

Statistical Modelling

Ralf Blöchliger

Contents

1	Refresher on Linear Algebra, Probability, and Statistics	3
1.1	Linear Algebra	3
1.2	Probability	8
1.3	Statistics	12
2	Estimating the Linear Regression Model	17
2.1	Variable transformations	17
2.2	Linear Regression as Orthogonal Projection	18
2.3	Simple vs Multiple Regression and Partial Correlation	20
2.4	Computational considerations for estimating linear regression	21
2.5	Categorical Variables and ANOVA	21
3	Properties and Efficiency of OLS Estimates	24
3.1	Gauss-Markov Theorem	25
3.2	Distribution of $\hat{\beta}$	29
3.3	Testing	33
3.4	F-Test	35
3.5	Residual Analysis	39
4	Model Selection	40
4.1	Approaches for model selection	41
4.2	Derivation of Mallows' C_p	42
4.3	Relation between C_p and AIC for Linear Regression	44
4.4	General Search Strategies	45

5	Non i.i.d. errors	46
5.1	Known covariance matrix – Generalised Least Squares	46
5.2	Unknown covariance matrix, known structure	48
5.3	Mixed Models	51
5.4	Comparison between different approaches for non iid errors	54
6	Generalized Linear Models	55
6.1	Logistic Regression	56
6.2	Poisson Regression	59
6.3	Gamma Regression	61
6.4	Overview: Generalized Linear Models	62
7	Extensions	63
7.1	Non-Linear Regression	63
7.2	Non-parametric regression	67
7.3	High-dimensional Regression	72
A	Mathematical Statistics	76
A.1	Fisher Information and Fisher Scoring	76
A.2	Convergence	77

1 Refresher on Linear Algebra, Probability, and Statistics

1.1 Linear Algebra

1.1.1 Vectors

We confine ourselves to the real numbers.

Definition 1.1 Vector

A *vector* is an element $x \in \mathbb{R}^n$.

A *linear combination of vectors* is given by

$$b_1 \cdot x_1 + \dots + b_p \cdot x_p$$

for some $b_i \in \mathbb{R}$ and vectors in \mathbb{R}^n .

Definition 1.2 Span

The *span* of a nonempty set of vectors S consists of all linear combinations of the vectors in S .

Definition 1.3 Basis

A set of vectors S is a *basis* for a space V if

$$\text{span}(S) = V$$

We are interested in assigning properties such as lengths and angles to vectors for which we introduce the scalar product.

Definition 1.4 Scalar product

The *scalar* product of two vectors $a, b \in \mathbb{R}^n$ is defined as

$$\langle a, b \rangle = b^T a = \sum_{i=1}^n a_i b_i$$

b^T denotes the transpose of vector b . Geometrically, the scalar product is the product of the Euclidean magnitudes of the two vectors and the cosine of the angle between them. We therefore define

Definition 1.5 Length of a vector

If $a \in \mathbb{R}^n$, then the *length* of a is defined as

$$|a|^2 = \langle a, a \rangle$$

An often useful result, which can be compared to the triangle inequality for real numbers, is the following:

Proposition 1.1 Cauchy-Schwartz Inequality

Let $a, b \in \mathbb{R}^n$. Then

$$|\langle a, b \rangle|^2 \leq \langle a, a \rangle \cdot \langle b, b \rangle$$

We give perpendicular vectors a special name. We will often encounter such vectors.

Definition 1.6 Orthogonal vectors

We call two vectors $a, b \in \mathbb{R}^n$ *orthogonal* if

$$\langle a, b \rangle = 0$$

One can see that this corresponds to the vectors being perpendicular by noting that

$$\cos \theta \cdot (\|a\| \|b\|) = a \cdot b$$

where $\|x\| = \sqrt{|\langle x, x \rangle|}$, and observing that $\cos \frac{\pi}{2} = 0$.

Definition 1.7 Orthonormal vectors

We call two vectors $a, b \in \mathbb{R}^n$ *orthonormal* if they are orthogonal and each vector has unit length.

A useful result, whose proof follows almost 'by inspection', is the following:

Proposition 1.2

If e_1, \dots, e_n are an orthogonal basis of \mathbb{R}^n and

$$y = b_1 \cdot e_1 + \dots + b_n \cdot e_n$$

then $\langle y, e_1 \rangle = b_1$.

1.1.2 Matrices

A *matrix* is a rectangular array of numbers, symbols, or expressions, arranged in rows and columns, which is used to represent a mathematical object.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} = (a_{ij}) \in \mathbb{R}^{m \times n}$$

A $m \times n$ matrix consists of m rows and n columns.

Whenever we multiply a matrix X with a vector b , assuming that dimensions are conformable, we obtain a linear function

$$y = Xb$$

whereby y is a linear combination of columns of X , using weights in b .

Definition 1.8 Identity Matrix

The *identity matrix* is an $n \times n$ matrix whose diagonal elements all consists of 1 and all off-diagonal elements are 0. We denote the identity matrix by

$$I = I_{n \times n} = \mathbb{I} = \text{diag}(1)_{n \times n}$$

Definition 1.9 Transpose Matrix

The *transpose* of a matrix is a matrix ‘mirrored’ on the diagonal.

$$A^T = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix}$$

Under some circumstances, a square matrix is *invertible*. An invertible matrix can *undo* the linear transformation of the vector space created by the original matrix.

Definition 1.10 Invertible Matrix and Inverse

A matrix A is **invertible**, if $\exists B$, s.t.:

$$AB = BA = I_n$$

We call B the **inverse matrix** of A and denote it by A^{-1}

A result we wil often use is the following:

Proposition 1.3

The matrix $A = (X^T X)^{-1}$ is symmetric.

1.1.3 Orthogonal projections

Let $n \geq p \in \mathbb{N}$. Consider a point $y \in \mathbb{R}^n$ and a subspace S spanned by linearly independent vectors x_1, \dots, x_p . We are interested in finding the orthogonal projection of y onto S . An orthogonal projection gives us the vector in S that lies ‘closest’ to y in the sense that the distance between the y and all vectors of the subspace S is minimised.

Definition 1.11

Let $X \in \mathbb{R}^{n \times p}$ be the matrix with columns being the linear independent vectors x_1, \dots, x_p each of length n , spanning a vector space S . The matrix

$$P_S = X \left(X^T X \right)^{-1} X^T \tag{1.1}$$

is called the **orthogonal projection matrix** onto S .

Proposition 1.4

For any $y \in \mathbb{R}^n$ and set of linear independent vectors x_1, \dots, x_p , each of length n , spanning a vector space S , the vector defined by

$$v = P_s y$$

is the vector in the space S which is closest to y .

We can think of the orthogonal projection matrix as applying two operations:

1. *Analysis*: $(X^T X)^{-1} X^T$ we find 'coordinates' or 'coefficients'
2. *Synthesis*: We use our above found coordinates, to form a linear combination of the columns in X .

We now discuss some important properties of orthogonal projections.

Proposition 1.5

The projection matrix P is idempotent and symmetric. I.e.,

$$P = P^2 = P^T$$

This proposition can be intuitively understood as follows: if we project a vector into a subspace S and then once again apply the same projection, the resulting vector will not change since it already lies in the subspace.

Proposition 1.6

Each eigenvalue of P is either 0 or 1.

Proof. If A is idempotent, λ is an eigenvalue and v a corresponding eigenvector then

$$\lambda v = Av = AA v = \lambda Av = \lambda^2 v$$

Since $v \neq 0$ we find $\lambda - \lambda^2 = \lambda(1 - \lambda) = 0$ so either $\lambda = 0$ or $\lambda = 1$.

□

Proposition 1.7

The rank of P is equal to the trace of $P = p$.

Proof. One can show that projection matrices are *diagonalizable*. I.e., there exists a diagonal matrix of eigenvalues E and matrix C such that

$$P = CEC^{-1}$$

Noting the cyclical property of the trace,

$$\text{tr}(P) = \text{tr}(CEC^{-1}) = \text{tr}(C^{-1}CE)$$

Thus the trace of P is equal to p .

This is equal to the rank of P , since the rank is equal to the number of non-zero eigenvalues for a diagonalizable matrix. \square

Proposition 1.8

The matrix $1 - P$ is also a projection and $P(1 - P) = 0$

1.1.4 Quadratic Forms and positive definite matrices

For any $n \times n$ matrix A and vector of length n \mathbf{x} , the quadratic form is given by

$$\mathbf{x}^T A \mathbf{x}$$

Definition 1.12 Positive (semi-)definite Matrix

A $n \times n$ matrix A is **positive definite** if

$$x^T A x > 0 \text{ for all } x \in \mathbb{R}^n \setminus \{0\}$$

A is **positive semi-definite** if

$$x^T A x \geq 0 \text{ for all } x \in \mathbb{R}^n$$

Proposition 1.9

All eigenvalues of a positive-definite matrix A are positive.

Proof. Assume for contradiction that some eigenvalue λ of A is zero or negative. Let x be the associated eigenvector. Then

$$x^T A x = \lambda x^T x = \lambda |x|^2 \leq 0$$

which contradicts that A is PD. \square

Using similar arguments relating traces and determinants to eigenvalues as above, we have the following corollary:

Corollary 1.1

The trace and determinant of a positive-definite matrix A are positive.

The following result will at times be useful:

Proposition 1.10

Let \mathbf{B} be an $n \times p$ -matrix. The matrix $\mathbf{B}^T \mathbf{B}$ is then symmetric and nonnegative definite. If $\text{rk}(\mathbf{B}) = p$, then $\mathbf{B}^T \mathbf{B}$ is positive definite.

One way of seeing this is to note that if x is any vector of appropriate dimensions then

$$\begin{aligned} x^T \mathbf{B}^T \mathbf{B} x &= (\mathbf{B}x)^T \mathbf{B}x \\ &= \sum (\mathbf{B}x)_{ii}^2 \geq 0 \end{aligned}$$

On occasion, we will use the following two results, relating positive-definite matrices to other matrices with nice properties:

Proposition 1.11

If A is psd, then there exists a unique symmetric square root matrix B such that

$$A = BB = BB^T$$

The proof uses the Eigendecomposition.

Proposition 1.12 Choleksy decomposition

For symmetric positive-definite matrix A , there exists a real lower triangular matrix with positive diagonal entries L such that

$$A = LL^T$$

1.2 Probability

1.2.1 Expectations

$$\mathbb{E}(X) = \int_{\Omega} x f(x)$$

Properties:

- Linearity: $\mathbb{E}[a\mathbf{X} + b\mathbf{Y}] = a\mathbb{E}[\mathbf{X}] + b\mathbb{E}[\mathbf{Y}]$
- For independent components: $\mathbb{E}[X_i X_j] = \mathbb{E}[X_i] \mathbb{E}[X_j]$

Estimator:

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}$$

Covariance matrix:

$$\Sigma = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T]$$

With elements:

$$\Sigma_{ij} = \text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$$

Univariate case:

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}(X)) \cdot (Y - \mathbb{E}(Y))] = \mathbb{E}[X \cdot Y] - \mathbb{E}(X)\mathbb{E}(Y)$$

Properties:

- Symmetry: $\Sigma = \Sigma^T$
- Positive semi-definite: $\mathbf{v}^T \Sigma \mathbf{v} \geq 0$ for all vectors \mathbf{v}
- For linear transformation $\mathbf{Y} = A\mathbf{X}$: $\text{Cov}(\mathbf{Y}) = A\Sigma A^T$

Estimator:

$$\hat{\Sigma} = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}})(\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}})^T$$

Note: Independence between X, Y implies $\text{Cov}(X, Y) = 0$ but the inverse does not hold except for if X, Y are jointly normal.

1.2.2 Normal Distribution

The normal distribution in one dimension is characterised by the density

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right\}$$

- Note: if $Z \sim \mathcal{N}(0, 1)$, then $Y = \mu + \sigma Z \sim \mathcal{N}(\mu, \sigma^2)$.
- We usually denote the density ϕ and the cdf by Φ .

We can extend the univariate normal to multiple dimensions. The multi-variate density is given by

$$f_{\mathbf{X}}(x_1, \dots, x_k) = \frac{\exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}}$$

If $Z = (Z_1, \dots, Z_n)$, $Z_i \sim \mathcal{N}(0, 1)$ i.i.d. , then $Z \sim N(0, I_n)$. If $Y = \mu + AZ$, then

$$Y \sim \mathcal{N}(\mu, \Sigma) \text{ with } \Sigma = AA^T$$

The reason, why we often encounter the normal distribution is due to the central limit theorem.

Theorem 1.1 Classical Central Limit Theorem

If X_i are i.i.d. with $\mathbb{E}(X_i) = \mu$, $\text{Var}(X_i) = \sigma^2 < \infty$ then

$$\frac{1}{\sqrt{n}\sigma} \sum (X_i - \mu) \rightarrow \mathcal{N}(0, 1) \text{ in distr as } n \rightarrow \infty$$

A more general statement is given by the Lindenberg CLT, which does not require identical but still independent distribution of our variables X_i :

Theorem 1.2 Lindenberg Central Limit Theorem

Suppose X_i are independently distributed random variables with $\mathbb{E}(X_i) = \mu_i$, $\text{Var}(X_i) = \sigma_i^2 < \infty$. Let $s_n^2 = \sum_{i=1}^n \sigma_i^2$. If

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{i=1}^n \mathbb{E} \left[(X_i - \mu_i)^2 \cdot \mathbf{1}_{\{|X_i - \mu_i| > s_n\}} \right] = 0,$$

then

$$\frac{1}{s_n} \sum_{i=1}^n (X_i - \mu_i) \rightarrow \mathcal{N}(0, 1) \text{ (in distr as } n \rightarrow \infty)$$

Note that if $X_1, \dots, X_n \sim N(\mu, \sigma_X^2)$ i.i.d. then $\bar{X}_n \sim N\left(\mu, \frac{\sigma_X^2}{n}\right)$ exactly.

We next discuss important properties of the normal distribution.

- If Y_i, Y_j are jointly normally distributed random variables, then

$$\text{Cov}(Y_i, Y_j) = 0 \iff Y_i \perp Y_j$$

I.e., zero covariance is equivalent to independence for normally distributed RVs¹

- Linear transformations preserve normality. If $Y \sim \mathcal{N}(\mu, \Sigma)$ and A is a matrix of appropriate dimensions then

$$AY \sim \mathcal{N}(A\mu, A\Sigma A^T)$$

We also have the following more technical result:

Lemma 1.1 Independent linear transformations of a normal vector

Suppose $Y \sim \mathcal{N}(0, \sigma^2 I)$. Assume $A \in \mathbb{R}^{p \times n}, B \in \mathbb{R}^{q \times n}$ with $p + q \leq n$. Define $U = AY, V = BY$. Then

$$U \perp V \iff AB^T = 0$$

Proof. Write

$$\begin{bmatrix} U \\ V \end{bmatrix} = \begin{bmatrix} A \\ B \end{bmatrix} Y = CY$$

¹Note that X, Y each being normal (but not jointly normal), is not sufficient. Counterexample: Let Z be standard normal and let $X = Z$ and $Y = Z$ if $|Z| \leq 1$, and $Y = -Z$ if $|Z| > 1$. Then X and Y are each normal, $\text{Cov}(X, Y) = 0$, but they are not independent.

then $CY \sim N(c_\mu, \Sigma)$. From normality, it follows that $U \perp V \iff \text{Cov}(U, V) = 0$. We have

$$\Sigma = \sigma^2 \begin{bmatrix} AA^T & AB^T \\ BA^T & BB^T \end{bmatrix}$$

Therefore, the covariance is zero if $AB^T = 0$ and since $\begin{bmatrix} U & V \end{bmatrix}^\top$ is multivariate normal, zero covariance implies independence. \square

1.2.3 Distributions related to the normal

Chi-Squared.

- If $Z \sim \mathcal{N}(0, 1)$ then $Z^2 \sim \chi_1^2$
- If $Z_i \perp Z_j$ and each $Z_i \sim \mathcal{N}(0, 1)$ then $\sum_{i=1}^n Z_i^2 \sim \chi_n^2$
- If $Y \sim \chi_n^2$ then $\mathbb{E}[Y] = n$, $\text{Var}(Y) = 2n$
- If $Y \sim \mathcal{N}(\mu, \Sigma)$ with $\dim(Y) = n$, then $(Y - \mu)^T \Sigma^{-1} (Y - \mu) \sim \chi_n^2$

The following shows that if we have a $n \times n$ matrix with rank smaller than n (i.e., degenerate columns/rows), then only the actual rank ‘matters’ for the χ^2 distribution:

Lemma 1.2

Let $\mathbf{e} \sim \mathcal{N}(0, M)$, where M is an idempotent matrix of rank r . Then $\mathbf{e}^\top M \mathbf{e}$ has a central chi-square distribution with r degrees of freedom.

This follows from the fact that for a symmetric, idempotent matrix there exists an orthogonal matrix Γ such that $M = \Gamma D_\lambda \Gamma'$, where $D_\lambda = \text{diag}(1, \dots, 1, 0, \dots, 0)$ with the number of ones being equal to the rank of the matrix M . Thus $\mathbf{e}' M \mathbf{e} = \mathbf{e}' \Gamma D_\lambda \Gamma' \mathbf{e} = \mathbf{z}' D_\lambda \mathbf{z} = \sum_{i=1}^r z_i^2$, where

$$\mathbf{z} = (z_1, \dots, z_n)' = \Gamma' \mathbf{e} \sim N(0, \Gamma M \Gamma') = N(0, D_\lambda).$$

I.e., the vectors \mathbf{z} are multivariate standard normal with the truncated identity matrix which is $D_\lambda = \text{diag}(1, \dots, 1, 0, \dots, 0)$. The matrix is not full rank, however the truncated matrix would be full-rank and is identical to the identity and thus identical to its inverse. We can thus apply our definition of the multi-variate χ^2 matrix from above.

F-distribution

- If $X \sim \chi_m^2, Y \sim \chi_n^2$ with $X \perp Y$, then

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

- If $Y \sim F_{m,n}$ then $\mathbb{E}[Y] = n/(n-2)$

t-distribution

- If $Z \sim \mathcal{N}(0, 1)$, $X \sim \chi_k^2$ with $Z \perp X$ then

$$T = \frac{Z}{\sqrt{X/k}} \sim t_k$$

- The t-distribution is related to the F-distribution in a similar way as \mathcal{N} is related to χ^2 :

$$T \sim t_k \rightarrow T^2 \sim F_{1,k}$$

1.3 Statistics

1.3.1 MLE

In most cases, we want to use maximum likelihood estimators. MLEs have great asymptotic properties

- Consistent: $\hat{\theta} \rightarrow \theta$ in distribution as $n \rightarrow \infty$
- Asymptotically normal: $\sqrt{n}(\hat{\theta} - \theta) \rightarrow N(0, I^{-1})$ in distribution as $n \rightarrow \infty$ where I denotes Fischer Information
- Efficient: MLE is asymptotically unbiased and has the smallest possible variance among reasonable estimators

1.3.2 Hypothesis Testing

z-Test

If we have normally distributed variables with known σ_x (or alternatively, if we can invoke a CLT to obtain asymptotic normality), then

$$\bar{X} \sim \mathcal{N}(\mu, \sigma_{\bar{x}}^2), \text{ with } \sigma_{\bar{x}} = \sigma_x / \sqrt{n}$$

We obtain our test and test result in five steps:

1. Specify null and alternative hypothesis:

$$H_0 : \mu = \mu_0, H_A : \mu \neq \mu_0$$

2. Test statistic:

$$Z = \frac{\bar{X}_n - \mu_0}{\sigma_{\bar{X}_n}} = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\sigma_X}$$

Distribution of test statistics under $H_0 : Z \sim \mathcal{N}(0, 1)$

3. Choose level of significance α (e.g. 0.05) and calculate rejection region (“tails” of null distribution):

$$K = \left(-\infty, -\Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\right] \cup \left[\Phi^{-1}\left(1 - \frac{\alpha}{2}\right), \infty\right) \text{ if } H_A : \mu \neq \mu_0$$

4. Finally, take samples, calculate test and reject or fail to reject

t-test

We usually do not know σ_x and need to estimate it from data. We proceed as before but modify the test statistic and rejection region based on the new distr. of the test statistic:

$$T = \frac{\bar{X}_n - \mu_0}{\hat{\sigma}_{\bar{X}_n}} = \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\hat{\sigma}_X}$$

where $\hat{\sigma}_X^2 = \frac{1}{n-1} \sum (X_i - \bar{X}_n)^2$.

- Distribution of test statistics under $H_0 : T \sim t_{n-1}$.
- Rejection region for test statistics:

$$K = \left(-\infty, -qt \left(n-1; 1 - \frac{\alpha}{2} \right) \right] \cup \left[qt \left(n-1; 1 - \frac{\alpha}{2} \right), \infty \right) \text{ if } H_A : \mu \neq \mu_0$$

Likelihood Ratio Test

Due to the Neyman-Pearson lemma, we know that likelihood ratio tests have largest power for deciding between two concrete single hypothesis. Asymptotically, LR tests are χ^2 distributed.

Definition 1.13 Likelihood ratio test

Define λ by

$$\lambda = \frac{L(\hat{\Omega}_0)}{L(\hat{\Omega})} = \frac{\sup_{\Theta \in \Omega_0} L(\Theta)}{\sup_{\Theta \in \Omega} L(\Theta)}.$$

A likelihood ratio test of $H_0 : \Theta \in \Omega_0$ versus $H_a : \Theta \in \Omega_a$ employs λ as a test statistic, and the rejection region is determined by $\lambda \leq k$ where k is chosen to achieve the desired significance level α .

