

Statistical Theory & Applied Statistics

Part III Courses, Michaelmas 2009

Revision Notes

Daniel Guetta

guetta@cantab.net

Preliminary Computing

Directory, Linux & Miscellaneous operations

- Changing password: **passwd**
- Creating a new directory: **mkdir**
- Changing to directory **newDir**: **cd newDir**. The command **cd ..** moves up one level in the directory tree.
- List all files in the current directory: **ls**. Can be followed by ***.tex**, for example, to only list tex files.
- Delete **File1**: **rm File1**
- Move **File1** to directory **newDir**: **mv File1 newDir**
- Copy **File1** to a file with name **File2**: **cp File1 File2**
- Examine the content of **File1**: **more File1**
- Leaving emacs: **ctrl x, ctrl c**

Printing

- Looking at the print queue: **lpq**
- Looking at my jobs: **ps x**

General Splus7 stuff

- Opening Splus for the first time in a given directory: **Splus7 CHAPTER**
- Opening Splus7: **Splus7 -e** The *-e* allows the use of arrows
- Leaving Splus7: **q()**
- Listing all objects created: **objects()**
- Removing an object: **rm(objectName)**
- Getting help with function blah: **?blah**
- **rep(1:5,9)** repeats the sequence "1 2 3 4 5" 9 times. **rep(1:5,9)** produces 9 1s, followed by 9 2s, etc...

Imputing and handling data

- Reading **FileName** into table **dataVar** with column titles **x** and **y**:
dataVar <- scan("FileName", list(x=0,y=0)) The "list" function

simply creates an empty table with two columns; x containing value 0 and y containing value 0. Without this argument, the entire file is read left-to-right into a single list.

- Reading FileName into table dataVar, when the column titles are already included in the file FileName: **dataVar <- read.table("tumour", header=T)**
- Inputting a vector directly into x: **x <- scan()** *End input with a blank line*
- Listing the columns of a table tableName: **names(tableName)**
- Extracting the x from the table dataVar: **dataVector <- dataVar\$x**
- **x[-11]** returns vector x without observation 11
- Saving each column in table dataVar as a separate vector, with the same name as the column title: **attach(dataVar)**
- Joining several columns into a matrix: **cbind(...)**
- Many functions (glm, plot, etc...) can be applied to subsets of the data by adding an argument to the function. For example, the run a function only on data in which var1 takes the value A, add the argument **subset=(var1=="A")** to the function.
- **table(var1,var2)** will create a contingency table, listing the number of times each possible pair of values occurs.

Simple graphical stuff

- Plotting two vectors against each other: **plot(xVector, yVector)**
- Follow the above by: **identify(xVector, yVector)** to be able to identify different points on the plot.
- Adding a straight line to the current plot: **lines(xVector, yVector, type^{Optional})** *"Type" takes values "p" for points, "l" for a line or "b" for both*
- Fitting two-by-two plots on each page: **par(mfrow=c(2,2))**

Outputting stuff to a file

- To output graphics, begin with **postscript(file='fileName.ps')** and create all graph using **plot** (none of them will appear on screen). When done, use **dev.off()**.
- To save output to a file, start with **sink('fileName')** – all subsequent output will go to the file. End with **sink()**.

Simple statistical stuff

- The upper α point of a t_n distribution is given by: **qt(1- α , n)**
- The upper α of an $F_{m,n}$ distribution is given by: **qf(1- α , m, n)**
- Generating vectors containing a sample of n simulated $N(0,1)$ observations: **rnorm(n)**

Preliminary Mathematics

This section contains preliminary mathematics and proofs that are assumed in these notes.

Projections

- A **idempotent** matrix P on an underlying vector space W satisfies $P^2 = P$ and is a **projection**. The eigenvalues of P can only be 0 or 1, with the following eigenspaces

$$U = \text{Range/Image of } P = \{P\mathbf{y} : \mathbf{y} \in W\}$$

$$V = \text{Null space/kernel of } P = \{\mathbf{x} : P\mathbf{x} = 0, \mathbf{x} \in W\}$$

We say P projects W onto U .

- $W = U \oplus V$; for any $\mathbf{x} \in W$ we can write $\mathbf{x} = P\mathbf{x} + (\mathbf{x} - P\mathbf{x})$, where the first vector is in U and the second is in V .

If U and V are orthogonal (see next point), we can show that such a decomposition is unique. Consider

$$\mathbf{x} = \mathbf{u} + \mathbf{v} = \mathbf{u}' + \mathbf{v}'$$

Then

$$0 = \|\mathbf{x} - \mathbf{x}\|^2 = \|(\mathbf{u} - \mathbf{u}') + (\mathbf{v} - \mathbf{v}')\|^2 = \|(\mathbf{u} - \mathbf{u}')\|^2 + \|(\mathbf{v} - \mathbf{v}')\|^2$$

Since $\{\mathbf{u} - \mathbf{u}' \in U\} \perp \{\mathbf{v} - \mathbf{v}' \in V\}$. Thus, we must have $\mathbf{u} = \mathbf{u}', \mathbf{v} = \mathbf{v}'$.

- A projection is **orthogonal** if $V = U^\perp = \{\mathbf{y} \in W : \mathbf{u}^T \mathbf{y} = 0, \forall \mathbf{u} \in U\}$.

Consider $\mathbf{x}, \mathbf{y} \in W$ and $\mathbf{u} = P\mathbf{x} \in U$, $\mathbf{v} = \mathbf{y} - P\mathbf{y} \in V$. We then have

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{x}^T P^T (\mathbf{y} - P\mathbf{y}) = \mathbf{x}^T (P^T - P) \mathbf{y}$$

We therefore have that a projection is orthogonal if and only if $P^T = P$.

- If P and Q are projections and $PQ = 0$, then the range of each of these vectors are orthonormal.
- For the orthogonal projection onto a vector \mathbf{u} , $P = \mathbf{u}\mathbf{u}^T$.

Generating Functions

- Let X be a real-valued random variable. The **moment generating function** of X is given by

$$M(t) = \mathbb{E}(e^{tX}) \quad \forall t \text{ s.t. } \mathbb{E}(e^{tX}) < \infty$$

- Note that, by an extension of linearity

$$M(t) = \mathbb{E}\left(\sum_{k=0}^{\infty} \frac{t^k X^k}{k!}\right) = \sum_{k=0}^{\infty} \mathbb{E}\left(\frac{t^k X^k}{k!}\right) = \sum_{k=0}^{\infty} \frac{t^k}{k!} \mathbb{E}(X^k)$$

$$\boxed{\left. \frac{d^n M}{dt^n} \right|_{t=0} = \mathbb{E}(X^n)}$$

- If $X \sim N(\mu, \sigma^2)$, then

$$M(t) = \exp\left(\mu t + \frac{1}{2} \sigma^2 t^2\right)$$

- The **cumulant generating function** is given by

$$K(t) = \log M(t)$$

We can write K as a power series

$$K(t) = \log M(t) = \sum_{n=1}^{\infty} \kappa_n \frac{t^n}{n!}$$

κ_n is called the n^{th} **cumulant**. Clearly

$$\kappa_r = \left. \frac{d^r K_Y(t)}{dt^r} \right|_{t=0}$$

Note also that we can use the expansion $\log(1+z) = z - \frac{1}{2}z^2 + \dots$ on the series expansion $M(t) = 1 + t\mathbb{E}(X) + \frac{1}{2}t^2\mathbb{E}(X^2) + \dots$ to get

$$\begin{aligned} K(t) &= \left\{ t\mathbb{E}(X) + \frac{1}{2}t^2\mathbb{E}(X^2) \right\} - \frac{1}{2} \left\{ t\mathbb{E}(X) + \frac{1}{2}t^2\mathbb{E}(X^2) \right\}^2 \\ &= \left\{ t\mathbb{E}(X) + \frac{1}{2}t^2\mathbb{E}(X^2) \right\} - \frac{1}{2} \left\{ t^2 [\mathbb{E}(X)]^2 + \dots \right\} \\ &= \left\{ \mathbb{E}(X) \right\} t + \left\{ \mathbb{E}(X^2) - [\mathbb{E}(X)]^2 \right\} \frac{t^2}{2} \\ &= \mu t + \sigma^2 \frac{t^2}{2} \end{aligned}$$

This confirms that $\kappa_1 = \mu$ and $\kappa_2 = \sigma^2$.

Distributions

- If $Z_i \sim N(0,1)$ where $i = 1, \dots, n$

$$\sum_{i=1}^n Z_i^2 \sim \chi_n^2$$

Where χ_n^2 is a χ^2 distribution on n degrees of freedom.

$$\mathbb{E}(\chi_n^2) = n$$

- If $Z \sim N(0,1)$ and $Y \sim \chi_n^2$ and Z and Y are independent, then

$$\frac{Z}{\sqrt{Y/n}} \sim t_n$$

Where t_n is a t distribution on n degrees of freedom, which is symmetric with heavier tails than the normal distribution, and converges in distribution to $N(0,1)$ as $n \rightarrow \infty$.

$$\mathbb{E}(t_n) = 0$$

- If $X \sim \chi_m^2$ and $Y \sim \chi_n^2$ and X and Y are independent

$$\frac{X/m}{Y/n} \sim F_{m,n}$$

Where $F_{m,n}$ is an F distribution on m and n degrees of freedom. Note that

$$\mathbb{E}(F_{m,n}) = \frac{n}{n-2}$$

- If $\mathbf{Z} \sim N_n(\boldsymbol{\mu}, I)$, then $U = \mathbf{Z}^T \mathbf{Z}$ has a non-central χ^2 distribution $\chi_n^2(\delta)$ with n degrees of freedom and non-centrality parameter $\delta = \boldsymbol{\mu}^T \boldsymbol{\mu}$.

$$\mathbb{E}(\chi_n^2(\delta)) = n + \delta$$

- If $W_1 \sim \chi_{n_1}^2(\delta)$ and $W_2 \sim \chi_{n_2}^2(\delta)$ and W_1 and W_2 are independent, then $\frac{W_1/n_1}{W_2/n_2}$ has a non-central F distribution $F_{n_1, n_2}(\delta)$ with n_1 and n_2 degrees of freedom and non-centrality parameter δ .

$$\mathbb{E}(F_{n_1, n_2}(\delta)) = \frac{n_2(n_1 + \delta)}{n_1(n_2 - 2)}$$

The Multivariate Normal

Y has a univariate normal distribution with mean μ and variance $\sigma^2 \in (0, \infty)$, and we write $Y \sim N(\mu, \sigma^2)$ if its density is

$$f_y(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y - \mu)^2\right\} \quad y \in \mathbb{R}$$

It's also convenient to define a degenerate normal distribution $Y \sim N(\mu, 0)$ in which $\mathbb{P}(y = \mu) = 1$.

Definition: We say a random vector $\mathbf{Y} \in \mathbb{R}^n$ has an **n -variate normal distribution** if for every vector $\mathbf{t} \in \mathbb{R}^n$, the random variable $\mathbf{t}^T \mathbf{Y}$ has a (possibly degenerate) univariate distribution. We write $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

We write $\boldsymbol{\mu} = \mathbb{E}(\mathbf{Y})$ for the mean vector and $\Sigma = \text{cov}(\mathbf{Y})$ for the covariance matrix¹ of \mathbf{Y} , so that

$$\Sigma_{ij} = \mathbb{E}\{(Y_i - \mu_i)(Y_j - \mu_j)\} = \mathbb{E}(Y_i Y_j) - \mu_i \mu_j$$

The entire dependence structure of \mathbf{Y} is determined by Σ . Components of \mathbf{Y} are independent if and only if Σ is diagonal².

In general, Σ is **symmetric** and **non-negative definite**, because $\mathbf{t}^T \Sigma \mathbf{t} = \text{var}(\mathbf{t}^T \mathbf{Y})$ and variances are non-negative.

If, in addition, Σ is **positive definite**, then none of its eigenvalues are 0 and it is **invertible**. We then have

$$f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^n \|\Sigma\|}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu})\right\} \quad \mathbf{y} \in \mathbb{R}^n$$

Even if Σ is only non-negative definite, its MGF is

$$M_{\mathbf{y}}(\mathbf{t}) = \mathbb{E}(e^{\mathbf{t}^T \mathbf{Y}}) = e^{\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \Sigma \mathbf{t}}$$

We finally note that if $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \Sigma)$ and A is a $p \times n$ matrix, then $A\mathbf{Y} \sim N_p(A\boldsymbol{\mu}, A\Sigma A^T)$ ³.

Laws of Large Numbers

Chebyshev's Inequality: Let X be a random variable with $\mathbb{E}(X) = \mu$ and $\text{Var}(X) = \sigma^2$. Then

$$\mathbb{P}(|X - \mu| \geq \alpha) \leq \frac{\sigma^2}{\alpha^2}$$

¹ The covariance matrix is given by $\text{cov}(\mathbf{X}) = \Sigma = \mathbb{E}\left\{(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\right\}$, and note that $\text{cov}(A\mathbf{X}) = A \text{cov}(\mathbf{X}) A^T$.

