

# Theory of Inference - Reviewed Notes

Dom Hutchinson

March 10, 2020

## Contents

<b>1</b>	<b>General</b>	<b>2</b>
1.1	Approaches to Inference . . . . .	2
1.2	Models . . . . .	2
1.3	Inference by Mathematical Manipulation . . . . .	3
1.4	Hypothesis Testing . . . . .	4
1.5	Intervals . . . . .	5
1.6	Causality . . . . .	5
1.6.1	Controlled Experiments . . . . .	5
1.6.2	Instrumental Variables . . . . .	5
1.7	Regularity Conditions . . . . .	5
<b>2</b>	<b>Linear Models</b>	<b>5</b>
2.1	Frequentist Approach . . . . .	6
2.1.1	Estimation . . . . .	6
2.1.2	Checking . . . . .	7
2.1.3	Evaluating . . . . .	8
2.1.4	Hypothesis Testing & Intervals . . . . .	8
2.2	Bayesian Approach . . . . .	10
2.3	Beyond . . . . .	10
<b>3</b>	<b>Maximum Likelihood Estimation</b>	<b>10</b>
3.1	By Calculus . . . . .	10
3.2	Numerical Optimisation . . . . .	11
3.3	Cramer-Rao Bound . . . . .	11
3.4	Hypothesis Testing . . . . .	11
3.5	Intervals . . . . .	11
<b>0</b>	<b>Appendix</b>	<b>12</b>
0.1	Definitions . . . . .	12
0.2	Theorems . . . . .	13

# 1 General

## 1.1 Approaches to Inference

### Definition 1.1 - *Statistical Inference*

*Statistical Inference* is the process of taking some data and inferring a property of the world from it. This is done by theorising a *Statistical Model* which may have generated the data and then calculating parameters for it from the data.

### Definition 1.2 - *Statistical Model*

*Statistical Models* are a, simplified, mathematical description for how a set of data could have been generated. In particular, a *Statistical Model* describes the random variability in the data generating process.

### Definition 1.3 - *Frequentist Inference*

The *Frequentist Approach to Statistical Inference* treats model unknowns (parameters or functions) as fixed states of nature whose values we want to estimate.

There is no modelling of random variability and thus any that occurs during data collection will be inherited by the model.

### Remark 1.1 - *Frequentist Inference*

Often in *Frequentist Inference* we use *asymptotic results* which only become exact as the sample size tends to infinity. This has practical drawbacks.

### Definition 1.4 - *Bayesian Inference*

The *Bayesian Approach to Statistical Inference* treats unknown model parameters as random variables. We define our initial uncertainty about parameter values (the *Prior Distribution*,  $\mathbb{P}(\Theta)$ ), observed data is used to update these distributions in order to reach a *Posterior Distribution*,  $\mathbb{P}(\Theta|X)$ .

*N.B.* This is done by using *Bayes' Theorem*.

### Remark 1.2 - *Bayesian Inference*

Often in *Bayesian Inference* we use *simulation methods*, which only become exact as the sample size tends to infinity. Again, there are practical drawbacks to this.

### Remark 1.3 - *Statistical Design*

When trying to infer a model from data there are a few common questions we ask

- i) What range of parameter values are consistent with the data?
- ii) Which of several alternative models could most plausibly have generated the data?
- iii) Could our model have generated the data at all?
- iv) How could we better arrange the data gathering process to improve the answers to the preceding questions?

## 1.2 Models

### Definition 1.5 - *Predictor Variables*

*Predictor Variables* are the dependent variables of a system, whose values we observe.

*N.B.* Typically denoted  $\mathbf{x}$  or  $\mathbf{X}$ .

**Definition 1.6 - Metric**

*Metrics* are *Predictor Variables* which measure an explicit quantity.

**Definition 1.7 - Factor**

*Factors* are *Predictor Variables* which act as labels to whether an observation belongs in a particular class due a property which cannot be explicitly quantified. (e.g. Male or Female).

**Definition 1.8 - Response Variables**

*Response Variables* are the independent variables of a system, whose value we observe. *N.B.* Typically denoted  $y$  or  $\mathbf{y}$ .

**Definition 1.9 - Fitted Values,  $\hat{y}$** 

*Fitted Values* are our estimated values for the *Response Variable*.

$$\hat{y}_i := f(\mathbf{x}_i)$$

**Definition 1.10 - Residual**

The *Residual* is the difference between the true value of the *Response Variables* & our *Fitted Values*.

$$\epsilon := |y_i - \hat{y}_i|$$

**Definition 1.11 - Residual Sum of Squares**

The *Residual Sum of Squares* is the sum of the squared value of the *Residuals* for each observation.