Where we use the following notation: null hypothesis specifies that Θ lies in a particular set of possible values Ω_0 and the alternative hypothesis specifies that Θ lies in another set of possible values Ω_a , which does not overlap Ω_0 . Let $L(\hat{\Omega}_0) = \sup_{\Theta \in \Omega_0} L(\Theta)$ and $L(\hat{\Omega}) = \sup_{\Theta \in \Omega} L(\Theta)$.

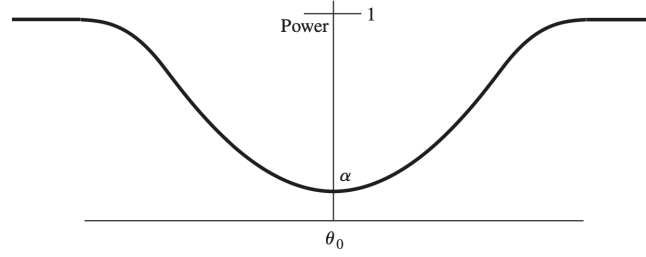
Type I, Type II error and power

Definition 1.14

A **type I error** is made if H_0 is rejected when H_0 is true. The probability of a type I error is denoted by α . The value of α is called the **level** of the test.

A **type II error** is made if H_0 is accepted when H_a is true. The probability of a type II error is denoted by β .

Figure 1: A typical power curve



I.e.,

$$\begin{aligned}\beta &= P(\text{type II error}) = P(\text{accepting } H_0 \text{ when } H_a \text{ is true}) \\ &= P(\text{value of the test statistic is not in RR when } H_a \text{ is true}).\end{aligned}$$

Clearly, α, β are inversely related as the higher α , the smaller the rejection region and thus the higher the possibility of not rejecting the null even though H_0 is wrong.

The trade-off between the two errors is closely related to a concept called power.

Definition 1.15 Power of a test

Suppose that W is the test statistic and RR is the rejection region for a test of a hypothesis involving the value of a parameter θ . Then the power of the test, denoted by $\text{power}(\theta)$, is the probability that the test will lead to rejection of H_0 when the actual parameter value is θ . That is, $\text{power}(\theta) = P(W \text{ in RR when the parameter value is } \theta)$.

If we write a decision function for our test such that $d(W) = 1$ if $W \in \text{RR}$ and 0 otherwise, then we can write

$$\text{power}(\theta) = \mathbb{E}_\theta(d(W)) \quad (1.2)$$

If θ_a is some value in H_a , i.e. a value not equal to the null hypothesis, then

$$\text{power}(\theta_a) = 1 - \beta(\theta_a)$$

P-value

Definition 1.16 P-Value

Consider an observed test-statistic t from unknown distribution T . Then the p -value p is what the prior probability would be of observing a test-statistic value at least as “extreme” as t if null hypothesis H_0 were true.

Note that under the null hypothesis, the p -value is uniformly distributed over $[0, 1]$. For simplicity, consider a one-sided t -test. Then $p = \Pr(T \geq t) = 1 - F(t)$ where F is the cdf of the distribution of the test statistic under H_0 . The *Probability Integral Transform* tells us

that the random variable $Z_t := F(t)$ is uniformly distributed over $[0, 1]$ for any cumulative distribution function.

Confidence Intervals

CI show a set of “compatible” parameter values. A $(1 - \alpha)$ -confidence interval contain all parameter values for which a test at significance level α would fail to reject.

We can often solve for CIs analytically: Find all values of μ , so that the test does *not* reject, i.e.:

$$\left| \frac{\sqrt{n}(\bar{x}_n - \mu)}{\hat{\sigma}_X} \right| \leq qt\left(n - 1; 1 - \frac{\alpha}{2}\right) \rightarrow \text{solve for } \mu$$

Thus, $(1 - \alpha)$ -confidence interval is:

$$\left(\bar{X}_n - qt\left(n - 1; 1 - \frac{\alpha}{2}\right) \cdot \frac{\hat{\sigma}_X}{\sqrt{n}}; \bar{X}_n + qt\left(n - 1; 1 - \frac{\alpha}{2}\right) \cdot \frac{\hat{\sigma}_X}{\sqrt{n}} \right)$$

- Rule of thumb for 95%-CI: $\left(\bar{X}_n - 2 \cdot \frac{\hat{\sigma}_X}{\sqrt{n}}; \bar{X}_n + 2 \cdot \frac{\hat{\sigma}_X}{\sqrt{n}} \right)$
- A $(1 - \alpha)$ -confidence interval covers the true value with probability $(1 - \alpha)$.

Multiple Testing

Suppose we consider performing 100 T-tests at a significance level $\alpha = 0.05$. Even if there is no effect, we expect to find about 5 significant results. One way of correcting for this is considering a conservative test level such that the probability of at least one type-1 error is $\leq \alpha$. This is called a **Familywise Error Rate**.

One simple and conservative approach is the **Bonferroni correction**: test at level α/m if m tests are performed.

1.3.3 Measuring association

Pearson correlation is a standardized measure of the covariance and measures linear association:

$$\rho_{X,Y} = \text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y} = \frac{\sigma_{XY}}{\sigma_X \cdot \sigma_Y}$$

For inference, we may use the *Fisher z-transform* which is given by

$$Z := \tanh^{-1}(\hat{\rho}) = \frac{1}{2} \log \left(\frac{1 + \hat{\rho}}{1 - \hat{\rho}} \right) \quad (1.3)$$

Proposition 1.13

If \mathbf{X}, \mathbf{Y} follows a multivariate normal distribution, then Z as defined in (1.3) is ap-

proximately

$$Z \sim \mathcal{N}\left(\tanh^{-1}(\rho), \frac{1}{n-3}\right)$$

A 95%-CI for ρ is thus given by

$$\tanh\left(z \pm 1.96 \cdot \sqrt{\frac{1}{n-3}}\right)$$

If we want to measure nonlinear associations as well, one possibility is **Spearman correlation** which measures monotone relationships. r_S can be obtained by computing ranks and then applying Pearson correlation, or using

$$r_S = 1 - \frac{6 \sum_{i=1}^n D_i^2}{n(n^2 - 1)}, \quad D_i := \text{rank}(x_i) - \text{rank}(y_i)$$

Simple Linear Regression

Consider estimating a linear relation between pairs of variables (x_i, Y_i) : I.e., we want to find a line of the form $Y_i = \beta_0 + \beta_1 x_i$ that is as close as possible to y_i . We can fit our model by minimizing the sum of squared residuals, i.e.,

$$\widehat{\beta}_0, \widehat{\beta}_1 = \arg \min \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

The solution is given by

$$\hat{\beta}_1 = \hat{\rho}_{XY} \cdot \frac{\hat{\sigma}_Y}{\hat{\sigma}_X}, \quad \hat{\beta}_0 = \bar{y}_n - \hat{\beta}_1 \bar{x}_n$$

2 Estimating the Linear Regression Model

Assume X are fixed and $\text{rank}(X) = p$ where p is the number of predictors (including the intercept). Then

$$Y_i = \beta_1 + \beta_2 X_i^{(2)} + \dots + \beta_p X_i^{(p)} + \epsilon_i$$

We may write our regression equation in vector or matrix form as

$$\begin{aligned} Y_i &= X_i^\top \beta + \epsilon_i, \\ Y &= X\beta + \epsilon \end{aligned}$$

with $\mathbb{E}(\epsilon_i) = 0, \text{Var}(\epsilon_i) < \infty$.

One way to estimate $\hat{\beta}$ from data is by minimizing the residual sum of squares:

$$\min \sum_{i=1}^n \left(Y_i - X_i^\top \beta \right)^2$$

This problem is equivalent to MLE with $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{1})$.

In matrix notation, we may write the objective function as

$$(Y - X\beta)^\top (Y - X\beta)$$

Taking gradient and setting equal to zero yields the normal equations:

$$X^T(Y - X\hat{\beta}) = 0 \rightarrow X^T X \hat{\beta} = X^T Y$$

The closed-form solution is then given by

$$\hat{\beta} = \left(X^T X \right)^{-1} X^T Y \quad (2.1)$$

To see equivalence between OLS and MLE under $\epsilon \sim \mathcal{N}$, note that the log likelihood of our data is given by

$$\log L = -n \cdot \log(\sigma \sqrt{2\pi}) - \frac{1}{2} \frac{\left(\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1,i} - \dots - \beta_{p-1} x_{p-1,i})^2 \right)}{\sigma^2}$$

clearly this quantity is maximized when the sum of squared residuals is minimized. One can also see from the quadratic function that we have a convex (or concave) optimization problem.

2.1 Variable transformations

Linear regression only requires linearity in parameters. We can for instance apply complex transformations to our X columns and our target variable Y .

We often encounter the following variable transformations in our regression models:

1. Exponential type:

$$y = \exp(\beta_0 + \beta_1 \cdot x) \cdot \tilde{\epsilon} \rightarrow \tilde{y} = \beta_0 + \beta_1 \cdot x + \epsilon \text{ with } \tilde{y} = \log(y)$$

2. Power type:

$$y = \beta_0 \cdot x^{\beta_1} \cdot \tilde{\varepsilon} \rightarrow \tilde{y} = \widetilde{\beta_0} + \beta_1 \cdot \tilde{x} + \varepsilon$$

$$\text{with } \tilde{y} = \log(y), \tilde{x} = \log(x), \widetilde{\beta_0} = \log(\beta_0)$$

In both case, we initially have a multiplicative error which we transform into an additive error.

Note, in general $g^{-1}(\mathbb{E}(g(Y))) \neq \mathbb{E}(Y)$. In particular, $\exp(\mathbb{E}(\log(Y))) \neq \mathbb{E}(Y)$. We therefore need to be careful when interpreting the results of our transformed equation back on the original scale.

However, if our distribution on the transformed scale is symmetric and since logarithms preserve ordering, $\mathbb{E}(\log(Y)) = \text{median}(\log(Y)) = \log(\text{median}(Y))$. Thus, $\exp(\mathbb{E}(\log(Y))) = \text{median}(Y)$. The effect interpretation of our coefficient is then

$$\frac{\exp(E(\log(Y) | x_i + 1))}{\exp(E(\log(Y) | x_i))} = \exp(\tilde{\beta}_1) = \frac{\text{Median}(Y | x_i + 1)}{\text{Median}(Y | x_i)}$$

I.e., “It is estimated that the median of Y given $x_i + 1$ is $\exp(\tilde{\beta}_1)$ times as large as the median of Y given x_i .”

2.2 Linear Regression as Orthogonal Projection

If we interpret our data in X ‘column-wise’, our columns are vectors spanning a (p -dimensional) subspace of \mathbb{R}^n .² I.e., the vector Y of observations is a single point in the n -dimensional space \mathbb{R}^n . If we vary the value of the parameter β , the product $X\beta$ describes a p -dimensional hyperplane through the origin.

We can obtain our OLS estimates by trying to find the element $X\beta$ on the hyperplane which lies closest to the point Y . I.e.,

$$\text{Choose } \beta \text{ to minimize } L_2\text{-norm: } |Y - X\beta|_2^2$$

Since we obtain the orthogonal projection of Y onto hyperplane spanned by columns of X , the residual vector $\hat{\varepsilon} = y - X\hat{\beta}$ must be orthogonal to any vector in the hyperplane, see Figure 2 for an illustration. With this, we recover the normal equations $X^T(Y - X\hat{\beta}) = 0$. Solving for $\hat{\beta}$ yields

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

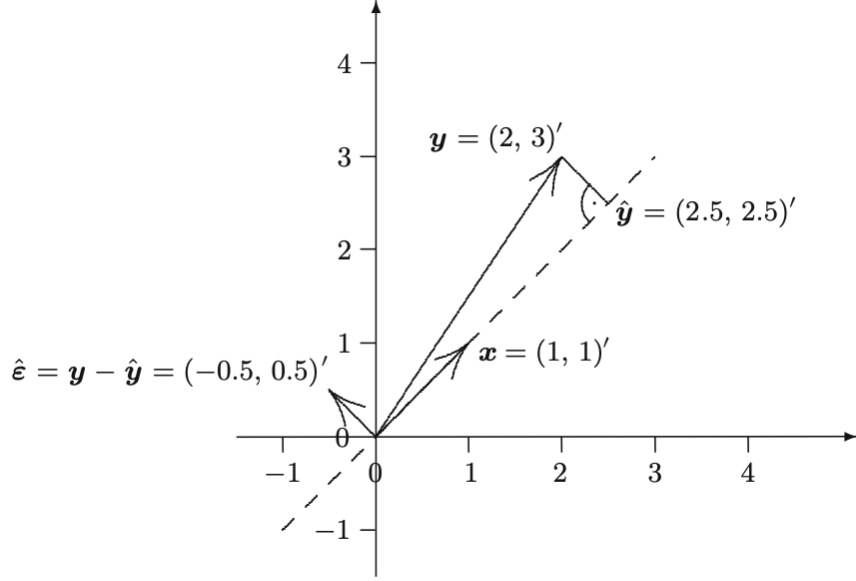
Our projection of Y onto this subspace is therefore

$$\hat{Y} = X\hat{\beta} = X \underbrace{(X^T X)^{-1} X^T}_{:=H} Y = HY$$

Note that H has the exact form of the orthogonal projection matrix in (1.1).

²Alternatively, we can interpret our data and regression ‘row-wise’ in which our features and target are each a different basis of \mathbb{R}^{p+1} . We then want to fit a hyperplane that is as close as possible to all of our points.

Figure 2: Graphical Intuition for Linear Regression as Projection



If we want to obtain our residuals, we can do so by using the matrix $M = 1 - H$, since

$$MY = (1 - H)Y = Y - \hat{Y} = \hat{\epsilon}$$

This is also a projection matrix. We also know (assuming $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$)

$$\hat{Y} \perp \hat{\epsilon} \text{ since } H \cdot M = 0 \quad (2.2)$$

which follows from Lemma 1.1. Also,

$$X^T \hat{\epsilon} = 0. \quad (2.3)$$

Intuitively, we can see (2.3) by noting that $\hat{\epsilon}$ is orthogonal to the column space of X . Formally, we can show this as

$$\begin{aligned} X^T \hat{\epsilon} &= X^T (y - \hat{y}) \\ &= X^T (1 - H)y \\ &= X^T (1 - X(X^T X)^{-1} X^T)y \\ &= X^T y - (X^T X)(X^T X)^{-1} X^T y \\ &= 0 \end{aligned}$$

This result leads to the following useful consequences:

- $\hat{\beta} \perp \hat{\epsilon}$
- $\hat{\epsilon} = MY = M(X\beta + \epsilon) = M\epsilon$

- If our model includes an intercept, one column of X is $E = (1, 1, \dots, 1)^T$. Thus

$$E \cdot \hat{\varepsilon} = \sum \hat{\varepsilon}_i = 0$$

- $\frac{1}{n} \sum_{i=1}^n \hat{y}_i = \bar{y}$.
- The regression hyperplane runs through the average of the data, i.e.,

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \dots + \hat{\beta}_k \bar{x}_k.$$

For the first of the above results, note

$$\hat{\beta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (\hat{Y} + \hat{\varepsilon}) = (X^T X)^{-1} X^T \hat{Y}$$

so $\hat{\beta}$ is only a function of X and \hat{Y} and $\hat{Y} \perp \hat{\varepsilon}$.

Our projection interpretation allows us to perform the following Pythagoras decomposition:

$$R^2 = 1 - \frac{RSS}{TSS} = \frac{ESS}{TSS}$$

- Called *coefficient of determination*
- Measure of goodness of fit
- Note that $|\hat{Y} - \bar{Y}|$ lies in the column space of X and $|Y - \hat{Y}|$ lies in the orthogonal complement. Thus these two vectors are two sides to a right-angled triangle so

$$TSS = ESS + RSS$$

where $TSS = |Y - \bar{Y}|^2$

- One can show that $R^2 = [\text{cor}(Y, \hat{Y})]^2$.

2.3 Simple vs Multiple Regression and Partial Correlation

Unless predictors are orthogonal, many simple regressions will yield different coefficient from multiple regression. Interpretation of multiple regression: Effect of x_k on Y when keeping x_{-k} fixed.

Special case: Orthogonal design

- Orthogonal predictors: $x_{\cdot,j}^\top x_{\cdot,k} = 0$ for all $j \neq k$
- $X^T X = \text{diag}(\sum_{i=1}^n x_{i,1}^2, \dots, \sum_{i=1}^n x_{i,p}^2)$ Thus:

$$\hat{\beta}_j = \left((X^T X)^{-1} X^T Y \right)_j = \frac{\sum_{i=1}^n x_{i,j} y_i}{\sum_{i=1}^n x_{i,j}^2} = \frac{\hat{\sigma}_{x_j, y}}{\hat{\sigma}_{x_j}^2}$$

which is the same result as simple regressions

We can therefore view multiple regression as measuring *partial correlation*. An informal definition of partial correlation is:

Partial correlation $\rho_{XY|Z}$: strength and direction of the linear dependence between X and Y after accounting for the linear dependence of X and Y on Z

One can estimate partial correlation in the following ways:

1. A recursive formula
2. Let r_x be residuals when regressing x on z , r_y residuals when regressing y on z , then $\rho_{XY|Z} = \text{cor}(r_x, r_y)$
3. Via precision matrix, $\rho_{YX^j|X^{-j}} = -\frac{K_{p+1,j}}{\sqrt{K_{p+1,p+1}K_{j,j}}}$ where $K = \Sigma^{-1}$ and $\Sigma = \text{Cov}\left((X_{*,1}, \dots, X_{*,p}, Y)^T\right)$

In multiple linear regression, we can write our estimated coefficients using partial correlation as

$$\beta_j = -\frac{K_{p+1,j}}{K_{p+1,p+1}}; \quad \rho_{YX^j|X^{-j}} = \beta_j \frac{\sqrt{K_{p+1,p+1}}}{\sqrt{K_{j,j}}} \quad (2.4)$$

2.4 Computational considerations for estimating linear regression

Instead of directly using the closed-form solution in (2.1), we can estimate the OLS coefficients in different ways. From (2.4) we can see that we can estimate them using partial correlations. The numerically preferred way is using a QR decomposition of X . I.e., let $X = QR$ where Q is $n \times p$ with orthonormal columns and R is $p \times p$ upper triangle. From the normal equations, we have

$$\begin{aligned} X^T X \hat{\beta} &= X^T y \\ (QR)^T (QR) \hat{\beta} &= (QR)^T y \\ R^T Q^T QR \hat{\beta} &= R^T Q^T y \\ R^T R \hat{\beta} &= R^T Q^T y \\ R \hat{\beta} &= Q^T y \\ \implies \hat{\beta} &= R^{-1} Q^T y \end{aligned}$$

2.5 Categorical Variables and ANOVA

We can include categorical variables as predictors by using dummy encoding (i.e., factors).

- E.g., we can allow for a different intercept between two groups, where x_2 is a binary variable:

$$Y_i = \beta_0 + \beta_1 x_1 + \beta_2 \mathbf{1}(x_2 = 1) + \epsilon_i$$

- One reference level and $k - 1$ binary dummies if we have k categories.

- Contrasts to derive the intercept in different groups. Let $\ell = \begin{pmatrix} 1 & 0 & 1 \end{pmatrix}^\top$. Then $\ell^\top \beta$ is the intercept for the group with $x_2 = 1$.
- Interactions to allow for different slopes between groups. E.g.,

$$Y_i = \beta_0 + \beta_1 x_1 + \beta_2 \mathbf{1}(x_2 = 1) + \beta_3 x_1 \cdot \mathbf{1}(x_2 = 1) + \epsilon_i$$

Analysis of Variance (ANOVA) is a special case of linear regression which is closely related to categorical variables. We illustrate the idea behind ANOVA using a simple example.

Let $g = 4$ be the number of groups (or ‘treatments’) and p denote the number of observations per group which we assume to be the same in each group. Our null hypothesis is

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$$

To test this, we calculate:

- **Variation between groups:**

$$SS_B = p \cdot \sum_{i=1}^g \left(\bar{Y}_{i.} - \bar{Y}_{..} \right)^2$$

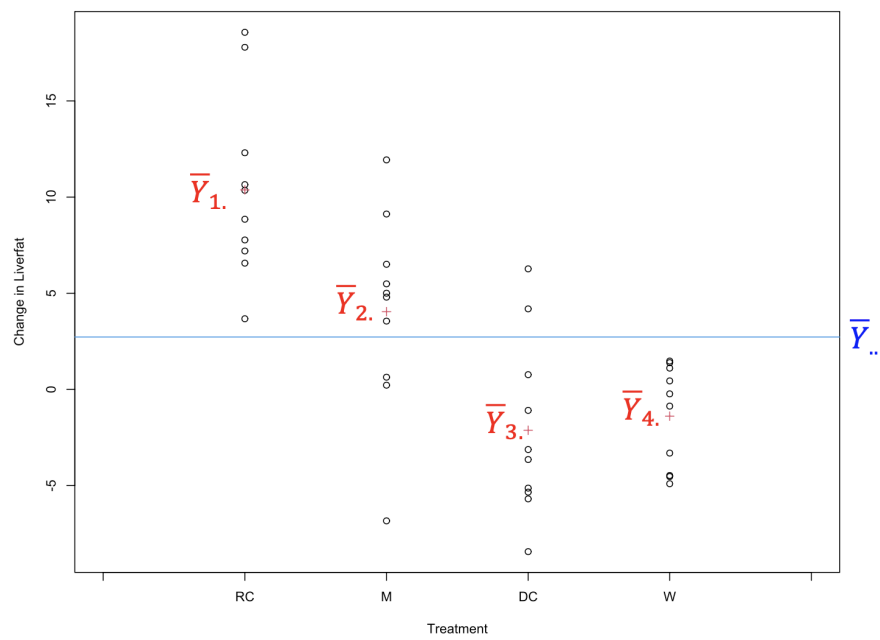
where $\bar{Y}_{..}$ is the overall mean of all of observations.

