_1 \\ \hat{\theta}_2(\theta_1) \end{bmatrix} \right)] \leq \chi_{k_1, \alpha}^2 \right\} \approx 1 - \alpha \Rightarrow$
 $Pr \left\{ l \left(\begin{bmatrix} \theta_1 \\ \hat{\theta}_2(\theta_1) \end{bmatrix} \right) \geq l(\hat{\theta}) - \frac{1}{2} \chi_{k_1, \alpha}^2 \right\} \approx 1 - \alpha$
- Thus, the set of values of θ_1 that when maximizing over θ_2 , yield a maximized likelihood within $\frac{1}{2} \chi_{k_1, \alpha}^2$ of the likelihood maximized over all θ , form a $100(1 - \alpha)\%$ confidence region for θ_1
- The ultimate distributional dependence is the asymptotic N or χ^2 distribution.
 - Only holds for infinitely large sample sizes
 - But, may be an appropriate approximate for practical but large sample sizes
 - may be inappropriate if sample sizes are too small
- Often, the above never worried about.

Why worry about a second ci method?

- When $\mathbf{Y} \sim N(\theta, \sigma^2 \mathbf{I})$, $\mathbf{C}\beta \sim N$ in any size sample. When inference uses $\hat{\sigma}^2$, the T and F distributions are exact for any size sample. No need for LR ci's or tests.
- Now, let $\mathbf{Y}_{indep} \sim F(\beta)$, some arbitrary, non-normal distribution. Imagine that the sample size is sufficiently large that $\hat{\beta}_{approx.} \sim N$. Now, you want inference on $\tau = \exp \beta_1$ or $\eta = \beta_2/\beta_1$.
- MLE of $\exp \beta$ is $\exp \hat{\beta}$, MLE of β_2/β_1 is $\hat{\beta}_2/\hat{\beta}_1$.
- se of $\exp \beta$ or β_2/β_1 from Delta method.
- ci? If $\beta_{approx} \sim N$, then $\exp \hat{\beta}$ and $\hat{\beta}_2/\hat{\beta}_1$ are most definitely not normal.

- If sample size increased, eventually will converge to asymptotic normal distribution, but can't use Wald methods on current sample.
- log-likelihood invariant to reparameterization:
 - $l(\log \tau) = l(\beta_1)$.
 - so LR ci is τ s.t. $2 \left[l(\hat{\beta}_1) - l(\log \tau) \right] < \chi^2_{1-\alpha, k}$
- for functions of multiple parameters, need an additional maximization:
 - so LR ci for η is η s.t. $2 \left[l(\hat{\beta}_1, \hat{\beta}_2) - \max_{\beta_1} l(\beta_1, \eta \beta_2) \right] < \chi^2_{1-\alpha, k}$
 - Just a computing exercise.
- In my experience, the χ^2 approximation for a LR statistic is usable at much smaller sample sizes than the N approximation for a $\hat{\beta}$, unless you're lucky in your choice of β .
- I know of no cases where the N approximation is usable at smaller sample sizes than the χ^2 approximation, unless $Y \sim N$.

Two/three notes for caution

- The regularity conditions do not hold if the true parameter or the value specified by the null hypothesis falls on the boundary of the parameter space.
- We will see examples of this soon, when we discuss mixed models
- One example, if a model has two random components (error and something else, call it $u \sim N(0, \sigma_u^2)$), testing $H_0 : \sigma_u^2 = 0$ is on the boundary of the parameter space
- None of the usual theory applies for this situation.
- because the regularity conditions are not met.

Generalized Linear Models

- Consider the normal theory Gauss-Markov linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I})$.
- Does not have to be written as function + error
- Could specify distribution and model(s) for its parameters
- i.e., $y_i \sim N(\mu_i, \sigma^2)$, where $\mu_i = \mathbf{X}'_i \boldsymbol{\beta}$ for all $i = 1, \dots, n$ and y_1, \dots, y_n independent.
- This is one example of a generalized linear model.
- Here is another example of a GLM:
 $y_i \sim \text{Bernoulli}(\pi_i)$, where $\pi_i = \frac{\exp(\mathbf{X}'_i \boldsymbol{\beta})}{1 + \exp(\mathbf{X}'_i \boldsymbol{\beta})}$ for all $i = 1, \dots, n$ and y_1, \dots, y_n independent.
- In both examples, all responses are independent and each response is a draw from one type of distribution whose parameters may depend on explanatory variables through a known function of a linear predictor $\mathbf{X}'_i \boldsymbol{\beta}$.

- The normal and Bernoulli models (and many others) are special cases of a generalized linear model.
- These are models where:
 - The parameters are specified functions of $\mathbf{X}\beta$
 - The distribution is in the exponential scale family
 - i.e., y_i has a probability density function (or probability mass function, for discrete distribution)

$$\exp \left(\frac{\eta(\theta) \mathbf{T}(y_i) - \mathbf{b}(\theta)}{a(\phi)} + c(y_i, \phi) \right) \quad (2)$$

where $\eta()$, $\mathbf{T}()$, $a()$, $b()$, and $c()$ are known functions and θ is a vector of unknown parameters depending on $\mathbf{X}\beta$ and ϕ is either a known or unknown parameter.

- Exponential family / exponential class is (2) without the $a(\phi)$
- $a(\phi)$ includes “overdispersed” distributions in the family

- For example, the pdf for a normal distribution can be written as:

$$\exp \left(\frac{-1}{2\sigma^2} y_i^2 + \frac{\mu}{2\sigma^2} y_i - \frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right)$$

- from which:

- $\eta(\theta) = \left(\frac{\mu}{\sigma^2}, \frac{-1}{2\sigma^2} \right)$
- and $\mathbf{T}(y_i) = (y_i, y_i^2)$