The *RSS* is used as a measure for how well our model fits the data

$$RSS := \sum_{i=1}^N \epsilon_i^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

**1.3 Inference by Mathematical Manipulation****Remark 1.4 - Inference by Mathematical Manipulation**

*Bayesian* and *Frequentist Inference* use mathematical computation to make inferences about parameter values & their uncertainty. An alternative approach is to use mathematical manipulation, rather than computation.

**Definition 1.12 - Bootstrapping**

Let  $X$  be a set of observed data.

*Bootstrapping* is a *simulation* of the data gathering process.

In *Bootstrapping* we uniformly sample values from  $X$  with replacement until we reach a desired sample size (often  $|X|$ ).

**Proposition 1.1 - Inference by Resampling**

Once we have used *Bootstrapping* to generate a set of new data-sets we can use *Bayesian* & *Frequentist Inference* techniques in order to estimate parameter values.

**Remark 1.5 - Bootstrap Interval**

An *Interval* generated from *Bootstrapping* datasets are generally narrower than those produced by *Bayesian* or *Frequentist Approaches*.

*N.B.* This discrepancy is reduced as sample size increases.

**Proposition 1.2 - Bootstrap Percentiles**

When wishing to create an interval for  $\theta$  using *Bootstrapped* data, we treat the  $\hat{\theta}$  values as if they came from  $\mathbb{P}(\mu|X)$ .

## 1.4 Hypothesis Testing

### Definition 1.13 - Simple Hypothesis

A *Simple Hypothesis* states that a parameter takes an exact value.

i.e.  $\theta = \theta_0$  for  $\theta_0 \in \Theta$ .

### Definition 1.14 - Composite Hypothesis

A *Composite Hypothesis* states that a parameter takes a value from a set.

i.e.  $\theta \in \Theta_0$  for  $\Theta_0 \subseteq \Theta$ .

### Definition 1.15 - Test Statistic

A *Test Statistic* is a random variable whose value depends on the observed set of data.

*Test Statistics* are used to assess the likelihood of observing a certain data set under a given *Null Hypothesis* in *Hypothesis Testing*.

### Definition 1.16 - p-Value

The *p-Value* of a *Test Statistic* is the probability of observing a value more extreme than the one produced by the observed data,  $\mathbf{x}$ , under the *Null Hypothesis*.

$$p(\mathbf{x}) := \sup_{\theta \in \Theta_0} \mathbb{P}(T(\mathbf{X}) \geq T(\mathbf{x}); \theta)$$

i.e. Under the null-hypothesis what is the probability of observing a more extreme test statistic value.

**Remark 1.6** -  $p(\mathbf{x})$  is the smallest Significance Level at which we would reject the Null Hypothesis

### Definition 1.17 - Hypothesis Testing

*Hypothesis Testing* is the process of determining which of two hypotheses about model parameters is more consistent with the data.

We define a *Null Hypothesis* and an *Alternative Hypothesis*. The *Null Hypothesis* acts as our default position, and we only reject it if the observed data is too extreme (given that it is true).

N.B. *Null* and *Alternative Hypothesis* are mutually exclusive.

### Proposition 1.3 - Process for Hypothesis Testing

Let  $\mathbf{x}$  be a realisation of  $\mathbf{X}$

- i) Choose a model  $f(\cdot; \theta)$  st  $\mathbf{X} \sim f(\cdot; \theta)$  for  $\theta \in \Theta$ .
- ii) Define a *Null Hypothesis*,  $H_0$ , and an *Alternative Hypothesis*,  $H_1$ .
- iii) Define a *Test Statistic*,  $T(\cdot)$ .
- iv) Choose a *Significance Level*,  $\alpha$ , and calculate the equivalent *Critical Value*,  $c$ , for the *Test Statistic*.
- v) Calculate value of the *Test Statistic* under the observed data,  $t_{\text{obs}} = T(\mathbf{x})$ .
- vi) If  $t_{\text{obs}} \geq c$  then reject  $H_0$  in favour of  $H_1$ , otherwise accept  $H_0$ .