- **Variation within groups:**

$$SS_W = \sum_{i=1}^g \sum_{j=1}^p \left(Y_{ij} - \bar{Y}_{i.} \right)^2$$

Our test statistic is then approximately $\frac{SS_B}{SS_W}$ which follows an F-distribution under suitable assumptions. We reject our null if the variation between group averages is much larger than within group variation. The setup is illustrated in Figure 3.

Figure 3: Example ANOVA setup



3 Properties and Efficiency of OLS Estimates

Note that our coefficients depend on Y and, as Y is random, are random variables themselves. We must therefore think about precision and distribution of estimated parameters. We assume the *Gauss-Markov Conditions*:

$$Y = X\beta + \varepsilon$$

$$E(\varepsilon) = 0, \text{Cov}(\varepsilon) = \sigma^2 I$$

We obtain the following moments:

- $\hat{\beta}$:
 - $\mathbb{E}[\hat{\beta}] = \beta$
 - $\text{Cov}(\hat{\beta}) = \sigma^2(X^T X)^{-1}$
- \hat{Y} :
 - $\mathbb{E}[\hat{Y}] = \mathbb{E}[Y] = X\beta$
 - $\text{Cov}(\hat{Y}) = \sigma^2 H$
- $\hat{\varepsilon}$:
 - $\mathbb{E}[\hat{\varepsilon}] = 0$
 - $\text{Cov}(\hat{\varepsilon}) = \sigma^2 M$ which implies $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - H_{ii})$
 - $\text{Cov}(\hat{\varepsilon}, \hat{Y}) = 0$

For the variance of a single element of our coefficient vector, one can obtain the following formula

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{(1 - R_j^2) \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2},$$

where R_j^2 is the coefficient of determination from regressing x_j as the response on all other variables in X . We note:

- The smaller the model variance σ^2 , the smaller the variance of $\hat{\beta}_j$ and thus the more accurate the estimation.
- The smaller the linear dependence between x_j and the other explanatory variables (measured through R_j^2), the smaller is the variance of $\hat{\beta}_j$. The variance, $\text{Var}(\hat{\beta}_j)$, is minimized for $R_j^2 = 0$, i.e., when the covariates are uncorrelated.
- The larger the variability of covariate x_j around its average, the smaller is the variance of $\hat{\beta}_j$.

Further, we can show that $1/(n-p) \sum_{i=1}^n \hat{\epsilon}_i^2$ is an unbiased estimate of σ^2 . To this end, note that

$$\begin{aligned}\mathbb{E}[\hat{\sigma}^2] &= \frac{1}{n-p} \sum_{i=1}^n \mathbb{E}[\hat{\epsilon}_i^2] = \frac{1}{n-p} \sum_{i=1}^n \text{Var}[\hat{\epsilon}_i^2] \\ &= \frac{1}{n-p} \sum_{i=1}^n \sigma^2(1 - H_{ii}) = \frac{1}{n-p} \sigma^2 (n - \text{tr}(H)) \\ &= \frac{1}{n-p} \sigma^2 (n - \text{tr}(X(X^T X)^{-1} X)) \\ &= \frac{1}{n-p} \sigma^2 (n - \text{tr}(\mathbf{1}_{p \times p})) = \sigma^2\end{aligned}$$

Note here also the following insight into the relation between degrees of freedom and projection matrices:

- We have n observations and estimate p parameters to obtain residuals. Thus residual df are $n - p$
- The trace of a projection matrix measures the rank of the subspace onto which it projects a vector. I.e., how many independent dimensions there are in a subspace. We have $\text{tr}(M) = \text{tr}(\mathbf{1} - H) = n - p$.

3.1 Gauss-Markov Theorem

Given data Y , a parameter estimate in general is some function $\hat{\beta} = f(Y)$. We want to understand goodness properties of OLS and how (or, whether) there are ways to improve upon OLS.

Note that OLS has the following properties

1. $\hat{\beta} = (X^T X)^{-1} X^T y = Ay$ is *linear*
2. $\mathbb{E}(\hat{\beta}) = \beta$, i.e., *unbiased*

We will show that, depending on more or less restrictive assumptions, OLS can be

- Optimal among all *unbiased estimators*, and
- Optimal among all *linear and unbiased estimators*.

In both cases we may obtain an estimator with lower MSE by allowing for small bias if we can achieve a large reduction of variance simultaneously, see equation (3.4).

These results and the underlying assumptions are summarised in the following two theorems.

Theorem 3.1 Gauss-Markov – Version 1

Let $Y = X\beta + \varepsilon$, $E(\varepsilon) = 0$, $\text{Cov}(\varepsilon) = \sigma^2 I$, $\text{rank}(X) = p$. Furthermore, let $\ell \in \mathbb{R}^p$, and $\hat{\beta}$ the OLS estimator.

Then, for all $c \in \mathbb{R}^n$ such that $E(c^T y) = \ell^T \beta$, we have:

$$\text{Var}(\ell^T \hat{\beta}) \leq \text{Var}(c^T y).$$

Thus, the OLS estimator $\ell^T \hat{\beta}$ has minimal variance among all linear unbiased estimators of $\ell^T \beta$.

Proof. Let $y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}, \beta \in \mathbb{R}^p$. The OLS estimator is given by

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ \implies \ell^T \hat{\beta} &= \ell^T (X^T X)^{-1} X^T y \end{aligned}$$

We know that $\mathbb{E}[\ell^T \hat{\beta}] = \ell^T \beta$. Now let $c^T y$ be any other linear unbiased estimator of $\ell^T \beta$. Since this estimator is unbiased,

$$\ell^T \beta = \mathbb{E}[c^T y] = c^T X \beta$$

so that we must have $\ell^T = c^T X$.

We now consider variances.

$$\begin{aligned} \text{Var}(c^T y) &= \sigma^2 c^T c \\ \text{Var}(\ell^T \hat{\beta}) &= \sigma^2 \ell^T (X^T X)^{-1} \ell \\ &= \sigma^2 c^T X (X^T X)^{-1} X^T c \end{aligned}$$

Subtracting both variances, we obtain

$$\begin{aligned} &\text{Var}(c^T y) - \text{Var}(\ell^T \hat{\beta}) \\ &= \sigma^2 c^T \left(1 - X^T (X^T X)^{-1} X^T \right) c \\ &= \sigma^2 c^T M c \geq 0 \end{aligned}$$

since M is positive semi-definite.³

□

If we are further willing to assume normality of the error term, we arrive at a stronger result.

Theorem 3.2 Gauss-Markov – Version 2

Let furthermore ε be normally distributed. Then $\ell^T \hat{\beta}$ has minimal variance among

³Since M is a orthogonal projection, $M^T = M = M^2$. Let z be any vector of appropriate dimensions, then

$$z^T M z = z^T M^T M z = (M z)^T M z = |M z|^2 \geq 0$$

all unbiased estimators of $\ell^T \beta$.

To proof the above theorem, we will use the Cramer-Rao bound result:

Let $(f_\eta(\mathbf{Y}))$ be a parametric family of strictly positive densities in \mathbb{R}^n . Let η be a variable parameter with values in an open subset of \mathbb{R}^k , and let $f_\eta(\mathbf{Y})$ be differentiable wrt η . Our parameter of interest is $g(\eta)$, where g is an arbitrary real-valued function of η . Then we obtain the following result:

Theorem 3.3 Cramér-Rao for unbiased estimators

If $T(\mathbf{Y})$ is an arbitrary unbiased estimate of $g(\eta)$, i.e.

$$\mathbb{E}_\eta[T(\mathbf{Y})] = g(\eta) \quad \forall \eta,$$

then g is differentiable and

$$\text{Var}_\eta(T(\mathbf{Y})) \geq \frac{\partial g^T}{\partial \eta} I(\eta)^{-1} \frac{\partial g}{\partial \eta} \quad (3.1)$$

where $I(\eta)$ denotes the Fisher information matrix:

$$I(\eta) = \mathbb{E}_\eta \left[\frac{\partial \log f_\eta(\mathbf{Y})}{\partial \eta} \frac{\partial \log f_\eta(\mathbf{Y})^T}{\partial \eta} \right].$$

Proof. We know that

$$\text{Var}(\ell^T \beta) = \ell^T \sigma^2 (X^T X)^{-1} \ell \quad (3.2)$$

We now derive the Cramer-Rao lower bound. We have $\eta = (\sigma^2, \beta^T)^T$ and $g(\eta) = \ell^T \beta$. Then the Fisher information matrix is given by

$$I(\eta) = \mathbb{E}_\eta [\nabla_\eta \log L \cdot \nabla_\eta \log L^T]$$

where we use $\log L$ to denote the log likelihood function.

Our log likelihood function is

$$\log L = \sum_{i=1}^N -\frac{1}{2} \log(2\pi\sigma^2) - \frac{(y_i - x_i^T \beta)^2}{2\sigma^2}$$

Taking derivatives,

$$\nabla_\eta \log L = \begin{bmatrix} \sum_{i=1}^N \left(-\frac{1}{2\sigma^2} + \frac{(y_i - x_i^T \beta)^2}{2\sigma^4} \right) \\ \sum_{i=1}^N \left(\frac{(y_i - x_i^T \beta) x_i}{\sigma^2} \right) \end{bmatrix} \quad (3.3)$$

Since observations are independent, the Fisher information matrix is the sum of the individual Fisher information matrices for single observations. We compute,

$$\begin{aligned} I_{11}^{(i)} &= \mathbb{E} \left[\left(-\frac{1}{2\sigma^2} + \frac{(y_i - x_i^T \beta)^2}{2\sigma^4} \right)^2 \right] = \frac{1}{2\sigma^4} \\ I_{jk}^{(i)} &= \mathbb{E} \left[\frac{(y_i - x_i^T \beta)^2}{\sigma^4} x_j x_k \right] = \frac{x_j x_k}{\sigma^2} \\ I_{1j}^{(i)} &= I_{j1}^{(i)} = \mathbb{E}_\eta \left[\left(-\frac{1}{2\sigma^2} + \frac{(y - x^T \beta)^2}{2\sigma^4} \right) \cdot \frac{(y - x^T \beta) \cdot x_j}{\sigma^2} \right] = 0 \end{aligned}$$

For the first result, we use $\mathbb{E} \left[(y_i - x_i^T \beta)^4 \right] = 3\sigma^4$ for the normal distribution. In the last result, after extending and simplifying, we again use this result and the unbiasedness of our estimate.

Summing over all observations, we arrive at

$$I(\eta) = \begin{bmatrix} \frac{n}{2\sigma^4} & 0 \\ 0 & \frac{1}{\sigma^2} X^T X \end{bmatrix}$$

Taking the inverse of this block matrix, we get

$$I(\eta)^{-1} = \begin{bmatrix} \frac{2}{n}\sigma^4 & 0 \\ 0 & \sigma^2 (X^T X)^{-1} \end{bmatrix}$$

We compute the gradient vectors:

$$\frac{\partial g}{\partial \eta} = \begin{bmatrix} \partial \ell^T \beta / \partial \sigma^2 \\ \partial \ell^T \beta / \partial \beta_1 \\ \vdots \\ \partial \ell^T \beta / \partial \beta_p \end{bmatrix} = \begin{bmatrix} 0 \\ \ell_1 \\ \vdots \\ \ell_p \end{bmatrix} = \begin{bmatrix} 0 \\ \ell \end{bmatrix}$$

Finally, multiplying out, we obtain

$$\frac{\partial g^T}{\partial \eta} I(\eta)^{-1} \frac{\partial g}{\partial \eta} = \ell^T \sigma^2 (X^T X)^{-1} \ell$$

so that our OLS estimator's variance in (3.2) achieves the Cramer-Rao lower bound in (3.1). \square

Note that there may be biased estimators which yield a lower mean squared error. In particular, the MSE can be decomposed as

$$\mathbb{E} \left((\hat{\theta} - \theta)^2 \right) = (\mathbb{E}(\hat{\theta}) - \theta)^2 + \text{Var}(\hat{\theta}) \quad (3.4)$$

When may our OLS estimates perform badly?

- If our normality assumption holds, one may obtain slightly biased estimators with far lower variance
- If normality assumption is violated, OLS and MLE no longer coincide so we no longer have asymptotic efficiency. Non-linear estimates may outperform OLS.

3.2 Distribution of $\hat{\beta}$

Normal errors

Under normal errors, MLE coincides with OLS, so we should expect asymptotic normality of our estimated coefficients. We can, however, show that this already holds in finite samples.

Proposition 3.1

Assume that $Y = X\beta + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$. Then

1. $\hat{\beta} \sim \mathcal{N}_p(\beta, \sigma^2 (X^T X)^{-1})$
2. $\hat{Y} \sim \mathcal{N}_n(X\beta, \sigma^2 H)$, $\hat{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 M)$
3. \hat{Y} and $\hat{\varepsilon}$ are independent
4. $\frac{\sum_{i=1}^n \hat{\varepsilon}_i^2}{\sigma^2} \sim \chi_{n-p}^2$
5. $\hat{\sigma}^2$ is independent of $\hat{\beta} = (X^T X)^{-1} X^T Y$

Note that result 2) implies that even though ε are i.i.d., the covariance structure of $\hat{\varepsilon}$ can be highly dependent, including heteroskedasticity.

Proof. 1.) Since $Y \sim N(X\beta, \sigma^2 \mathbb{I})$ and $\hat{\beta} = (X^T X)^{-1} X^T y$, this implies that $\hat{\beta} \sim N(., .)$. Note that

$$\mathbb{E}[\hat{\beta}] = (X^T X)^{-1} X^T \mathbb{E}[y] = (X^T X)^{-1} X^T (X\beta + \mathbb{E}[\varepsilon]) = \beta$$

and

$$\text{Var}(\hat{\beta}) = [(X^T X)^{-1} X^T] \text{Var}(y) [(X^T X)^{-1} X^T]^T = \sigma^2 (X^T X)^{-1}$$

2.)

$\hat{Y} = HY$ where $H = X(X^T X)^{-1} X^T$ with $\text{rank}(H) = p$. Thus $\hat{Y} \sim N(., .)$. We have

$$\mathbb{E}[\hat{Y}] = H \mathbb{E}[Y] = HX\beta = X\beta = \mathbb{E}[Y]$$

This follows since H is a projection onto the column space of X and $X\beta$ evidently already lies in the col. space of X . Also,

$$\text{Var}(\hat{Y}) = H \text{Var}(Y) H^T = \sigma^2 H H^T = \sigma^2 H$$

$\hat{\varepsilon} = MY$ where $M = \mathbb{I} - H$.

$$\mathbb{E}[\hat{\varepsilon}] = M \mathbb{E}[y] = MX\beta = (\mathbb{I} - H)X\beta = 0$$

$$\text{Var}(\hat{\varepsilon}) = \text{Var}(MY) = \sigma^2 M$$

3.)

$\hat{Y} = HY, \hat{\varepsilon} = MY$. Note that

$$HM = 0$$

Thus, by Lemma 1.1, we have $\hat{Y} \perp \hat{\varepsilon}$

4.)

Let $e_i = \hat{\varepsilon}_i / \sigma$. Then

$$\frac{\sum \hat{\varepsilon}_i}{\sigma^2} = e^T e = e^T M e$$

We can multiply by the matrix M since e lies in the orthogonal complement of the column space of X and is therefore unaffected by multiplying with the orthogonal projection matrix M .

We know that $e_i \sim N, \mathbb{E}[e] = 0, \text{Var}(e) = \frac{1}{\sigma^2} \text{Var}(\hat{\varepsilon}_i) = M$. Therefore, $e \sim \mathcal{N}(0, M)$. We also know that M is idempotent and

$$\text{rank}(M) = \text{rank}(\mathbb{I} - H) = n - p$$

Then, using Lemma 1.2, $e^T M e \sim \chi_{n-p}^2$.

5.)

Since $X^T \hat{\varepsilon} = 0$.

$$\hat{\beta} = (X^T X)^{-1} X^T (\hat{Y} + \hat{\varepsilon}) = (X^T X)^{-1} X^T \hat{Y}.$$

Thus, $\hat{\sigma}$ is a function of $\hat{\varepsilon}$ and $\hat{\beta}$ is a function of \hat{Y} . Since $\hat{Y} \perp \hat{\varepsilon}$, we also have $\hat{\beta} \perp \hat{\sigma}$.

□

Gauss-Markov Assumptions

We now turn to the case where the errors are iid but not normally distributed. In this case, OLS no longer coincides with MLE so we must consider the use of some central limit theorem to derive asymptotic normality. In particular, we will use the Lindenberg-Feller CLT:

Theorem 3.4 Lindenberg-Feller Central Limit Theorem

Suppose $\{x_{ni}\}$ is a triangular array of $p \times 1$ random vectors such that $z_n = \frac{1}{n} \sum_{i=1}^n x_{ni}$ and

$$V_n = \frac{1}{n} \sum_{i=1}^n \text{Var}(X_{ni}) \rightarrow V, \text{ where } V \text{ is psd.}$$

If for every $\varepsilon > 0$ we have

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} \left\{ \|x_{ni}\|^2 \mathbf{1}(\|x_{ni}\| \geq \varepsilon \sqrt{n}) \right\} \rightarrow 0,$$

then $\sqrt{n} z_n \xrightarrow{d} \mathcal{N}(0, V)$. Or equivalently, $\sqrt{n} V_n^{-1/2} z_n \xrightarrow{d} \mathcal{N}(0, \mathbb{I})$.

For the asymptotic approximation to hold, we need some weak conditions on the explanatory variables \mathbf{x}_i :

- The smallest eigenvalue of $X^T X = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$, namely $\lambda_{\min, n}$, converges to ∞ .
- $\max_j P_{jj} = \max_j \mathbf{x}_j^T (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T)^{-1} \mathbf{x}_j$ converges to zero.

The first condition states that increasing n always yields more information, while the second condition prohibits any \mathbf{x}_j from dominating the others.

Theorem 3.5 Asymptotic normality of coefficients under Gauss-Markov

If the errors ε_i are i.i.d. with mean 0 and variance σ^2 , and if (\mathbf{x}_i) satisfies the conditions just given, then the LS estimators $\hat{\beta}$ are consistent (for β), and the distribution of

$$(X^T X)^{1/2} (\hat{\beta} - \beta)$$

converges weakly to $\mathcal{N}_p(\mathbf{0}, \sigma^2 I)$.

We will proof a slightly easier version of the above statement, by assuming that $x_i^T x_i$ is bonded instead of the second assumption.

Proof. Consistency.

The i -th component $\hat{\beta}_i$ is unbiased and has variance $\sigma^2 ((X^T X)^{-1})_{ii}$, which converges to zero by the first assumption. Then consistency follows from Chebyshev's inequality.

Asymptotic normality.

Suppose $y = X\beta + \omega$ where $\omega_i \sim F$ i.i.d with $\mathbb{E}[\omega_i] = 0$, $\text{Var}(\omega_i) = \sigma^2$. Suppose further that

$$\frac{1}{n} X^T X \rightarrow \Sigma, \text{ where } \Sigma \text{ is psd}$$

Let x_i denote the fixed $p \times 1$ vector of covariates of observation i and assume there exists $M \in \mathbb{R}$ such that $\|x_i\| < M < \infty, \forall i$.

Step 1: Rearrange OLS estimate

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X\beta + \omega) \\ &= \beta + (X^T X)^{-1} X^T \omega \end{aligned}$$

We can write

$$X^T \omega = \sum_{i=1}^n x_i \omega_i$$

That is, our estimate $\hat{\beta}$ is the true value plus some function of the noise term. While ω_i are iid, premultiplying with x_i yields independently but not identically distributed RVs. We wish to show that this noise term is asymptotically normal.

Step 2: Check Lindenberg-Feller Condition

We now consider the Lindenberg-Feller CLT.

A triangular array is given by random variables in the form

$$\begin{array}{cccc} x_{11} & \curvearrowright & x_i \omega_i & \\ x_{21} & & x_{22} & \\ x_{31} & & x_{32} & x_{33} \\ \vdots & & \dots & \end{array}$$

where the RVs in each row are independent, have expected value zero, and have finite variances.

We consider each element to be an element of $x_i \omega_i$. We can see that the conditions hold as $x_i \omega_i \perp x_j \omega_j$ since ω_i iid. and x_i, x_j are assumed to be fixed, $\mathbb{E}(x_i \omega_i) = x_i \mathbb{E}(\omega_i) = 0$ and $\text{Var}(x_i \omega_i) = \sigma^2 x_i x_i^T < \infty$ by assumption.

Thus, $x_i \omega_i$ is a triangular array with

- row-wise means: $z_n = \frac{1}{n} \sum x_i \omega_i = \frac{1}{n} X^T \omega$.
- row-wise average variance:

$$V_n = \frac{1}{n} \sum_i \text{Var}(x_i \omega_i) = \frac{1}{n} \sum \sigma^2 x_i x_i^T = \frac{\sigma^2}{n} X^T X \rightarrow \sigma^2 \Sigma \quad (3.5)$$

where convergence follows by assumption.