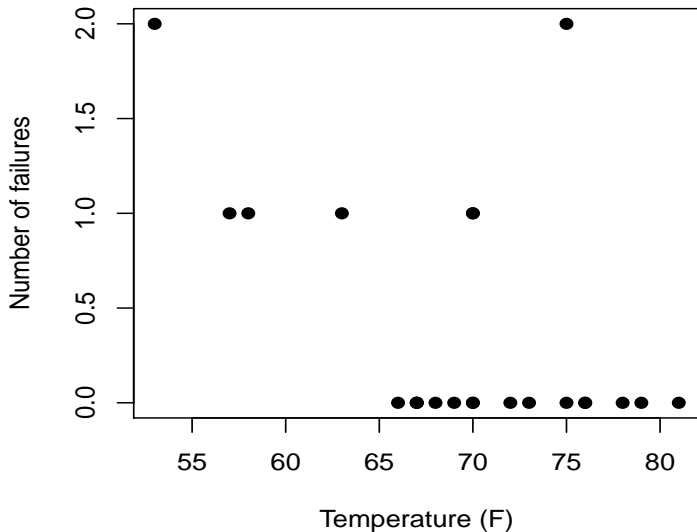
- This family includes many common distributions:

| Distribution | $\eta(\theta)'$ | $T(y_i)'$ | $a(\phi)$ |
|-------------------|--|-------------------|-----------|
| Normal | $(\frac{\mu}{\sigma^2}, \frac{-1}{2\sigma^2})$ | (y_i, y_i^2) | 1 |
| Bernoulli | $\log\left(\frac{\pi}{1-\pi}\right)$ | y_i | 1 |
| Poisson | $\log \lambda$ | y_i | 1 |
| Overdisp. Poisson | $\log \lambda$ | y_i | ϕ |
| Gamma | $(\frac{-1}{\theta}, (k-1))$ | $(y_i, \log y_i)$ | 1 |

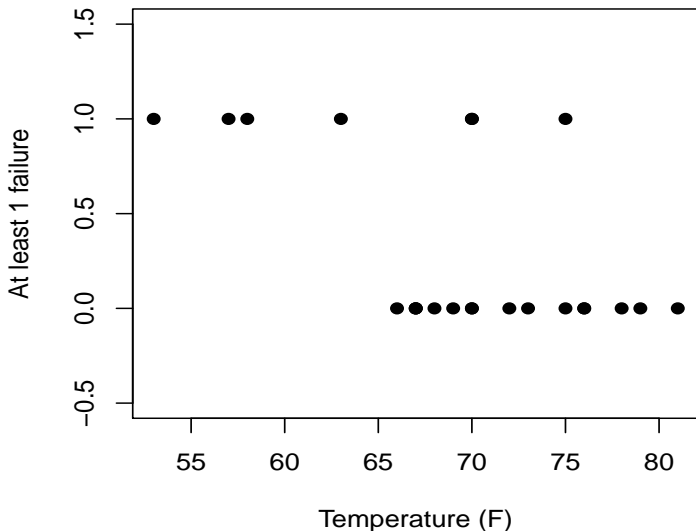
- A lot of stat theory results follow immediately from the exponential form
- e.g. $T(y_i)$ is the vector of sufficient statistics
- A lot of theory is much nicer when the distribution is parameterized in terms of $\eta(\theta)$ instead of θ
- e.g. use $\log\left(\frac{\pi}{1-\pi}\right)$ as the parameter of a Bernoulli distribution instead of π .
- This is called the canonical or natural form

- A generalized linear model has:
 - A probability distribution in the exponential scale family
 - A linear predictor, $\eta = \mathbf{X}\beta$
 - A link function that connects $E(\mathbf{y})$ and the linear predictor:
 $g(E(\mathbf{Y})) = g(\mu) = \eta$

- Example: Challenger O-ring data, with flight temperature



- Q: does $P[\text{one or more failures}]$ depend on temperature?
If so, what is a reasonable model?



- Need a reasonable model for independent observations taking values of 0 or 1.
- $y \sim \text{Bernoulli}(\pi)$ has probability mass function

$$f(y) = \begin{cases} \pi^y(1 - \pi)^{1-y} & \text{for } y \in \{0, 1\} \\ 0 & \text{otherwise} \end{cases}$$

- Thus, $Pr(y = 0) = \pi^0(1 - \pi)^{1-0} = 1 - \pi$ and $Pr(y = 1) = \pi^1(1 - \pi)^{1-1} = \pi$.
 $E(y) = \sum_y yf(y) = 0(1 - \pi) + 1 * \pi = \pi$
 $Var(y) = E(y) - \{E(y)\}^2 = \pi - \pi^2 = \pi(1 - \pi)$
 Note that $Var(y)$ is a function of $E(y)$.

The Logistic Regression Model

- For $i = 1, \dots, n$; $y_i \sim \text{Bernoulli}(\pi_i)$, where $\pi_i = \frac{\exp(\mathbf{X}'_i\beta)}{1 + \exp(\mathbf{X}'_i\beta)}$ and y_1, \dots, y_n are independent.
- logit function: $g(\pi) = \log\left(\frac{\pi}{1-\pi}\right)$.
maps the interval $(0,1)$ to the real line $(-\infty, \infty)$.
- Odds of event $A \equiv \frac{\Pr(A)}{1-\Pr(A)}$, so $\log\left(\frac{\pi}{1-\pi}\right)$ is the log("odds").
- Note that:

$$\begin{aligned}g(\pi_i) &= \log\left(\frac{\pi}{1-\pi}\right) \\&= \log\left[\frac{\exp(\mathbf{X}'_i\beta)}{1 + \exp(\mathbf{X}'_i\beta)} / \frac{1}{1 + \exp(\mathbf{X}'_i\beta)}\right] \\&= \log[\exp(\mathbf{X}'_i\beta)] \\&= \mathbf{X}'_i\beta\end{aligned}$$

- Thus, the logistic regression model says that $y_i \sim \text{Bernoulli}(\pi_i)$ where $\log\left(\frac{\pi_i}{1-\pi_i}\right) = \mathbf{X}'_i\beta$

- An initial model for the Challenger data is $\log\left(\frac{\pi_i}{1-\pi_i}\right) = \beta_0 + \beta_1 \text{temp}_i$
- The logistic regression model is one example of a Generalized Linear Model
- Reminder: a generalized linear model has:
 - A probability distribution in the exponential scale family: Bernoulli is.
 - A linear predictor, $\eta = \mathbf{X}\beta$: $\mathbf{X}\beta = \beta_0 + \beta_1 \text{temp}_i$
 - A link function that connects $E(\mathbf{y})$ and the linear predictor:

$$g(E(\mathbf{Y})) = g(\mu) = \eta$$
- In this model, the logit is the canonical link function.
- Interpretation of parameters especially easy because of the connection between change in log odds and β
- But, some other link function may fit the data better