### Definition 1.18 - Power Function, $\pi$

The *Power Function*,  $\pi(\cdot)$ , measures the probability of rejecting the *Null-Hypothesis* given that another set of parameter values is true (usually test with the *Alternative Hypothesis*).

Let  $\mathbf{X} \sim f(\cdot; \theta)$ ,  $T(\cdot)$  be a *Test Statistic* and  $c$  be the *Critical Value* of  $T$ . Then

$$\pi(\theta_1; T, c) = \mathbb{P}(T(\mathbf{X}) \geq c; \theta_1)$$

## 1.5 Intervals

### Definition 1.19 - Random Interval

Let  $\mathbf{X} \sim f_n(\cdot; \theta^*)$  for  $\theta^* \in \Theta$  and  $L, U : \mathcal{X}^n \rightarrow \Theta$  st  $\forall \mathbf{x} \in \mathcal{X}^n$   $L(\mathbf{x}) < U(\mathbf{x})$ .

A *Random Interval* is an *Interval* whose bounds depends on a *Random Variable*.

Here  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  is a *Random Interval*.

N.B.  $L(\cdot)$  &  $U(\cdot)$  are maps from observed data to parameter values.

### Definition 1.20 - Coverage of an Interval

Let  $\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})]$  be a *Random Interval* for  $\theta$  with true value  $\theta^*$ .

The *Coverage of an Interval* is the probability that the true value of the parameter it is estimating lies in the interval.

$$C_{\mathcal{I}} = \mathbb{P}(\theta^* \in \mathcal{I}(\mathbf{X}); \theta^*)$$

### Definition 1.21 - Confidence Interval

A  $1 - \alpha$  *Confidence Interval* for a parameter is an interval with *Coverage* at least  $1 - \alpha$ .

$$\mathcal{I}(\mathbf{X}) := [L(\mathbf{X}), U(\mathbf{X})] \text{ is a } 1 - \alpha \text{ Confidence Interval if } \mathbb{P}(\theta^* \in \mathcal{I}) \geq 1 - \alpha$$

N.B. If  $\mathbb{P}(\theta^* \in \mathcal{I}(\mathbf{X})) = 1 - \alpha$  then  $\mathcal{I}$  is an Exact *Confidence Interval*.

**Remark 1.7** - *Confidence Intervals are a part of Frequentist Statistics, not Bayesian*

### Definition 1.22 - Credible Interval

**Remark 1.8** - *Credible Intervals are a part of Bayesian Statistics, not Frequentist*

## 1.6 Causality

### Proposition 1.4 - Causality v Causation

### Definition 1.23 - Confounding

### Definition 1.24 - Counfounding Variable

### 1.6.1 Controlled Experiments

### Definition 1.25 - Randomisation

### 1.6.2 Instrumental Variables

### Definition 1.26 - Instrumental Variables

## 1.7 Regularity Conditions

## 2 Linear Models

### Proposition 2.1 - Implementing Factors

Suppose a model has *Factor Variable*  $g_i$  which separates observations into  $n$  categories.

In the function for  $y_i$ ,  $g_i$  would be represented by a single term  $\gamma_{g_i}$  whose value depends on the

value of  $g_i$ . (*i.e.* A different weight is assigned to each group).

We want to find the  $n$  values which  $\gamma_{g_i}$  can take.

We can express this in terms of matrices, as below, with each row on the LHS denoting which category each observation belongs to

$$\begin{pmatrix} 0 & 1 & 0 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_n \end{pmatrix}$$

*N.B.* Each row has a single 1 element and  $n - 1$  zero elements.

**Remark 2.1** - We want to find the parameters to the Predictor Variables which produce accurate values for the Response Variables.

**Definition 2.1** - *Model Matrix*

Each element in a *Model Matrix* is a function of the *Predictor Variables*.

Each row depends on a different set of observations.

**Definition 2.2** - *Linear Model*

A *Linear Model* is a *Statistical Model* whose response vector,  $\mathbf{y}$ , is linear wrt its *Model Matrix*,  $\mathbf{X}$ , and some zero-mean random error,  $\boldsymbol{\varepsilon}$ .

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$\mathbf{X} \in \mathbb{R}^{n \times p}$ ,  $\boldsymbol{\beta} \in \mathbb{R}^p$  and  $\boldsymbol{\varepsilon} \sim \text{Normal}(\mathbf{0}, \sigma^2 I)$ .