We now turn to verify the LF-condition. Let $\varepsilon > 0$, recall that we assume existence of M such that $\|x_i\| < M$, $\forall i$.

$$\begin{aligned} & \frac{1}{n} \sum_i E \left\{ \|x_i \omega_i\|^2 \mathbf{1} \left(\|x_i \omega_i\| \geq \varepsilon \sqrt{n} \right) \right\} \\ &= \frac{1}{n} \sum_i E \left\{ \omega_i^2 \|x_i\|^2 \mathbf{1} \left(|\omega_i| \|x_i\| \geq \varepsilon \sqrt{n} \right) \right\} \\ &< \frac{1}{n} \sum_i E \left\{ \omega_i^2 M^2 \mathbf{1} \left(|\omega_i| M \geq \varepsilon \sqrt{n} \right) \right\} \\ &= E \left\{ \omega^2 M^2 \mathbf{1}(|\omega| M \geq \varepsilon \sqrt{n}) \right\} \\ &=: E \{T_n\} \end{aligned}$$

Note that $T_n \xrightarrow{p} 0$ as $\varepsilon \sqrt{n} M^{-1} \rightarrow \infty$ as ω are iid with finite variance.⁴

Define $Z = M^2 \omega^2$, $\mathbb{E}[Z] = M^2 \sigma^2 < \infty$. Note that $\forall n \in \mathbb{N} : T_n \leq Z$. Since T_n converges pointwise to the zero function and we have an integrable function which dominates T_n we can apply the *Dominated convergence theorem* to deduce convergence of $\mathbb{E}(T_n)$:

$$T_n \xrightarrow{p} 0 \implies T_n \xrightarrow{d} 0 \implies \mathbb{E}[T_n] \rightarrow 0$$

⁴If ω_i are iid with finite second moment, we have from Chebyshev's inequality that

$$\mathbb{E} \left[\omega_i^2 \mathbf{1}_{[|\omega_i| > d]} \right] = \mathbb{E} \left[\omega_1^2 \mathbf{1}_{[|\omega_1| > d]} \right] \xrightarrow{d \rightarrow \infty} 0.$$

3: Putting it all together

Recall:

$$\begin{aligned}\hat{\beta} &= \beta + (X^T X)^{-1} X^T \omega \\ z_n &= \frac{1}{n} \sum x_i \omega_i = \frac{1}{n} X^T \omega \\ V_n &= \frac{1}{n} \sum \text{Var}(x_i \omega_i) \rightarrow \sigma^2 \Sigma\end{aligned}$$

Then the LF-CLT says

$$\sqrt{n} z_n \xrightarrow{d} \mathcal{N}(0, V)$$

Hence,

$$\frac{1}{\sqrt{n}} X^T \omega \xrightarrow{d} \mathcal{N}(0, \sigma^2 \Sigma)$$

We know that $(X^T X)(\hat{\beta} - \beta) = X^T \omega$, thus (dividing by the square root of the variance from (3.5)), we normalise and obtain

$$\frac{1}{\sigma} (X^T X)^{-1/2} (X^T X)(\hat{\beta} - \beta) = \frac{1}{\sigma} (X^T X)^{1/2} (\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(0, \mathbb{I})$$

□

3.3 Testing

We next prove the distribution of various potential test statistics. Note that we always assume $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 \mathbb{I}_{n \times n})$. In particular, p-values in summary outputs in R rely on this normality assumption. If instead only Gauss-Markov assumptions hold, p values and confidence intervals are only valid asymptotically.

Confidence interval for β

Proposition 3.2

Assume $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 \mathbb{I}_{n \times n})$ then $\hat{\beta} \sim N_p(\beta, \sigma^2 (X^T X)^{-1})$. Further, let $\hat{\sigma}^2 = \frac{\sum \hat{\varepsilon}_i^2}{n-p}$. We have $\frac{\sum \hat{\varepsilon}_i^2}{\sigma^2} \sim \chi_{n-p}^2$ and $\hat{\sigma} \perp \hat{\beta}$. Then:

$$\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma} \sqrt{\left((X^T X)^{-1}\right)_{ii}}} \sim t_{n-p}$$

Proof. Recall, by definition of the t-distribution that if $Z \sim \mathcal{N}(0, 1)$, $V \sim \chi_n^2$, $Z \perp V$, then

$$T = \frac{Z}{\sqrt{V/n}} \sim t_n$$

Since $\hat{\beta}_i \sim \mathcal{N}(\beta_i, \sigma^2 ((X^T X)^{-1})_{ii})$, we have

$$Z = \frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{((X^T X)^{-1})_{ii}}} \sim \mathcal{N}(0, 1)$$

$$V = \frac{\sum \hat{\varepsilon}_i^2}{\sigma^2} \sim \chi_{n-p}^2 \text{ and } V \perp Z$$

Thus

$$T = \frac{\frac{\hat{\beta}_i - \beta_i}{\sigma \sqrt{((X^T X)^{-1})_{ii}}}}{\sqrt{\frac{\sum \hat{\varepsilon}_i^2}{\sigma^2 (n-p)}}} = \frac{\hat{\beta}_i - \beta_i}{\hat{\sigma} \sqrt{((X^T X)^{-1})_{ii}}} \sim t_{n-p}.$$

Here we used $\sum \hat{\varepsilon}_i^2 / (n-p) = \hat{\sigma}^2$. □

Using our typical ‘Wald-type-CI’, we can achieve a α level CI by rewriting

$$\hat{\beta}_i - q_{t_{\alpha/2, n-p}} \hat{\sigma}_{\hat{\beta}_i} \leq \beta_i \leq \hat{\beta}_i + q_{t_{\alpha/2, n-p}} \hat{\sigma}_{\hat{\beta}_i}$$

where $\hat{\sigma}_{\hat{\beta}_i} = \hat{\sigma} \sqrt{((X^T X)^{-1})_{ii}}$.

Confidence Interval for $\mathbb{E}(Y)$

The confidence interval for the expected value of Y (i.e., confidence interval for the level of the regression line at a point), is entirely similar to above. Noting that

$$\hat{Y}_i \sim \mathcal{N}(\mu_i, \sigma^2 H_{ii})$$

where H is the matrix used in $\hat{y} = Hy$, we obtain the test statistic

$$T = \frac{\hat{y}_i - \mu_i}{\hat{\sigma} \sqrt{H_{ii}}} \sim t_{n-p}.$$

The proof of the distribution follows our approach above.

Prediction interval

If we want to have a confidence interval for the predicted interval of accuracy single observation, i.e., a **prediction interval**, we need to consider not only the source of uncertainty from a wrong estimate of β but also the underlying noise in the data. Thus, prediction intervals are larger than the corresponding CI for $\mathbb{E}(\hat{Y})$.

$y_0 \perp \hat{y}_0$, since we consider a new observation $y_0 = x_0^T \beta + \varepsilon_0$ which was not used in the estimation of $\hat{\beta}$.

Consider $D = y_0 - \hat{y}_0 \sim \mathcal{N}(\mu_D, \sigma_D^2)$ where normality follows as this is a difference between two independent normal variables. By unbiasedness and $\mathbb{E}(\varepsilon) = 0$, $\mu_D = 0$. By independence⁵

$$\begin{aligned} \sigma_D^2 &= \text{Var}(y_0) + \text{Var}(\hat{y}_0) \\ &= \sigma^2(1 + x_0^T (X^T X)^{-1} x_0) \end{aligned}$$

⁵ As $n \rightarrow \infty$, the second term goes to zero while the first term will always remain non-zero.

Continuing as in our derivation above, we arrive at the test statistic

$$T = \frac{y_0 - \hat{y}_0}{\hat{\sigma} \sqrt{1 + x_0^T (X^T X)^{-1} x_0}} \sim t_{n-p}$$

Note that a prediction interval is interested in *one new observation*. If we want to arrive at an interval that contains $P\%$ of *all future observations* (with a certain level of confidence), we must instead consider a **tolerance interval**.

3.4 F-Test

Colinearity

Highly colinear variables lead to two problems:

- Interpretation is difficult since these variables tend to move together. Holding one fixed is not a natural phenomenon.
- Reduced accuracy in parameter estimates

We may measure colinearity using **Variance Inflation Factor (VIF)**:

$$\text{VIF} = \frac{1}{1-R_j^2} \text{ where } R_j^2 \text{ is the } R^2 \text{ value of the regression } x_j \sim x_{-j}$$

I.e., how well can x_j be explained by the other covariates? In practice, t-tests tend not to be significant but become significant when other variables are dropped. This will motivate the use of tests for whether at least one of multiple coefficients is non-zero.

Recap: Different likelihood-based tests

Recall that we can perform inference test based on the likelihood function in three different ways:

- *Wald type*: differences between $\hat{\alpha}$ and α_0 (closely related to t-test)⁶
- *LR*: Difference in likelihood $\log L(\hat{\alpha})$ and $\log L(\alpha_0)$ (closely related to F-test)⁷
- *Score*: Slope of derivative of likelihood w.r.t. α_0

⁶Tests for individual coefficients. *t*-Test uses a t-distribution, square root of Wald test uses asymptotically normal distribution. Equivalent as $n \rightarrow \infty$.

⁷LR uses ratio of likelihood functions, F-test uses ratio of RSS. It turns out that for linear regression, the (partial) F-test is equivalent to the likelihood ratio test as $\text{RSS} \propto \text{likelihood}$.

General form of F-tests

We focus on F-Tests. We are interest in asking questions such as are two coefficients identical, or are multiple coefficients equal to 0? These questions take the general form:⁸

$$H_0 : B\beta = b \text{ for some } ((p - q) \times p) \text{-matrix } B \text{ and some } (p - q)\text{-vector } b \text{ (often } b = 0)$$

We can perform one **partial F-Test** or multiple t-tests (one per row, with multiple testing adjustments). Partial F-tests compare nested models whereby the large model consists of p parameters, thh small of q parameters and therea are $p - q$ side constraints.

Example 3.1

Suppose we want to test the null hypothesis $\beta_3, \beta_4, \beta_5 = 0$. Then we may write this in matrix form as follows:

$$\begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \\ \beta_5 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

Example 3.2 (Global F-Test)

We test the null hypothesis that all coefficients other than the intercept are 0. We have

$$B = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix}; \quad \mathbf{b} = \mathbf{0}$$

The test statistic is then $\sim F_{p-1, n-p}$ since we have $q = 1$ (intercept in restricted models).

As under the null, the simple model always predicts \bar{Y} , we can rewrite this test statistic using the relation between the F-test and SSE as follows:

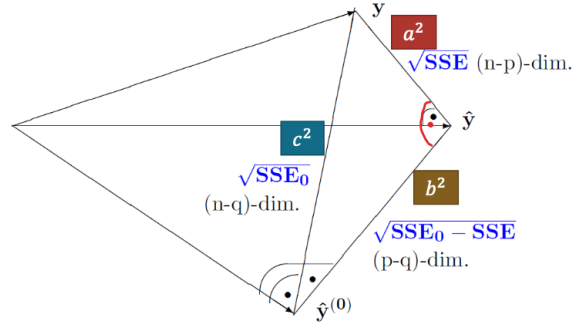
$$\begin{aligned} F &= \frac{n-p}{k} \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum \hat{\varepsilon}_i^2} \\ &= \frac{n-p}{k} \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2 - \sum (\hat{y}_i - \bar{y})^2} \\ &= \frac{n-p}{k} \frac{R^2}{1 - R^2}. \end{aligned}$$

Example 3.3 (Exact F-Test for specific value of β)

If we specify a complete set of values for the vector β under the null hypothesis, our

⁸A particularly common F-test is the global F-test which tests whether at least one coefficient other than the intercept is different from 0 (vs all coefficients being zero).

Figure 4: Graphical Intuition for F-Test



test statistic reduces to

$$\frac{(\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta)}{p\hat{\sigma}^2} \sim F_{p,n-p}$$

Intuition Partial F-Test

Suppose we have a point y and we first project it onto a p -dimensional subspace to obtain \hat{y} . Next consider adding additional side-constraints ($p - q$ such constraints), and projecting onto that subspace to obtain \hat{y}^0 . Note that \hat{y}^0 is also in the column space of X . The line connecting \hat{y}^0 and \hat{y} must thus be orthogonal to the line from y to \hat{y} . We therefore have a right-angled triangle in Figure 4 and can use Pythagoras to obtain a relation between the different line lengths. An intuitive test statistic would be to say that the simple model is true if b^2 is small relative to a^2 .

Under the assumption that the simple model is true, we have

$$\frac{(SSE_0 - SSE)/(p - q)}{SSE/(n - p)} = \frac{\|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}^{(0)}\|^2/(p - q)}{\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2/(n - p)} \sim F_{p-q, n-p}$$

To show this, we use the following lemma which can be obtained using Lagrange multipliers:

Lemma 3.1

The least squares estimator $\hat{\beta}_{(0)}$ under the supplementary condition $B\beta = \mathbf{b}$ is

$$\hat{\beta}_{(0)} = \hat{\beta} - (X^T X)^{-1} B^T \left(B (X^T X)^{-1} B^T \right)^{-1} (B\hat{\beta} - \mathbf{b})$$

Furthermore,^a

$$\underbrace{SSE_0}_{c^2} = \underbrace{SSE}_{a^2} + \underbrace{(B\hat{\beta} - \mathbf{b})^T \left(B (X^T X)^{-1} B^T \right)^{-1} (B\hat{\beta} - \mathbf{b})}_{b^2}.$$

^aTo obtain this, one may use

$$b^2 = \|\hat{Y} - \hat{Y}_0\|^2 = (X(\hat{\beta} - \hat{\beta}_0))^T X(\hat{\beta} - \hat{\beta}_0) = \dots$$

We use this, and the fact $\hat{\sigma}^2 = \epsilon^T \epsilon / (n - p)$, to rewrite:

$$\frac{(SSE_0 - SSE) / (p - q)}{SSE / (n - p)} = \frac{(B\hat{\beta} - \mathbf{b})^T (B(X^T X)^{-1} B^T)^{-1} (B\hat{\beta} - \mathbf{b})}{(p - q)\hat{\sigma}^2} \sim F_{p-q, n-p}$$

One can show that this is a likelihood ratio test. We now show that this follows an F -distribution if $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{1})$.

Proof. Recall that if $Y \sim \mathcal{N}(\mu, \Sigma)$, then $(y - \mu)^T \Sigma^{-1} (y - \mu) \sim \chi_n^2$. Also, $X = \frac{S_1/d_1}{S_2/d_2} \sim F_{d_1, d_2}$ if $S_i \sim \chi_{d_i}^2$ for $i = 1, 2$ and $S_1 \perp S_2$.

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

$$B\hat{\beta} \sim N(\underbrace{B\beta}_{=b}, \sigma^2 B(X^T X)^{-1} B^T)$$

Thus

$$S_1 := (B\hat{\beta} - b)^T (\sigma^2 B(X^T X)^{-1} B^T)^{-1} (B\hat{\beta} - b) \sim \chi_{p-q}^2.$$

We obtain the degrees of freedom by noting that b is a $(p - q)$ -dimensional vector.

We have previously shown in Proposition 3.1 (which relies on $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{1})$) that

$$\frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{\sigma^2} \sim \chi_{n-p}^2$$

A simple rearrangement, using $\hat{\sigma}^2 = \sum \hat{\epsilon}_i^2 / (n - p)$ yields

$$S_2 := \frac{(n - p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2.$$

Since $\hat{\epsilon} \perp \hat{\beta}$, we have $S_1 \perp S_2$.

Finally, after cancelling out some terms,

$$\frac{S_1 / (p - q)}{S_2 / (n - p)} = \frac{(B\hat{\beta} - b)^T (B(X^T X)^{-1} B^T)^{-1} (B\hat{\beta} - b)}{(p - q)\hat{\sigma}^2} \sim F_{p-q, n-p}.$$

□

In R: Fit two nested models (e.g. fm and fm2) and compare with anova(fm, fm2).

3.5 Residual Analysis

Recall that errors are not identical to residuals. In particular,

- **Errors:** $\varepsilon \sim N(0, \sigma^2 \mathbb{I}) \rightarrow$ Errors are uncorrelated and have constant variance
- **Residuals:** $\hat{\varepsilon} \sim N(0, \sigma^2 M) \rightarrow$ Residuals are correlated and have different variance $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - H_{ii})$
- By design: \hat{y} and $\hat{\varepsilon}$ are uncorrelated
- We sometimes instead consider standardized residuals: $\hat{\varepsilon}_i^S = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1-H_{ii}}} \rightarrow$ unit variance if the true model is correct.⁹

There are four main assumptions to check:

1. Independent samples:
 - Plot residuals vs recording time or order of observation: is there serial correlation or clustering?
2. Functional relationship and constant variance
 - Scatter plot $x - y$ in simple regression
 - Tukey-Anscombe plot (residual vs fitted value $\hat{\varepsilon} - \hat{y}$)¹⁰
 - Scale Location plot (standardised residual vs fitted value with horizontal smoother), see Figure 5
3. Normal distribution of errors
 - QQ-Plot of residuals

Outliers

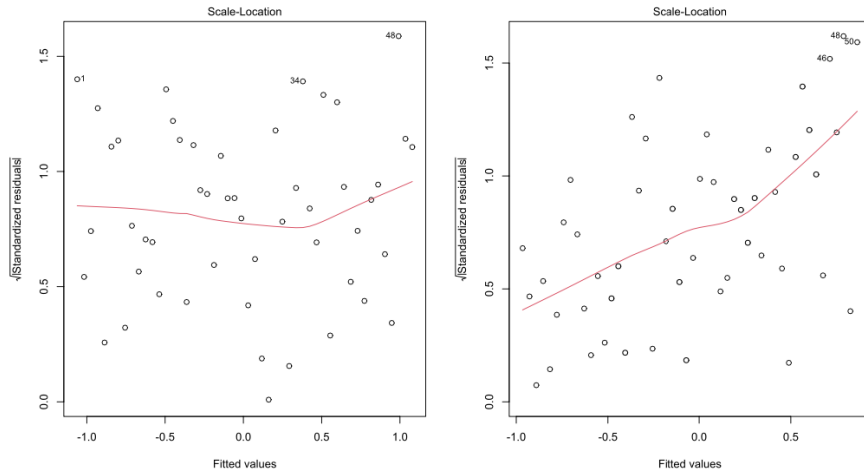
Outliers can occur in the y or x direction. Some ways of measuring outliers are:

- **Leverage H_{ii} :** How influential is y_i on prediction \hat{y}_i ?
 - $\hat{Y} = HY \rightarrow \frac{d\hat{y}_i}{dy_i} = H_{ii}$ depends on X but not on y
- **Cook's Distance D_i :** How influential is y_i on whole fit?
 - Note: $\hat{y}_{j(i)}$ is the fitted response value for observation j obtained when excluding observation i

⁹If for instance $\text{Var}(\varepsilon) = \Sigma \neq \sigma^2 \mathbb{I}$, then the standardization will not yield a unit variance since we will not have $\hat{\varepsilon} \sim N(0, \sigma^2 M)$. We can spot this in adequate plots.

¹⁰Note: cannot be too picky with varying variances since residuals do not have constant variances.

Figure 5: Examples of scale-location plots. The left plot is okay, the right hand plot indicates problematic violations of assumptions as the variance seems to increase with the value of \hat{y}



- $D_i = \sum_{j=1}^n \frac{(\hat{y}_j - \hat{y}_{j(i)})^2}{p\hat{\sigma}^2}$ can be transformed to $D_i = \frac{\hat{\epsilon}_i^2}{p\hat{\sigma}^2} \cdot \left(\frac{H_{ii}}{(1-H_{ii})^2} \right)$ ¹¹
- Rule of thumb: $D_i > 1$ is problematic

We can investigate outliers in standard residual plots (e.g., residual vs leverage).

4 Model Selection

If we aim at **inference**, we should postulate one model as this allows for valid p-values.

If we instead explore multiple models, p-values etc. are no longer valid, we have a multiple testing issue. This is a research topic, post-selection inference.

Consequences of too many or too few variables in the model may be:

1. *Too few* (sparse model): Bias, but reduced variance¹²
2. *Too many*: Unbiased, large variance

In particular, overfitting is an issue which can reduce prediction accuracy when considering bigger models. This is the fundamental issue behind the bias-variance trade-off. We have

$$\text{MSE} = \text{Var} + \text{Bias}^2.$$

¹¹Similarly, there is a transformation which can be applied to calculate LOOCV errors: $CV = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_{i(i)}^2$ can be transformed to $CV = \frac{1}{n} \sum_{i=1}^n \left(\frac{\hat{\epsilon}_i}{1-H_{ii}} \right)^2$

¹²The model is unbiased if the missing variables are not correlated with the included Xs. E.g., if we have a completely orthogonal design matrix.

For this, note that when z is an estimator of the unknown constant c ,

$$\begin{aligned} E[(z - c)^2] &= E[(z - E[z]) + (E[z] - c)^2] \\ &= \underbrace{E[(z - E(z))^2]}_{\text{Var}(z)} + 2E[\underbrace{(z - E[z])}_{\rightarrow 0} \cdot \underbrace{(E[z] - c)}_{\text{constant}}] + \underbrace{E[(E[z] - c)^2]}_{\text{bias}^2} \end{aligned}$$

Note that $E[z] - c$ is constant since we take expectation with respect to z , so c is a constant and $E(E(z))$ is also just the expectation of a constant real number.

4.1 Approaches for model selection

Local approaches using F and t-tests

Search strategies based on partial F-Test:

- Stepwise forward
- Stepwise backward (only possible if $n > p$)
- Add "most significant" / remove "least significant" variable according to partial F-test

Drawbacks:

- p-values are used but not strictly valid anymore
- Only nested models can be compared

Global Approaches using Cp and AIC

Often we do not want to use data splitting or cross-validation due to computational concerns. In this case, we can **estimate Test MSE** by using a formula that approximates the additional error in test MSE out-of-sample. That is, we want to minimize a quantity like

$$\min RSS + c \cdot p$$

where we have a penalty for complex models.

- Mallows's C_p

$$C_p = \frac{RSS}{\hat{\sigma}^2} + 2|M| - n$$

- Akaike's Information Criterion:

$$AIC = -2 \cdot l(\hat{\theta}_M) + 2 \cdot p$$

- Bayesian Information Criterion:

$$BIC = -2 \cdot l(\hat{\theta}_M) + \log(n) \cdot p$$

The last two criteria are more general than Mallows's C_p as they are valid for all models which admit a likelihood.

4.2 Derivation of Mallow's C_p

We want to derive Mallow's C_p .