² Note that in general, $\text{cov}[u, v] \neq 0$ does **not** imply u and v are independent unless they are jointly normal.

³ To prove it's p -variate normal, choose $\mathbf{t} \in \mathbb{R}^p$, and note that $\mathbf{t}^T A\mathbf{Y} = (A^T \mathbf{t})^T \mathbf{Y}$ where $(A^T \mathbf{t})^T \in \mathbb{R}^n$. For the mean, note that $\mathbb{E}(A\mathbf{Y}) = A\mathbb{E}(\mathbf{Y}) = A\boldsymbol{\mu}$ and for the variance

$$\text{cov}\left((A\mathbf{Y})_i, (A\mathbf{Y})_j\right) = \text{cov}\left(A_{i\alpha} Y_\alpha, A_{j\beta} Y_\beta\right) = A_{i\alpha} A_{j\beta} \text{cov}(Y_\alpha, Y_\beta) = A_{i\alpha} A_{j\beta} \Sigma_{\alpha\beta} = (A\Sigma A^T)_{ij}$$

Proof:

$$\begin{aligned}
 \mathbb{P}\left(|X - \mu| \geq \alpha\right) &= \mathbb{E}\left(\mathbb{I}\{|X - \mu| \geq \alpha\}\right) \\
 &= \mathbb{E}\left(\mathbb{I}\left\{\left(\frac{X - \mu}{\alpha}\right)^2 \geq 1\right\}\right) \\
 &\leq \mathbb{E}\left(\left(\frac{X - \mu}{\alpha}\right)^2\right) \\
 &= \frac{1}{\alpha^2} \mathbb{E}\{(X - \mu)^2\} \\
 &= \frac{\sigma^2}{\alpha^2}
 \end{aligned}$$

As required. ■

We define the following modes of convergence

We say that a sequence of random vectors (Y_n) converges **almost surely** to Y if

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} (Y_n) = Y\right) = 1$$

or equivalently, for every $\varepsilon > 0$

$$\mathbb{P}\left(\sup_{m \geq n} |Y_m - Y| > \varepsilon\right) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

And we write $Y_n \xrightarrow{a.s.} Y$.

We say that a sequence of random vectors (Y_n) converges **in probability** to Y if for every $\varepsilon > 0$

$$\mathbb{P}\left(|Y_n - Y| > \varepsilon\right) \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

And we write $Y_n \xrightarrow{p} Y$.

We say that a sequence of random vectors (Y_n) converges **in distribution** to Y if

$$\mathbb{E}\{f(Y_n)\} \rightarrow \mathbb{E}\{f(Y)\}$$

for all bounded, continuous, real-valued functions f . In fact, it is enough that the convergence occurs when f is bounded and Lipschitz (ie: $\exists L > 0$ such that $|f(x) - f(y)| \leq L\|x - y\|$). Equivalently, if and only if

$$\mathbb{P}(Y_n \leq y) \rightarrow \mathbb{P}(Y \leq y)$$

at all points where the distribution function of Y is continuous.

We write $Y_n \xrightarrow{d} Y$.

$$Y_n \xrightarrow{a.s.} Y \Rightarrow Y_n \xrightarrow{p} Y \Rightarrow Y_n \xrightarrow{d} Y$$

For any h continuous and real-valued:

$$Y_n \xrightarrow{?} Y \Rightarrow f(Y_n) \xrightarrow{?} f(Y)$$

In fact, h may have a set of discontinuities D_h provided that $\mathbb{P}(Y \in D_h) = 0$.

We are now ready to prove the *weak law of large numbers*

Weak Law of Large Numbers: Let X_1, X_2, \dots be an infinite sequence of IID random variables with $\mathbb{E}(X_1) = \mathbb{E}(X_2) = \dots = \mu < \infty$. Define

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n)$$

Then assuming $\text{Var}(X_1) = \text{Var}(X_2) = \dots = \sigma^2 < \infty$ ⁴

$$\bar{X}_n \xrightarrow{p} \mu \quad \text{as } n \rightarrow \infty$$

Proof: By linearity and independence of the X_i , we have that

$$\begin{aligned} \mathbb{E}(\bar{X}_n) &= \frac{1}{n} \mathbb{E}(X_1 + \dots + X_n) = \mu \\ \text{Var}(\bar{X}_n) &= \frac{1}{n^2} \text{Var}(X_1 + \dots + X_n) = \sigma^2 / n \end{aligned}$$

Then by Chebyshev's Inequality

$$\mathbb{P}\left(|\bar{X}_n - \mu| \geq \varepsilon\right) \leq \frac{\sigma^2}{n\varepsilon^2} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

Thus, by the definition of convergence in probability, we have proved our theorem. ■

Strong Law of Large Numbers: With all the definitions above,

$$\bar{X}_n \xrightarrow{a.s.} \mu \quad \text{as } n \rightarrow \infty$$

⁴ The theorem is still true if this condition fails to hold, but the proof is more involved.

Uniform Law of Large Numbers: Suppose $f(x, \theta)$ is a function defined for $\theta \in \Theta$. Further suppose that

1. Θ is compact (closed and bounded).
2. $f(x, \theta)$ is continuous at each $\theta \in \Theta$ for almost all x .
3. There exists a dominating function $g(x)$ such that $\mathbb{E}\{g(X)\} < \infty$ and

$$|f(x, \theta)| \leq g(x) \quad \forall \theta \in \Theta$$

Then $\mathbb{E}\{f(X, \theta)\}$ is continuous in θ and

$$\sup_{\theta \in \Theta} \left\| \left\{ \frac{1}{n} \sum_{i=1}^n f(X_i, \theta) \right\} - \mathbb{E}\{f(X, \theta)\} \right\| \xrightarrow{p} 0$$

Slutsky's Theorem

Slutsky's Theorem: If \mathbf{Y}_n and \mathbf{Z}_n be sequences of random vectors with $\mathbf{Y}_n \xrightarrow{d} \mathbf{Y}$ and $\mathbf{Z}_n \xrightarrow{d} \mathbf{c}$ where \mathbf{c} is constant, then

$$(\mathbf{Y}_n, \mathbf{Z}_n) \xrightarrow{d} (\mathbf{Y}, \mathbf{c})$$

and for any continuous real-valued function g

$$g(\mathbf{Y}_n, \mathbf{Z}_n) \xrightarrow{d} g(\mathbf{Y}, \mathbf{c})$$

Proof: Let f be a bounded, Lipschitz function with Lipschitz constant L .⁵

Given $\varepsilon > 0$, let $\delta = \varepsilon / \{3(L+1)\}$.

Since $\mathbf{Z}_n \xrightarrow{d} \mathbf{c}$ we can choose $n_0 \in \mathbb{N}$ such that⁶

$$\mathbb{P}(\|\mathbf{Z}_n - \mathbf{c}\| > \delta) \leq \frac{\varepsilon}{6(\|f\|_\infty + 1)} \text{ for all } n \geq n_0$$

Since $\mathbf{Y}_n \xrightarrow{d} \mathbf{Y}$ and $f(\cdot, \mathbf{c})$ is bounded and continuous, we can also choose $n_1 \in \mathbb{N}$ such that

$$\left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{c})\} - \mathbb{E}\{f(\mathbf{Y}, \mathbf{c})\} \right| < \frac{\varepsilon}{3} \text{ for all } n \geq n_1$$

We then have that for $n = \max\{n_0, n_1\}$

$$\begin{aligned} & \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{Z}_n)\} - \mathbb{E}\{f(\mathbf{Y}, \mathbf{c})\} \right| \\ &= \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{Z}_n)\} - \mathbb{E}\{f(\mathbf{Y}, \mathbf{c})\} - \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{c})\} + \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{c})\} \right| \\ &= \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{Z}_n) - f(\mathbf{Y}_n, \mathbf{c})\} - \mathbb{E}\{f(\mathbf{Y}, \mathbf{c})\} + \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{c})\} \right| \\ &\leq \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{Z}_n) - f(\mathbf{Y}_n, \mathbf{c})\} \mathbb{I}_{\{\|\mathbf{Z}_n - \mathbf{c}\| \leq \delta\}} \right| \\ &\quad + \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{Z}_n) - f(\mathbf{Y}_n, \mathbf{c})\} \mathbb{I}_{\{\|\mathbf{Z}_n - \mathbf{c}\| > \delta\}} \right| \\ &\quad + \left| \mathbb{E}\{f(\mathbf{Y}_n, \mathbf{c})\} - \mathbb{E}\{f(\mathbf{Y}, \mathbf{c})\} \right| \end{aligned}$$

As required. ■

⁵ Lipschitz continuity is a smoothness condition for functions which is stronger than regular continuity; it limits how fast a function can change. If a function $f: \mathbb{R} \rightarrow \mathbb{R}$ is Lipschitz continuous with Lipschitz constant L , then

$$\frac{|f(x_1) - f(x_2)|}{|x_1 - x_2|} \leq L \quad \forall x_1 \neq x_2$$

Intuitively, a line joining any two points on the graph of f will never have a slope steeper than L ; so there is a double white cone whose vertex can be translated along the graph so that the graph always remains entirely outside the cone.

⁶ For a bounded function f defined over the set S , $\|f\|_\infty = \|f\|_{\infty, S} = \sup\{|f(x)| : x \in S\}$.

Linear Models

Introduction and MLEs

A general linear model in which we make n **observations** and we want to fit p **parameters** takes the following form:

$$\boxed{\underbrace{\mathbf{Y}}_{n \times 1} = \underbrace{X}_{n \times p} \underbrace{\boldsymbol{\beta}}_{p \times 1} + \underbrace{\boldsymbol{\varepsilon}}_{n \times 1}}$$

Where

- X is an $n \times p$ matrix of known **covariates** or **explanatory variables** (with the **first column** consisting of 1s, for the intercept). Also called the **design matrix**.

It is usual to assume that this matrix has **full rank** so that $X^T X$ is **positive definite**⁷.

- \mathbf{Y} is an $n \times 1$ vector of **observations** or **response variables**.
- $\boldsymbol{\varepsilon}$ is an $n \times 1$ vector of IID **errors** $\boldsymbol{\varepsilon} \sim N_n(0, \sigma^2)$.
- $\boldsymbol{\beta}$ is an unknown **vector** of **regression coefficients**.

\mathbf{Y} is a linear combination of constants and normal variables, and in fact

$$\boxed{\mathbf{Y} \sim N_n(X\boldsymbol{\beta}, \sigma^2 \mathbb{I})}$$

In this parametric model, the unknown parameter is $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2)^T$. The **maximum likelihood estimators** of $\boldsymbol{\beta}$ and σ can be worked out⁸

⁷ If X had linearly dependent rows, there would be a vector \mathbf{z} such that $\mathbf{z}X^T X \mathbf{z} = \|\mathbf{Xz}\|^2 = 0$, which would imply $X^T X$ were not positive definite.

⁸ ...by maximising, with respect to both $\boldsymbol{\beta}$ and σ , the loglikelihood

$$\ell(\mathbf{Y}; \boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2 + \text{constants}$$

Specifically,

- In the $\boldsymbol{\beta}$ case, this amounts to minimising the following, with respect to $\boldsymbol{\beta}$

$$\|\mathbf{Y} - X\boldsymbol{\beta}\|^2 = (\mathbf{Y} - X\boldsymbol{\beta})^T (\mathbf{Y} - X\boldsymbol{\beta}) = \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T X \boldsymbol{\beta} - \boldsymbol{\beta}^T X^T \mathbf{Y} + \boldsymbol{\beta}^T X^T X \boldsymbol{\beta}$$

Differentiate using $\frac{\partial(\mathbf{x}\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{\partial(\boldsymbol{\beta}\mathbf{x})}{\partial \boldsymbol{\beta}} = \mathbf{x}$, $\frac{\partial(\mathbf{x}\boldsymbol{\beta}^T)}{\partial \boldsymbol{\beta}} = \frac{\partial(\boldsymbol{\beta}^T \mathbf{x})}{\partial \boldsymbol{\beta}} = \mathbf{x}^T$ and $\frac{\partial(\boldsymbol{\beta}^T A \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta}^T (A^T + A)$,

then transpose the lot. For distribution, note that $\boldsymbol{\beta}$ is a linear transformation of \mathbf{Y} .