## 2.1 Frequentist Approach

**Proposition 2.2** - *Frequentist Approach to Linear Models*

In the *Frequentist Approach to Linear Models* we treat  $\boldsymbol{\beta}$  and  $\sigma^2$  as fixed (but unknown) states of nature.

Thus all random variability from the data will be inherited into the model.

### 2.1.1 Estimation

**Proposition 2.3** - *Point Value Estimates*

We can make *Point Value Estimates* of parameter values by finding the set of parameters  $\boldsymbol{\beta}$  which minimises the *Residual Sum of Squares*.

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{LSE}} &:= \underset{\boldsymbol{\beta}}{\text{argmin}} \sum_{i=1}^N (y_i - (\mathbf{X}\boldsymbol{\beta})_i)^2 \\ &= \underset{\boldsymbol{\beta}}{\text{argmin}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \end{aligned}$$

*N.B.* This is the *Least Squares Estimate* of  $\boldsymbol{\beta}$ .

**Remark 2.2** -  $\hat{\boldsymbol{\beta}}_{\text{LSE}} = R^{-1}Q^T \mathbf{y}$

where  $Q, R$  are from the decomposition of  $\mathbf{X}$  st  $\mathbf{X} = QR$  with  $Q$  being *Orthogonal* and  $R$  being *Upper-Triangle*.

**Proposition 2.4** - *Deriving Least Squares Estimate for  $\boldsymbol{\beta}$*

Let  $\mathbf{X}, \mathbf{y}$  be  $n$  observed data points &  $\boldsymbol{\beta} \in \mathbb{R}^p$  be a parameter vector we are fitting to our model.

We want to find  $\hat{\beta}_{\text{LSE}} = \operatorname{argmin}_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2$ .

Since  $\mathbf{X}$  is a real-valued matrix it can be decomposed into

$$\mathbf{X} = \mathcal{Q} \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} = \mathcal{Q} R^1$$

where  $R \in \mathbb{R}^{p \times p}$  is an upper triangle matrix,  $\mathcal{Q} \in \mathbb{R}^{n \times n}$  is orthogonal &  $Q \in \mathbb{R}^{n \times p}$  is the first  $p$  columns of  $\mathcal{Q}$ .

Note that  $\mathcal{Q}^T$  is *Orthogonal*.

Thus

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\beta\|^2 &= \left\| \mathbf{y} - \mathcal{Q} \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \beta \right\|^2 \\ &= \left\| \mathcal{Q}^T \mathbf{y} - \mathcal{Q}^T \mathcal{Q} \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \beta \right\|^2 \quad \text{since } \mathcal{Q}^T \text{ is orthogonal} \\ &= \left\| \mathcal{Q}^T \mathbf{y} - \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \beta \right\|^2 \end{aligned}$$

Decompose  $\mathcal{Q}^T \mathbf{y} = \begin{pmatrix} \mathbf{f} \\ \mathbf{r} \end{pmatrix}$  with  $\mathbf{f} \in \mathbb{R}^p$  &  $\mathbf{r} \in \mathbb{R}^{n-p}$ .

Note that  $\mathbf{f} = Q^T \mathbf{y}$ .

$\mathbf{f}$  is the first  $p$  rows of  $\mathcal{Q}^T \mathbf{y}$  and  $\mathbf{r}$  is the last  $n - p$  rows.

Thus

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\beta\|^2 &= \left\| \begin{pmatrix} \mathbf{f} \\ \mathbf{r} \end{pmatrix} - \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \beta \right\|^2 \\ &= \|\mathbf{f} - R\beta\|^2 + \|\mathbf{r}\|^2 \end{aligned}$$

$\|\mathbf{r}\|^2$  is indepdent of  $\beta$  and thus irreducible.

This final expression is minimised when  $\|\mathbf{f} - R\beta\|^2 = 0$  (Meaning  $\|\mathbf{r}\|^2 = \|\mathbf{y} - \mathbf{X}\beta\|^2$ ).

Thus

$$\hat{\beta}_{\text{LSE}} = R^{-1} \mathbf{f} = R^{-1} Q^T \mathbf{y}$$

This requires that  $R$  is full rank, in order for its inverse to exist.

Further,  $\mathbf{X}$  has to have full rank, which we can ensure by our design of the model.