Parameter estimation and inference

- Done using likelihood
- The likelihood function for logistic regression is

$$\begin{aligned}l(\boldsymbol{\beta}|\mathbf{y}) &= \sum_{i=1}^n \log[\pi_i^{y_i}(1 - \pi_i)^{1-y_i}] \\&= \sum_{i=1}^n [y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i)] \\&= \sum_{i=1}^n [y_i \{\log(\pi_i) - \log(1 - \pi_i)\} + \log(1 - \pi_i)] \\&= \sum_{i=1}^n [y_i \log\left(\frac{\pi}{1 - \pi}\right) + \log(1 - \pi_i)] \\&= \sum_{i=1}^n [y_i \mathbf{X}'_i \boldsymbol{\beta} - (1 + \exp \{ \mathbf{X}_i \boldsymbol{\beta} \})]\end{aligned}$$

- The likelihood function has to be maximized numerically.
- Fisher's Scoring algorithm commonly used for likelihoods for Generalized Linear Models.
 - uses the expected value of the matrix of second derivatives (-Fisher Information matrix)
 - For Generalized Linear Models, Fisher' Scoring Method results in an iteratively reweighted least squares procedure.
 - The algorithm is presented for the general case in Section 2.5 of Generalized Linear Models 2nd Edition (1989) by McCullagh and Nelder.
- Inference on β
 - For sufficiently large samples, $\hat{\beta}$ is approximately normal with mean β and a variance-covariance matrix that can be approximated by the estimated inverse of Fisher information matrix.
 - Use either Wald tests/intervals or likelihood ratio tests/profile likelihood intervals for inference

Interpretation of Logistic Regression Parameters:

- Let $\tilde{\mathbf{x}}' = (x_1, x_2, \dots, x_{j-1}, x_j + 1, x_{j+1}, \dots, x_p)'$

In other words, $\tilde{\mathbf{X}}$ is the same as \mathbf{x} except that the j^{th} explanatory variable has been increased by one unit.

Let $\pi = \frac{\exp(\mathbf{x}'\beta)}{1 + \exp(\mathbf{x}'\beta)}$ and $\tilde{\pi} = \frac{\exp(\tilde{\mathbf{x}}'\beta)}{1 + \exp(\tilde{\mathbf{x}}'\beta)}$

- The odds ratio

$$\begin{aligned}\frac{\tilde{\pi}}{1 - \tilde{\pi}} / \frac{\pi}{1 - \pi} &= \exp \left\{ \log\left(\frac{\tilde{\pi}}{1 - \tilde{\pi}}\right) - \log\left(\frac{\pi}{1 - \pi}\right) \right\} \\ &= \exp \left\{ \tilde{\mathbf{x}}'\beta - \mathbf{x}'\beta \right\} \\ &= \exp \left\{ (x_j + 1)\beta_j - x_j\beta_j \right\} \\ &= \exp \left\{ \beta_j \right\}\end{aligned}$$

- All other explanatory variables held constant, the odds of success at $x_j + 1$ are $\exp(\beta_j)$ times the odds of success at x_j .

- A 1 unit increase in the j^{th} explanatory variable (with all other explanatory variables held constant) is associated with a multiplicative change in the odds of success by the factor $\exp(\beta_j)$.
- This is true regardless of the initial value x_j .
- Effect on probability of event does depend on the initial odds

| x_j | $P[Y_j = 1 x_j]$ | Odds $ x_j$ | Odds $ x_j + 1$ | $P[Y_j = 1 x_j + 1]$ |
|-------|------------------|-------------|-----------------|----------------------|
| -5 | 0.0067 | -5 | -4 | 0.018 |
| -1 | 0.268 | -1 | 0 | 0.5 |
| 0 | 0.5 | 0 | 1 | 0.731 |
| 3 | 0.952 | 3 | 4 | 0.982 |

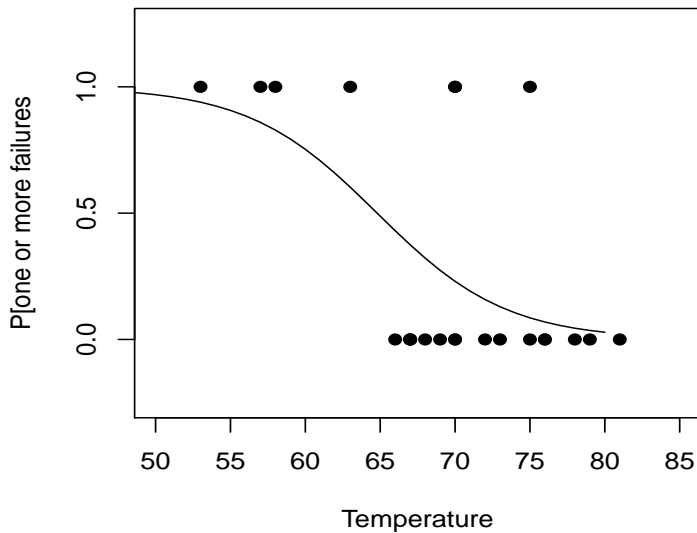
- If (L_j, U_j) is a $100(1 - \alpha)\%$ confidence interval for β_j , then $(\exp\{L_j\}, \exp\{U_j\})$ is a $100(1 - \alpha)\%$ confidence interval for $\exp\{\beta_j\}$.
- Also note that $\pi = \frac{\exp(\mathbf{x}'\beta)}{1+\exp(\mathbf{x}'\beta)} = \frac{1}{\frac{1}{\exp(\mathbf{x}'\beta)} + 1} = \frac{1}{\exp(-\mathbf{x}'\beta)}$