- In the σ^2 case, it's a straight differentiation with respect to σ^2 . σ is **not** a linear transformation of \mathbf{Y} , and so the distribution requires Cochran's Theorem – start with

$$\boxed{\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y} \sim N_p \left(\beta, (X^T X)^{-1} \sigma^2 \right)} \quad \boxed{\hat{\sigma}^2 = \frac{1}{n} \|\mathbf{Y} - X\hat{\beta}\|^2 \sim \frac{\sigma^2}{n} \chi_{n-p}^2}$$

And these are **independent of each other**⁹. (Note that the $\|\mathbf{Y} - X\hat{\beta}\|^2$ term in the expression for $\hat{\sigma}^2$ is just the residual sum of squares, $\text{RSS} = \hat{\epsilon}^T \hat{\epsilon}$, and that it can also be written as $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 = \|(I - H)\mathbf{Y}\|^2$).

Note that

- $\hat{\beta}$ is an **unbiased estimator**, but its different components may be correlated, since the covariance matrix is not, in general, diagonal.

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \|\mathbf{Y} - X\hat{\beta}\|^2 = \frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 = \frac{1}{n} \|(I - H)\mathbf{Y}\|^2 = \frac{1}{n} \mathbf{Y}^T (I - H)^T (I - H) \mathbf{Y} \\ &= \frac{1}{n} \mathbf{Y}^T (I - H) \mathbf{Y} = \frac{1}{n} (\mathbf{Y} - X\beta)^T (I - H) (\mathbf{Y} - X\beta) \end{aligned}$$

(The last step is non-trivial and requires some working – start with the result and work backwards). We also have that $I = H + (I - H)$ with both H and $I - H$ symmetric. All we now need to show is that $\text{rank}(H) + \text{rank}(I - H) = n$. We do that by noting that $H^2 = H$ which means that all the eigenvalues of H are 0 or 1, and we can find Q such that $Q^T Q = I$ and $Q^T H Q = D$ where D is diagonal and contains only 0s or 1s. Thus

$$\text{rank}(H) = \text{rank}(D) = \text{tr}(D) = \text{tr}(DQ^T Q) = \text{tr}(QDQ^T) = \text{tr}(H) = \text{tr} \left(\overbrace{(X^T X)^{-1} X^T X}^{I_p} \right) = p$$

And $I_n - H = I_n - QDQ^T = Q(I_n - D)Q^T$, so

$$\text{rank}(I_n - H) = \text{rank}(I_n - D) = n - p$$

⁹ From the penultimate line in the previous footnote,

$$(I_n - H) = Q \begin{pmatrix} I_{n-p} & 0 \\ 0 & 0 \end{pmatrix} Q^T = LL^T \quad L = Q \begin{pmatrix} I_{n-p} \\ 0_{p, n-p} \end{pmatrix}$$

Now consider the vector $\mathbf{T} = \begin{pmatrix} B \\ L^T \end{pmatrix} \mathbf{Y}$ with $B = (X^T X)^{-1} X^T$. Note that $B\mathbf{Y} = \hat{\beta}$ and that $\hat{\sigma}^2 = \mathbf{Y}^T (I - H) \mathbf{Y} = \mathbf{Y}^T LL^T \mathbf{Y} = \|L^T \mathbf{Y}\|^2$. So if we can show the two components of the vector are independent, we have shown what we wanted. Consider the covariance matrix of \mathbf{T} .

$$\Sigma = \sigma^2 \begin{pmatrix} B \\ L^T \end{pmatrix} \begin{pmatrix} B^T & L \end{pmatrix} = \sigma^2 \begin{pmatrix} BB^T & BL \\ L^T B^T & L^T L \end{pmatrix}$$

And finally, note that

$$\begin{aligned} L^T L &= \begin{pmatrix} I_{n-p} & 0_{n-p,p} \end{pmatrix} Q^T Q \begin{pmatrix} I_{n-p} \\ 0_{p, n-p} \end{pmatrix} = I_{n-p} \\ BL &= BL^T L = B(I - H)L = X(X^T X)^{-1} X^T \left(I - X(X^T X)^{-1} X^T \right) L = 0 \end{aligned}$$

The covariance matrix is therefore diagonal, and so the components are independent.

- $\mathbb{E}(\hat{\sigma}^2) = \frac{n-p}{n}\sigma^2$, and so $\hat{\sigma}^2$ is a slightly biased estimator. An unbiased estimator for σ^2 is

$$s^2 = \frac{1}{n-p} \|\mathbf{Y} - X\hat{\boldsymbol{\beta}}\|^2 \sim \frac{\sigma^2}{n-p} \chi_{n-p}^2$$

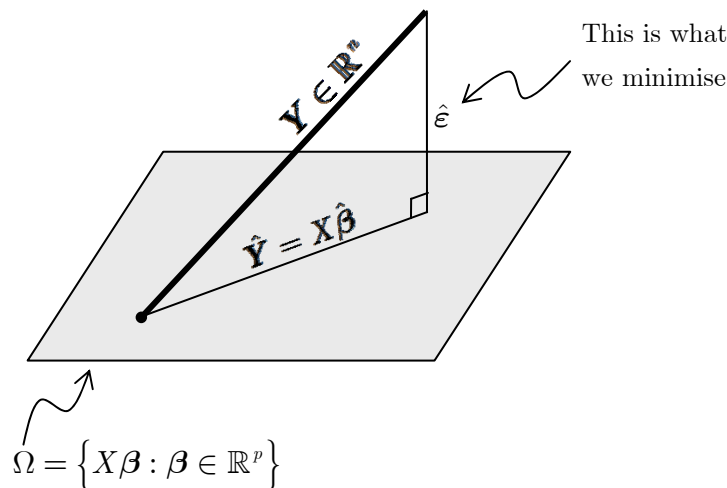
This is the estimator to use in working out the covariance of $\boldsymbol{\beta}$, which is therefore given by $\text{cov}(\hat{\boldsymbol{\beta}}) = (X^T X)^{-1} s^2$.

Geometrical interpretation

The **vector of fitted values** contains our best guess of what \mathbf{Y} should be and is given by

$$\hat{\mathbf{Y}} = X\hat{\boldsymbol{\beta}} = X(X^T X)^{-1} X^T \mathbf{Y} = H\mathbf{Y}$$

We note that H is **symmetric** and **idempotent** ($HH = H$, which also implies that $(I - H)(I - H) = I - H$). These are the properties of a **projection matrix**; indeed H **projects** \mathbf{Y} onto the space spanned by the columns of X , which contains all $X\boldsymbol{\beta}$. Thus, H effectively finds the closest fit in $\Omega = \{X\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ for \mathbf{Y} . (The residuals are given by $\hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \hat{\mathbf{Y}} = (I - H)\mathbf{Y} \sim N_n(0, (I - H)\sigma^2)$):



Hypothesis tests on $\boldsymbol{\beta}$

For each component of $\hat{\boldsymbol{\beta}}$, we have that

- $\text{var}(\hat{\beta}_i) = \sigma^2 (X^T X)^{-1}_{ii}$
- Our best estimator of this is $\text{se}(\hat{\beta}_i) = \sqrt{s^2 (X^T X)^{-1}_{ii}}$

So:

$$\begin{aligned} \frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)} &= \frac{\hat{\beta}_i}{\sqrt{s^2(X^T X)_{ii}^{-1}}} \\ &= \frac{\hat{\beta}_i / \sqrt{\sigma^2(X^T X)_{ii}^{-1}}}{\sqrt{\frac{n-p}{n-p} \sqrt{s^2(X^T X)_{ii}^{-1}} / \sqrt{\sigma^2(X^T X)_{ii}^{-1}}}} \\ &= \frac{\hat{\beta}_i / \sqrt{\sigma^2(X^T X)_{ii}^{-1}}}{\sqrt{\frac{1}{n-p} \frac{s^2(n-p)}{\sigma^2}}} \end{aligned}$$

Since $\hat{\beta}_i \sim N(\beta_i, \sigma^2(X^T X)_{ii}^{-1})$ and $s^2 = \frac{\sigma^2}{n-p} \chi_{n-p}^2$, we can write

$$\sim \frac{N(\beta_i, 1)}{\sqrt{\frac{1}{n-p} \chi_{n-p}^2}}$$

If $\beta_i = 0$, then this ratio is a t distribution with $n - p$ degrees of freedom:

$$H_0 : \mathbb{P}_{\beta_i=0, \sigma^2} \left(\frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)} \leq t_{n-p}(\alpha) \right) = 1 - \alpha$$

So a $(1-\alpha)$ -level confidence set for β_i **assuming** $\beta_i = 0$ is

$$\left\{ \beta_i \in \mathbb{R} : \frac{\hat{\beta}_i}{\text{se}(\hat{\beta}_i)} \leq t_{n-p}(\alpha) \right\}$$

We can also find a confidence set for β . First, notice that

$$(\hat{\beta} - \beta) \sim N_p(0, (X^T X)^{-1} \sigma^2)$$

this implies that¹⁰

$$(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \sim \sigma^2 \chi_p^2$$

we also know that

$$s^2 = \frac{\sigma^2}{n-p} \chi_{n-p}^2$$

This implies that the ratio of these two quantities has an F distribution:

$$\mathbb{P}_{\beta, \sigma^2} \left(\frac{\frac{1}{p} (\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta)}{s^2} \leq F_{p, n-p}(\alpha) \right) = 1 - \alpha$$

So a $(1-\alpha)$ -level confidence set for β is

¹⁰ If $\mathbf{X} \sim N(0, \Sigma)$, then $\Sigma^{-1/2} \mathbf{X} \sim N(0, I)$ and so $\mathbf{X}^T \Sigma^{-1} \mathbf{X} \sim \chi^2$. In this case, $\Sigma = \sigma^2 (X^T X)^{-1}$.

$$\boxed{\left\{ \beta \in \mathbb{R}^p : \frac{\frac{1}{p}(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta)}{s^2} \leq F_{p, n-p}(\alpha) \right\}}$$

[Note: imagine $p = 1$, so that we only have 1 β to test, and that we're testing whether $\beta = 0$, this becomes

$$\left\{ \beta \in \mathbb{R}^p : \frac{X^T X \hat{\beta}^2}{s^2} \leq F_{1, n-1}(\alpha) \right\} \Rightarrow \left\{ \beta \in \mathbb{R}^p : \frac{\hat{\beta}^2}{s^2 (X^T X)^{-1}} \leq F_{1, n-1}(\alpha) \right\}$$

it so happens that $F_{1, n-1}(\alpha) = t_{n-1}^2(\alpha)$, so this is precisely a t -test].

Diagnostic Tests

A number of diagnostic tests can be carried out to check whether a model we have fitted is appropriate:

- We can check whether a larger model is needed (ie: more columns of X) by adding extra terms (like x^2) and seeing whether they're needed using the t -test above.
- We can use the fact that $\hat{\varepsilon}$ and $\hat{\mathbf{Y}}$ are uncorrelated under the MLE fit¹¹, as are $\hat{\varepsilon}$ and $\hat{\beta}$. This means that a plot of errors against fitted values should show no systematic relationship.

If it did, then a transformation of the Y_i might make sense. To understand why, imagine a situation in which the variance and mean of Y are correlated, and a “well-behaved” transformation $f(\cdot)$

$$Y_i \sim N(\mu, \sigma^2(\mu))$$

$$f(Y) = f(\mu) + (Y - \mu)f'(\mu) \stackrel{\text{approx}}{\sim} N\left(f(\mu), \sigma^2(\mu)[f'(\mu)]^2\right)$$

We can then choose an f such that $\sigma^2(\mu)[f'(\mu)]^2$ is approximately constant. We try to choose f to **stabilise the variance**. A standard example is to try $\log Y_i = \alpha + \beta x_i + \varepsilon$.

¹¹ This follows from the fact that

$$\text{cov}(\hat{\varepsilon}, \hat{\mathbf{Y}}) = \text{cov}((I - H)\mathbf{Y}, H\mathbf{Y}) = (I - H)H \text{cov}(\mathbf{Y}, \mathbf{Y}) = 0$$

and the fact that the variables are normally distributed.

Another example is to use a **Box-Cox transformation**, in which each y is transformed to

$$y^{(\lambda)} = \begin{cases} (y^\lambda - 1) / \lambda & \lambda \neq 0 \\ \log y & \lambda = 0 \end{cases}$$

This family of transformations combines power and log transformations, is parameterised by λ and is continuous in λ . We wish to find the “best λ ”. If we have $\mathbb{E}(Y^{(\lambda)}) = X\beta$, then the likelihood of λ given the observed Y is¹²

$$\frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y}^{(\lambda)} - X\beta)^T(\mathbf{Y}^{(\lambda)} - X\beta)\right) \prod_{i=1}^n Y_i^{\lambda-1}$$

Fitting the linear model gives $\hat{\beta}(\lambda) = (X^T X)^{-1} X^T \mathbf{Y}^{(\lambda)}$ and $\hat{\sigma}^2 = \frac{1}{n}(\mathbf{Y}^{(\lambda)})^T (I - H) \mathbf{Y}^{(\lambda)} = RSS(\lambda) / n$, and the log-likelihood becomes

$$L(\lambda) = \text{const} - \frac{n}{2} \log \sigma^2 + (\lambda - 1) \sum_{i=1}^n \log(Y_i)$$

$$L(\lambda) = \text{const} - \frac{n}{2} \log(RSS(\lambda) / n) + (\lambda - 1) \sum_{i=1}^n \log(Y_i)$$

We can plot this function and pick an MLE $\hat{\lambda}$ that is closest to its maximum. The corresponding likelihood ratio test rejects the hypothesis $H_0 : \lambda = \lambda_0$ if

$$2(L(\hat{\lambda}) - L(\lambda_0)) > \chi_1^2(\alpha)$$

Where $\chi_1^2(\alpha)$ is the upper 100 α % point of a χ_1^2 -distribution.