**Proposition 2.5** - *Least Squares Estimate of Parameter Vector is Unbiased*

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(R^{-1} Q^T \mathbf{y}) = R^{-1} Q^T \mathbb{E}(\mathbf{y}) = R^{-1} Q^T \mathbf{X} \beta = R^{-1} Q^T Q R \beta = \beta$$

**Proposition 2.6** - *Variance of Least Squares Estimate of Parameter Vector*

$$\begin{aligned} \implies \quad \begin{aligned} \operatorname{Cov}(\mathbf{y}) &= I\sigma^2 \\ \operatorname{Cov}(\mathbf{f}) &= Q^T \mathbf{y} \\ &= Q^T Q \sigma^2 \\ &= I\sigma^2 \end{aligned} \\ \implies \quad \begin{aligned} \operatorname{Cov}(\hat{\beta}_{\text{LSE}}) &= \operatorname{Cov}(R^{-1} \mathbf{f}) \\ &= R^{-1} \operatorname{Cov}(\mathbf{f}) R^{-T} \\ &= R^{-1} I\sigma^2 R^{-T} \\ &= R^{-1} R^{-T} \sigma^2 \end{aligned} \end{aligned}$$

### 2.1.2 Checking

**Remark 2.3** - *Assumptions*

We assume that each  $\varepsilon_i$  is independent & has constant variance (we also assume they are normally distributed but this generally holds due to CLT).

---

<sup>1</sup>Known as *QR Decomposition* and can be performed in R using `qr.Q(qr(X), complete = TRUE)` & `qr.R(qr(X))`

We need a way to check this assumption holds in order for inferences (beyond point estimates) to be sound.

**Proposition 2.7 - Graphical Checks**

Plotting  $\hat{\epsilon} = y_i - (\mathbf{X}\hat{\boldsymbol{\beta}})_i$  on a graph tends to indicate whether an assumption has been broken, and if so, how it was broken.

- Systematic patterns in the mean indicate independence assumption is broken.
- Systematic patterns in the variability indicate the constant variance assumption is broken.

### 2.1.3 Evaluating

**Remark 2.4 - Choice of measure to minimise?**

Was choosing to minimise *Residual Sum of Squares* a good one?

*N.B.* Choosing  $\sum_i |\epsilon_i|$ ,  $\sum_i \epsilon_i^4, \dots$  could have worked.

**Remark 2.5 - Problem with  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|$  as measure**

Suppose our data has a lot of information for estimating  $\beta_i$  but not mch for  $\beta_j$ , should we weight them equally?

**Remark 2.6 - Preferred Estimators**

We require estimators to be *Unbiased*, and then we shall choose the estimator with the least variance among those which are *Unbiased*.

*N.B.* Least variance means smallest covariance matrix (in a way which accounts for weighting individual parameters).

**Theorem 2.1 - Gauss Markov Theorem**

Let  $\mathbf{X}, \mathbf{y}$  be some observed data.

Consider a model where  $\boldsymbol{\mu} := \mathbb{E}(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$  and  $\Sigma_y^2 = \sigma^2 I$ .

Let  $\tilde{\theta} := \mathbf{c}^T \mathbf{y}$  be any *Unbiased Linear Estimator* of  $\theta = \mathbf{t}^T \boldsymbol{\beta}$  for some arbitrary vector,  $\mathbf{t}$ .

Then

$$\text{Var}(\tilde{\theta}) \geq \text{Var}(\hat{\theta})$$

where  $\hat{\theta} = \mathbf{t}^T \hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\beta}} = R^{-1} Q^T \mathbf{y}$  where  $\mathbf{X} = QR$ .

Thus each element of  $\hat{\boldsymbol{\beta}}$  is a *minimum variance unbiased estimator*, since  $\mathbf{t}$  is arbitrary.

### 2.1.4 Hypothesis Testing & Intervals

**Remark 2.7 - Populat Hypothesis Test**

Often we want to test whether any  $\beta_i = 0$  as this would indicate that those predictors do not affect the model accuracy.

**Proposition 2.8 - Distribution of  $\hat{\boldsymbol{\beta}}$**

We assume that  $\epsilon_i \sim \text{Normal}(0, \sigma^2)$ . Thus

$$\begin{aligned} \mathbf{y} &\sim \text{Normal}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 I) \\ \implies \hat{\boldsymbol{\beta}} &\sim \text{Normal}(\boldsymbol{\beta}, R^{-1} R^{-T} \sigma^2) \end{aligned}$$

Note that  $\boldsymbol{\beta}$  and  $\sigma^2$  are unknown.

*N.B.*  $\mathbf{X} = QR$  where  $Q$  is orthogonal &  $R$  upper-triangle.