- Thus, if (L_j, U_j) is a $100(1 - \alpha)\%$ confidence interval for $\mathbf{x}'\beta$, then a $100(1 - \alpha)\%$ confidence interval for π is $(\frac{1}{1+\exp(-L_j)}, \frac{1}{1+\exp(-U_j)})$
- Results for Challenger data

| Param. | estimate | se | 95% ci (profile) | p-value | |
|--------|----------|------|---------------------|---------|--------|
| | | | | Wald | LRT |
| Int. | 15.04 | 7.38 | (3.33, 34.34) | 0.042 | |
| Temp. | -0.23 | 0.11 | (-0.515, -0.061) | 0.0320 | 0.0048 |

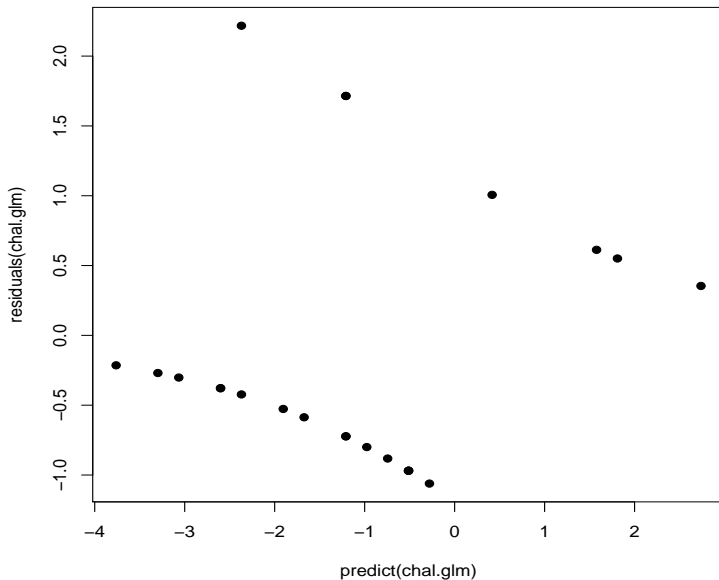
- Predictions:

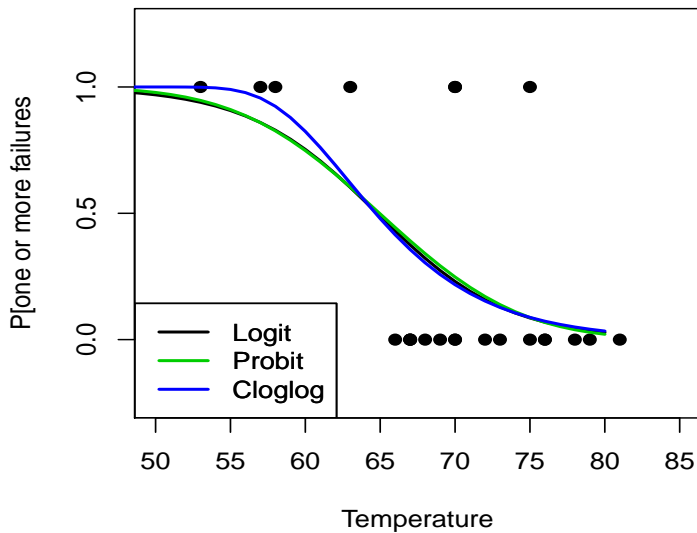
| Temp | Xb | P[one or more failure] | se |
|------|-------|------------------------|-------|
| 50 | 3.43 | 0.968 | 0.061 |
| 60 | 1.11 | 0.753 | 0.191 |
| 70 | -1.21 | 0.230 | 0.105 |
| 80 | -3.53 | 0.028 | 0.039 |



Model assessment / diagnostics

- Harder with Bernoulli data because residuals not very informative
- If sufficient observations, define groups of obs with similar x_i
compute $P[Y = 1]$ for each group, compare to model predictions
- Consider more complicated models:
logit $\pi_i = \beta_0 + \beta_1 x_i$: AIC = 24.38
logit $\pi_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2$: AIC = 25.3887
- Consider different link functions:
 - Probit: $\Phi^{-1}(\pi_i) = \beta_0 + \beta_1 x_i$: AIC = 24.38
where Φ^{-1} is inverse cdf of a normal distribution
 - Comp. log log: $\log(-\log \pi) = \beta_0 + \beta_1 x_i$: AIC = 23.53
- If $Y \sim \text{Poisson}(\lambda)$, $P[Y = 0] = \exp(-\lambda)$.
 $\log(-\log P[Y = 0]) = \log \lambda$.





R code for 0/1 logistic regression

```
# fit a logistic regression to the Challenger data
#   coded as 1 = 1 or more failures

chal <- read.table('challenger2.txt', as.is=T, header=
plot(chal$temp, chal$fail)

# fit a bernoulli logistic regression using
#   default link (logit if binomial)

chal.glm <- glm(fail~temp, data=chal, family=binomial)
# family= specifies which exponential family dn to use
#   this sets the link function and the variance funct
#   can override those defaults if needed, see ?glm
```

```
# all the usual helper functions:
```

```
# print(), summary()
```

```
# coef(), vcov()
```

```
# anova()
```

```
# predict(), residual()
```

```
# most behave in the way you might expect
```

```
# anova(), predict() and residual() are tricky
```

```
anova(chal.glm)    # change in deviance only
```

```
# to get a test, need to specify the appropriate dn
```

```
anova(chal.glm, test='Chi')    # Chi-square test
```

```
anova(chal.glm, test='F')      # if est. overdispersion
```



```
# predict() can return different types of predictions
# ?predict.glm() gives the full story
predict(chal.glm)
# returns predictions on the Xb (linear predictor scale)
predict(chal.glm,type='response')
# returns predictions on the Y scale

# residuals() can return 5 different types of residuals
residuals(chal.glm,type='pearson')
# pearson Chi-square residuals
residuals(chal.glm,type='deviance')
# deviance residuals

plot(predict(chal.glm),
      residuals(chal.glm,type='deviance'))
```