Setup:

- Assume $y_i, i = 1, \dots, n$ independent.
- $\mathbb{E}[y_i] = \mu_i, \text{Var}(y_i) = \sigma^2$
- Regressors $1, x_1, \dots, x_k \implies X$ is a $n \times (k+1)$ design matrix
- Consider a subset $M \subseteq \{0, 1, \dots, k\}$, denote the design matrix of this reduced set of regressors by X_M which is $n \times |M|$

Then

$$\hat{\beta}_M = (x_M^T X_M)^{-1} X_M^T y, \quad \hat{y}_M = X_M \hat{\beta}_M = H_M y$$

Facts:

1. $E(\hat{Y}_M) = H_M E(y) = X_M (X_M^T X_M)^{-1} X_M^T X \beta$.
 - Note that we may have bias as the X_M and X term do not cancel each other out
2. $\text{Cov}(\hat{y}_M) = \sigma^2 H_M$
3. $\sum_{i=1}^n \text{Var}(\hat{Y}_{iM}) = \text{tr}(\sigma^2 H_M) = \sigma^2 \text{tr}(X_M (X_M^T X_M)^{-1} X_M^T) = \sigma^2 |M|$
4. Sum of mean squared errors:¹³
Use $\mu_{iM} = \mathbb{E}[\hat{y}_{iM}]$. Then

$$\begin{aligned} \text{SMSE} &= \sum_{i=1}^n \mathbb{E}[(\hat{y}_{iM} - \mu_i)^2] \\ &= \sum_{i=1}^n E((\hat{y}_{iM} - \mu_{iM}) + (\mu_{iM} - \mu_i))^2 \\ &= \sum \text{Var}(\hat{y}_{iM}) + 2 \sum E(\underbrace{(\hat{y}_{iM} - \mu_{iM})}_{\rightarrow 0} (\mu_{iM} - \mu_i)) \\ &\quad + \sum \underbrace{E(\mu_{iM} - \mu_i)^2}_{\text{constant, can remove } E} = \\ &= \sigma^2 |M| + \sum (\mu_{iM} - \mu_i)^2 \end{aligned} \tag{4.1}$$

We can then define

$$\gamma(M) = \frac{\text{SMSE}}{\sigma^2} = |M| + \frac{\text{bias}^2}{\sigma^2} \tag{4.2}$$

We note that $\gamma(M) \geq M$ with equality if the model is unbiased.

¹³Note, this is a bit of a misnomer. Instead it is the sum of expected squared deviations between predictions and expected value for y_i .

5. Future observations:

Assume that we record again observations y for all $i = 1, \dots, n$ with the same regressors x_{i1}, \dots, x_{ik} . Then

$$Y_{n+i} = \mu_i + \underbrace{\varepsilon_{n+i}}_{\text{independent of } \varepsilon_1, \dots, \varepsilon_n}$$

Since regressors are identical, our prediction $\hat{Y}_{n+i} = \hat{y}_n$.

We define the *expected squared prediction error (SPSE)*¹⁴

$$\begin{aligned} \text{SPSE} &= \sum_{i=1}^n \mathbb{E} \left[(Y_{n+i} - \hat{Y}_{iM})^2 \right] \\ &= \sum E \left((Y_{n+i} - \mu_i) - (\hat{Y}_{iM} - \mu_i) \right)^2 \\ &= \sum (E(Y_{n+i} - \mu_i)^2 - 2E \underbrace{(Y_{n+i} - \mu_i)(\hat{Y}_{iM} - \mu_i)}_{\rightarrow 0} + E(\hat{Y}_{iM} - \mu_i)^2) \quad (4.3) \\ &= n\sigma^2 + \text{SMSE} \\ &= \underbrace{n\sigma^2}_{(1)} + \underbrace{|M|\sigma^2}_{(2)} + \underbrace{\sum (\mu_{iM} - \mu_i)^2}_{(3)} \end{aligned}$$

Note in the third line that y_{n+i} and \hat{y}_{iM} are independent. Therefore the expectation of the product is the product of expectations and we know that $E(Y_{n+i} - \mu_i) = 0$.

We have the following three elements: (1) Irreducible error, (2) Variance, (3) Bias.

Before deriving Mallows's C_p , we show the following relation:

Lemma 4.1

$$\mathbb{E}(\text{RSS}) = \text{SPSE} - 2|M|\sigma^2$$

For this, note that $\mathbb{E}(\text{RSS})$ is the sum of expected squared residuals in sample (as opposed to SPSE which considers out of sample residuals).

Proof.

$$\begin{aligned} E(\text{RSS}) &= E \left(\sum (y_i - \hat{y}_i^M)^2 \right) = \sum E \left((y_i - \hat{y}_i^M)^2 \right) \\ &= \sum \underbrace{\text{Var} (y_i - \hat{y}_i^M)}_{(1)} + \sum \underbrace{\left(E (y_i - \hat{y}_i^M) \right)^2}_{(2)} \end{aligned}$$

which follows from the usual formula for the variance, $\text{Var}(X) = \mathbb{E}(X^2) - [\mathbb{E}(X)]^2$.

ad (1):

¹⁴Note that we only consider squared errors on new observations.

$$y - \hat{y}^M = y - H^M y = (\mathbf{1} - H^M) y$$

$$\text{Cov}(\hat{y} - \hat{y}^M) = (\mathbf{1} - H^M) \sigma^2$$

Therefore,

$$\begin{aligned} \sum \text{Var}(y_i - \hat{y}_i^M) &= \sigma^2 \text{tr}(\mathbf{1} - H^M) \\ &= \sigma^2 (n - \text{tr}(H^M)) = \sigma^2 (n - |M|) \end{aligned}$$

ad (2):

$$\begin{aligned} \sum (E(y_i - \hat{y}_i^M))^2 &= \sum (E(\hat{y}_i^M) - \mu_i)^2 = \\ &= \text{SMSE} - |M| \sigma^2 \\ &= \text{SPSE} - n \sigma^2 - |M| \sigma^2 \end{aligned}$$

In the first step we use $\mathbb{E}(y_i) = \mu_i$. Then we used equations (4.1). In particular, we use $E(\hat{y}_{iM}) = \mu_{iM}$ and the last line in that block of equations. Finally, we use (4.3) to get the result.

Continuing with what we want to show, we have

$$\begin{aligned} \mathbb{E}[RSS] &= (n - |M|) \sigma^2 + \text{SPSE} - n \sigma^2 - |M| \sigma^2 = \\ &= \text{SPSE} - 2|M| \sigma^2 \end{aligned}$$

□

Mallow's C_p now wants to estimate $\gamma(M) = \frac{\text{SMSE}}{\sigma^2}$ as defined in (4.2). From the above,

$$\begin{aligned} \text{SPSE} &= n \sigma^2 + \text{SMSE} = E(RSS) + 2|M| \sigma^2 \\ \Rightarrow \text{SMSE} &= E(RSS) + 2|M| \sigma^2 - n \cdot \sigma^2 \end{aligned}$$

We now replace unobserved aspects with estimates quantities,

$$C_p = \frac{RSS}{\hat{\sigma}^2} + 2|M| - n$$

4.3 Relation between C_p and AIC for Linear Regression

In linear regression,¹⁵

$$\text{AIC} = -2 \cdot l(\hat{\beta}_M, \hat{\sigma}^2) + 2(|M| + 1)$$

The term of the log-likelihood is given by

$$\begin{aligned} -2l(\hat{\beta}_M, \hat{\sigma}^2) &= n \log(\hat{\sigma}^2) + \frac{1}{\hat{\sigma}^2} (\mathbf{y} - \mathbf{X}_M \hat{\beta}_M)' (\mathbf{y} - \mathbf{X}_M \hat{\beta}_M) \\ &= n \log(\hat{\sigma}^2) + \frac{n \hat{\sigma}^2}{\hat{\sigma}^2} \\ &= n \log(\hat{\sigma}^2) + n. \end{aligned}$$

¹⁵Note that we have $|M| + 1$ parameters since we estimate $|M|$ coefficients and the residual standard error.

We get

$$AIC = n \cdot \log(\hat{\sigma}^2) + 2(|M| + 1) + n \quad (4.4)$$

This looks quite different to C_p , but we can show similarity by applying a Taylor expansion of $\log \hat{\sigma}^2$ around σ^2 to obtain

$$\log \hat{\sigma}^2 \approx \log \sigma^2 + \frac{1}{\sigma^2} (\hat{\sigma}^2 - \sigma^2) = \log \sigma^2 + \frac{SSE}{n\sigma^2} - 1$$

Thus, ignoring the constant $2 + n$ in (4.4),

$$AIC \approx n \log(\sigma^2) + \underbrace{\frac{SSE}{\sigma^2} - n + 2|M|}_{\approx C_p}$$

4.4 General Search Strategies

• Exhaustive Search

- For given $p = 1, \dots, p_{\max}$: Find best model according to RSS¹⁶
- We obtain p_{\max} different models
- We then choose the best one with C_p or AIC,
 - * Pros: Finds global optimum
 - * Cons: Computational cost p variables $\rightarrow 2^p$ subsets

• Forward / Backward Selection

- Forward: Start with empty model, add one variable at a time
- Backward: Start with full model, remove one variable at a time
- E.g., in forward selection we calculate C_p for each potential $p' + 1$ variable model with first p variables fixed after having found best p' model. If there is no model with a better C_p value, we stop.¹⁷
 - * Pros: Fast
 - * Cons: Only local optimum \rightarrow Forward and Backward often get to different final models

In R:

```
regsubsets(y ~ ., data, method = "exhaustive", nvmax = 15)
```

¹⁶We do not require an information criteria for a fixed number of parameters p since in this case the model with lowest RSS automatically has the lowest value for the information criteria.

¹⁷The R implementation in `regsubsets` finds the best model according to RSS for any potential number of variables and only in the end chooses according to information criteria between models with different number of parameters.

5 Non i.i.d. errors

Assume $\text{Cov}(\varepsilon) = \sigma^2 \mathbb{I}$. Then we can show that

$$\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

Now, instead assume $\text{Cov}(\varepsilon) = \sigma^2 W^{-1}$. Then

$$\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} X^T W^{-1} X (X^T X)^{-1} \neq \sigma^2 (X^T X)^{-1}$$

\implies Using 'normal' OLS leads to wrong p-values and CIs.

Often times we can detect deviations from $\text{Cov}(\varepsilon) = \sigma^2 \mathbb{I}$ by investigating TA-plots or from context knowledge (e.g., grouped measurements). In particular if we have additional information about the form of the covariance matrix, we can often use such information to estimate the covariance matrix and obtain correct CIs.

5.1 Known covariance matrix – Generalised Least Squares

Suppose we know the covariance matrix of our errors up to a factor. E.g.,

$$\mathbf{Y} = X\beta + \varepsilon \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma)$$

where Σ is known and PD. By Proposition 1.11 (and due to PD, thus invertibility), there exists an invertible matrix A such that

$$\Sigma^{-1} = (AA^T)^{-1}$$

We can thus formulate our OLS problem in a 'tilde' dimension where we use

$$\tilde{\mathbf{Y}} := A^{-1}\mathbf{Y} = A^{-1}(X\beta + \varepsilon) = \underbrace{A^{-1}X}_{\tilde{X}}\beta + \underbrace{A^{-1}\varepsilon}_{\tilde{\varepsilon}} = \tilde{X}\beta + \tilde{\varepsilon}$$

We have

$$\begin{aligned} \mathbb{E}\tilde{\varepsilon} &= \mathbb{E}A^{-1}\varepsilon = A^{-1}\mathbb{E}\varepsilon = \mathbf{0} \\ \text{Cov}[\tilde{\varepsilon}] &= \text{Cov}[A^{-1}\varepsilon] = A^{-1} \text{Cov}[\varepsilon] (A^{-1})^T \\ &= A^{-1} \sigma^2 (AA^T) (A^{-1})^T = \sigma^2 I. \end{aligned}$$

We can thus apply our usual OLS results to estimate $\hat{\beta}$ from the 'tilde model'

$$\|\tilde{\mathbf{Y}} - \tilde{X}\beta\|^2 = (\mathbf{Y} - X\beta)^T A^{-T} A^{-1} (\mathbf{Y} - X\beta) = (\mathbf{Y} - X\beta)^T \Sigma^{-1} (\mathbf{Y} - X\beta).$$

We can think of this as rescaling our residuals such that observations with larger variance (more uncertainty) receive a lower weight. This will be more obvious in the special case where Σ is diagonal which we treat below.

Solving the above for $\hat{\beta}$ gives us our usual result in the tilde dimension:

$$\hat{\beta} = \left(\tilde{X}^T \tilde{X} \right)^{-1} \tilde{X}^T \tilde{\mathbf{Y}} = \left(X^T \Sigma^{-1} X \right)^{-1} X^T \Sigma^{-1} \mathbf{Y}.$$

In the last step we plug back in our definitions of $\tilde{X}, \tilde{\mathbf{Y}}$. It is then easy to see that

$$\hat{\beta} \sim \mathcal{N}_p \left(\beta, \sigma^2 \left(X^T \Sigma^{-1} X \right)^{-1} \right), \quad \hat{\sigma}^2 = \frac{1}{n-p} \varepsilon^T \Sigma^{-1} \varepsilon$$

If we know that each y_i is an average of w_i -many i.i.d. measurements, we obtain the special case of **weighted least squares**.

$$\text{Var}(\varepsilon_i) = \frac{\sigma^2}{w_i}; \quad \text{Cov}(\varepsilon) = \sigma^2 \text{diag} \left(\frac{1}{w_1}, \dots, \frac{1}{w_n} \right) = \sigma^2 \Sigma$$

Thus

$$\begin{aligned} |\tilde{\mathbf{Y}} - \tilde{X}\beta|^2 &= \left[A^{-1}(\mathbf{Y} - X\beta) \right]^T A^{-1} [\mathbf{Y} - X\beta] \\ &= (\mathbf{Y} - X\beta)^T \Sigma^{-1} (\mathbf{Y} - X\beta) = \sum_{i=1}^n w_i \left(y_i - x_i^T \beta \right)^2 \end{aligned}$$

Each observation in our objective function is weighted proportionally to the **number of samples** or **inverse variance of error**. Note that WLS in general encapsulates cases where we have heteroskedastic but uncorrelated errors and know the heteroskedasticity structure. In this case, we have $\Sigma^{-1} = W^{1/2} W^{1/2}$ where $W = \text{diag}(w_1, \dots, w_n)$, so $A^{-1} = A^{-\top} = W^{1/2}$.

In R: `fm <- lm(y ~ x, data = df, weights = nreps), weights argument`

From the fact that we can rescale our model to a tilde dimension with diagonal covariance matrix, we can apply many of our previous results, in particular:

Proposition 5.1 Gauss-Markov for WLS

Among all linear and unbiased estimators $\hat{\beta}^L = \mathbf{A}\mathbf{y}$, the WLS estimator has minimal variance, i.e.,

$$\text{Var}(\hat{\beta}_j) \leq \text{Var}(\hat{\beta}_j^L), \quad j = 0, \dots, k$$

A similar extension holds with normally distributed errors and WLS being the minimal variance unbiased estimator.

Note that we could use our normal OLS estimate and the fact that

$$\text{Var}(\hat{\beta}_{\text{OLS}}) = \sigma^2 (X^T X)^{-1} X^T \Sigma X (X^T X)^{-1}$$

while plugging in Σ to obtain correct coverage probabilities. However, the above application of the Gauss-Markov theorem shows that instead using Generalized Least Squares leads to more efficient estimates.

5.2 Unknown covariance matrix, known structure

MLE and two-stage OLS

If our covariance matrix is unknown, but we can specify the structure, we can usually use this information to model the covariance matrix. There are two approaches:

1. Two-stage procedure
 - (a) Fit OLS
 - (b) Use OLS residual estimates and information about Covariance structure to model coefficients of covariance
 - (c) Use above to derive weights and fit GLS
 - (d) Iterate
 - (e) Example: Serial correlation of AR(1) form

$$\text{Cov}(\varepsilon) = \sigma^2 W^{-1} = \frac{\sigma_u^2}{1 - \rho^2} \begin{pmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & 1 & \rho & \dots & \rho^{n-2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \rho^{n-1} & \rho^{n-2} & \rho^{n-3} & \dots & 1 \end{pmatrix}$$

We estimate residuals using OLS, then estimate ρ_1, \dots, ρ_k using residuals and use ρ to define weights to be used in GLS.

2. Maximum-Likelihood
 - (a) Make assumption about distribution of ε and write out likelihood using known information about covariance matrix
 - (b) Estimate covariance matrix using MLE

Huber-White HC (Sandwich) Estimator

Consider the model

$$Y = X\beta + \varepsilon, \varepsilon \sim \mathcal{N}_n(0, D)$$
$$D = \text{diag}(\sigma_1^2, \dots, \sigma_n^2), \sigma_i^2 = \text{Var}(\varepsilon_i).$$

We know that

$$\hat{\beta} \sim \mathcal{N}_p\left(\beta, (X^T X)^{-1} X^T D X (X^T X)^{-1}\right).$$

For statistical testing or confidence intervals, we need to estimate the covariance matrix. An easy but powerful approach is to use

$$\hat{D} = \text{diag}(r_1^2, \dots, r_n^2)$$

where $r = Y - X\hat{\beta}$ are the residuals (from the least squares estimator). The estimated covariance matrix is then

$$\hat{V} = \widehat{\text{Cov}}(\hat{\beta}) = (X^T X)^{-1} X^T \hat{D} X (X^T X)^{-1}.$$

Note that it may seem implausible to estimate this, as we have n observation and the matrix \hat{D} consists of n parameters. However, we can show that we can estimate the term $X^T \hat{D} X$ consistently. In particular, since we know $V^{-1/2}(\hat{\beta} - \beta) \sim \mathcal{N}_p(0, I)$, we now want to show:

$$\hat{V}^{-1/2}(\hat{\beta} - \beta) \implies \mathcal{N}_p(0, I) \text{ as } n \rightarrow \infty$$

It follows from Slutsky's theorem (Theorem A.1) that this is the case if we can consistently estimate V .¹⁸ Thus, we want to show that $\hat{\beta} \xrightarrow{p} \beta$ as $n \rightarrow \infty$ and

$$n(\widehat{\text{Cov}}(\hat{\beta}) - \text{Cov}(\hat{\beta})) \rightarrow 0 \quad (5.1)$$

We multiply by n since Slutsky's theorem requires sufficiently fast convergence of $\hat{V} \rightarrow V$. Since $\text{Cov}(\hat{\beta})$ itself shrinks at n^{-1} , we must have convergence of $\hat{V} \rightarrow V$ at a faster rate which we ensure by checking whether $n(\hat{V} - V) \rightarrow 0$.

Lemma 5.1

Assume D diagonal and $Y = X\beta + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, D)$. Let $\hat{\beta} = (X^T X)^{-1} X^T y$. Estimate $\hat{D} = \text{diag}(\hat{\varepsilon}_1^2, \dots, \hat{\varepsilon}_n^2)$. Assume further that

1. $\frac{X^T X}{n} \rightarrow C$ as $n \rightarrow \infty$ where C is psd
2. $\sigma_i^2 < C_1 < \infty$ for all i , and $1/n \sum X_{ij}^2 \leq C_2 < \infty$ for all j

Then

$$\hat{\beta} \xrightarrow{p} \beta \text{ as } n \rightarrow \infty$$

Proof. We know that our estimated $\hat{\beta}$ is unbiased. We therefore only need to show that the covariance converges to 0.

We write

$$\text{Cov}(\hat{\beta}) = \frac{1}{n} \underbrace{\left(\frac{X^T X}{n} \right)}_{(1)}^{-1} \underbrace{\frac{X^T D X}{n}}_{(2)} \underbrace{\left(\frac{X^T X}{n} \right)^{-1}}_{(3)}.$$

We know that terms (1) and (3) are bounded by assumption 1. Further, term (2) is bounded by assumption 2 since any element of $\frac{X^T D X}{n}$ is weakly smaller than $C_1 \cdot C_2$. Therefore,

$$\text{Var}(\hat{\beta}_j) \rightarrow 0 \text{ as } n \rightarrow \infty, \forall j.$$

By Tchebshev's inequality,

$$\lim_{n \rightarrow \infty} P(|\hat{\beta}_j - \beta_j| > \varepsilon) = 0, \forall \varepsilon > 0$$

□

¹⁸Normality of $V^{-1/2}(\hat{\beta} - \beta)$ follows from our usual arguments of either the distribution of ε or a central limit theorem. The part of interest is now whether we can replace $V^{-1/2}$ with our estimated covariance matrix $\hat{V}^{-1/2}$.

Lemma 5.2

In addition, assume $\max_{i,j} |X_{ij}| \leq K < \infty$. Then

$$\frac{X^T \hat{D} X}{n} - \frac{X^T D X}{n} \xrightarrow{p} 0 \text{ as } n \rightarrow \infty$$

Proof. Consider an arbitrary elements (r, s) . We can write

$$\frac{1}{n} (X^T \hat{D} X)_{rs} = \frac{1}{n} \sum_i X_{ir} X_{is} \hat{\varepsilon}_i^2 \quad (5.2)$$

Then,

$$\begin{aligned} \frac{1}{n} (X^T \hat{D} X)_{rs} &= \frac{1}{n} \sum_i X_{ir} X_{is} \hat{\varepsilon}_i^2 \\ &= \frac{1}{n} \sum_i X_{ir} X_{is} (y_i - \hat{y}_i)^2 \\ &= \frac{1}{n} \sum_i X_{ir} X_{is} \left(\varepsilon_i - \underbrace{x_i^T (\hat{\beta} - \beta)}_{\rightarrow 0} \right)^2 \\ &= \frac{1}{n} \sum_i X_{ir} X_{is} (\varepsilon_i^2 + o_p(1)) \end{aligned}$$

In the second last line we use our above result that $\hat{\beta} \rightarrow \beta$ in probability as $n \rightarrow \infty$.