**Proposition 2.9 -  $\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_{\hat{\beta}_i}} \sim t_{n-p}$**

Note that we can produce a decomposition  $\mathbf{X} = Q \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix}$  where  $Q$  is orthogonal &  $R$  is upper



triangular.

We have

$$\text{Cov}(\mathcal{Q}^T \mathbf{y}) = \mathcal{Q}^T \text{Cov}(\mathbf{y}) \mathcal{Q}^{-T} = \mathcal{Q}^T \text{Cov}(\mathbf{y}) \mathcal{Q} = \mathcal{Q}^T I \sigma^2 \mathcal{Q} = I \sigma^2$$

This implies that elements of  $\mathcal{Q}^T \mathbf{y}$  are independent, due to their assumed normal distribution. Note that

$$\mathbb{E}(\mathcal{Q}^T \mathbf{y}) = \mathbb{E} \left( \begin{pmatrix} \mathbf{f} \\ \mathbf{r} \end{pmatrix} \right) \quad \text{and} \quad \mathbb{E}(\mathcal{Q}^T \mathbf{y}) = \mathcal{Q}^T \mathbb{E}(\mathbf{y}) = \mathcal{Q}^T \mathbf{X} \boldsymbol{\beta} = \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \boldsymbol{\beta}$$

Thus

$$\mathbb{E} \left( \begin{pmatrix} \mathbf{f} \\ \mathbf{r} \end{pmatrix} \right) = \begin{pmatrix} R \\ \mathbf{0} \end{pmatrix} \boldsymbol{\beta} \implies \mathbb{E}(\mathbf{f}) = R \boldsymbol{\beta} \ \& \ \mathbb{E}(\mathbf{r}) = \mathbf{0}$$

Further

$$\mathbf{f} \sim \text{Normal}(R \boldsymbol{\beta}, I_p \sigma^2) \quad \text{and} \quad \mathbf{r} \sim \text{Normal}(\mathbf{0}, I_{n-p} \sigma^2)$$

and  $\mathbf{f}$  &  $\mathbf{r}$  are independent.

Thus  $\hat{\boldsymbol{\beta}}$  &  $\hat{\sigma}^2$  are independent.

Since each  $r_i \sim \text{Normal}(0, \sigma^2)$

$$\frac{\|\mathbf{r}\|^2}{\sigma^2} = \frac{1}{\sigma^2} \sum_{i=1}^{n-p} r_i^2 \sim \chi_{n-p}^2$$

Since  $\mathbb{E}(\chi_{n-p}^2) = n - p \implies \hat{\sigma}^2 := \frac{1}{n-p} \|\mathbf{r}\|^2$  is an unbiased estimator of  $\sigma^2$ .

$\hat{\Sigma}_{\hat{\boldsymbol{\beta}}} := \Sigma_{\hat{\boldsymbol{\beta}}} \frac{\hat{\sigma}^2}{\sigma^2} = R^{-1} R^{-T} \hat{\sigma}^2$  is an unbiased estimator of  $\Sigma_{\hat{\boldsymbol{\beta}}}$ .

Thus  $\hat{\sigma}_{\hat{\beta}_i} := \sqrt{[\hat{\Sigma}_{\hat{\boldsymbol{\beta}}}]_{ii}} = \sigma_{\hat{\beta}_i} \frac{\hat{\sigma}}{\sigma}$ .

Finally

$$\frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_{\hat{\beta}_i}} = \frac{\hat{\beta}_i - \beta_i}{\sigma_{\hat{\beta}_i} \frac{\hat{\sigma}}{\sigma}} = \frac{\frac{1}{\sigma_{\hat{\beta}_i}} (\hat{\beta}_i - \beta_i)}{\sqrt{\hat{\sigma}^2 / \sigma^2}} = \frac{\frac{1}{\sigma_{\hat{\beta}_i}} (\hat{\beta}_i - \beta_i)}{\sqrt{\frac{1}{\sigma^2} \frac{1}{n-p} \|\mathbf{r}\|^2}} \sim \frac{\text{Normal}(0, 1)}{\sqrt{\frac{1}{n-p} \chi_{n-p}^2}} \sim t_{n-p}$$

N.B.  $\|\mathbf{r}\|^2 = \|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2$  by the results in **Proposition 2.4**.