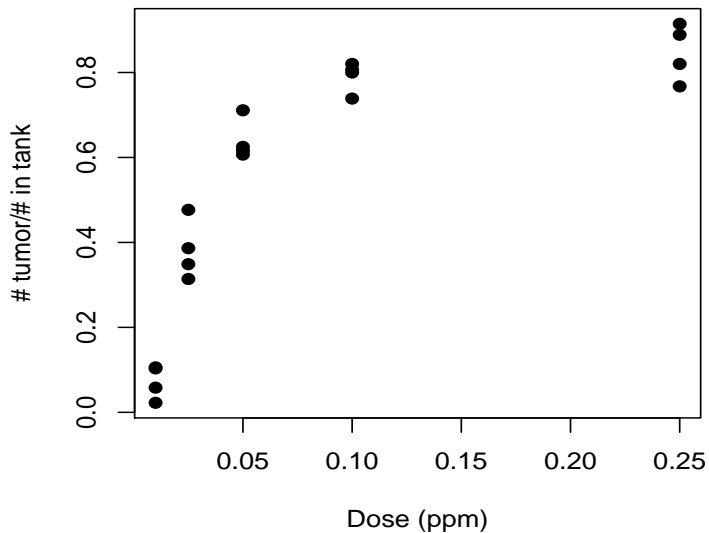
```
# How similar are the two types of resid?
plot(residuals(chal.glm, type='pearson'),
     residuals(chal.glm,type='deviance'))

# to overlay data and curve, need type='response'
plot(chal$temp, chal$fail, pch=19)
lines(30:85, predict(chal.glm, newdata=data.frame(temp=
  type='response'))),

# change link functions
chal.glm2 <- glm(fail~temp, data=chal,
  family=binomial(link=probit))
chal.glm3 <- glm(fail~temp, data=chal,
  family=binomial(link=cloglog))
lines(30:85,
  predict(chal.glm2, newdata=data.frame(temp=30:85),
    type='response'), col=3)
lines(30:85,
```

Logistic Regr. Model for Binomial Count Data

- Bernoulli model appropriate for 0/1 response on an individual
- What if data are # events out of # trials per subject?
- Example: Toxicology study of the carcinogenicity of aflatoxicol.
 - (from Ramsey and Schaefer, *The Statistical Sleuth*, p 641)
 - Tank of trout randomly assigned to dose of aflatoxicol
 - 5 doses. CRD. 4 replicate tanks per dose
 - 86-90 trout per tank
 - Response is # trout with liver tumor
- Could use Bernoulli model for each individual fish
- But, all fish in a tank have the same covariate values (dose)
- easier to analyze data in summarized form (# with tumor, # in tank)



- Each response is # “events” out of # trials.
- $y_i \sim \text{Binomial}(m_i, \pi_i)$, $i = 1, \dots, n$, where m_i is a known number of trials for observation i .

•

$$\pi_i = \frac{\exp(\mathbf{x}'_i \beta)}{1 + \exp(\mathbf{x}'_i \beta)}$$

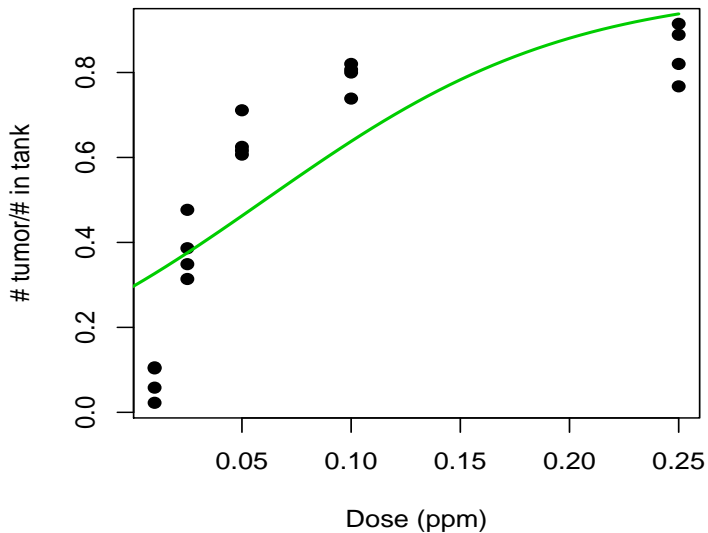
- y_1, \dots, y_n are independent.
- Note: two levels of independence assumed in this model
 - each response, y_i , is independent
 - trials within each each response are independent $|\pi_i$
- $y_i \sim \text{Binomial}(m_i, \pi_i)$ when
 - $y_i = \sum_{j=1}^{m_i} Z_{ij}$, where $Z_{ij} \sim \text{Bernoulli}(\pi_i)$
 - **And** Z_{ij} independent
- May be an issue for both the Challenger and trout data sets
- Binomial model assumes no flight (tank) effects
- Data are same as ≈ 360 fish raised individually, or the same as one tank with ≈ 360 fish

- For now, assume Binomial model reasonable
- Facts about Binomial distributions: If $y_i \sim \text{Binomial}(m_i, \pi_i)$
 - $E(y_i) = m_i \pi_i$
 - $\text{Var}(y_i) = m_i \pi_i (1 - \pi_i)$
 - $f(y_i) = \binom{m_i}{y_i} \pi_i^{y_i} (1 - \pi_i)^{m_i - y_i}$ for $y_i \in \{0, \dots, m_i\}$
 - $l(\beta|\mathbf{y}) = \sum_{i=1}^n [y_i \log(\frac{\pi_i}{1-\pi_i}) + m_i \log(1 - \pi_i)] + \text{const} = \sum_{i=1}^n [y_i \mathbf{x}_i' \beta - m_i \log(1 + \exp\{-\mathbf{x}_i' \beta\})] + \text{const}.$
- $l(\beta|\mathbf{y}, \mathbf{m}, \mathbf{x})$ for $(y_1, m_1, x_1), (y_2, m_2, x_2), \dots, (y_n, m_n, x_n)$ same (apart from constant) as Bernoulli lnL: $l(\beta|\mathbf{z})$ for $(z_{ij}, x_{ij}) = (0, x_1), (0, x_1), \dots, (1, x_1), \dots, (0, x_2), \dots, (1, x_2), \dots, \dots, (0, x_n), \dots, (1, x_n), \dots$
- MLE's $\hat{\beta}$ obtained by numerically maximizing $l(\beta|\mathbf{y})$ over $\beta \in \mathbb{R}^p$

- Results for trout data:

| Coefficient | Estimate | se | z | p |
|-------------|----------|-------|--------|---------|
| Intercept | -0.867 | 0.076 | -11.30 | <0.0001 |
| dose | 14.33 | 0.937 | 15.30 | <0.0001 |

- Looks impressive, but ...