To show consistency, we first show asymptotic unbiasedness.

$$\mathbb{E} \left[\frac{1}{n} \sum_i X_{ir} X_{is} \varepsilon_i^2 \right] = \frac{1}{n} \sum_i X_{ir} X_{is} \sigma_i^2 = \frac{1}{n} (X^T D X)_{rs}$$

where we use (5.2) for population values (and in reverse).

Next, we show that the variance converges to 0 as $n \rightarrow \infty$,

$$\begin{aligned} \text{Var} \left(\frac{1}{n} \sum_i X_{ir} X_{is} \varepsilon_i^2 \right) &= \frac{1}{n} \frac{1}{n} \sum_{i=1}^n X_{ir}^2 X_{is}^2 \text{Var}(\varepsilon_i^2) \\ &= \frac{1}{n} \frac{1}{n} n \cdot \text{constant} \rightarrow 0 \text{ (} n \rightarrow p \text{)} \end{aligned}$$

For this, we may use that if $\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$ then $\text{Var}(\varepsilon_i^2) = 3\sigma_i^4 < \infty$.

Thus, $\frac{X^T \hat{D} X}{n} - \frac{X^T D X}{n} \xrightarrow{p} 0$ as $n \rightarrow \infty$. □

This allows us to proof (5.1) since

$$n \left[\widehat{\text{Cov}}(\hat{\beta}) - \text{Cov}(\hat{\beta}) \right] = \left(\frac{X^T X}{n} \right)^{-1} \left[\frac{X^T \hat{D} X}{n} - \frac{X^T D X}{n} \right] \left(\frac{X^T X}{n} \right)^{-1}$$

where the mid term converges to zero and the other terms converge to some constant matrix C .

What is the risk that we “badly” estimate at least one element of the matrix if we also allow the number of covariates to grow? We can upper-bound this error by $p^2 \sup \Pr(\text{error}_i)$. Note that the probability of an error goes down at rate n^{-1} , as shown. If the number of covariates p grows at a smaller rate, e.g., $\log N$, this error probability still converges to 0.

One fundamental application of the HC estimator is when using linear models to approximate complex models. Even if we have $Y_i = f(X_i) + \varepsilon_i$ with ε_i independent of X_i , and $\text{Var}(\varepsilon_i) = \sigma^2$ when we use a linear approximation, we will have a heteroskedastic error.

In R: `coeftest(fm, vcov = vcovHC(fm, type = "HC0"))` from sandwich package.

5.3 Mixed Models

Another approach of modelling more complex error structures is to use mixed models. Suppose we have i individuals or clusters, and n_i observations per individual. Coefficients are likely to be different across people. One approach would be to have a block factor:

$$y_{ij} = (\beta_0 + \beta_{0,i}) + \beta_1 x_{ij} + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma^2) \text{ i.i.d}$$

Here we only consider *fixed effects* and include a custom level-shift per cluster. With this we can easily make statements about individual intercepts but less easily derive statements about the underlying distribution.

A different approach uses so called random effects.

Random Intercept Model

$$\begin{aligned} y_{ij} &= (\beta_0 + u_i) + \beta_1 x_{ij} + \varepsilon_{ij}, \\ \varepsilon_{ij} &\sim N(0, \sigma^2), \quad u_i \sim N(0, \sigma_1^2) \text{ i.i.d} \\ \varepsilon_{ij} &\perp u_k, \forall i, j, k \end{aligned}$$

I.e., we have a cluster-specific random shift (u_i).

This is a **mixed model** since we consider both fixed and random effects. We can easily study statements about the dispersion of initial slopes in the population but not about individuals.

It turns out that we implicitly model the correlation among samples of the same cluster.

Within block i :

$$\begin{aligned} \text{Cov}(Y_{ij}, Y_{ik}) &= \text{Cov}(u_i + \varepsilon_{ij}, u_i + \varepsilon_{ik}) \\ &= \text{Cov}(u_i, u_i) + \text{Cov}(\varepsilon_{ij}, \varepsilon_{ik}) \\ &= \sigma_1^2 + \sigma^2 \mathbf{1}\{j = k\} \end{aligned}$$

Between blocks, we have zero covariance. Let $m \neq i$,

$$\text{Cov}(Y_{ij}, Y_{mk}) = \underbrace{\text{Cov}(\mu_i, \mu_m)}_{=0} + \underbrace{\text{Cov}(\varepsilon_{ij}, \varepsilon_{mk})}_{=0}$$

We therefore obtain a block-diagonal covariance matrix.

- Constant correlation structure
- Relative strength of within-subject correlation depends on relative magnitudes of σ^2 and σ_1^2

Random Intercept Random Slope Model

We can naturally extend this model to allow for varying slopes between clusters. We now not only consider a random shift in the intercept but also a random shift for the slope parameter.

$$\begin{aligned}
 y_{ij} &= (\beta_0 + u_{1,i}) + (\beta_1 + u_{2,i}) x_{ij} + \varepsilon_{ij}, \\
 \varepsilon_{ij} &\sim N(0, \sigma^2) \text{ i.i.d} \\
 u_{1,i} &\sim N(0, \sigma_1^2), \\
 u_{2,i} &\sim N(0, \sigma_2^2), \\
 \text{cor}(u_1, u_2) &= \rho
 \end{aligned}$$

We initially estimate models using MLE to compare models. To obtain the final model we use restricted MLE (RMLE), as this removes the biased estimation of standard deviations.¹⁹

Note: The standard deviation of our slope and intercept estimates will approach the population values as $n \rightarrow \infty$, however, the standard error of the estimates will go to 0.

In R: lmer function in lme4 library.

Comparing Fixed Effects and Random Effects

Note that the question of fixed vs random effects depends fundamentally on the question one wants to answer. If one is interested in the specific sample, fixed effects are more useful, while random effects allow inference about population attributes.

In terms of estimation, fixed effects are modelled by explicitly including a dummy variable for each cluster in the design matrix X and estimating coefficients for each cluster, e.g., α_i , through OLS. On the other hand, to estimate random effects we use MLE and specify that the random effects follow a normal distribution from which we then find $\hat{\sigma}_1^2$.²⁰ Since we estimate fewer parameters, this allows scaling to larger number of random effects.

¹⁹We initially use MLE due to better properties of MLE and allowing comparison of non-nested models.

²⁰Some libraries such as lme4 can still extract actual random effects by using a technique called *best linear unbiased prediction (BLUP)*. After estimating $\hat{\sigma}_1^2$, BLUPs are calculated for each random effect by taking conditional expectations of random effects given data and estimated coefficients. The formula in a RI model is given by:

$$\hat{b}_j = \frac{\sigma_1^2}{\sigma_1^2 + \sigma^2/n_j} (\bar{y}_j - \beta_0 - \beta_1 \bar{x}_j)$$

I.e., instead of estimating each coefficient directly, we estimate distribution parameters and then obtain random effects by taking conditional expectations.

General Linear Mixed Models

We can greatly increase complexity of such models by having more complicated random effects and more complex covariance structures. This gives rise to *General Linear Mixed Models*.

Consider y_i to be an n -dimensional vector for cluster i :

$$y_i = x_i\beta + u_i\gamma_i + \varepsilon_i$$

$$\gamma_i \sim \mathcal{N}(0, Q), \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2 \Sigma_{n_i}), \quad \gamma_i, \varepsilon_i \text{ indep.}$$

We obtain the *conditional formulation* of our model by conditioning on random effects (i.e., we fix cluster i):

$$y_i \mid \gamma_i \sim \mathcal{N}(x_i\beta + u_i\gamma_i, \sigma^2 \Sigma_{n_i})$$

The *marginal formulation* subsums random effects and other noise terms:

$$y_i = X_i\beta + \varepsilon_i^*; \quad E(y_i) = X_i\beta$$

$$\text{Cov}(y_i) = \text{Cov}(u_i\gamma_i) + \text{Cov}(\varepsilon_i) = u_i Q u_i^\top + \sigma^2 \Sigma_{n_i}$$

$$\Rightarrow Y_i \sim \mathcal{N}(x_i\beta, \sigma^2 \Sigma_{n_i} + u_i Q u_i^\top)$$

Example 5.1

To recover the RI model from the above, use

$$u_i = (1, \dots, 1)^\top, \quad Q = \sigma_1^2, \quad \Sigma = \mathbb{I}$$

We would then recover the exact same structure we had above.

Residual analysis is more difficult, as the complicated error structure leads to many technicalities. We can investigate Turkey-Anscombe plots but should only consider major deviations. We can also use QQ-plots on the estimated random effects to investigate whether the normality assumption holds.

To summarise: A general linear mixed model is given by

$$y = X\beta + U\gamma + \varepsilon$$

with

$$\begin{pmatrix} \gamma \\ \varepsilon \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} G & 0 \\ 0 & R \end{pmatrix} \right).$$

In this model, \mathbf{X} and \mathbf{U} are design matrices, β is a vector of fixed effects, and γ is a vector of random effects. The covariance matrices for γ and ε are assumed to be nonsingular, and therefore positive definite, and γ and ε are independent.

$$\mathbf{U} = \text{blockdiag}(\mathbf{U}_1, \dots, \mathbf{U}_i, \dots, \mathbf{U}_m) = \begin{bmatrix} \mathbf{U}_1 & & & & \mathbf{0} \\ & \ddots & & & \\ & & \mathbf{U}_i & & \\ & & & \ddots & \\ \mathbf{0} & & & & \mathbf{U}_m \end{bmatrix}$$

$\varepsilon \sim N(\mathbf{0}, \mathbf{R}), \gamma \sim N(\mathbf{0}, \mathbf{G})$ with block diagonal covariance matrices:

$$\mathbf{R} = \text{blockdiag} \left(\sigma^2 \mathbf{\Sigma}_{n_1}, \dots, \sigma^2 \mathbf{\Sigma}_{n_i}, \dots, \sigma^2 \mathbf{\Sigma}_{n_m} \right)$$

$$\mathbf{G} = \text{blockdiag}(\mathbf{Q}, \dots, \mathbf{Q}, \dots, \mathbf{Q})$$

Note: dimensions of above elements

- \mathbf{y}_i is an $n_i \times 1$ response vector.
- \mathbf{X}_i is an $n_i \times p$ fixed-effects design matrix.
- $\boldsymbol{\beta}$ is a $p \times 1$ vector of fixed-effect coefficients (common across groups).
- \mathbf{U}_i is an $n_i \times q_i$ random-effects design matrix. E.g., this could consist of one column of ones and one column with time observations in a typical RIRS example.
- $\boldsymbol{\gamma}_i$ is a $q_i \times 1$ vector of random-effect coefficients. E.g., we assume that there is some randomness in initial reaction time and response to days of sleep deprivation between people. $\boldsymbol{\gamma}_i$ is the “drawn” value from this random distribution which specifies initial reaction and reaction to sleep deprivation for individual i .
- ε_i is an $n_i \times 1$ vector of random errors.

5.4 Comparison between different approaches for non iid errors

- Generalised Least Squares:
 - applicable to heteroskedastic and correlated errors but requires knowledge of covariance matrix up to multiplicative factor.
 - Different coefficient estimates.
- Two-stage or MLE:
 - applicable to heteroskedastic and correlated errors but requires a parametric form of covariance structure (e.g., AR process) (do not require exact covariance matrix ex ante).
 - Different coefficient estimates.
- Huber-White:
 - Applicable to heteroskedastic but uncorrelated errors.
 - Note: same coefficient estimates as OLS
 - Heteroskedasticity and autocorrelation robust estimators exist.
- Mixed-models:
 - Implicitly models error structure through random effect structure, can include both heteroskedastic and correlated errors but requires knowledge of grouping of data.
 - Different coefficient estimates.
 - Higher variance compared to other approaches but less detailed knowledge about covariance matrix required.

6 Generalized Linear Models

Generalized Linear Models (GLMs) model relationships between explanatory variables and the parameter of the distribution. That is, we assume our target variable follows some distribution $F(\theta(X))$ whose parameters depend on X .

$$\begin{aligned} S : Y &\sim F(\theta(X)) \\ D : g(\theta(X)) &= \beta_0 + \beta_1 x \end{aligned}$$

Note that while Y is stochastic, we view the link between the parameter and our covariates as being deterministic (no error term).²¹

Three key components:

- Link function $g(\cdot)$
- Linear (in β) predictor
- Distribution of Y

We will see that GLM usually assumes that

- $F(\theta)$ is in exponential family
- θ is expected value
- Every $F(\theta)$ then has a canonical link function with ‘nice’ properties

Example 6.1 (Linear Regression as GLM)

Linear regression is a special case of GLM where we model Y to be normally distributed, with the mean depending on X .

$$\begin{aligned} S : Y &\sim N(\mu(x), \sigma^2) \\ D : \mu(x) &= \beta_0 + \beta_1 x \end{aligned}$$

Example 6.2 (Logistic Regression as GLM)

Assume we have a binary outcome variable. We then naturally model this as a binomial distribution:

$$S : Y \sim \text{Bin}(1, p(x))$$

We further assume that p depends on explanatory variables X . We may model the relationship between X and p as:

$$p(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$$

²¹I.e., we do not explicitly include or model an error term ε . Previously, randomness in $Y \mid X$ only depended on the error, we now instead model Y to follow the distribution $F(\theta(X))$ directly.

where the right hand term is the *Logistic function*. Then

$$D : \log \left(\frac{p(x)}{1 - p(x)} \right) = \beta_0 + \beta_1 x.$$

We call D the *link function*.

We could use any *cdf* since this maps $X \rightarrow [0, 1]$ Logistic regression is popular since it models *Log-Odds*.

Recall that odds of an event A are $P(A)/(1 - P(A))$. We note that log-odds and odds grow monotonically with probability. Log-odds are directly modelled in Logistic regression.

A related concept is the risk ratio. If we are intersted in some event A given B , the risk ratio is $P(A|B)/P(A|B^C)$. The odds of $A|B$ are $P(A|B)/P(A^C|B) = P(A|B)/[1 - P(A|B)]$.

6.1 Logistic Regression

We now study logistic regression in more detail. Note that other than Logistic regression, any CDF can be used to model the link between p and predictors X as it ensures that the predicted probability is in the interval $[0, 1]$.

Logistic regression as latent variable model

Logistic regression can be thought of as a latent variable model. Assume $Z_i = x_i^T \beta + \varepsilon_i$, where ε_i i.i.d. and symmetric around 0. But we only observe

$$\begin{aligned} &\rightarrow Y_i = 1, \text{ if } Z_i > 0 \\ &\rightarrow Y_i = 0, \text{ if } Z_i \leq 0 \end{aligned}$$

We can then see

$$\begin{aligned} P(Y_i = 1) &= P(x_i^T \beta + \varepsilon_i > 0) \\ &= P(\varepsilon_i > -x_i^T \beta) \\ &= P(\varepsilon_i < x_i^T \beta) \text{ (by symmetry)} \end{aligned}$$

If we assume $\varepsilon_i \sim \text{Logistic}(\mu = 0, s = 1)$,²² then

$$\begin{aligned} \underbrace{P(Y_i = 1)}_{=:\pi} &= P(\varepsilon_i < x_i^T \beta) = \frac{\exp\{x_i^T \beta\}}{1 + \exp\{x_i^T \beta\}} \\ \Rightarrow \log \left(\frac{\pi}{1 - \pi} \right) &= x_i^T \beta \end{aligned}$$

²²Note that the Logistic distribution has pdf $f(x; \mu, s) = \frac{e^{-(x-\mu)/s}}{s(1+e^{-(x-\mu)/s})^2}$, mean is given by μ , standard deviation $s^2\pi^2/3$. s is called the scale parameter. The CDF for the standard logistic distribution is $P(\varepsilon_i \leq x) = F(x) = \frac{e^x}{1+e^x}$.

We can instead also assume that $\varepsilon_i \sim \mathcal{N}(0, 1)$ and arrive at **Probit regression**.²³

If we have more than two *ordered* categories, we can use proportional odds logistic regression.

Parameter Estimation

We use maximum likelihood estimation, given a vector x_i of explanatory variables and binary outcome y_i . I.e., we find $\hat{\beta}$ to maximize

$$l(\hat{\beta}) = \prod_{i:y_i=1} p(x_i) \prod_{j:y_j=0} (1 - p(x_j)) = \prod_{i:y_i=1} \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \prod_{j:y_j=0} \frac{1}{1 + \exp(x_j^T \beta)}.$$

We now assume just one regressor and more succinctly write this as

$$L(y_1, \dots, y_n; \beta) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{(1-y_i)}$$

with $\pi_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}$ and $\eta_i = x_i^T \beta$. Taking logs:²⁴

$$\begin{aligned} \ell(y_1, \dots, y_n; \beta) &= \sum y_i \log \pi_i + (1 - y_i) \log (1 - \pi_i) \\ &= \sum_i \left(y_i x_i^T \beta - \log (1 + e^{x_i^T \beta}) \right) \end{aligned}$$

Consider the partial derivatives,

$$\frac{\partial \pi_i}{\partial \beta_0} = \frac{\partial \pi_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_0} = \pi_i (1 - \pi_i); \quad \frac{\partial \pi_i}{\partial \beta_1} = x_i \pi_i (1 - \pi_i)$$

Further, let $n\ell = -\ell$. We aim to minimize the negative log-likelihood,

$$\begin{aligned} \frac{\partial n\ell}{\partial \beta_0} &= - \sum y_i \frac{1}{\pi_i} \pi_i (1 - \pi_i) + (1 - y_i) \frac{(-1)}{1 - \pi_i} \pi_i (1 - \pi_i) \\ &= - \sum y_i (1 - \pi_i) - (1 - y_i) \pi_i = - \sum y_i - \pi_i \\ \frac{\partial n\ell}{\partial \beta_1} &= - \sum (y_i - \pi_i) x_i \end{aligned}$$

Thus,

$$\nabla_{\beta} n\ell = \begin{pmatrix} \partial n\ell / \partial \beta_0 \\ \partial n\ell / \partial \beta_1 \end{pmatrix} = \begin{pmatrix} \sum (\pi_i - y_i) \\ \sum x_i (\pi_i - y_i) \end{pmatrix}$$

We will not be able to set these equations equal to zero and solve them analytically. We therefore consider Newton-Raphson instead.

²³The link function affects units of estimated parameters. Due to different variances, estimated parameters for Probit and Logit regression will be different. If we choose scale parameters such that $\text{Var}(\epsilon_i)$ is identical for Logit and Probit, then we will obtain very similar results. Also, no matter the scale, estimated probabilities will be similar.

²⁴Use $\ln(1 - \pi_i) = \ln \left(1 - \frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right) = \ln \left(\frac{1}{1 + \exp(\eta_i)} \right)$.

Second-order derivatives:

$$\begin{aligned}\frac{\partial^2 nl}{\partial \beta_0^2} &= \frac{\partial}{\partial \beta_0} \left(\sum (\pi_i - y_i) \right) = \sum \frac{\partial}{\partial \beta_0} \pi_i = \sum \pi_i (1 - \pi_i) \\ \frac{\partial^2 nl}{\partial \beta_1^2} &= \dots = \sum x_i^2 \pi_i (1 - \pi_i) \\ \frac{\partial^2 nl}{\partial \beta_1 \partial \beta_1} &= \dots = \sum x_i \pi_i (1 - \pi_i)\end{aligned}$$

The Hessian matrix is therefore given by

$$H_{nl} = \begin{bmatrix} \sum \pi_i (1 - \pi_i) & \sum x_i \pi_i (1 - \pi_i) \\ \sum x_i \pi_i (1 - \pi_i) & \sum x_i^2 \pi_i (1 - \pi_i) \end{bmatrix} = \sum_{i=1}^N \begin{bmatrix} 1 \\ x_i \end{bmatrix} \pi_i (1 - \pi_i) \begin{bmatrix} 1 & x_i \end{bmatrix} \quad (6.1)$$

We can use Newton's method to estimate $\hat{\beta}$, using some initial value:

$$\beta_{n+1} = \beta_n - H_{nl}^{-1} \nabla_{nl}$$

Note that y does not appear in Eq. (6.1). Therefore, taking expectation is identical to (6.1) itself. Recalling the definition of Fisher information (Definition A.1 and Eq. (A.1)), the Fisher information matrix is given by

$$J_{ij}(\beta) = -\mathbb{E} \left(\frac{\partial^2}{\partial \beta_i \partial \beta_j} \ell(y_1, \dots, y_n, \beta) \right)$$

Thus, Newton-Method is identical to Fisher scoring. This is always the case when canonical links are used. Further, the optimisation problem is convex and therefore has a global minimum.

Inference

Due to MLE estimation, we have asymptotic normality. In particular,

$$\hat{\theta} \stackrel{a}{\sim} \mathcal{N}(\theta, V(\theta)),$$

where the asymptotic covariance matrix $V(\theta)$ of $\hat{\theta}$ is the inverse of the Fisher information which we calculated for the simple case of one predictor in (6.1). The general form is given by

$$V(\theta)^{-1} = I(\theta) = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \mathbb{E} \left[(y_i - P_\theta[Y_i = 1])^2 \right] = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \frac{\exp(\mathbf{x}_i^T \theta)}{(1 + \exp(\mathbf{x}_i^T \theta))^2}.$$

We can test nested models based on the *Likelihood Ratio test*. If we consider two nested models with $q < p$ dimension, then

$$2 \left(\ell(\hat{\theta}^{(p)}) - \ell(\hat{\theta}^{(q)}) \right)$$

is known to be asymptotically χ^2 with $(p - q)$ degrees of freedom.

In particular, it is often useful to compare our model to the worst model with only an intercept, $\ell(\hat{\beta}^0)$ and the best possible models with as many parameters as observations, $\ell(\hat{\beta}^S)$ (*saturated model*).