**Proposition 2.10** - *Confidence Interval for  $\beta_i$*

Using the result in **Proposition 2.9** we can construct the following  $1 - \alpha$  confidence interval

$$\mathbb{P} \left( -t_{n-p, \frac{\alpha}{2}} < \frac{\hat{\beta}_i - \beta_i}{\hat{\sigma}_{\hat{\beta}_i}} < t_{n-p, \frac{\alpha}{2}} \right) = \mathbb{P} \left( \hat{\beta}_i - t_{n-p, \frac{\alpha}{2}} \sigma_{\hat{\beta}_i} < \beta_i < \hat{\beta}_i + t_{n-p, \frac{\alpha}{2}} \sigma_{\hat{\beta}_i} \right) = 1 - \alpha$$

**Proposition 2.11** - *Hypothesis Testing on  $\beta_i$*

Suppose we want to test  $H_0 : \beta_i = \beta_{i0}$  against  $H_1 : \beta_i \neq \beta_{i0}$ .

We use test statistic

$$T = \frac{\hat{\beta}_i - \beta_{i0}}{\hat{\sigma}_{\hat{\beta}_i}}$$

under  $H_0$   $T \sim t_{n-p}$  where  $n$  is the number of observations &  $p$  the number of parameters.

Thus we can assess the test using  $p = \mathbb{P}(|T| \geq |t_{obs}|)$ .

**Proposition 2.12** - *Testing Multiple Variables in a Model*

This can be expressed as the test of  $H_0 : \mathbf{C} \boldsymbol{\beta} = \mathbf{d}$  against  $H_1 : \mathbf{C} \boldsymbol{\beta} \neq \mathbf{d}$  where  $\mathbf{C} \in \mathbb{R}^{q \times p}$  &  $\mathbf{d} \in \mathbb{R}^q$  with  $q < p$ .

Under  $H_0$  we have  $(\mathbf{C} \hat{\boldsymbol{\beta}} - \mathbf{d}) \sim \text{Normal}(\mathbf{0}, \mathbf{C} \Sigma_{\hat{\boldsymbol{\beta}}} \mathbf{C}^T)$ .

We can produce a *Cholesky Decomposition*  $\mathbf{L}^T \mathbf{L} = \mathbf{C} \Sigma_{\hat{\beta}} \mathbf{C}^T$ .

Thus

$$\begin{aligned} \mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d}) &\sim \text{Normal}(0, I) \\ \Rightarrow (\mathbf{C}\hat{\beta} - \mathbf{d})^T (\mathbf{C} \Sigma_{\hat{\beta}} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\beta} - \mathbf{d}) &= (\mathbf{C}\hat{\beta} - \mathbf{d})^T \mathbf{L}^{-1} \mathbf{L}^{-T} (\mathbf{C}\hat{\beta} - \mathbf{d}) \\ &= \|\mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d})\|^2 \\ &\sim \sum_{i=1}^q \text{Normal}(0, 1)^2 \\ &\sim \chi_q^2 \end{aligned}$$

Setting  $\hat{\Sigma}_{\hat{\beta}} := \frac{\hat{\sigma}^2}{\sigma^2} \Sigma_{\hat{\beta}}$  we can produce a test statistic

$$F := \frac{1}{q} (\mathbf{C}\hat{\beta} - \mathbf{d})^T (\mathbf{C} \Sigma_{\hat{\beta}} \mathbf{C}^T)^{-1} (\mathbf{C}\hat{\beta} - \mathbf{d})$$

Which has the distribution

$$\begin{aligned} F &= \frac{1}{q} \|\mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d})\|^2 \\ &= \frac{\sigma^2}{q \hat{\sigma}^2} \|\mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d})\|^2 \\ &= \frac{\frac{1}{q} \|\mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d})\|^2}{\hat{\sigma}^2 / \sigma^2} \\ &= \frac{\frac{1}{q} \|\mathbf{L}^{-T}(\mathbf{C}\hat{\beta} - \mathbf{d})\|^2}{\frac{1}{\sigma^2} \frac{1}{n-p} \|\mathbf{r}\|^2} \\ &\sim \frac{\frac{1}{q} \chi_q^2}{\frac{1}{n-p} \chi_{n-p}^2} \\ &\sim F_{q, n-p} \end{aligned}$$

**Proposition 2.13** -  $F = \frac{\frac{1}{q}(RSS_0 - RSS_q)}{\frac{1}{n-p}RSS_1}$

Where  $RSS_0$  is the residual sum of squares when  $H_0$  is true and  $RSS_1$  is the residual sum of squares when  $H_1$  is true.