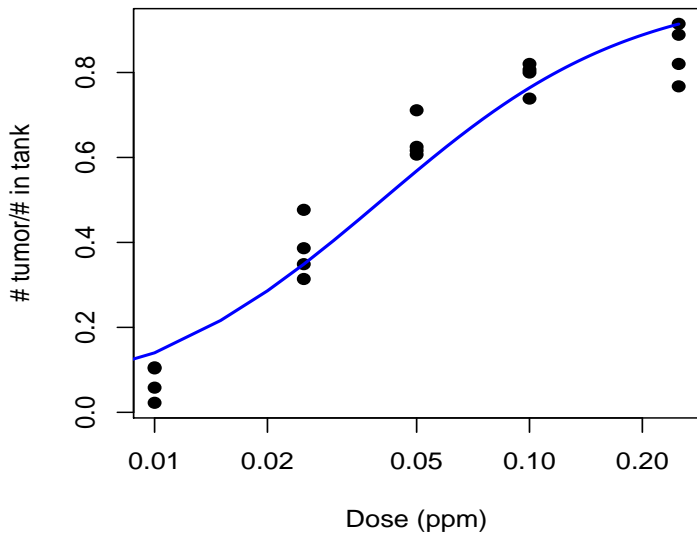


Testing lack of fit

- Remember ANOVA lack of fit test from 500
 - Compare fit of regression to fit of a means model.
 - Required replicated observations so you can estimate $\text{Var } y$
- Can use same ideas in logistic regression:

| Model | df | Deviance | p-value |
|-----------|----|----------|---------|
| lin. reg. | 18 | 277.047 | |
| means | 15 | 25.961 | |
| diff. | 3 | 251.09 | <0.0001 |

- Better fit when $X=\log(\text{dose})$, but still large LOF
- Problem is that $P[\text{tumor} | \text{dose}]$ seems to asymptote at ≈ 0.8
- Standard logistic model asymptotes at 1.0



Residual deviance as a LOF test

- When obs. \sim Binomial, there is another way to assess model fit
- Remember, don't need to estimate σ^2 .
- $\text{Var } y_i = m_i E y_i(1 - E y_i)$
- Can fit (and use in sensible ways) a model with a separate parameter for each observation.
- Called a “saturated” model. Has one π_i parameter for each y_i observation.

- Logistic Regression Model

$$y_i \sim \text{Binomial}(m_i, \pi_i)$$

y_i, \dots, y_n independent

$$\pi_i = \frac{\exp(\mathbf{x}_i' \beta)}{1 + \exp(\mathbf{x}_i' \beta)}$$

for some $\beta \in \mathbb{R}^p$

$1, \dots, n$ with no other restrictions

p parameters

- Saturated Model

$$y_i \sim \text{Binomial}(m_i, \pi_i)$$

y_i, \dots, y_n independent

$$\pi_i \in [0, 1]$$

for $i =$

n parameters

- Can compare the fit of the two models using $D = -2 \ln L(\hat{\beta}|\mathbf{y}) - (-2 \ln L(\hat{\pi}|\mathbf{y}))$
- This statistic is sometimes called the Deviance Statistic, the Residual Deviance, or just the Deviance.
- Under H_0 : regression model fits the data, i.e.

$$\pi_i = \frac{\exp(\mathbf{x}_i' \beta)}{1 + \exp(\mathbf{x}_i' \beta)}, \quad D \sim \chi_{n-p}^2$$
- To test lack of fit, compare D to χ_{n-p}^2
- This is an asymptotic result.
- The χ^2 approximation to the null distribution works reasonably well if $m_i \geq 5$ for most i .
- Do not have to fit two models to calculate the Deviance statistic
 - Let $\hat{\pi}_i = \frac{\exp(\mathbf{x}_i' \beta)}{1 + \exp(\mathbf{x}_i' \beta)}$ for all $i = 1, \dots, n$.
 - Then, the likelihood ratio statistic for testing the logistic regression model as the reduced model vs. the saturated model as the full model is $D = \sum_{i=1}^n 2[y_i \log(\frac{y_i}{m_i \hat{\pi}_i}) + (m_i - y_i) \log(\frac{m_i - y_i}{m_i - m_i \hat{\pi}_i})]$

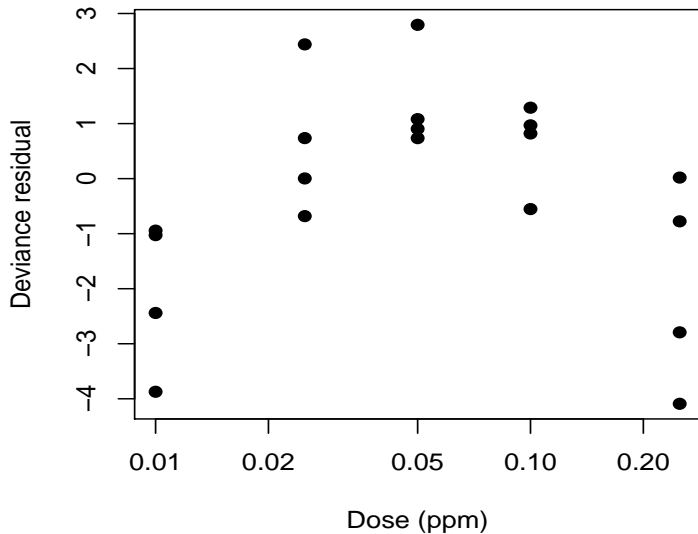
- For the trout data:

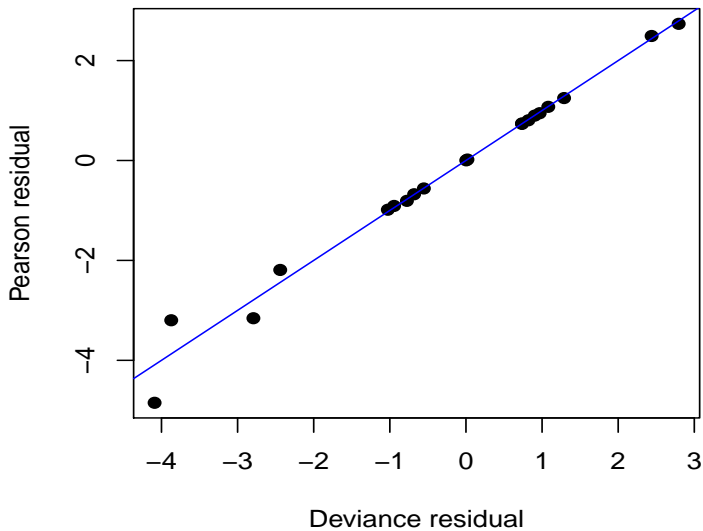
| X | Residual Deviance | d.f. | p for LOF |
|----------|-------------------|------|-----------|
| dose | 277.05 | 18 | <0.0001 |
| log dose | 68.90 | 18 | <0.0001 |

- Plotting data and curve works for one X variable. Harder for two X 's and impossible for many X 's
- Goal: a quantity like the residual, $y_i - \hat{y}_i$, in a linear regression that we can use to diagnose problems with the generalized linear model

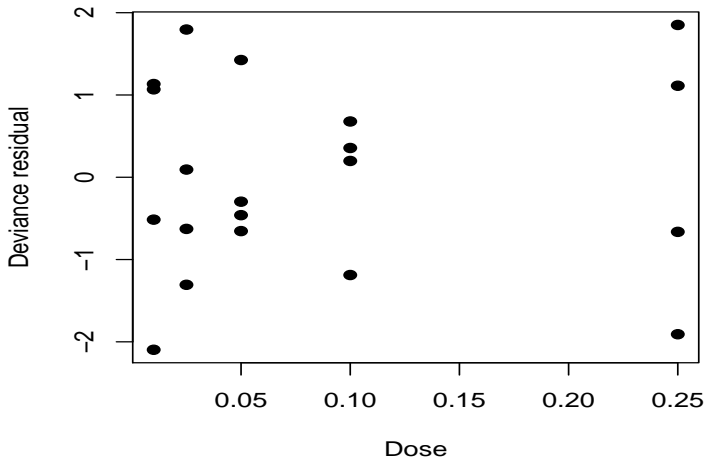
Defining a residual in a logistic regression

- in usual GM model, y_i have constant variance: $y_i \sim (\mathbf{X}\beta, \sigma^2)$
- in a logistic regression $\text{Var } y_i$ depends on π_i
- Two common definitions of residuals for logistic regression:
 - 1 Deviance residual:
 - $d_i = \text{sign}(y_i - m_i \hat{\pi}_i) \sqrt{2[y_i \log(\frac{y_i}{m_i \hat{\pi}_i}) + (m_i - y_i) \log(\frac{m_i - y_i}{m_i - m_i \hat{\pi}_i})]}$
 - The residual deviance statistic $D = \sum_{i=1}^n d_i^2$.
 - 2 Pearson χ^2 residual:
 - $r_i = \frac{y_i - m_i \hat{\pi}_i}{\sqrt{m_i \hat{\pi}_i (1 - \hat{\pi}_i)}}$
 - Because $\sum_{i=1}^n r_i^2 = \sum_{i=1}^n (\frac{y_i - m_i \hat{\pi}_i}{\sqrt{m_i \hat{\pi}_i (1 - \hat{\pi}_i)}})^2 = \sum_{i=1}^n (\frac{y_i - \hat{E}(y_i)}{\sqrt{\hat{\text{Var}}(y_i)}})^2$
- For large m_i 's, both d_i and r_i should behave like standard normal random variables if the logistic regression model is correct.





- For fun, let's evaluate LOF of the means model
5 parameters, one π_j for each dose



- But Residual deviance = 25.96 with 15 df. $p = 0.038$
- Test says that the means model doesn't fit the data!

Overdispersion:

- The LOF test is evaluating all aspects of the model.
- One is the fit of the model for π_i
- A second is the implied variance
 - For many GLM distributions, $\text{Var}(y_i)$ is a function of $E(y_i)$.
 - Logistic regression:
$$\text{Var}(y_i) = m_i \pi_i (1 - \pi_i) = m_i \pi_i - \frac{(m_i \pi_i)^2}{m_i} = E(y_i) - \frac{[E(y_i)]^2}{m}$$
 - So estimating π_i provides estimates of $\text{Var } y_i$
- If the variability of our response is greater than we should expect based on our estimates of the mean, we say that there is overdispersion.
- That is the problem with the means model for the trout data

Overdispersion

- Account for overdispersion by introducing an additional scale parameter, ϕ : $\text{Var}(y_i) = \phi m_i \pi_i (1 - \pi_i)$
- Observations no longer have Binomial distributions, but have the variance pattern characteristic (apart from ϕ) of a Binomial distribution.
- Called a quasiBinomial distribution
- Can also use for underdispersion, but that rarely happens.
- The dispersion parameter ϕ can be estimated from
 - The residual deviance: $\frac{\sum_{i=1}^n d_i^2}{n-p}$
 - Or, the Pearson Chi-square statistic: $\frac{\sum_{i=1}^n r_i^2}{n-p}$
- Beware: can not distinguish between model lack of fit and overdispersion.
- Make sure model is reasonable before estimating ϕ

- For the trout data, using means model:
 - Residual deviance: 25.96 with 15 df.
 - $\hat{\phi} = 1.73$
- Because we have replicate tanks, can estimate $E y_i/m_i$ and $\text{Var } y_i/m_i$ for each dose. Compare $\text{Var } y_i/m_i$ to Binomial implied variance. Approx. calculation because m_i not constant. Use $m_i = 88$, so implied $\text{Var } y_i/m_i \approx \hat{\pi}_i(1 - \hat{\pi}_i)/88$

| Dose | ave. y_i/N_i | implied Var | sample Var | ratio |
|-------|----------------|-------------|------------|-------|
| 0.010 | 0.072 | 0.000764 | 0.00159 | 2.09 |
| 0.025 | 0.381 | 0.00268 | 0.00491 | 1.83 |
| 0.05 | 0.640 | 0.00262 | 0.00232 | 0.88 |
| 0.10 | 0.791 | 0.00187 | 0.00131 | 0.70 |
| 0.25 | 0.847 | 0.00147 | 0.00446 | 3.04 |
| ave. | | | | 1.71 |