- Our model assumes that $\mathbb{P}(\varepsilon_i \leq x) = \Phi(x / \sigma)$. We can approximate this by the **empirical distribution function** of the residuals

$$\hat{F}_n = \frac{\#\{i : \hat{\varepsilon}_i \leq x\}}{n}$$

If our assumptions are acceptable, then

$$\Phi^{-1}(\hat{F}_n(x)) \approx \frac{x}{\sigma}$$

A **Q-Q Plot** $\hat{\varepsilon}_i$ against $\Phi^{-1}(\hat{F}_n(\varepsilon_i))$. If everything is as it should be, the points should, roughly, lie on a straight line.

¹² Note that if $u = u(x)$, then $P_u du = P_x dx$, and so $P_u = P_x dx/du$, where the last term is called the **Jacobian**. This is where the last term comes from.

We note that since $\text{var}(\hat{\varepsilon}_i) = (1 - H_{ii})\sigma^2$, we could also use $\hat{\varepsilon}_i / \sigma\sqrt{1 - H_{ii}}$ or $\hat{\varepsilon}_i / s\sqrt{1 - H_{ii}}$, to ensure all variables has variance approximately equal to 1.

- We can check if any points appear to be **influential points** (pints whose deletion causes a large change in the analysis) by calculating the **Cook Distance**, D_i for each point i . Let $\hat{\mathbf{Y}}_{(i)}$ be the fitted values from a model omitting point i . Then

$$D_i = \frac{(\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)})^T (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)})}{s^2 p}$$

Where p is the number of parameters in β . Large values could indicate an influential observation.

Nested Models

Sometimes, we want to test whether *a number of* the β in our model are redundant. In other, we want to test the following hypotheses

$$\Omega : \mathbf{Y} = X\beta + \varepsilon \qquad \omega_1 : \mathbf{Y} = X_1\beta_1 + \varepsilon_1 \Rightarrow \beta_2 = \mathbf{0}$$

Where

$$\begin{array}{ccc}
 & \begin{array}{c} \nearrow \\ \nearrow \\ \nearrow \end{array} & \\
 X = \begin{pmatrix} X_1 & | & X_2 \end{pmatrix} & & \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \\
 \begin{array}{c} \nwarrow \\ \nwarrow \\ \nwarrow \end{array} & & \begin{array}{c} \nwarrow \\ \nwarrow \\ \nwarrow \end{array} \\
 n \times \tilde{p} & \begin{array}{c} n \times p_1 \quad n \times p_2 \end{array} & \begin{array}{c} \tilde{p} \times 1 \\ p_1 \times 1 \\ p_2 \times 1 \end{array}
 \end{array}$$

[Note even though β partitions into β_1 and β_2 , it is **not** always the case that $\hat{\beta}$ partitions into $\hat{\beta}_1$ and $\hat{\beta}_2$. Similarly, $\hat{\sigma}$ is not necessarily equal to $\hat{\sigma}_1$].

The alternative, model, in fact, corresponds to a projection of \mathbf{Y} onto a smaller subspace spanned by the columns of X_1 only, $\omega_1 = \{X_1\beta_1 : \beta_1 \in \mathbb{R}^{p_1}\}$. $= \{H_{\omega_1} \mathbf{Y} : \mathbf{Y} \in \mathbb{R}^n\}$. The projection matrices are¹³

¹³ It is useful to note that

$$\mathbb{R}^n = \underbrace{\omega_1 \oplus (\Omega \cap \omega_1^\perp)}_{\text{Partition } \mathbb{R}^n} \oplus \Omega^\perp$$

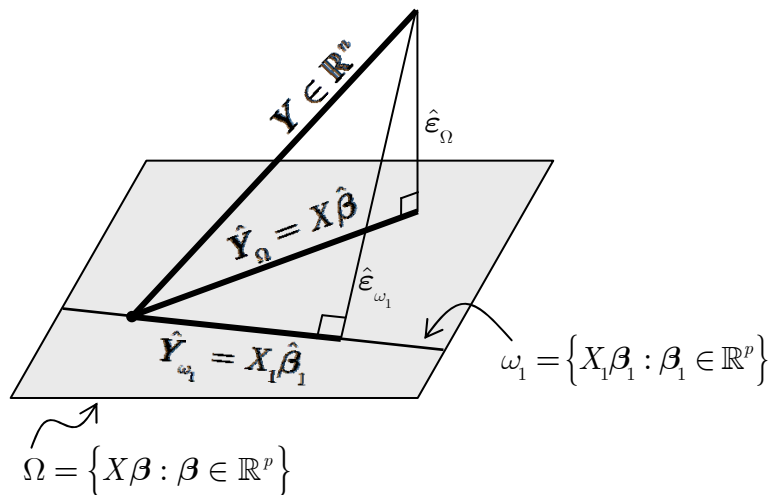
To show this is true, consider the fact that any vector $\mathbf{y} \in \mathbb{R}^n$ may be written as

$$\mathbf{y} = I\mathbf{y} = (H_{\omega_1} + (H_\Omega - H_{\omega_1}) + (I - H_\Omega))\mathbf{y} = H_{\omega_1}\mathbf{y} + (H_\Omega - H_{\omega_1})\mathbf{y} + (I - H_\Omega)\mathbf{y}$$

And that

$$H_{\omega_1} = X_1(X_1^T X_1)^{-1} X_1^T \qquad H_{\Omega} = X(X^T X)^{-1} X^T$$

$\omega_1 = \{X_1 \beta_1 : \beta_1 \in \mathbb{R}^{p_1}\} \subseteq \{X \beta : \beta \in \mathbb{R}^p\} = \Omega$ because any vector in the LHS set can be obtained by setting some entries of β in the RHS set to 0.



Clearly, $\hat{e}_{\omega_1} > \hat{e}_{\Omega}$. What will determine whether to accept or reject ω_1 is exactly how much bigger \hat{e}_{ω_1} is. If it's not much bigger, then we might as well drop all the extra superfluous parameters. We can formalise this in terms of a likelihood-ratio test. The log-likelihood is as usual in the linear model

- $H_{\omega_1} \mathbf{y} \in \omega_1$ (obviously!)
- $(H_{\Omega} - H_{\omega_1}) \mathbf{y} \in (\Omega \cap \omega_1^{\perp})$. To see why, we must show it is both in Ω and ω_1^{\perp}
 - Clearly $(H_{\Omega} - H_{\omega_1}) \mathbf{y} \in \Omega$ because $H_{\Omega} \mathbf{y} \in \Omega$ and $H_{\omega_1} \mathbf{y} \in \omega_1 \subseteq \Omega$.
 - Consider a $\mathbf{z} \in \omega_1$. Then $[(H_{\Omega} - H_{\omega_1}) \mathbf{y}]^T \cdot \mathbf{z} = \mathbf{y}^T (H_{\Omega} - H_{\omega_1}) \mathbf{z} = \mathbf{y}^T (H_{\Omega} \mathbf{z} - H_{\omega_1} \mathbf{z}) = 0$, since $\mathbf{z} \in \omega_1 \subseteq \Omega \Rightarrow H_{\Omega} \mathbf{z} = H_{\omega_1} \mathbf{z} = \mathbf{z}$. Thus, $(H_{\Omega} - H_{\omega_1}) \mathbf{y} \in \omega_1^{\perp}$.
- $(I - H_{\Omega}) \mathbf{y} \in \Omega^{\perp}$. Consider a $\mathbf{z} \in \Omega$. Then $[(I - H_{\Omega}) \mathbf{y}]^T \cdot \mathbf{z} = \mathbf{y}^T (\mathbf{z} - H_{\Omega} \mathbf{z}) = 0$, since $\mathbf{z} \in \Omega \Rightarrow H_{\Omega} \mathbf{z} = \mathbf{z}$.

Furthermore, we note that

- $\omega_1 \perp (\Omega \cap \omega_1^{\perp})$, because $(\Omega \cap \omega_1^{\perp}) \subseteq \omega_1^{\perp}$
- $(\Omega \cap \omega_1^{\perp}) \perp \Omega^{\perp}$, because $(\Omega \cap \omega_1^{\perp}) \subseteq \Omega$
- $\omega_1 \perp \Omega^{\perp}$, because $\omega_1 \subseteq \Omega$

$$\begin{aligned}
\ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2) &= -\frac{n}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} \|\mathbf{Y} - X\hat{\boldsymbol{\beta}}\|^2 \\
&= -\frac{n}{2} \log \left(\frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 \right) - \frac{1}{2\hat{\sigma}^2} n\hat{\sigma}^2 \\
&= -\frac{n}{2} \log \left(\frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 \right) - \frac{n}{2}
\end{aligned}$$

And the log-likelihood ratio statistic for ω_1 is

$$\begin{aligned}
W(\omega_1) &= 2 \left\{ \ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}_\Omega^2) - \ell(\hat{\boldsymbol{\beta}}, \hat{\sigma}_{\omega_1}^2) \right\} \\
&= 2 \left\{ -\frac{n}{2} \log \left(\frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}_\Omega\|^2 \right) - \frac{n}{2} + \frac{n}{2} \log \left(\frac{1}{n} \|\mathbf{Y} - \hat{\mathbf{Y}}_{\omega_1}\|^2 \right) + \frac{n}{2} \right\} \\
&= n \log \left(\varepsilon_{\omega_1}^2 / \varepsilon_\Omega^2 \right)
\end{aligned}$$

By looking at the diagram above and using Pythagoras' Theorem¹⁴ we get

$$\begin{aligned}
\varepsilon_{\omega_1}^2 &= \varepsilon_\Omega^2 + \|\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}\|^2 \\
\Rightarrow \|\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}\|^2 &= \varepsilon_{\omega_1}^2 - \varepsilon_\Omega^2
\end{aligned}$$

This allows us to write

$$\frac{\varepsilon_{\omega_1}^2}{\varepsilon_\Omega^2} = 1 - \frac{\|\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}\|^2}{\varepsilon_\Omega^2}$$

Now, by the extension of Cochran's Theorem¹⁵:

¹⁴ And the fact that ε_Ω^T and $\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}$ are perpendicular, because

$$\begin{aligned}
\varepsilon_\Omega^T (\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}) &= (\mathbf{Y} - \hat{\mathbf{Y}}_\Omega)^T (\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}) = \mathbf{Y}^T (I - H_\Omega) (H_\Omega - H_{\omega_1}) \mathbf{Y} \\
&= \mathbf{Y}^T (H_\Omega - H_{\omega_1} - H_\Omega + H_\Omega H_{\omega_1}) \mathbf{Y} = \mathbf{Y}^T (H_\Omega H_{\omega_1} \mathbf{Y} - H_{\omega_1} \mathbf{Y}) = \mathbf{0}
\end{aligned}$$

In the last step, we use the fact that $H_{\omega_1} \mathbf{Y} \in \omega_1 \subseteq \Omega \Rightarrow H_\Omega H_{\omega_1} \mathbf{Y} = H_{\omega_1} \mathbf{Y}$.

¹⁵ First note that $\mathbf{Y} \sim N(X\boldsymbol{\beta}, \sigma^2 I)$. Consider that

$$\begin{aligned}
\|\mathbf{Y}\|^2 &= \|H_{\omega_1} \mathbf{Y}\|^2 + \|(H_\Omega - H_{\omega_1}) \mathbf{Y}\|^2 + \|(I - H_\Omega) \mathbf{Y}\|^2 \\
&= \|H_{\omega_1} \mathbf{Y}\|^2 + \|\hat{\mathbf{Y}}_\Omega - \hat{\mathbf{Y}}_{\omega_1}\|^2 + \varepsilon_\Omega^2 \\
&= \mathbf{Y}^T H_{\omega_1} \mathbf{Y} + \mathbf{Y}^T (H_\Omega - H_{\omega_1}) \mathbf{Y} + \mathbf{Y}^T (I - H_\Omega) \mathbf{Y}
\end{aligned}$$

Note that:

- $I = H_{\omega_1} + (H_\Omega - H_{\omega_1}) + (I - H_\Omega)$
- $\dim(\omega_1) = p_1$ and $\dim(\Omega) = \tilde{p}$, so $\dim(\Omega^\perp) = n - \tilde{p}$ and $\dim(\Omega \cap \omega_1^\perp) = \tilde{p} - p_1$.