**Proposition 2.14** - *Testing whether a Factor Variable belongs in a Model*

*Factor Variables* have multiple parameters associated to them in a model and thus to test whether the *Factor Variable* should be in the model requires testing whether all of these parameters should equal 0.

This can be tested using the results in **Proposition 2.12** with  $\mathbf{d} = \mathbf{0}$  and  $\mathbf{C}$  is the rows of the  $I_p$  which indicate the parameters we wish to test.

In this case  $q$  is the number of parameters we wish to test.

## 2.2 Bayesian Approach

## 2.3 Beyond

**Definition 2.3** - *Linear Mixed Models*

# 3 Maximum Likelihood Estimation

## 3.1 By Calculus

**Definition 3.1** - *Likelihood*

**Definition 3.2** - *Log-Likelihood*

**Definition 3.3** - *Maximum Likelihood Estimator*

**Proposition 3.1** - *Consistency of MLE*

**Proposition 3.2** - *Large Sample Distribution of MLE*

## 3.2 Numerical Optimisation

**Definition 3.4** - *Numerical Optimisation*

**Definition 3.5** - *Objective Function*

**Definition 3.6** - *Newton's Method*

## 3.3 Cramer-Rao Bound

**Definition 3.7** - *Fisher Information*

**Definition 3.8** - *Fisher Information Matrix*

## 3.4 Hypothesis Testing

**Definition 3.9** - *Neyman-Pearson Lemma*

**Definition 3.10** - *Generalised Likelihood Ratio Test Statistic*

## 3.5 Intervals

## 0 Appendix

### 0.1 Definitions

#### Definition 0.1 - Parametric Models

*Parameteric Models* are *Statistical Models* whose only unknowns are parameters.

#### Definition 0.2 - Semi-Parametric Models

*Parameteric Models* are *Statistical Models* which contain unknown parameters and unknown functions.

#### Definition 0.3 - Non-Parametric Models

*Non-Parametric Models* make *few* prior assumptions about how data was generated and instead depend mainly on the observed data.

We cannot simulate data from *Non-Parameteric Models*.

#### Definition 0.4 - Orthogonal Matrix

A matrix  $\mathbf{X}$  is *Orthogonal* if

$$\mathbf{X}^T \mathbf{X} = \mathbf{X} \mathbf{X}^T = I \implies \mathbf{X}^T = \mathbf{X}^{-1}$$

*Orthogonal Matrices* rotate & reflect vectors without changing their magnitude.

N.B.  $\mathbf{X}^T$  is *Orthogonal*.

#### Definition 0.5 - Full Rank Matrix

Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$ .

If  $m > n$  then  $\mathbf{X}$  has *Full Rank* iff all its columns are linearly independent.

If  $n > m$  then  $\mathbf{X}$  has *Full Rank* iff all its rows are linearly independent.

N.B. In statistics the number of  $m > n$  always as we should have more observations than fields.

#### Definition 0.6 - Upper Triangle Matrix

A matrix  $X$  is an *Upper Triangle Matrix* if  $X_{i,j} = 0$  for  $i > j$ .

#### Definition 0.7 - Unbiased Estimator

An *Estimator* of a parameter,  $\hat{\theta}$ , is unbiased if its expected value is the true value of the parameter for all possible parameter values

$$\mathbb{E}(\hat{\theta};) = \theta^*$$

#### Definition 0.8 - Conjugacy

#### Definition 0.9 - Fisher Information

#### Definition 0.10 - Correlation

#### Definition 0.11 - Covariance

#### Definition 0.12 - Expected Value

#### Definition 0.13 - Variance

#### Definition 0.14 - Positive Semi-Definite Matrix

#### Definition 0.15 - Taylor's Theorem

## 0.2 Theorems

### Theorem 0.1 - Bayes' Theorem

Suppose  $X \sim f(\cdot; \Theta)$ . Then

$$\underbrace{\mathbb{P}(\Theta|X)}_{\text{Posterior}} = \frac{\overbrace{\mathbb{P}(X|\Theta)}^{\text{Likelihood}} \overbrace{\mathbb{P}(\Theta)}^{\text{Prior}}}{\underbrace{\mathbb{P}(X)}_{\text{Evidence}}}$$

### Theorem 0.2 - Euclidean Distance Identities

$$\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \sum_{i=1}^n x_i^2$$