Adjusting inference for overdispersion

- Estimation done by quasilielihood, similar to a likelihood but only depends on mean and variance of the “distribution”
- Overdispersed Logistic regression: $\text{Var } y_i/m_i = \phi \pi_i(1 - \pi_i)/m_i$
- Adjustments to inference
 - 1 The estimated variance of $\hat{\beta}$ is multiplied by $\hat{\phi}$.
 - 2 For Wald type inferences (tests, ci's), the standard normal null distribution is replaced by t with $n - p$ degrees of freedom.
 - 3 A test statistic T that was assumed X_q^2 under H_0 is replaced with $\frac{T}{q\hat{\phi}}$ and compared to an F distribution with q and $n - p$ degrees of freedom.
- Above are analogous to changes to inference for normal theory Gauss-Markov linear models if we switched from assuming $\sigma^2 = 1$ to assuming σ^2 was unknown and estimating it with MSE.
(Here ϕ is like σ^2 and $\hat{\phi}$ is like MSE.)

Poisson regression

- What if y_i is a count?
- Often modelled by a Poisson distribution
- We begin with a review of the basics of the Poisson distribution.
- $y \sim \text{Poisson}(\mu) \Rightarrow$

$$f(y) = \begin{cases} \frac{\mu^y e^{-\mu}}{y!} & \text{for } y \in \{0, 1, 2, \dots\} \\ 0 & \text{otherwise} \end{cases}$$

$$E(y) = \mu$$

$$\text{Var}(y) = \mu$$

- μ must be ≥ 0

- The usual Poisson regression model:
 $y_i \sim \text{Poisson}(\mu_i)$ for $i = 1, \dots, n$
- $\mu_i \equiv \exp(\mathbf{x}'_i \boldsymbol{\beta})$
- y_1, \dots, y_n are independent.
- Note that $\mu_i \equiv \exp(\mathbf{x}'_i \boldsymbol{\beta}) \Leftrightarrow \log(\mu_i) = \mathbf{x}'_i \boldsymbol{\beta}$
- Thus, log is the link function in this case.
- Both the Binomial and Poisson distributions are in the exponential family
- Both Logistic and Poisson regression are Generalized Linear Models.
- Estimation, inference, and model diagnosis analogous to those for logistic regression
- Differences: $\text{Var } y_i = \mu_i$ and log link instead of logit link

Interpretation of parameters in a Poisson model with a log link

- Consider two \mathbf{X}_i vectors:

$\mathbf{X}_1 = (\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_{j-1}, \mathbf{X}_j, \mathbf{X}_{j+1}, \dots, \mathbf{X}_p)$, and

$\mathbf{X}_2 = (\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_{j-1}, 1 + \mathbf{X}_j, \mathbf{X}_{j+1}, \dots, \mathbf{X}_p)$



$$\mu_2 = \exp(\mathbf{X}_2\beta) = \exp(\mathbf{X}_1\beta + \beta_j) = \mu_1 \exp(\beta_j)$$

- If X_j increases by one unit, the mean is multiplied by e^{β_j}
- Or, e^{β_j} is the multiplicative effect of a one unit change in X_j .

Using counts to model rates

- Sometimes the response is a count, but the real interest is a rate.
e.g. Number of accidents per man-hours worked on a construction job
- Not Binomial because # man-hours may not be integer
 - define: $y_i = \#$ accidents, $o_i = \#$ man-hours,
 - Interested in $r_i = y_i/o_i$ as a function of covariates
- Model: $y_i \sim \text{Poisson}(o_i r_i) = \text{Poisson}(o_i \exp[f(\mathbf{X}_i, \beta)])$
- $\log \mu_i = \log o_i + \mathbf{X}_i \beta$
- $\log o_i$ is like another X variable, but associated β is exactly 1
- $\log o_i$ is called an offset
- NB: need $\log o_i$, not o_i , as the offset

R code for Binomial and Poisson models

```
# analysis of the trout aflatoxicol data

trout <- read.table('trout.txt',as.is=T,header=T)

# for binomial data, the response is a 2 col. matrix
#   (# event, # not event)
#   Often n is the number tested!

plot(trout$dose, trout$y/trout$n)

trout$f <- cbind(trout$y,trout$n-trout$y)

trout.glm <- glm(f~dose,data=trout,
  family=binomial)
```

```
# lack of fit test relative to means model
trout$dose.f <- factor(trout$dose)
trout.glm2 <- glm(f~dose.f,data=trout,
  family=binomial)
```

```
anova(trout.glm,trout.glm2,test='Chi')
```

```
# or could get in one fit from sequential SS
trout.glm2b <- glm(f~dose+dose.f,data=trout,
  family=binomial)
anova(trout.glm2b,test='Chi')
```

```
trout.glm3 <- glm(f~log(dose),data=trout,
  family=binomial)
plot(trout$dose,residuals(trout.glm3,log='x',
  type='deviance'))

# estimate and account for overdispersion
trout.glm4 <- glm(f~log(dose),data=trout,
  family=quasibinomial)
summary(trout.glm4)
# glm() uses Pearson Chi2/residual df as
#   estimate of overdisp.

anova(trout.glm4,test='F')
# use F instead of Chi2 when estimating over disp.
```

```
# Here's how to fit a Poisson model - use trout
#   data, but assume count ~Poisson

trout.pgml <- glm(y~log(dose),data=trout,
  family=poisson)
summary(trout.pgml)
trout.pgml2 <- glm(y~log(dose),data=trout,
  family=quasipoisson)

# include offset as a vector, remember needs
#   to be a log scale value, # so to consider
#   # tested (trout$n) as the basis for a rate,
trout.pgml3 <- glm(y~log(dose),data=trout,
  family=quasipoisson, offset=log(n))
```