Therefore, $\text{rank}(H_{\omega_1}) = p_1$, $\text{rank}(H_\Omega - H_{\omega_1}) = \tilde{p} - p_1$ and $\text{rank}(I - H_\Omega) = n - \tilde{p}$.

We can therefore apply the extension to Cochran's Theorem to find the second two terms in the expansion above. The relevant non-centrality parameters δ are (note that $X\boldsymbol{\beta} \in \Omega$):

$$\sigma^2 \delta_{\varepsilon_\Omega^2} = \boldsymbol{\beta}^T X^T (I - H_\Omega) X \boldsymbol{\beta} = \boldsymbol{\beta}^T X^T (X\boldsymbol{\beta} - X\boldsymbol{\beta}) = 0$$

- $\epsilon_{\Omega}^2 \sim \sigma^2 \chi_{n-\bar{p}}^2$
- $\left\| \hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega_1} \right\|^2 \sim \sigma^2 \chi_{\bar{p}-p_1}^2 \left(\frac{1}{\sigma^2} \left\| (I - H_{\omega_1}) X \beta \right\|^2 \right)$.

Therefore

$$F = \frac{\frac{1}{\bar{p}-p_1} \left\| \hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega_1} \right\|^2}{\frac{1}{n-\bar{p}} \epsilon_{\Omega}^2} = \frac{\frac{1}{\bar{p}-p_1} (\epsilon_{\omega_1}^2 - \epsilon_{\Omega}^2)}{\frac{1}{n-\bar{p}} \epsilon_{\Omega}^2} \sim F_{\bar{p}-p_1, n-\bar{p}} \left(\frac{1}{\sigma^2} \left\| (I - H_{\omega_1}) X \beta \right\|^2 \right)$$

Under ω_1 , $X\beta$ lies in ω_1 already and so $(I - H_{\omega_1})X\beta = (X\beta - X\beta) = 0$ and so we expect a standard **F-distribution**.

Intuitively, the F statistic is as follows

$$F = \frac{\frac{1}{\bar{p}-p_1} (\epsilon_{\omega_1}^2 - \epsilon_{\Omega}^2)}{\frac{1}{n-\bar{p}} \epsilon_{\Omega}^2}$$

Error per extra degree of freedom added by going from $\omega_1 \rightarrow \Omega$

Error per degree of freedom present in Ω anyway

Under Ω , this will be large, because all the components of β are important; indeed, we get a non-central F distribution, which is stochastically increasing in the non-centrality parameter. Thus, we **reject ω_1 for large values of F** .

The data used to carry out F tests is often summarised in a table. Consider the following nested models:

$$\begin{aligned} \omega &: \mathbf{Y} = \mu \mathbf{1} && [\text{ie: } \beta_1 = \beta_2 = 0] && 1 \\ \omega_1 &: \mathbf{Y} = \mu \mathbf{1} + X_1 \beta_1 + \epsilon && [\text{ie: } \beta_2 = 0] && p_1 + 1 \\ \Omega &: \mathbf{Y} = \mu \mathbf{1} + X_1 \beta_1 + X_2 \beta_2 + \epsilon && && p_1 + p_2 + 1 \end{aligned}$$

Our table is then

Source	Degrees of freedom	Sum of squares	Mean sum of squares	F-value
Due to β_1 [ie: $\omega \rightarrow \omega_1$]	p_1	$SS_{\beta_1 \mu} = \epsilon_{\omega}^2 - \epsilon_{\omega_1}^2$	$S_{\beta_1 \mu} / p_1$	$\frac{SS_{\beta_1 \mu} / p_1}{\epsilon_{\Omega}^2 / (n - p_1 - p_2 - 1)}$
Due to β_2 [ie: $\omega_1 \rightarrow \Omega$]	p_2	$SS_{\beta_2 \mu, \beta_1} = \epsilon_{\omega_1}^2 - \epsilon_{\Omega}^2$	$S_{\beta_2 \mu, \beta_1} / p_2$	$\frac{SS_{\beta_2 \mu, \beta_1} / p_2}{\epsilon_{\Omega}^2 / (n - p_1 - p_2 - 1)}$
Residual [in Ω anyway]	$n - p_1 - p_2 - 1$	ϵ_{Ω}^2	$\frac{\epsilon_{\Omega}^2}{n - p_1 - p_2 - 1}$	
Total	$n - 1$	$\epsilon_{\omega}^2 = \sum (Y_i - \bar{Y})^2$		

$$\sigma^2 \delta_{\left\| \hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega_1} \right\|^2} = \beta^T X^T (H_{\Omega} - H_{\omega_1}) X \beta = \beta^T X^T (I - H_{\omega_1}) X \beta = \left\| (I - H_{\omega_1}) X \beta \right\|^2$$

Note that the expected value of each of the F values in the table above is

$$\frac{n - p_1 - p_2 - 1}{(n - p_1 - p_2 - 1) - 2}$$

What we're effectively doing with this table is

- Taking the **total variation in \mathbf{Y}** , given by $\varepsilon_\omega^2 = \sum (Y_i - \bar{Y})^2$
- Seeing how much of it remains unexplained even in our full model ε_Ω^2
- Seeing how much we “lose out” each time we reduce the size of our model.

[Note, of course, that unless only factor is present, each of these F values will be different to the F value testing the drop in variance from $\omega \rightarrow \Omega$].

As we mentioned above, even though β partitions into β_1 and β_2 , it is **not** always the case that $\hat{\beta}$ partitions into $\hat{\beta}_1$ and $\hat{\beta}_2$. This adds extra complications when calculating sum-of-squares, in that $SS_{\beta_2|\beta_1} = SS_{\beta_2}$.

To investigate, consider the following models

$$\begin{aligned} \omega &: \mathbf{Y} = \varepsilon \\ \omega_1 &: \mathbf{Y} = X_1\beta_1 + \varepsilon \\ \omega_2 &: \mathbf{Y} = X_2\beta_2 + \varepsilon \\ \Omega &: \mathbf{Y} = X\beta + \varepsilon \end{aligned} \quad X = \begin{pmatrix} X_1 & X_2 \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

We say that the models are **orthogonal** if $X_1^T X_2 = 0$.

Let's first consider the MLEs $\hat{\beta}_i$. If the models are orthogonal, then

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T \mathbf{Y} = \left(\begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} \begin{pmatrix} X_1 & X_2 \end{pmatrix} \right)^{-1} \begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} \mathbf{Y} = \begin{pmatrix} X_1^T X_1 & 0 \\ 0 & X_2^T X_2 \end{pmatrix}^{-1} \begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} \mathbf{Y} \\ &= \begin{pmatrix} (X_1^T X_1)^{-1} & 0 \\ 0 & (X_2^T X_2)^{-1} \end{pmatrix} \begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} \mathbf{Y} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \end{aligned}$$

Thus, the coefficients partition neatly if and only if the models are orthogonal.

Next, consider that

$$\begin{aligned}
SS_{\beta_2|\beta_1} &= \epsilon_{\omega_1}^2 - \epsilon_{\Omega}^2 \\
&= \mathbf{Y}^T (I - H_1) \mathbf{Y} - \mathbf{Y}^T (I - H) \mathbf{Y} \\
&= \mathbf{Y}^T (H - H_1) \mathbf{Y} \\
SS_{\beta_2} &= \epsilon_{\omega}^2 - \epsilon_{\omega_2}^2 \\
&= \mathbf{Y}^T \mathbf{Y} - \mathbf{Y}^T (I - H_2) \mathbf{Y} \\
&= \mathbf{Y}^T H_2 \mathbf{Y}
\end{aligned}$$

If the models are orthogonal, then by the calculation above

$$H = X(X^T X)^{-1} X^T = \begin{pmatrix} X_1 & X_2 \end{pmatrix} \begin{pmatrix} (X_1^T X_1)^{-1} & 0 \\ 0 & (X_2^T X_2)^{-1} \end{pmatrix} \begin{pmatrix} X_1^T \\ X_2^T \end{pmatrix} = H_1 + H_2$$

So the above becomes

$$SS_{\beta_2|\beta_1} = \mathbf{Y}^T H_2 \mathbf{Y} = SS_{\beta_2}$$

Thus, the order in which parameters are added to a model doesn't matter if and only if the models are orthogonal.

Coefficient of Determination

Consider the two models

$$\begin{aligned}
\omega &: \mathbf{Y} = \mu \mathbf{1} \\
\Omega &: \mathbf{Y} = \mu \mathbf{1} + X\boldsymbol{\beta} + \epsilon
\end{aligned}$$

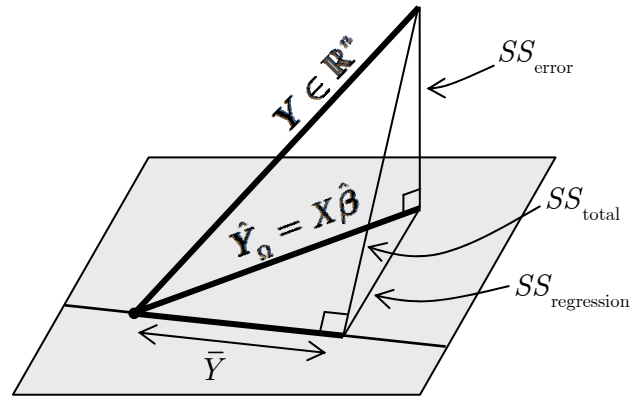
Clearly,

- $\epsilon_{\omega}^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}_{\omega}\|^2 = \|\mathbf{Y} - \bar{Y}\mathbf{1}\|^2 = SS_{\text{total}}$ is the **total variability in the sample**.
- $\epsilon_{\Omega}^2 = \|\hat{\mathbf{Y}}_{\Omega} - \mathbf{Y}\|^2 = SS_{\text{error}}$ is the **variability not explained by Ω** .
- $\|\hat{\mathbf{Y}}_{\Omega} - \bar{Y}\mathbf{1}\|^2 = \|\hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega}\|^2 = SS_{\text{regression}}$ is the variability that **is explained by Ω** .

We know, however, from the previous section that $\|\hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega}\|^2 = \epsilon_{\omega}^2 - \epsilon_{\Omega}^2$, and this implies that¹⁶ the total variability can be split as follows:

$$SS_{\text{total}} = SS_{\text{(Explained by) regression}} + SS_{\text{(Explained by) error}}$$

¹⁶ Reminder: because $SS_{\text{total}} = \|\mathbf{Y} - \hat{\mathbf{Y}}_{\omega}\|^2 = \|(\hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega}) + (\mathbf{Y} - \hat{\mathbf{Y}}_{\Omega})\|^2 = \|(\hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega}) + \epsilon_{\Omega}\|^2$, the first term squared is $SS_{\text{regression}}$, and see previous footnote for proof $\epsilon_{\Omega}^T (\hat{\mathbf{Y}}_{\Omega} - \hat{\mathbf{Y}}_{\omega}) = 0$.



The **coefficient of determination**, R^2 is given by the **ratio** between the **error explained by regression** and the **total error**:

$$R^2 = \frac{SS_{\text{regression}}}{SS_{\text{total}}}$$

It gives us an idea of how well our model explains the variability in the \mathbf{Y} .

ANOVA

ANOVA is a model for the analysis of categorical data, in which our observation depends on one or more discrete factors.

We first begin with a simple example in which one factor is involved. If Y_{ij} is the j^{th} observation at level i of the factor, our model is

$$Y_{ij} = \mu + \alpha_i + \varepsilon_{ij} \quad \varepsilon_{ij} = \text{NID}(0, \sigma^2)$$

This model can be expressed in least squares form. Consider the example in which the factor can take 3 values, and each group contains 2 observations:

$$\beta = (\mu \quad \alpha_1 \quad \alpha_2 \quad \alpha_3)^T$$

$$X = \begin{pmatrix} 1 & 1 & & \\ 1 & 1 & & \\ 1 & & 1 & \\ 1 & & 1 & \\ 1 & & & 1 \\ 1 & & & 1 \end{pmatrix}$$

The easiest way to find the MLE is to use the equation $X^T X \beta = X^T \mathbf{Y}$ directly. So in the case above, we get

$$\begin{pmatrix} 10 & 2 & 2 & 2 \\ 2 & 2 & & \\ 2 & & 2 & \\ 2 & & & 2 \end{pmatrix} \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} = \begin{pmatrix} Y_{++} \\ Y_{1+} \\ Y_{2+} \\ Y_{3+} \end{pmatrix}$$

This model is over-defined, because we can add a fixed amount to μ and remove that amount from every α . We therefore need to impose an additional constraint on the parameterisation. Two common choices are:

- **Sum-to-zero constraint:** $\boxed{\sum \alpha_i = 0}$. The MLE is then

$$\hat{\mu} = \bar{Y}_{..} \quad \hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{..}$$

Using these constraints, every score is compared to the *overall* mean.