- **Null deviance:** $D_0 = 2 \left(\ell(\hat{\beta}^S) - \ell(\hat{\beta}^0) \right)$
- **Residual deviance:**

$$2 \left(\ell(\hat{\beta}^S) - \ell(\hat{\beta}) \right) = \dots = \sum 2 \left(y_i \log \left(\frac{\hat{\pi}_i^S}{\hat{\pi}_i} \right) + (1 - y_i) \log \left(\frac{1 - \hat{\pi}_i^S}{1 - \hat{\pi}_i} \right) \right)$$

Note that we can evaluate the contribution of each observation to the residual deviance. We can use this to analyse residuals which is much more tricky than in linear regression.

We can compare nested models using AIC.

Interpretation

The estimated β_i coefficient can be interpreted as follows.

- **Odds-scale:** If x_i increases by one unit, *odds* increase by **factor** $\exp\{\beta_i\}$.
- **Log-odds scale:** If x_i increases by one unit, *log-odds* increase by β_i .
- Note: it is difficult to make simple and compact statement in terms of the probability scale.

In R:

- `fm1 <- glm(y ~ x, data = df, family = binomial())`
- Predict response either in terms of log-odds or on response scale, by specifying `predict.glm(..., type = "response")` or `type = "link"`

6.2 Poisson Regression

Suppose our Y_i is count data. In this case we may assume $Y_i \sim \text{Poisson}(\lambda_i)$. Recall that both variance and expectation of Y_i are given by λ_i .

Poisson regression assumes that the parameter λ_i can be explained by certain variables. We model

$$\begin{aligned} \text{D: } \log(\mu_i) &= \log(\lambda_i) = x_i^T \beta \\ \text{S: } Y_i &\sim \text{Pois}(\lambda_i) \end{aligned}$$

Note that we use the $\log(\cdot)$ function to ensure non-negativity of our estimated λ_i .

Estimation

We estimate our model using MLE.

Log-likelihood:

$$L(\beta_0, \beta_1) = \prod_{i=1}^n \frac{\lambda_i^{y_i}}{y_i!} e^{-\lambda_i} \text{ where } \lambda_i = \exp\{\beta_0 + \beta_1 x_i\}$$

$$l(\beta_0, \beta_1) = \sum_{i=1}^n y_i \log \lambda_i - \log(y_i!) - \lambda_i =$$

$$= \sum_{i=1}^n y_i (\beta_0 + \beta_1 x_i) - \log(y_i!) - e^{\beta_0 + \beta_1 x_i}$$

Score equations:

$$s(\hat{\beta}_0, \hat{\beta}_1) = \begin{pmatrix} \frac{\partial l}{\partial \beta_0} \\ \frac{\partial l}{\partial \beta_1} \end{pmatrix} \stackrel{!}{=} 0 \Rightarrow \hat{\beta}_0, \hat{\beta}_1$$

where

$$\frac{\partial l}{\partial \beta_0} = \sum y_i - \underbrace{e^{\beta_0 + \beta_1 x_i}}_{\lambda_i} \stackrel{!}{=} 0$$

$$\frac{\partial l}{\partial \beta_1} = \sum y_i x_i - x_i e^{\beta_0 + \beta_1 x_i} = \sum x_i (y_i - e^{\beta_0 + \beta_1 x_i}) \stackrel{!}{=} 0$$

Note that this is similar in the form to normal equations in linear regression where $X^T(y - \hat{y}) = 0$. This system of equations is non-linear, we use numerical optimization (Fisher Scoring). We will later show that this is a convex optimization problem.

Hessian matrix

$$\frac{\partial^2 l}{\partial \beta_0^2} = -\sum \lambda_i; \quad \frac{\partial^2 l}{\partial \beta_0 \partial \beta_1} = -\sum x_i \lambda_i; \quad \frac{\partial^2 l}{\partial \beta_1^2} = -\sum x_i^2 \lambda_i$$

Therefore,

$$F(\beta_0, \beta_1) = -\sum \lambda_i \begin{bmatrix} 1 & x_i \\ x_i & x_i^2 \end{bmatrix} = -\sum \begin{bmatrix} 1 \\ x_i \end{bmatrix} \begin{bmatrix} 1 & x_i \end{bmatrix} \lambda_i$$

Fisher Scoring. Note that our hessian matrix does not depend on y . Thus taking expectation is identical to directly taking the matrix. Let $\mathcal{I}(\beta_0, \beta_1) = -F(\beta_0, \beta_1)$ be the Fisher matrix, we have

$$\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + \mathcal{I}(\hat{\beta}^{(t)})^{-1} s(\hat{\beta}^{(t)})$$

This optimization process can be rewritten as iteratively re-weighted least squares. This holds true for all GLMs (i.e., all exponential distributions) when using the canonical link.

Inference

Using asymptotic properties of the MLE, we have $\hat{\beta} \stackrel{a}{\sim} \mathcal{N}(\beta, \mathcal{I}^{-1}(\hat{\beta}))$.

Model Comparison

- Null model: worst model with only intercept

$$\Rightarrow l(\hat{\beta}_0) = \sum y_i \hat{\beta}_0 - \log(y!) - e^{\hat{\beta}_0}$$

- Saturated model (best): #pars = #obs

$$\Rightarrow y_i = \lambda_i$$

$$\Rightarrow l(\hat{\beta}^s) = \begin{cases} \sum y_i \log(y_i) - \log(y_i!) - y_i & \text{if } y_i \neq 0 \\ 0 & \text{else} \end{cases}$$

We then calculate

- Null deviance: $-2(l(\hat{\beta}^s) - l(\hat{\beta}^0))$
- Residual deviance: $-2(l(\hat{\beta}^s) - l(\hat{\beta}))$

We can write the residual deviance as a sum of **squared deviance residuals** of each observation.

Caveat: Poisson assumes linear relation between expectation and variance with slope 1. One may use quasi-poisson to loosen the relationship to a linear relationship with slope $\neq 1$. Using quasi-poisson leads to same point estimates but different estimates for standard errors.²⁵ Large overdispersion should lead to use of other models, e.g., negative binomial regression.

6.3 Gamma Regression

Recall that if we have positive continuous data we may want to log-transform the data. The issue in interpreting effect size is that

$$g^{-1}(E(g(Y))) \neq E(Y)$$

However, if we have a *symmetric* distribution, as $\log(\cdot)$ preserves ordering, the expectation equals median so we can make statements about changes in median.

GLM offers an alternative (for the case of skewed distributions) by considering gamma regressions. For each level of the explanatory variable, we find μ_i based on x_i and model a separate gamma distribution with term μ_i .

- Gamma distribution $\Gamma(k, \theta)$ where k measures shape and θ scale

$$E(X) = k\theta, \text{Var}(X) = k\theta^2$$

- Gamma regression then models

$$\text{D} : \log(\mu_i) = x_i^T \beta$$

$$\text{S} : Y_i \sim \Gamma(\mu_i, v)$$

²⁵Since we estimate the quasi-dispersion parameter, we will obtain t instead of z values.

6.4 Overview: Generalized Linear Models

All previous cases are examples of Generalized Linear Models which we now summarize.

- Three components:
 - Distribution from exponential family
 - link function (preferably use canonical link)
 - linear predictor
- $D : g(\mu(x)) = x_i^T \beta$ (or $\mu(x) = h(x_i^T \beta)$, thus $g = h^{-1}$)
- $S : Y \sim F(\mu)$

Recall: Exponential family is given by distributions with density of the form

$$f(y | \theta) = \exp \left(\frac{y\theta - b(\theta)}{\phi} w + c(y, \phi, w) \right)$$

The log-density is given by

$$\begin{aligned} \log f(y | \theta) &= \frac{y\theta - b(\theta)}{\phi} w + c(y, \phi, w). \\ \mathbb{E}(y) &= \mu = b'(\theta), \quad \text{Var}(y) = \phi b''(\theta)/w. \end{aligned}$$

here θ is related to the linear predictor, ϕ is a dispersion parameter, and c a nuisance parameter.

Example 6.3 (Bernoulli is in exp. family)

We have $P(X = x) = \pi^x (1 - \pi)^{(1-x)}$. We can rewrite this by taking log and exponential:

$$\exp \left\{ x \log \frac{\pi}{1 - \pi} + \log(1 - \pi) \right\}$$

we can then set $\theta = \log \frac{\pi}{1 - \pi}$, $-b(\theta) = \log(1 - \pi)$ and $w = \phi = c = 1$.

MLE is done in the same spirit as above, by finding roots of score functions using Fisher scoring. Note: Fisher scoring can be written as **iteratively re-weighted least squares (IRLS)** estimates whereby the weights depend on the derivative of the link function, the variance function of the distribution, and the current estimate.

One can easily show that for exponential family densities, the log likelihood is concave, so one can find a global optimum:

$$\begin{aligned} l(\theta) &= \frac{y\theta - b(\theta)}{\phi} \omega + c(y, \phi, \omega) \\ l'(\theta) &= \frac{\partial l}{\partial \theta} = \frac{y - b'(\theta)}{\phi} \omega \\ l''(\theta) &= \frac{\partial^2 l}{\partial \theta^2} = -\frac{b''(\theta)}{\phi} \omega = -\frac{\text{Var}(Y)\omega^2}{\phi^2} \leq 0 \end{aligned}$$

We can also extend GLMs to include mixed effects which leads to Generalized Linear Mixed Models.

7 Extensions

7.1 Non-Linear Regression

Consider

$$y_i = f(x_i, \theta) + \varepsilon_i,$$

- p : Number of parameters (θ),
- m : Number of explanatory variables (x)

where potentially $m \neq p$.

Assumptions on residual term: Usually one of the two below

- $\varepsilon_i \sim N(0, \sigma^2)$ iid (then OLS = MLE \rightarrow properties of MLEs), or
- $E(\varepsilon_i) = 0$ and $\text{Cov}(\varepsilon_i) = \sigma^2 I \rightarrow$ asymptotic normality of pars under technical assumptions

Usually the form of $f(x, \theta)$ is given by *context*. That is, we assume to know $f(\cdot)$.

7.1.1 Linearizing

In many cases, we may be able to linearize a non-linear function and thus obtain coefficients and predictions in our usual OLS framework. Example:

Michaelis-Menten $y = \frac{\theta_1 x}{\theta_2 + x} \rightarrow \frac{1}{y} = \frac{\theta_2}{\theta_1} \cdot \frac{1}{x} + \frac{1}{\theta_1}$ can be linearized to $\tilde{y} = \frac{1}{y}, \tilde{x} = \frac{1}{x} \rightarrow \tilde{y} = \frac{\theta_2}{\theta_1} \cdot \tilde{x} + \frac{1}{\theta_1}$

- prediction is easily done
- however, mainly interested in parameters:
 - $\rightarrow \theta_1$ easy to get from $\frac{1}{\theta_1}$; θ_2 much harder to untangle from θ_2/θ_1
- Error structure might not be additive anymore
 - $y = \frac{\theta_1 x}{\theta_2 + x} + \varepsilon \rightarrow \tilde{y} = \frac{\theta_2}{\theta_1} \cdot \tilde{x} + \frac{1}{\theta_1} + ?$

Instead use GLM with inverse link (if error structure is additive):

$$\frac{1}{\mu} = \frac{\theta_2}{\theta_1} \cdot \frac{1}{x} + \frac{1}{\theta_1} \text{ and } Y \sim N\left(\mu(x), \sigma^2\right)$$

Error structure is important as violation leads to wrong coverage probabilities of CI!

We can investigate deviations from assumed error structure by investigating the Turkey-Anscombe plot. If we find deviations, we may want to instead consider fitting a nonlinear regression model.

7.1.2 Model fitting: Iterative Least Squares

If we have n observations, we generate a vector $\eta(\theta)$ which contains $f(x, \theta)_i$ for each observation i . Our vector of predictions $\eta(\theta)$ and our true y vector are (each) a point in \mathbb{R}^n . As we change θ , we generate a new vector and we trace out a model surface in \mathbb{R}^n . We then want to find the point on the model surface closest to actual observations. I.e., we want to find $\hat{\theta}$ to minimize

$$S(\theta) = \sum_{i=1}^n (y_i - \eta_i(\theta))^2. \quad (7.1)$$

We need to use *numerical optimization* as this generally does not allow for a closed-form solution.

Below we denote by $A(\theta)$ the Jacobian Matrix of $\eta(\theta)$.²⁶

Algorithm (Gauss-Newton Method)

1. Start with initial value $\hat{\theta}^{(0)}$
2. For $k = 1, 2, \dots$
 - Calculate tangent plane of $\eta(\theta)$ at $\hat{\theta}^{(k-1)}$:

$$\eta(\theta) \approx \eta(\hat{\theta}^{(k-1)}) + A(\hat{\theta}^{(k-1)}) \cdot (\theta - \hat{\theta}^{(k-1)})$$

- Project on tangent plane (linear regression problem):

Note: $Y = \eta(\theta) + \varepsilon$, thus

$$Y = \eta(\theta) + \varepsilon \approx \eta(\hat{\theta}^{(k-1)}) + A(\hat{\theta}^{(k-1)}) \cdot (\theta - \hat{\theta}^{(k-1)}) + \varepsilon$$

$$Y - \eta(\hat{\theta}^{(k-1)}) \approx A(\hat{\theta}^{(k-1)}) \cdot (\theta - \hat{\theta}^{(k-1)}) + \varepsilon$$

«Preliminary residuals»: $\tilde{Y}^{(k-1)} = Y - \eta(\hat{\theta}^{(k-1)})$

«correction term»: $\beta^{(k-1)} = (\theta - \hat{\theta}^{(k-1)})$

Solve:

$$\tilde{Y}^{(k-1)} = A(\hat{\theta}^{(k-1)}) \cdot \beta^{(k-1)} + \varepsilon$$

This is a linear regression which yields $\beta^{(k-1)}$.

- Use $\beta^{(k-1)}$ to update estimated coefficients:

$$\hat{\theta}^{(k)} = \hat{\theta}^{(k-1)} + \beta^{(k-1)}$$

3. Iterate until convergence

That is, at each iteration we're using a linear approximation by a tangent plane. We then solve a "linear regression problem" with the preliminary residuals as our target value and the Jacobian matrix at our parameter estimate as regressors. We use the estimated coefficients $\hat{\beta}$ as our correction term to update our coefficients $\hat{\theta}$.

Convergence depends on starting values since model surface might be arbitrarily complex (thus non-convex, local optima). Three methods for choosing starting values:

²⁶I.e., the i -th column of A is the gradient $d\eta(\theta)/d\theta_i$. If $\eta(\theta) = X\beta$, then $A(\theta) = X$ so we recover the usual OLS problem in the below algorithm.

1. From context or experience
2. From linearized functions (OLS fit)
3. From data (parameters have meaning.)

Self-starting functions automate this to find initial values.

Example 7.1 (Interpreting coefs in Biochemical Oxygen Demand (BOD))

Model:

$$f(x, \theta) = \theta_1 (1 - \exp(-\theta_2 x)); \text{ assume } \theta_1 > 0, \theta_2 > 0$$

Interpretation of parameters:

$$\begin{aligned} - \lim_{x \rightarrow \infty} f(x, \theta) &= \theta_1 \rightarrow \text{Asymptote} \\ - \frac{df}{dx}(x=0) &= \theta_1 \theta_2 \rightarrow \text{Slope at } x=0 \end{aligned}$$

In R: `fm2 <- nls(yObs ~ SSgompertz(x, Asym, b2, b3), data = df)`.

Manually:

`fm <- nls(yObs ~ t1 * exp(-t2 * t3^x), data = df, start = c(t1 = 12, t2 = 5, t3 = 0.5))`

7.1.3 Inference

In general, inference and assessing model fit is more difficult in nonlinear least squares.²⁷ Nevertheless, there are multiple useful approaches for inference in nonlinear problems.

Based on Linear Approximation

- From Gauss-Newton, we approximate the model surface by the tangential plane given by $A(\hat{\theta})$
- Fit standard linear regression: $Y = A(\hat{\theta}) \cdot \beta + \varepsilon$
- One can show that at optimum with $\hat{\beta} \approx 0$, $\text{Cov}(\hat{\theta}) \approx \text{Cov}(\hat{\beta}) = \sigma^2 \left(A(\theta_0)^T A(\theta_0) \right)^{-1}$. Therefore,

$$\hat{\theta} \stackrel{\text{asym.}}{\sim} \mathcal{N} \left(\theta_0, \sigma^2 \left(A(\theta_0)^T A(\theta_0) \right)^{-1} \right)$$

where we in practice replace θ_0 with $\hat{\theta}$.

- Thus we get for a $(1 - \alpha)$ -confidence interval:

$$\hat{\theta}_k \pm t_{n-p; 1-\alpha/2} \text{se}(\hat{\theta}_k)$$

$$\text{where } \text{se}(\hat{\theta}_k) = \hat{\sigma} \sqrt{\left(\left(A(\hat{\theta})^T A(\hat{\theta}) \right)^{-1} \right)_{kk}} \text{ and } \hat{\sigma}^2 = \frac{s(\hat{\theta})}{n-p} = \frac{1}{n-p} \sum_{i=1}^n \left(y_i - \eta_i(\hat{\theta}) \right)^2$$

²⁷E.g., is the assumption for Gompertz growth really satisfied? \rightarrow hard to answer from fit. But sometimes there are nested models which we can test.

- Our test statistic is:

$$\delta_k(\theta_k^*) = \frac{\hat{\theta}_k - \theta_k^*}{\text{s.e.}(\hat{\theta}_k)} \approx t_{n-p} \quad (7.2)$$

Note: CI symmetric by design

Based on Likelihood

Idea: Points where drop of in likelihood from optimum is small are still very likely. Use dropoff in likelihood ("vertical distance") to find confidence set.²⁸

- Assume $H_0 : \theta = \theta^*$ is true which imposes q restrictions²⁹
- Then it approximately holds:

$$T = \frac{n-p}{q} \cdot \frac{S(\theta^*) - S(\hat{\theta})}{S(\hat{\theta})} \approx F_{q, n-p}$$

- This is still only an approximate result; but in practice often much more accurate than linear approximation
- Construct confidence regions: Collect all vectors θ^* that are not rejected by the F-Test \rightarrow difficult to visualize and summarize.

Alternative, to obtain tests for single parameters, we can use **profiling**

- Assume $H_0 : \theta_k = \theta_k^*$
- Fix parameter of interest at some arbitrary value $\theta_k = \theta_k^*$ and minimize $S(\theta)$ wrt. $\theta_j, j \neq k$; denote the minimum by $\tilde{S}_k(\theta_k^*)$
 - Repeating for all possible values of θ_k , then the function $\tilde{S}_k(\theta_k)$ is the “*profile likelihood*”
- Under $H_0 : \theta_k = \theta_k^*$, it approximately holds: $\tilde{T}_k(\theta_k^*) = (n-p) \cdot \frac{\tilde{S}_k(\theta_k^*) - S(\hat{\theta})}{S(\hat{\theta})} \approx F_{1, n-p}$

²⁸Note that $S(\theta)$ is not a likelihood. We can, however, recover the likelihood and create a link between the two values: Under the assumption that errors are normally distributed ($\varepsilon_i \sim N(0, \sigma^2)$), the likelihood function is:

$$L(\theta, \sigma^2) = \prod \left(1/\sqrt{2\pi\sigma^2} \right) \exp \left(- (y_i - f(x_i, \theta))^2 / (2\sigma^2) \right)$$

Taking the log of this gives us the log-likelihood:

$$\ell(\theta, \sigma^2) = \text{constant} - (n/2) \log(\sigma^2) - (1/(2\sigma^2)) \sum (y_i - f(x_i, \theta))^2$$

Notice that the last term contains $S(\theta) = \sum (y_i - f(x_i, \theta))^2$, our sum of squared residuals. (For $n \rightarrow \infty$ the F -test is the same as the likelihood-ratio test, and the sum of squares is, up to a constant, equal to the negative log-likelihood).

²⁹Note that we have some inconsistent notation. In linear regression, we imposed $p - q$ restrictions in our submodel. Here we instead denote by q itself the number of restrictions in the submodel.

- Due to relation between $F_{1,n-p}$ and t_{n-p} this corresponds to:

$$T_k(\theta_k^*) = \text{sign}(\hat{\theta}_k - \theta_k^*) \cdot \frac{\sqrt{\tilde{S}_k(\theta_k^*) - S(\hat{\theta})}}{\hat{\sigma}} \approx t_{n-p} \quad (7.3)$$

- CI:

$$\left\{ \theta_k^* \mid \sqrt{\tilde{S}_k(\theta_k^*) - S(\hat{\theta})} \leq t_{n-p; 1-\alpha/2} \hat{\sigma} \right\}.$$

- CI based on profiling are usually not symmetric³⁰ and perform better than CI based on linear approximation
- We can apply a monotone transformation to θ_k and apply the same monotone transformation to the CI. This does not hold for the linear approximation.

7.1.4 Assessing non-linearity

We can test how well the linear approximation works. To do so, we compare the two tests based on linear approximation and profiling.

- If linear approximation is good, $\delta_k(\theta_k^*)$ and $T_k(\theta_k^*)$ should behave similarly
- Check in profile t-plot: Plot $T_k(\theta_k^*)$ against $\delta_k(\theta_k^*)$ (or $T_k(\theta_k^*)$ against θ_k^*) when varying θ_k^*
- In a linear setting, we would expect a diagonal straight line if $T_k(\theta_k^*)$ against θ_k^*
- In a strongly non-linear setting, we should see a clear deviation from this

Additionally, we can investigate likelihood profile traces. I.e., we fix arbitrary values for θ_1 and plot the optimal θ_2 (and vice versa). So for each point on the line, we're seeing what value of θ_2 gives us the best model fit when we force θ_1 to take a specific value. This creates a "trace" showing how the optimal values of the parameters move together.