- **Corner-point constraint:** $\boxed{\alpha_1 = 0}$. The MLE is then

$$\hat{\mu} = \bar{Y}_{1.} \quad \hat{\alpha}_i = \bar{Y}_{i.} - \bar{Y}_{1.}$$

Using these constraints, every score is compared to the mean of group 1. (Where the \bar{Y} imply an average of Y over \bullet indices)

In **two-way ANOVA**, two factors are involved. If Y_{ijk} is the k^{th} observation at level i of factor A and level j of factor B, we can fit two models:

- Model with **no interaction between factors**

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk} \quad (*)$$

- Model **with interaction between factors**

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk} \quad (\#)$$

The corner point constraints for this model are $\alpha_1 = \beta_1 = \gamma_{1j} = \gamma_{i1} = 0$.

In the second model, the effect of changing factor A depends on the precise value of factor B.

This can also be extended to models mixing “factor” and “non-factor” in the same model. An example might be the weight W of a child depending on their age A (non-factor) and sex S (factor). Two possible models can then be fit:

- Model with **no interaction**

$$W_{ij} = \mu + \alpha_i + \beta A_{ij} + \varepsilon_{ij} \quad \left[i \in \{M, F\}, \alpha_M = 0 \right]$$

- Model **with interaction**

$$W_{ij} = \mu + \alpha_i + (\beta + \gamma_i) A_{ij} + \varepsilon_{ij} \quad \left[i \in \{M, F\}, \alpha_M = \gamma_M = 0 \right]$$

{Proof of Cochran's Theorem}

I have been asked to remove all notes from the Part III Statistical Theory class from this website. Please email me if you have any questions.

Computer Code

Carrying out regression

- Regressing one variable on another: `linearModelName <- lm(yVariable~xVariable)`
- To carry out the regression on two covariates, use `linearModelName <- lm(yVariable~firstCovariate+secondCovariate)`
- To fit orthogonal polynomials such as $P_0 = a_0$, $P_1 = a_1 + b_1x$, $P_2 = a_2 + b_2x + c_2x^2$, etc..., use `linearModelName <- lm(yVariable~Poly(Covariate,2))`

Interpreting a regression

- Selected components of the variable `linearModelName`
 - `coefficients`
 - `residuals`
 - `fitted.values`

Each of those can be extracted using the `linearModelName$coefficients`, for example.
- Seeing a summary of the linear model: `summary(linearModelName)`.
 - The “residual standard error” given in the output is s .
 - Remember that only $|t|$ counts; the sign matters not.

Add `cor=F` as an argument to ommitt correlation of coefficients.
- `deviance(linearModelName)` gives the residual sum of squares. It can be used to carry out an F test directly.
- Carrying out successive F -tests on nested models in which each term in the regression is added sequentially: `anova(linearModelName)`
 - The last F statistic tested in the output tests the null model $Y = \mu\mathbf{1} + \varepsilon_i$ against the full model.

Model selection (*see later for theory*)

- *Backwards elimination*
 - Fit successive models, find t values using `summary` removing the covariate with lowest $|t|$ each time.
 - `library(MASS)` followed by `dropterm(linearModelName, test="F")` gives the effect of dropping each variable in turn, including F values. These, however, are just the square of the t values given by `summary`.

- *Forward selection*
 - `library(MASS)` followed by `addterm(linearModelName, covariate1+covariate2)` gives the effect of adding each variable in turn to `linearModelName`. Choose the one with the largest reduction in RSS each time.
- *Best subset regression*
 - `library(MASS)` followed by `stepAIC(linearModel)` successively removes the covariate that decreases the AIC the most, until there is no way to reduce the AIC further.
 - `library(MASS)` followed by


```
stepAIC(NulllinearModelName,
        list(upper=~1+Covt1+Covt2+Covt3
            lower=~1),
        direction="forward")
```

 starts with an empty model and successively adds the covariate that increases the AIC least.

Omitting the “direction” causes the program to go **stepwise** – at each step, it checks whether adding back a covariate that was previously removed could decrease the AIC

Diagnostic plots

- `plot(linearModelName, ask=T)` lists and shows all possible plots.
- `res <- linearModelName$residuals` followed by `qqnorm(res)` and `qqline(res)` gives a Q-Q plot.
- `library(MASS)` followed by `boxcox(linearModel)` gives a Box-Cox plot.

Categorical data – single factor

- Best illustrated by example

Control	Treat A	Treat B	
1	2	3	Should be
4	5	6	input as two
7	8	9	vectors

Result = 1 2 3 4 5 6 7 8 9

treat = 1 2 3 1 2 3 1 2 3

Followed by the command `Treat <- factor(treat)`

- A box-and-whisker plot of the data can be obtained using `plot(factorName, resultName)`

- To specify the ANOVA constraint, use
 - Corner-point constraints: `options(contrasts = c("contr.treatment", "contr.poly"))`
 - Sum-to-zero constraints: `options(contrasts = c("contr.sum", "contr.poly"))`
- The linear model is then fitted using `linearModelName <- lm(Result~Factor)`.
- Alternatively, `linearModelName <- aov(Result~Factor)` will also fit the model, but then the summary needs to be produced using `summary.lm(linearModelName)`.
- Note:
 - `tapply(Result, Factor, mean)` will return a table with the mean of each group. Alternatively, if an aov-style model was fitted, `model.tables(LinearModelName, type="means", se=T)` will automatically calculate means for every factor in the model (and the standard errors that appear in summary).
 - `tapply(Result, Factor, function(x) sqrt(var(x)/length(x)))` will return a table with the standard error of each group.

Categorical data – multiple factors

- Consider running


```
var1 <- c(1,2); var2 <- 1:5; var3 <- 1:4
grid <- expand.grid(var1,var2,var3)
```

 this creates a table with every possible permutation of the variables, in that order. We can then apply


```
Var1 <- factor(grid[,1])
Var2 <- factor(grid[,2])
```
- To fit a box-and-whisker diagram involving two factors


```
ourData <- data.frame(Var1, Var2, Observations)
plot.design(ourData)
```
- To create an interaction plot


```
interaction.plot(Var1,Var2,Observations)
```
- When fitting more than one factor
 - No interaction: `linearModelName <- lm(Result~Factor1+Factor2)`
 - Interaction: `linearModelName <- lm(Result~Factor1*Factor2)`

{Likelihood Theory}

I have been asked to remove all notes from the Part III Statistical Theory class from this website. Please email me if you have any questions.

Condition (7) allows us to write

$$\begin{aligned} 0 &= \bar{U}(\hat{\theta}_n) \mathbb{1}_{A_n} \\ &= \left(\bar{U}(\theta_0) - \bar{j}(\tilde{\theta}_n)^T (\hat{\theta}_n - \theta_0) \right) \mathbb{1}_{A_n} \end{aligned}$$

Where $\tilde{\theta}_n$ lies on the line segment between θ_0 and $\hat{\theta}_n$:

$$\tilde{\theta}_n \in \left\{ t\theta_0 + (1-t)\hat{\theta}_n : t \in [0,1] \right\}$$

Now, by the triangle rule:

$$\begin{aligned} \left\| \bar{j}(\tilde{\theta}_n) \mathbb{1}_{A_n} - i^{(1)}(\theta_0) \right\| &\leq \left\| \bar{j}(\tilde{\theta}_n) \mathbb{1}_{A_n} - i^{(1)}(\tilde{\theta}_n) \mathbb{1}_{A_n} \right\| \\ &\quad + \left\| i^{(1)}(\tilde{\theta}_n) \mathbb{1}_{A_n} - i^{(1)}(\theta_0) \mathbb{1}_{A_n} \right\| \\ &\quad + \left\| i^{(1)}(\theta_0) \mathbb{1}_{A_n} - i^{(1)}(\theta_0) \right\| \\ &\leq \sup_{\theta \in \Theta} \left\| \bar{j}(\tilde{\theta}_n) - i^{(1)}(\tilde{\theta}_n) \right\| \\ &\quad + \left\| i^{(1)}(\tilde{\theta}_n) \mathbb{1}_{A_n} - i^{(1)}(\theta_0) \mathbb{1}_{A_n} \right\| \\ &\quad + \left\| i^{(1)}(\theta_0) \mathbb{1}_{A_n^c} \right\| \end{aligned}$$

Generalised Linear Models

- Consider a family of distributions with respect to a σ -finite measure of the form²²

$$f(y; \mu, \sigma^2) = a(\sigma^2, y) \exp \left[\frac{1}{\sigma^2} \left\{ \theta(\mu)y - k(\theta(\mu)) \right\} \right] \quad \begin{array}{l} y \in E_1 \subseteq \mathbb{R} \\ \mu \in M \subseteq \mathbb{R} \\ \sigma^2 \in \Phi \subseteq (0, \infty) \end{array}$$

Where $a(\sigma^2, y)$ is a known positive function. Such a family is called an **exponential dispersion family** and σ^2 is called the **dispersion parameter**.

- If y is of exponential dispersion family form, then its **moment generating function** is

$$\begin{aligned} M(t) &= \int_{E_1} a(\sigma^2, y) \exp \left[t \cdot y + \frac{1}{\sigma^2} \left\{ \theta(\mu)y - k(\theta(\mu)) \right\} \right] d\mu(y) \\ &= \exp \left[\frac{1}{\sigma^2} \left\{ k(\sigma^2 t + \theta(\mu)) - k(\theta(\mu)) \right\} \right] \\ &\quad \int_{E_1} a(\sigma^2, y) \exp \left[\frac{1}{\sigma^2} \left\{ (\sigma^2 t + \theta(\mu)) \cdot y - k(\sigma^2 t + \theta(\mu)) \right\} \right] d\mu(y) \end{aligned}$$

²² Note: in applied statistics, the distribution is given in the form

$$f(y; \mu, \sigma^2) = \exp \left[\frac{y\mu - k(\mu)}{\sigma^2} + c(y, \sigma^2) \right]$$

In other words, it assumes $\theta(\cdot) = \cdot$ and $a(y, \sigma^2) = \exp(c(y, \sigma^2))$.

The integral is 1, because it is the sum over all space of another exponential distribution family. Thus

$$M(t) = \exp\left[\frac{1}{\sigma^2}\left\{k(\sigma^2 t + \theta(\mu)) - k(\theta(\mu))\right\}\right]$$

And the **cumulant generating function** (ie: the logarithm of the mgf) is

$$K_Y(t; \mu, \sigma^2) = \frac{1}{\sigma^2}\left\{k(\sigma^2 t + \theta(\mu)) - k(\theta(\mu))\right\}$$

- The discussion above allows to us to find

$$\begin{aligned}\mathbb{E}_{\mu, \sigma^2}(y) &= \kappa_1 = k'(\theta(\mu)) \\ \text{var}_{\mu, \sigma^2}(y) &= \kappa_2 = \sigma^2 k''(\theta(\mu)) = \sigma^2 V(\mu)\end{aligned}$$

The function $V(\mu)$ is called the **variance function**. (Note: to find it, differentiate k with respect to θ - not with respect to μ).

- A different derivation of these results uses the loglikelihood

$$\ell(y; \mu, \sigma^2) = \frac{\theta(\mu)y - k(\theta(\mu))}{\sigma^2} + \log\{a(\sigma^2, y)\}$$

Note that

$$\frac{\partial \ell}{\partial \theta} = \frac{y - k'(\theta(\mu))}{\sigma^2} \quad \frac{\partial^2 \ell}{\partial \theta^2} = -\frac{k''(\theta(\mu))}{\sigma^2}$$

It now suffices to use two results derived in the previous section:

$$\begin{aligned}\mathbb{E}\left(\frac{\partial \ell}{\partial \theta}\right) &= 0 \Rightarrow \mathbb{E}\left(\frac{y - k'(\theta(\mu))}{\sigma^2}\right) = 0 \\ &\quad \boxed{\mathbb{E}(y) = k'(\theta(\mu))} \\ \mathbb{E}_\theta\left\{\left(\frac{\partial \ell}{\partial \theta}\right)^2\right\} &= -\mathbb{E}_\theta\left\{\frac{\partial^2 \ell}{\partial \theta^2}\right\} \Rightarrow \frac{k''(\theta(\mu))}{\sigma^2} = \mathbb{E}\left(\frac{(y - \mathbb{E}(y))^2}{\sigma^4}\right) = \frac{\text{var}(y)}{\sigma^4} \\ &\quad \boxed{\text{var}(y) = \sigma^2 k''(\theta(\mu))}\end{aligned}$$

As above.