When the lines in a profile trace are nearly parallel, it suggests that the parameters are highly correlated or potentially redundant in the model.

7.2 Non-parametric regression

We are in a scenario where

$$Y = f(x) + \epsilon$$

with $\mathbb{E}(\epsilon) = 0$, $\text{Cov}(\epsilon) = \sigma^2 \mathbf{1}$ and f twice continuously differentiable.

Difference to non-linear regression: We do not know the form of the nonlinear function f . Often in this case polynomial regressions are unstable and nonparametric approaches produce better results.

³⁰Note that compared to the linear test statistic in (7.2), here $s(\theta)$ enters the test, not θ directly. I.e., we compare the "vertical distance" of the likelihood rather than the "horizontal distance" between parameter estimates.

7.2.1 Kernel Regression

- For each x , we average the observed y -values in neighborhood (e.g. $x \pm h$) where h is a fixed bandwidth
- We weight observations inversely to the distance from x (i.e., lower weights on more distant obs)
- Define weight via “Kernel function” (often a probability density function)

– Problem:

$$\hat{f}_n(x) = \operatorname{argmin}_{m_x \in \mathbb{R}} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) (Y_i - m_x)^2$$

I.e., at each point x , we are searching for the best local constant m_x such that the localized sum of squares is minimized

– Solution: Nadaraya-Watson kernel estimator

$$\hat{f}_n(x) = \sum_{i=1}^n w_i(x) Y_i \text{ where } w_i(x) = \frac{K\left(\frac{x-x_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{x-x_j}{h}\right)}$$

In practice, choice of kernel function does not have a strong influence on outcome, but choice of bandwidth h does. h governs *Bias-Variance tradeoff*.³¹

- bias $\sim h^2$;
- variance $\sim 1/h$

Rates of convergence:

- OLS (parametric): n^{-1}
- Non-parametric: $n^{-4/5}$

Derivation of convergence rates and h proportionality

OLS, $\hat{y} = \ell^T \beta = \hat{f}(\ell)$.

MSE equals variance since OLS is unbiased:

$$\operatorname{Var}(\hat{y}) = \ell^T \sigma^2 (X^T X)^{-1} \ell = \ell^T \sigma^2 \left(\underbrace{\frac{X^T X}{n}}_{\rightarrow \Sigma} \right)^{-1} \ell n^{-1}$$

Here we make our usual assumption that $X^T X n^{-1}$ converges to a constant, as for instance also in Lemma 5.1.³²

³¹MSE convergence and bias proportionality requires f to be twice continuously differentiable. Variance proportionality does not require this assumption.

³²For each entry (j, k) , we’re essentially assuming that the sample covariance between predictors j and k converges to some population value.

Kernel regression

We have

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) Y_i \bigg/ \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right)$$

We for now assume that $x_i = \frac{i}{n}$. i.e. the x values are on a grid with difference $\frac{1}{n}$. Then, using $\mathbb{E}(Y_i) = f(X_i)$,

$$\begin{aligned} \mathbb{E}[\hat{f}_h(x)] &= \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) f(i/n) \bigg/ \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-x_i}{h}\right) \\ &\approx \frac{\frac{1}{h} \int K\left(\frac{x-\mu}{h}\right) f(\mu) d\mu}{\frac{1}{h} \int K\left(\frac{x-\mu}{h}\right) d\mu} \end{aligned}$$

where we used that an integral is approximately a sum. Note that n was moved from outside the integral since the integration variable $\mu \approx \frac{1}{n}$ (loosely speaking).

We now perform a change of variables: $(x - \mu)/h = -w$. Then $d\mu = h dw$. Further, assume $K()$ is a symmetric density. Then we obtain

$$\frac{\frac{1}{h} \int K(-w) f(x + wh) h dw}{\frac{1}{h} \int K(-w) h dw}$$

Using symmetry and the fact that a density integrates to 1, simplifying $1/h \cdot h$ in numerator and denominator, we get

$$= \int K(w) f(x + wh) dw$$

Using a Taylor approximation, $f(x + hw) \approx f(x) + f'(x)hw + 1/2 f''(x)h^2 w^2$, noting that due to twice continuous differentiability the remaining terms are small,³³ we have

$$\mathbb{E}[\hat{f}_h(x)] = f(x) \underbrace{\int K(w) dw}_{=1} + f'(x)h \underbrace{\int K(w)w dw}_{=0 \text{ (symmetry)}} + f''(x)h^2 \underbrace{\frac{1}{2} \int K(w)w^2 dw}_{=: C_1(K)} + \dots$$

Our bias is therefore, approximately,

$$\mathbb{E}[\hat{f}_h(x)] - f(x) = f''(x)h^2 \frac{1}{2} \int K(w)w^2 dw$$

which is proportional to h^2 . Additionally, our bias depends on curvature.³⁴

Similarly, we can compute the variance. $\hat{f}_h(x)$ is a sum of independent terms since it is

³³This follows from the rest term in the Taylor expansion and continuous functions being bounded on closed intervals.

³⁴If $f''(x) < 0$, e.g. at a maximum, we underestimate, if $f''(x) > 0$, e.g. at a minimum, we overestimate. This is a general feature of Kernel regression, called *erosion*.

a sum of independent Y_i , so

$$\begin{aligned}\text{Var}\left(\hat{f}_h(x)\right) &= \frac{1}{n^2 h^2} \sum_{i=1}^n K^2\left(\frac{x-x_i}{h}\right) \sigma_\varepsilon^2 / \left(\frac{1}{nh} \sum K\left(\frac{x-x_i}{h}\right)\right)^2 \\ &\approx \frac{1}{nh^2} \int K^2(w) h dw \sigma_\varepsilon^2 / \underbrace{\left(\frac{1}{h} \int K(-w) h dw\right)^2}_{\rightarrow 1} \\ &= \frac{1}{nh} \underbrace{\int K^2(w) dw}_{=: C_2(K)} \sigma_\varepsilon^2\end{aligned}$$

which is proportional to h^{-1} . Notice that we have one n remaining when going to the Riemann sum (integral), since we first square the $1/nh$ term when taking it out of the variance.

Therefore, we conclude:

$$MSE(\hat{f}_h(x)) = h^4 (f''(x))^2 C_1(K)^2 + \frac{\sigma_\varepsilon^2}{nh} C_2(K)$$

Optimizing with respect to h gives³⁵

$$h_{opt} = n^{-1/5} \left(\frac{\sigma^2 C_2(K)}{4(f''(x))^2 C_1(K)} \right)^{1/5}$$

Plugging back into the MSE, we get

$$MSE(\hat{f}_{h,opt}(x)) = n^{-4/5} \cdot C$$

Inference

The kernel estimator at design points is actually a linear operator. To see this, assume $Y = m(x) + \varepsilon$, $\mathbb{E}(\varepsilon) = 0$, $\text{Cov}(\varepsilon) = \sigma^2 1$ and consider $(\hat{m}(x_1), \dots, \hat{m}(x_n)) =: \hat{m}(x_0) = \hat{y}$. Since

$$\hat{m}(X) = \frac{\sum K\left(\frac{X-X_i}{h}\right) y_i}{\sum K\left(\frac{X-X_i}{h}\right)} = \sum w_i(X) y_i$$

we can write $\hat{y} = SY$ with $[S]_{rs} = w_s(x_r)$. We have

- $\mathbb{E}(\hat{m}(x_0)) = Sm(x_0) \neq m(x_0)$
- $\text{Cov}(\hat{m}(x_0)) = \sigma_\varepsilon^2 SS^T$
- We may assume asymptotically $\hat{m}(x) \approx \mathcal{N}(E(\hat{m}(x)), \text{Var}(\hat{m}(x)))$
- We estimate $\widehat{\text{s.e.}}(\hat{m}(x_i)) = \hat{\sigma} \cdot \sqrt{(SS^T)_{ii}}$ and $\hat{\sigma}^2 = \sum_{i=1}^n \frac{(Y_i - \hat{m}(x_i))^2}{n-df}$ where $df = \text{tr}(S)$.
- Then $\hat{m}(x_i) \pm 1.96 \cdot \widehat{\text{s.e.}}(\hat{m}(x_i))$ is an approximate 95% CI for $\mathbb{E}(\hat{f}(x_i))$

³⁵Notice that the optimal value depends on the unknown function f . Thus we cannot use this practically to choose h . Instead we use cross-validation.

Local Polynomial Regression

Instead of calculating the average within a bandwidth (with reduced weights for far away points), we may also fit a weighted local polynomial regression, centered at x :

$$\hat{\beta}(x) = \arg \min_{\beta \in R^p} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right) \left(Y_i - \beta_1 - \beta_2(x_i - x) - \dots - \beta_p(x_i - x)^{p-1}\right)^2.$$

If we use a degree 1 polynomial, this is also called a *local linear estimator* or local linear regression.

- Low order polynomials used
- Have local regression function at a point x :

$$g_x(u) = \sum_{j=1}^p \hat{\beta}_j(x)(u - x)^{j-1} \text{ evaluated at } u = x$$

- Use **intercept only** for prediction:³⁶

$$\hat{f}_h(x) = \hat{\beta}_1(x)$$

I.e., for any new point x' , we refit the model around this point and take the newly fitted $\hat{\beta}_1(x')$ to be the prediction.

- Advantages:
 - Better performance at edges
 - Can estimate derivatives: Compute r -th derivative wrt. u and evaluate at $u = x$ ³⁷

$$\rightarrow \hat{f}^{(r)}(x) = r! \hat{\beta}_{r+1}(x), \text{ for } r = 0, 1, \dots, p-1$$

Smoothing Splines

Idea: Among all functions $m(x)$ with continuous second derivatives, find the one which minimizes the penalized residual sum of squares:

$$\sum_{i=1}^n (Y_i - m(x_i))^2 + \lambda \int m''(z)^2 dz$$

Solution: *Natural cubic spline* with knots at sorted observations x_1, \dots, x_n .³⁸

- Cubic polynomials around each point x_i

³⁶Note that if the design matrix is not orthogonal, this is not equal to just taking the average.

³⁷Note that x is fixed and we fitted a linear function near x and evaluate it at u . Hence, to find the correct slope, we also have to take the derivative w.r.t. u and evaluate this derivative at $u = x$.

³⁸Our estimated function is piecewise cubic polynomial and thus has $4(n-1)$ parameters. Requiring continuity for function and first two derivatives imposes $3(n-2)$ inner continuity constraints, additionally restrictions $\hat{m}''(x_1) = \hat{m}''(x_n) = 0$ leads to n free parameters.

- Smooth connections
- Curvature 0 at outer observations

Since the solution is a cubic spline, we can represent the solution using adequate basis functions B_j for natural splines:

$$m_\lambda(x) = \sum_{j=1}^n \beta_j B_j(x)$$

Optimize:

$$|Y - B\beta|^2 + \lambda\beta^T \Omega \beta,$$

where the design matrix B has columns $(B_j(x_1), \dots, B_j(x_n))^T$ and $\Omega_{jk} = \int B_j''(z) B_k''(z) dz$.

- Solution: $\hat{\beta} = (B^T B + \lambda \Omega)^{-1} B^T Y$
- Predictions: $\hat{Y} = B \hat{\beta} = B (B^T B + \lambda \Omega)^{-1} B^T Y = S_\lambda Y$
- Approximately similar to kernel regression with adaptive bandwidth³⁹

Generalized Additive Models

In higher dimensions, close neighbors become rare and distances between all points become large. Thus, kernel regressions which depend on “neighborhoods” will fail. GAMs reduce this problem by considering non-parametric regression per dimension (at price of assuming additivity).

$$g_{add}(x) = \mu + \sum_{j=1}^p g_j(x_j), \mu \in \mathbb{R}$$

$$g_j(\cdot) : \mathbb{R} \rightarrow \mathbb{R}, \mathbb{E}(g_j(X_j)) = 0 \quad (j = 1, \dots, p)$$

7.3 High-dimensional Regression

If we have $p > n$, we have a perfect fit \hat{y} and parameter estimates $\hat{\beta}$ are no longer unique. We instead impose regularisation penalties. E.g., we minimize

$$PLS(\beta) = (y - X\beta)^T (y - X\beta) + \lambda \cdot \text{pen}(\beta)$$

which we may equivalently write as

$$\min (y - X\beta)^T (y - X\beta) \quad \text{s.t.} \quad \text{pen}(\beta) \leq s$$

Even if $p \leq n$, such a setup may be useful due to the bias-variance tradeoff as a biased model may outperform OLS w.r.t. prediction accuracy (see (3.4)).

³⁹In regions where data suggests more curvature (higher second derivatives), the smoothing spline effectively decreases its bandwidth to capture these features

Ridge Regression

$$(y - X\beta)^T(y - X\beta) + \lambda \cdot \sum_{j=1}^p \beta_j^2$$

- Do not penalize intercept
- “Shrinkage penalty” λ is a tuning parameter, find with cross-validation
- Standardise predictors
- Closed form solution:

$$\hat{\beta}_{\text{Ridge}} = (X^T X + \lambda \mathbf{1})^{-1} X^T y$$

- Note that $X^T X + \lambda \mathbf{1}$ is always pd whenever $\lambda > 0$.
 - Thus, $\hat{\beta}_{\text{Ridge}}$ always exists and is always unique.
 - All p coefficients non-zero even if $p \gg n$
- Covariance matrix: $\text{Cov}(\hat{\beta}_{\text{ridge}}) = \sigma^2 (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$
- $\text{Cov}(\hat{\beta}_{\text{LS}}) - \text{Cov}(\hat{\beta}_{\text{Ridge}})$ is p.d. for $\lambda > 0$ so Ridge coefficients have a smaller variance
- As $\lambda \rightarrow \infty$, the L_2 norm of $\hat{\beta}$ decreases monotonically, but individual coefficients might have a bump and increase.⁴⁰
- If $p > n$, as $\lambda \rightarrow 0$, Ridge bias does not go to zero. If $n \leq p$, bias goes to zero since $\lim_{\lambda \rightarrow 0} \hat{\beta}_{\text{Ridge}}(\lambda) = (X^T X)^{-1} X^T Y = \hat{\beta}_{\text{OLS}}$.

Using penalisation introduces bias but may lead to better out-of-sample prediction performance due to bias-variance tradeoff.

LASSO

Issue with ridge: still have all variables in model. Advantage of LASSO: Shrink some coefficients to zero, thus perform variable selection.⁴¹

Objective function:

$$\min (y - X\beta)^T(y - X\beta) + \lambda \cdot \sum_{j=1}^p |\beta_j|$$

⁴⁰If the design matrix is orthogonal, the decrease is continuous and

$$\hat{\beta}_{\text{Ridge},j}(\lambda) = \frac{1}{1 + \lambda} \hat{\beta}_{\text{LS},j}.$$

⁴¹We can view these as approximations of the NP-hard problem:

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad \text{subject to} \quad \sum_{j=1}^p I(\beta_j \neq 0) \leq s.$$

- ℓ_1 norm not differentiable, no closed-form solution, non-linear estimator
- Still convex optimization problem
- At most $\min\{n, p\}$ non-zero coefficients
- Many theoretical advantages if true model is sparse
- In case of orthonormal design matrix, one can show that

$$\hat{\beta}_{\text{LASSO},j}(\lambda) = \text{sign}(\hat{\beta}_{\text{LS},j}) \left[\left| \hat{\beta}_{\text{LS},j} \right| - \frac{\lambda}{2} \right]_+$$

Theoretical Advantages of LASSO

Under technical assumptions (true model sparsity, relation between predictors, ...):

- If we knew the true non-zero coefficients β^0 , MSE of OLS would converge as:

$$\frac{|X(\hat{\beta} - \beta^0)|^2}{n} = O_P\left(\frac{s_0}{n}\right)$$

where s_0 is number of non-zero coefficients in true model.

- MSE convergence of LASSO:

$$\frac{|X(\hat{\beta} - \beta^0)|^2}{n} = O_P\left(\frac{s_0 \log(p)}{n}\right)$$

We thus pay a small price of using LASSO and not knowing true non-zero coefficients

- LASSO consistently estimates parameters:

$$|\hat{\beta} - \beta^0|_2 = O_P\left(\sqrt{\frac{s_0 \log(p)}{n}}\right)$$

where β^0 is the vector of active β s (padded with 0s).

Issues in High-Dimensional Statistics

If we have $p \geq n$, OLS overfits. Fit on training data will be exact while out-of-sample testing error gets worse the more variables are added. Cannot use BIC, AIC, etc. as we cannot estimate $\hat{\sigma}^2$.

Penalisation helps tame the harmful flexibility leading to better predictions in high dimensions.

- Inference:

- Cannot use traditional measures such as p-values, training MSE due to perfect fit in-sample
 - Instead focus on out-of-sample predictive performance
 - Advanced in-sample methods exist
- Model interpretation:
 - Perfect multicollinearity is guaranteed if $p > n$.
 - Any model could be replaced by other model with correlated predictors
 - Interpret model as one of many possible (good fitting) models

A Mathematical Statistics

A.1 Fisher Information and Fisher Scoring

(From DeGroot & Shervish)

Definition A.1 Fisher Information

Let X be a random variable whose distribution depends on a parameter θ that takes values in an open interval Ω of the real line. Let the p.f. or p.d.f. of X be $f(x | \theta)$. Assume that the set of x such that $f(x | \theta) > 0$ is the same for all θ and that $\lambda(x | \theta) = \log f(x | \theta)$ is twice differentiable as a function of θ . The **Fisher information** $I(\theta)$ in the random variable X is defined as

$$I(\theta) = E_{\theta} \left\{ [\lambda'(X | \theta)]^2 \right\}.$$

Thus, if $f(x | \theta)$ is a p.d.f., then

$$I(\theta) = \int_S [\lambda'(x | \theta)]^2 f(x | \theta) dx.$$

The idea behind Fischer information is that if the density f is sharply peaked with respect to θ , then it is easier to find the ‘correct’ value of θ . We therefore use the *score* (i.e., the derivative of the log of the density with respect to the parameter). A random variable carrying high Fisher information implies that the absolute value of the score is often high. High Fischer information indicates that the maximum is ‘sharp’ while low Fischer information means that the maximum is ‘blunt’.

When there are N parameters, so that θ is an $N \times 1$ vector $\theta = [\theta_1 \ \theta_2 \ \dots \ \theta_N]^\top$, the Fisher information takes the form of an $N \times N$ matrix. The typical element of this matrix is given by:

$$[\mathcal{I}(\theta)]_{i,j} = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log f(X; \theta) \right) \left(\frac{\partial}{\partial \theta_j} \log f(X; \theta) \right) \middle| \theta \right].$$

Under certain regularity conditions, the Fisher information matrix may also be written as

$$[\mathcal{I}(\theta)]_{i,j} = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(X; \theta) \middle| \theta \right]. \quad (\text{A.1})$$

We may use the Fischer information matrix to perform parameter estimation. This is known as **Fisher scoring**. Suppose we are given a log-likelihood function $\ell(\theta)$ and we want to maximize with respect to the parameter θ . We may use the scores (i.e., derivative of log-likelihood functions) and the Hessian matrix \mathcal{J} to optimize this using Newton-Raphson:

$$\theta_{t+1} = \theta_t - [\mathcal{J}(\theta_t)]^{-1} \nabla \ell(\theta_t)$$

If we instead use the expectation of the negative Hessian, which is, under certain regularity conditions the Fisher information matrix given in (A.1), we do Fischer scoring.

$$\theta_{t+1} = \theta_t + [\mathcal{I}(\theta_t)]^{-1} \nabla \ell(\theta_t)$$

A.2 Convergence

We introduce different notions of convergence of random variables.

Definition A.2 Convergence in Distribution

A sequence X_1, X_2, \dots of real-valued random variables, with cumulative distribution functions F_1, F_2, \dots , is said to **converge in distribution**, or *converge weakly*, or converge in law to a random variable X with cumulative distribution function F if

$$\lim_{n \rightarrow \infty} F_n(x) = F(x),$$

for every number $x \in \mathbb{R}$ at which F is continuous.

- Convergence in distribution implies that the probability for X_n to be in a given range is approximately equal to the probability of X being in that range for n sufficiently large
- Convergence in distribution does not generally imply that the sequence of densities will converge
- There exist many equivalent definitions

Definition A.3 Convergence in Probability

A sequence $\{X_n\}$ of random variables **converges in probability** towards the random variable X if for all $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0$$

More explicitly, let $P_n(\varepsilon)$ be the probability that X_n is outside the ball of radius ε centered at X . Then X_n is said to converge in probability to X if for any $\varepsilon > 0$ and any $\delta > 0$ there exists a number N (which may depend on ε and δ) such that for all $n \geq N$, $P_n(\varepsilon) < \delta$,

- Idea: probability of *unusual* outcomes becomes smaller as n increases.
- Convergence in probability implies convergence in distribution

Convergence in probability is closely related to an important characteristic of many estimators:

Definition A.4 Consistent estimator

An estimator T_n of parameter θ is said to be **weakly consistent**, if it converges in probability to the true value of the parameter:

$$\text{plim}_{n \rightarrow \infty} T_n = \theta,$$

i.e. if, for all $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|T_n - \theta| > \varepsilon) = 0.$$

A very powerful theorem is the following which tells us that usual operations such as addition of random variables has a nice “consistency” property. Part (3) of the theorem allows us to replace covariances with estimated covariances for asymptotic distributions of test statistics.

Theorem A.1 Slutsky's theorem

Let X_n, Y_n be sequences of scalar/vector/matrix random elements. If X_n converges in distribution to a random element X and Y_n converges in probability to a constant c , then

1. $X_n + Y_n \xrightarrow{d} X + c$;
2. $X_n Y_n \xrightarrow{d} Xc$;
3. $X_n/Y_n \xrightarrow{d} X/c$, provided that c is invertible,

where \xrightarrow{d} denotes convergence in distribution.