- An exponential dispersion family is completely characterised by $(V(\mu), M, \Phi)$. We may therefore write $Y \sim \text{ED}(\mu, \sigma^2 V(\mu))$ with $\mu \in M$ and $\sigma^2 \in \Phi$ to mean that Y is of exponential family form with mean μ and variance $\sigma^2 V(\mu)$.
- The linear we have thus far been studying can be specified as follows:
 - **Distribution:** $Y_i \sim N(\mu_i, \sigma^2)$
 - **Link function:** $\mathbb{E}(Y_i) = \mu_i = \beta^T \mathbf{X}_i$.

this, however, is only a subset of **general linear models** (GLMs). A GLM is a model for **independent** responses Y_1, \dots, Y_n in which

- **Distribution:** $Y_i \sim \text{ED}(\mu_i, \sigma_i^2 V(\mu_i))$, $\mu_i \in M$, $\sigma_i^2 = \sigma^2 a_i$ where σ^2 is an **unknown dispersion parameter** and a_1, \dots, a_n are **known constants**.

For example, if $Y_i \sim \frac{1}{n_i} \text{bin}(n_i, \mu_i)$, we have $\sigma_i^2 = \frac{1}{n_i}$, so we can take $a_i = \frac{1}{n_i}$ and $\sigma^2 = 1$.

- **Link function:** The i^{th} component of the **linear predictor** $\eta_i = \mathbf{X}_{i\bullet}^T \boldsymbol{\beta}$ are related through $g(\mu_i) = \eta_i$ where
 - $\mathbf{X}_{i\bullet}$ is a vector of known explanatory variables
 - $\boldsymbol{\beta}$ is an unknown vector of regression coefficients
 - g is a strictly increasing twice differentiable function called the **link function**.

(Note that for the standard linear model, $g(\mu) = \mu$).

- The choice $g(\mu) = \theta(\mu)$ is called the **canonical link function** and it simplifies calculations in certain cases.

For example, consider the density of $\mathbf{y} = (y_1, \dots, y_n)^T$

$$f_Y(\mathbf{y}; \mu, \sigma^2) = \left\{ \prod_{i=1}^n a(\sigma^2, y_i) \right\} \exp \left\{ \sum_{i=1}^n \left(\frac{\theta(\mu_i) y_i}{\sigma_i^2} - \frac{K(\theta(\mu_i))}{\sigma_i^2} \right) \right\}$$

However, if we use the canonical link function, $\mu_i = g^{-1}(x_{i\bullet}^T \boldsymbol{\beta}) = \theta^{-1}(x_{i\bullet}^T \boldsymbol{\beta})$

$$f_Y(\mathbf{y}; \mu, \sigma^2) = \left\{ \prod_{i=1}^n a(\sigma^2, y_i) \right\} \exp \left\{ \boldsymbol{\beta}^T \sum_{i=1}^n \left(\frac{x_{i\bullet} y_i}{\sigma_i^2} - \frac{K(x_{i\bullet}^T \boldsymbol{\beta})}{\sigma_i^2} \right) \right\}$$

In that case, we see that the vector

$$\sum_{i=1}^n \left(\frac{x_{i\bullet} y_i}{\sigma_i^2} \right) = \left(\sum_{i=1}^n \frac{x_{i1} y_i}{a_i}, \dots, \sum_{i=1}^n \frac{x_{ip} y_i}{a_i} \right)^T$$

is sufficient for $\boldsymbol{\beta}$, for each fixed value of σ^2 .

- In general, there is no closed-form expression for the MLE $\hat{\boldsymbol{\beta}}$, but we can use a Newton-Raphson type of algorithm (**Fisher Scoring**) to find a sequences converging to $\hat{\boldsymbol{\beta}}$. Moreover, under mild conditions on $\mathbf{X}_{i\bullet}$, we can apply the result of section 2.2 to deduce that

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N_p(0, \mathcal{I}^{(1)}(\boldsymbol{\beta})^{-1})$$

This result can be used to estimate the standard deviation of components of $\hat{\beta}$, or to test hypotheses about β .

- Alternatively, tests can be based on the **deviance**, which is closely related to the likelihood ratio statistic.

Let ω_s be the **saturated model** with as many parameters as observations. This achieves a perfect fit with $\hat{Y}_i^{(s)} = Y_i$. Let the corresponding maximised loglikelihood be $\ell_{\max}^{(s)}$. Let $\omega_f \subset \omega_s$ be a fitted model, with maximised loglikelihood $\ell_{\max}^{(f)} < \ell_{\max}^{(s)}$.

The **deviance** of the fitted model is then given by

$$D(\omega_f, \omega_s) = 2\phi \left(\ell_{\max}^{(s)} - \ell_{\max}^{(f)} \right)$$

The **scaled deviance** is given by

$$S(\omega_f, \omega_s) = 2 \left(\ell_{\max}^{(s)} - \ell_{\max}^{(f)} \right)$$

Good models will have a **small** scaled deviance.

To compare two models $\omega_1 \subset \omega_2 \subset \omega_s$, we consider the drop in deviance:

$$S(\omega_1, \omega_s) - S(\omega_2, \omega_s) = \frac{D(\omega_1, \omega_s) - D(\omega_2, \omega_s)}{\phi} \approx \chi_{p_2 - p_1}^2 \quad \text{if } \omega_1 \text{ true}$$

where $p_i = \dim(\omega_i)$.

This can be used in two contexts

- Checking whether it is worth adding an extra term in the model, by comparing the drop in deviance to $\chi_{p_2 - p_1}^2$.
- Testing whether a model is a “good fit” – in other words, whether it approximates the saturated model correctly. If it does, we would expect S for the model to be distributed as χ_{n-p}^2 . We are effectively using the test above with $\omega_1 =$ our model and $\omega_2 = \omega_s$, which implies that $S(\omega_2, \omega_s) = 0$.

When ϕ is known (for example, Binomial and Poisson), this can be used to test whether the drop in deviance is significant. Otherwise, we can use F -tests.

GLMs ~ Normal data

- If $Y \sim N(\mu, \sigma^2)$. The density is

$$\begin{aligned}
 f(y) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^2\right) \\
 &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{y^2}{2\sigma^2} + \frac{y\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\
 &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{y^2}{2\sigma^2}\right) \exp\left(\frac{1}{\sigma^2}(y\mu - \frac{1}{2}\mu^2)\right)
 \end{aligned}$$

This is of exponential distribution family form with

- $a(\sigma^2, y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{y^2}{2\sigma^2}\right)$
- $\theta(\mu) = \mu$
- $k(\mu) = \frac{1}{2}\mu^2$

We note that $V(\mu) = 1$.

- Now, consider variables $Y_1, \dots, Y_n \sim N(\mu_i, \sigma^2)$. The loglikelihood is given by

$$\ell = -\frac{1}{2\sigma^2} \|\mathbf{Y} - \boldsymbol{\mu}\|^2 - \frac{n}{2} \log(2\pi\sigma^2)$$

We then have

- ω_s has fitted values $\hat{\mu}_i^{(s)} = Y_i$
- ω_f has fitted values $\hat{\mu}_i^{(f)} = \hat{\boldsymbol{\beta}}^T \mathbf{X}_{\bullet i}$

The scaled deviance is then

$$\begin{aligned}
 S(\omega_f, \omega_s) &= 2(\ell_{\max}^{(s)} - \ell_{\max}^{(f)}) \\
 &= -\frac{1}{\sigma^2} \left(\|\mathbf{Y}_i - \mathbf{Y}_i\|^2 - \|\mathbf{Y}_i - X\boldsymbol{\beta}\|^2 \right) \\
 &= \frac{\text{RSS}_f}{\sigma^2} = \frac{\varepsilon_f^2}{\sigma^2} \\
 &\sim \chi_{n-p_f}^2 \quad \text{if } \omega_f \text{ is true}
 \end{aligned}$$

GLMs ~ Binomial Data

- If $Z \sim \text{bin}(n, p)$. Let $Y = Z/n$ be the proportion of successes. We note that $\mathbb{E}(Y) = p$. Now:

$$\begin{aligned}
 \mathbb{P}(Y = y) &= \mathbb{P}(Z = ny) \\
 &= {}^n C_{ny} p^{ny} (1-p)^{n-ny} \\
 &= {}^n C_{ny} \exp\left\{ \frac{y \log\left(\frac{p}{1-p}\right) + \log(1-p)}{1/n} \right\}
 \end{aligned}$$

This is, once again, of exponential distribution family form, with

- $\sigma^2 = 1/n$

- $\mu = p$ and $\theta(p) = \log\left(\frac{p}{1-p}\right) \Rightarrow p = \frac{e^\theta}{1+e^\theta}$
- $k(\theta(p)) = -\log(1-p)$
- $V(p) = p(1-p)$
- Now consider variables $Y_i \sim \text{bin}(n_i, p_i)$.
 - The canonical link is given by

$$g(p_i) = \log\left(\frac{p_i}{1-p_i}\right) = \boldsymbol{\beta}^T \mathbf{X}_{i\bullet} \quad p_i = p_i(\boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}}$$

This is known as a **logit link function**.

Note that it can be written as:

$$\text{Odds} = e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}$$

and so the interpretation of the model is that an increase in a given covariate i will result in the odds being increase by a factor of $\exp(\beta_i)$.

- $\sigma_i^2 = \frac{1}{n_i} \sigma^2$, with $\sigma^2 = 1$.
- The likelihood is given by

$$L(\mathbf{Y}) = \prod_{i=1}^n {}^{n_i}C_{y_i} p_i^{y_i} (1-p_i)^{n_i-y_i}$$

- The loglikelihood is given by

$$\begin{aligned} \ell(\mathbf{Y}) &= \sum_{i=1}^n \left\{ Y_i \log p_i + (n_i - Y_i) \log(1-p_i) + \log\left({}^{n_i}C_{y_i}\right) \right\} \\ &= \sum_{i=1}^n \left\{ Y_i \log \frac{p_i}{1-p_i} + n_i \log(1-p_i) + \log\left({}^{n_i}C_{y_i}\right) \right\} \\ \ell(\boldsymbol{\beta}) &= \sum_{i=1}^n \left\{ Y_i \boldsymbol{\beta}^T \mathbf{X}_{i\bullet} + n_i \log\left(1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}\right) + \log\left({}^{n_i}C_{y_i}\right) \right\} \\ &= \boldsymbol{\beta}^T \left(\sum_{i=1}^n Y_i \mathbf{X}_{i\bullet} \right) - \sum_{i=1}^n n_i \log\left(1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}\right) + \sum_{i=1}^n \log\left({}^{n_i}C_{y_i}\right) \end{aligned}$$

- As such, the MLE is found using

$$\begin{aligned} \frac{d\ell}{d\boldsymbol{\beta}} &= \sum_{i=1}^n Y_i \mathbf{X}_{i\bullet} - \sum_{i=1}^n n_i \mathbf{X}_{i\bullet} \frac{e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}} = 0 \\ \sum_{i=1}^n Y_i \mathbf{X}_{i\bullet} &= \sum_{i=1}^n n_i \mathbf{X}_{i\bullet} \underbrace{\frac{e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\bullet}}}}_{p_i} \end{aligned}$$

(The equation basically says that the observed value should be equal to the expected value, and can be solved using an iterative method).

- We also note that

$$\begin{aligned}
 -\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= \sum_{i=1}^n n_i \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \left(\frac{e^{\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}}} - \frac{e^{2\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}}}{(1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}})^2} \right) \\
 &= \sum_{i=1}^n n_i \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T \frac{e^{\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}}}{(1 + e^{\boldsymbol{\beta}^T \mathbf{X}_{i\cdot}})^2} \\
 &= \sum_{i=1}^n n_i \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T p_i (1 - p_i)
 \end{aligned}$$

So

$$\mathbb{E} \left(-\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \right) = V(\boldsymbol{\beta})^{-1} = \sum_{i=1}^n n_i \mathbf{X}_{i\cdot} \mathbf{X}_{i\cdot}^T p_i (1 - p_i) = i(\boldsymbol{\beta})$$

And so

$$\hat{\boldsymbol{\beta}} \approx N_p(\boldsymbol{\beta}, i(\hat{\boldsymbol{\beta}})^{-1})$$

- We now look at deviance. Note that

- For the saturated model ω_s , we fit $\hat{p}_i^{(s)} = Y_i / n_i$, so

$$\ell_{\max}^{(s)} = \sum_{i=1}^n \left\{ Y_i \log \frac{Y_i}{n_i} + (n_i - Y_i) \log \left(\frac{n_i - Y_i}{n_i} \right) + \log \binom{n_i}{C_{y_i}} \right\}$$

- For the fitted model ω_f , $\hat{p}_i^{(f)} = p_i(\hat{\boldsymbol{\beta}}) = e^{\hat{\boldsymbol{\beta}}^T \mathbf{X}_{i\cdot}} / (1 + e^{\hat{\boldsymbol{\beta}}^T \mathbf{X}_{i\cdot}})^{-1}$

$$\ell_{\max}^{(f)} = \sum_{i=1}^n \left\{ Y_i \log [p_i(\hat{\boldsymbol{\beta}})] + (n_i - Y_i) \log [1 - p_i(\hat{\boldsymbol{\beta}})] + \log \binom{n_i}{C_{y_i}} \right\}$$

The deviance (=scaled deviance since $\sigma^2 = 1$) is given by

$$D(\omega_f, \omega_s) = 2 \sum_{i=1}^n \left\{ Y_i \log \left(\frac{Y_i}{n_i \hat{p}_i} \right) + (n_i - Y_i) \log \left(\frac{n_i - Y_i}{n_i - n_i \hat{p}_i} \right) \right\}$$

As usual, we assess the fit by comparing D to χ_{n-p}^2 where $p = \dim \omega_f$. If ω_f is a bad fit, D will be **larger** than this.

- Expanding the logs, we find that that

$$D(\omega_f, \omega_s) = \sum_{i=1}^n \left\{ \frac{(Y_i - n_i \hat{p}_i)^2}{n_i \hat{p}_i} + \frac{\{(n_i - Y_i) - n_i(1 - \hat{p}_i)\}^2}{n_i(1 - \hat{p}_i)} \right\}$$

This is **Pearson's chi-squared statistic**²³!

- Logit is the most commonly used link functions. Other possibilities:

²³ It is of the form $\sum_{\text{all variables for each variable}} \sum_{\text{2 possibilities}} \frac{(O-E)^2}{E}$.

- **Probit link:** $g(p_i) = \Phi^{-1}(p_i)$ where Φ denotes the $N(0,1)$ CDF.
- **Complimentary log-log:** $g(p_i) = \log(-\log(1 - p_i))$

GLMs ~ Poisson Data

- If $Y \sim \text{poi}(\lambda)$, then

$$\mathbb{P}(Y = y) = \exp\{y \log \lambda - \lambda\} \frac{1}{y!}$$

Once again of exponential distribution family form with

- $\sigma^2 = 1$
- $\mu = \lambda$ and $\theta(\lambda) = \log \lambda$
- $k(\theta(\lambda)) = \lambda = e^\theta$
- $V(\lambda) = \lambda$
- Now consider variables $Y_i \sim \text{Po}(\lambda_i)$.

- The canonical link is given by

$$\log(\lambda) = \boldsymbol{\beta}^T \mathbf{X}_{\cdot i} \quad \lambda_i = e^{\boldsymbol{\beta}^T \mathbf{X}_{\cdot i}} = \lambda_i(\boldsymbol{\beta})$$

- The likelihood is given by

$$\begin{aligned} \ell(\mathbf{Y}) &= \sum_{i=1}^n \{-\lambda_i + Y_i \log \lambda_i - \log(Y_i!)\} \\ \ell(\boldsymbol{\beta}) &= -\sum_{i=1}^n e^{\boldsymbol{\beta}^T \mathbf{X}_{\cdot i}} + \boldsymbol{\beta}^T \sum_{i=1}^n Y_i \mathbf{X}_{\cdot i} - \sum_{i=1}^n \log(Y_i!) \end{aligned}$$

- As such, the MLE is found using

$$\begin{aligned} \frac{\partial \ell}{\partial \boldsymbol{\beta}} &= -\sum_{i=1}^n \mathbf{X}_{\cdot i} e^{\boldsymbol{\beta}^T \mathbf{X}_{\cdot i}} + \sum_{i=1}^n Y_i \mathbf{X}_{\cdot i} = 0 \\ \sum_{i=1}^n Y_i \mathbf{X}_{\cdot i} &= \sum_{i=1}^n \mathbf{X}_{\cdot i} e^{\boldsymbol{\beta}^T \mathbf{X}_{\cdot i}} \end{aligned}$$

- Also

$$-\frac{\partial^2 \ell}{\partial \boldsymbol{\beta}^T \partial \boldsymbol{\beta}} = -\sum_{i=1}^n \mathbf{X}_{\cdot i} \mathbf{X}_{\cdot i}^T e^{\boldsymbol{\beta}^T \mathbf{X}_{\cdot i}} = i(\boldsymbol{\beta})$$

and so

$$\hat{\boldsymbol{\beta}} \approx N_p(\boldsymbol{\beta}, i(\hat{\boldsymbol{\beta}})^{-1})$$

- We now look at deviance. Note that

- For the saturated model ω_s , we fit $\hat{\lambda}_i^{(s)} = Y_i$, so

$$\ell_{\max}^{(s)} = \sum_{i=1}^n \{Y_i \log Y_i - Y_i - \log(Y_i!)\}$$

- For the fitted model ω_f , $\hat{\lambda}_i^{(f)} = \lambda_i^{(f)}(\hat{\boldsymbol{\beta}}) = e^{\hat{\boldsymbol{\beta}}^T \mathbf{X}_{\cdot i}}$

$$\ell_{\max}^{(f)} = \sum_{i=1}^n \left\{ Y_i \log[\lambda_i(\hat{\boldsymbol{\beta}})] - \lambda_i(\hat{\boldsymbol{\beta}}) - \log(Y_i!) \right\}$$

The deviance (=scaled deviance since $\sigma^2 = 1$) is given by

$$D = 2 \left\{ \ell_{\max}^{(s)} - \ell_{\max}^{(f)} \right\} = 2 \left\{ \sum_{i=1}^n Y_i \log \left(\frac{Y_i}{\lambda_i(\hat{\boldsymbol{\beta}})} \right) - \sum_{i=1}^n (Y_i - \lambda_i(\hat{\boldsymbol{\beta}})) \right\}$$

As usual, we assess the fit by comparing D to χ_{n-p}^2 , where $p = \dim \omega_f$. If ω_f is a bad fit, D will be **larger** than this.

- Expanding the logs using $s \log \left(\frac{s}{t} \right) \approx (s-t) + \frac{(s-t)^2}{2t}$, we find that that

$$D \approx \sum_{i=1}^n \frac{(Y_i - \lambda_i(\hat{\boldsymbol{\beta}}))^2}{\lambda_i(\hat{\boldsymbol{\beta}})} \left(= \sum \frac{O - E}{E} \right)$$

Once again, this is **Pearson's chi-squared statistic!**

- Suppose Y_1, \dots, Y_n are counts for different “exposures” m_1, \dots, m_n , and

$$Y_i \sim \text{Po}(\lambda_i) \quad \lambda_i = m_i \theta_i$$

where θ_i is a **rate**, and the interest lies in modelling how this rate depends on the covariates. Our model is then

$$\log(\lambda_i) = \log(m_i \theta_i) = \log(m_i) + \underbrace{\log(\theta_i)}_{\beta^T \mathbf{X}_i = \log \theta_i}$$

In this type of situation, $\log(m_i)$ is called an **offset**, and its coefficient is forced to be 1.

- Note that if we fit a model of the kind $\log(\lambda_i) = \alpha + \beta x_i$, we are effectively saying that an increase of 1 in x results in an increase of $\exp \beta$ in the mean.

GLMs ~ Contingency Tables

- Consider an $r \times c$ contingency table containing n individuals, where the probability of being in cell (i, j) is p_{ij} and the number of items in cell (i, j) is Y_{ij} . Let $\mathbf{Y} = (Y_{11}, Y_{12}, \dots, Y_{rc})^T$. Then

$$\mathbb{P}(\mathbf{Y} = \mathbf{y}) = \frac{n!}{\prod_{i,j} (y_{ij}!)} \prod_{i,j} p_{ij}^{y_{ij}}$$

To see, why, consider that

- $n!$ is the total number of ways of arranging n items ignoring bins.

- Each bin can be arranged in $y_{ij}!$ ways, so we have $\prod(y_{ij}!)$ repeats.
- Thus, $n!/\prod(y_{ij}!)$ is the total number of ways of getting a given arrangement.
- $\prod p_{ij}^{y_{ij}}$ is the probability of any one such arrangement.

So $\mathbf{Y} \sim \text{Multinomial}(n, \mathbf{p})$, where $\mathbf{p} = (p_{11}, p_{12}, \dots, p_{rc})^T$

- Unfortunately, the GLM framework does not admit a multinomial distribution. However, it can be shown that if $Y_{ij} \sim \text{Po}(\mu_{ij})$

$$\mathbf{Y} \mid \left(\sum_{i,j} Y_{ij} = n \right) \sim \text{Multinomial} \left(n, \left(\frac{\mu_{11}}{\sum_{i,j} \mu_{ij}}, \dots, \frac{\mu_{rc}}{\sum_{i,j} \mu_{ij}} \right)^T \right)$$

So we can simply use a Poisson-like GLM, with the sum of all Y constrained to be n .

Computer Code

Generalised Linear Models

- The error estimates given in `summary(glmName)` come from the asymptotic normality of MLEs.
- `anova(glmName, test="Chisq")` will return an analysis of the deviance for the GLM, with relevant chi-squared values.

Binomial Data

- Best illustrated by example

Covariate	Num Trials	Num "yes"	
1	50	10	Should be input as three vectors
2	51	11	
3	52	10	

Covariate	=	1	2	3
numTrials	=	50	51	52
numYes	=	10	11	10

The model is then fit by setting `pYes <- numYes/numTrials` followed by

```
glmName <- glm(pYes~Covariate,family=binomial,weights=numTrials)
```

- To use a **probit** link function instead, use


```
glmName <- glm(pYes~Covariate,family=binomial(link=probit),weights=numTrials)
```
- It is worth noting that another way to find these models is


```
glmName <- glm(dataTable~Covariate,family=binomial)
```

where `dataTable` contains two columns; the first being the “successes” and the second being the “fails”. The table can be generated using, for example, `cbind`.

- For categorical data, every point made previously regarding design plots, interaction plots, and fitting applies here as well.
- Note that if our binomial data is binary (ie: only values are 0 and 1), there is no need to specify a weight in the `glm` function.
- In a model with only one observation for each combination of factors, the fit will be perfect. The maximum loglikelihood of the saturated model will therefore be 0.

Poisson Data

- To fit the original model, use `glmName <- glm(Observations~Covariates,poisson)`
- If **Observations** are counts for different amount of exposures, stored in **Exposures**, fit a rate model using `glmName <- glm(Observations~offset(log(Exposures)) + Covariates,poisson)`
- To display a different kind of residual, use `residuals(glmName, type="pearson")`

{High Dimensional Problem}

{Multiple Testing}

I have been asked to remove all notes from the Part III Statistical Theory class from this website. Please email me if you have any questions.

Non-Parametric Statistics

Two Independent Samples

- Consider two samples $X_1, \dots, X_m \sim \text{IID } F_X$ and $Y_1, \dots, Y_n \sim \text{IID } F_Y$
- Assume that $X_i = e_i$ and $Y_j = e_{m+j} + \Delta$, where the e are IID from the **same** continuous distribution.
- We are interested in testing

$$H_0 : \Delta = 0 \qquad H_1 : \Delta > 0$$

- The **Wilcoxon rank sum test** combines both samples, and assigns a rank S_1, \dots, S_n to each of the Y_1, \dots, Y_n . We then **reject H_0** if

$$T = \sum_{i=1}^n S_i > c$$

where c is such that

$$\mathbb{P}(T > c \mid H_0) = \alpha$$

- We can work c out because under H_0 , each of the ${}^{m+n}C_n$ arrangements of in the pooled sample are equally likely and have probability $1 / {}^{m+n}C_n$.
- For large samples, approximations need to be used.

Matched Pairs

- Consider a set of matched pairs $(X_1, Y_1), \dots, (X_n, Y_n)$.
- The differences are given by $Z_i = Y_i - X_i$, and we assume $X_i = \theta + e_i$, where the e are IID from a continuous distribution symmetric about 0.
- We are interested in testing

$$H_0 : \theta = 0 \qquad H_1 : \theta > 0$$

- The **Wilcoxon rank sum test** assigns a rank S_i to each of the $|Z_i|$ and then adds a sign to that rank; positive if Z_i is positive, negative otherwise. We then **reject H_0** if

$$W = \sum_{\substack{\text{only positive} \\ \text{signed ranks}}} S_i > c$$

where c is such that

$$\mathbb{P}(W > c \mid H_0) = \alpha$$

- We can work c out because under H_0 , each of the 2^n possible assignments of signs to the ranks is equally likely, with probability $1 / 2^n$.

- For large samples, approximations need to be used.

Computer Code

- For an unpaired Wilcoxon test on vectors \mathbf{x} and \mathbf{y} in which we want to test whether $\mathbf{x} < \mathbf{y}$, use `wilcox.test(x,y,alternative="less")`
- For a paired Wilcoxon test on vectors \mathbf{x} and \mathbf{y} , use `wilcox.test(x,y,paired=T)`. This function will try to run the test directly and will otherwise use an approximation and return a warning. Adding an extra argument `exact=F` uses an approximation directly and mutes the warning.