



Regression Analysis for Data Scientists (ZZSC5806)

Atefeh Zamani

University of New South Wales
School of Mathematics and Statistics

H2- 2023

Notations

- **Random Variables:** Y_1, Y_2, \dots, Y_n (not bold, CAPITAL LETTERS)
- **Realisations:** y_1, y_2, \dots, y_n (not bold, small letters)
- **Parameters:** α, β , etc. (Greek Letters)
- **Estimators:** $\hat{\alpha}, \hat{\beta}$, etc.

- **Vectors:** $\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$ (not italic/italic, small letters, **bold**)

- **Matrices:** $\mathbf{X} = \begin{bmatrix} X_{11} & \dots & X_{1p} \\ X_{21} & \dots & X_{2p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \dots & X_{np} \end{bmatrix}$ (not italic/italic, CAPITAL LETTERS, **bold**)

- **Transpose:** $\mathbf{y}^T, \mathbf{y}^T, \mathbf{X}^T$.

Week 1: Estimation and Inference for Regression Models

This week we have a quick review on introduction to regression modelling and key statistical concepts that we will use throughout the course. This includes:

- estimation methods,
- maximum likelihood estimation,
- least squares estimation,
- exponential family of distributions,
- prediction,
- inference.

1.1 Introduction to regression analysis

- **What is regression Analysis?**

Regression Analysis investigates the functional **relationship** between **statistical variables**. Data are usually **multiple observations** of a random vector (Y, \mathbf{x}) .

- $\mathbf{x} = (X_1, \dots, X_p)^\top$ is a p -vector of variables (explanatory variables, regressors, predictors, input variables or independent variables).
- Y may be continuous ($\in \mathbb{R}$), discrete ($\in \{1, \dots, K\}$) or ordinal (ordered discrete)(response variable, target variable, output variable, outcome variable or dependent variable).

Response variables are usually treated as **random variables**, while **predictors** are treated as **fixed observations**.

1.1 Introduction to regression analysis

- **Response variable**

- **nominal**: categories (Male/Female or Car/Bus/Bike)
- **ordinal**: classes with natural order or ranking (Young/Middle-Aged/Old)
- **continuous**: continuous scale, at least in theory (Time, weight and distance).

Nominal and ordinal data can be qualitative or quantitative.

- They are called categorical or discrete variables.
- The numbers of observations, counts or frequencies in each category are usually recorded.

Continuous data are quantitative.

- **Explanatory variables**

- quantitative (Covariate) or qualitative (Factor).

1.1 Introduction to regression analysis

Response	Explanatory	Method
Continuous	Binary Nominal, > 2categories Ordinal Continuous	t-test ANOVA ANOVA Multiple regression
Binary	Categorical Continuous	Contingency tables Logistic or probit regression
Nominal, > 2categories	Nominal Categorical & Continuous	Contingency tables Nominal logistic regression
Ordinal	Categorical & Continuous	Ordinal logistic regression
Counts	Categorical Categorical & Continuous	Log-linear models Poisson regression

1.1 Introduction to regression analysis

- Regression

Aim: Find a "good" functional relationship of the form $Y = f(\mathbf{x}) + \varepsilon$.

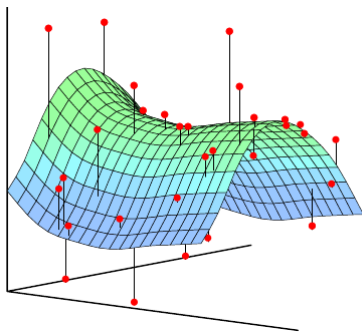


Figure: Regression of Y (vertical, continuous) on (X_1, X_2) (horizontal).

1.1 Introduction to regression analysis

- **General framework of statistical learning**

- **Statistical learning**: A vast set of tools for understanding data.
 - supervised
 - unsupervised
- **Regression**: statistical model for predicting or estimating an output based on one or more inputs (Supervised).

In contrast, unsupervised methods cover situations where there are inputs but no supervising output. In these type of analysis we learn about relationships and structure of data. Example of unsupervised analysis is cluster analysis.

1.3 Fundamental definitions

- **The purpose of regression analysis:** Prediction and Inference

- **Prediction**

In many situations, \mathbf{x} is available but the output Y cannot be easily obtained.

Since $E(\varepsilon) = 0$,

$$\hat{Y} = \hat{f}(X), \quad (1.3.1)$$

where

- \hat{f} : estimate for f
- \hat{Y} : the resulting prediction for Y .

The estimate \hat{f} is characterised by a **reducible** error and by an **irreducible** error.

$$\mathbb{E}[Y - \hat{Y}]^2 = \mathbb{E}[f(X) + \varepsilon - \hat{f}(X)]^2 = [f(X) - \hat{f}(X)]^2 + \mathbb{V}ar(\varepsilon) \quad (1.3.2)$$

1.3 Fundamental definitions

- **The purpose of regression analysis**

- **Inference**

Goal: Understanding the relationship between \mathbf{x} and Y . How Y changes as a function of X_1, X_2, \dots, X_p .

- Which predictors are associated with the response?
 - What is the relationship between the response and each predictor?
 - Is the relationship between Y and each predictor linear or more complex?

Depending on the purpose of the analysis, various methods of estimating f may be more appropriate.

- **Prediction:** Non-linear models (quite accurate predictions, but at the expense of a less interpretable model for which inference is more challenging. In some cases overfitting may cause problems)
 - **Inference:** linear models (relatively simple and interpretable inference, but may not yield accurate predictions)

1.3 Fundamental definitions

- **Method of estimation**

- **Parametric:**

- Step 1: Assumption on the functional form of f

$$f(\mathbf{x}) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \quad (1.3.3)$$

- Step 2: Fit the model (estimating $\beta_0, \beta_1, \dots, \beta_p$)

- **Nonparametric:** no explicit assumptions about the functional form of f at the price of not having a small and fixed amount of parameters (i.e. larger sample sizes are needed)

1.3 Fundamental definitions

- **Maximum likelihood estimation (MLE)**

- $\mathbf{y} = [Y_1, Y_2, \dots, Y_n]^\top$: a random vector.
- $f(\mathbf{y}; \theta)$: **joint** probability density function of Y_i , depends on $\theta = [\theta_1, \theta_2, \dots, \theta_p]^\top$.

The **likelihood function** $L(\theta; \mathbf{y})$ is the same as $f(\mathbf{y}; \theta)$ but the emphasis is on θ while \mathbf{y} stays fixed.

Note: L is a random variable!

The maximum likelihood estimator of θ

MLE of θ is the value $\hat{\theta}$ which maximizes the likelihood function:

$$L(\hat{\theta}, \mathbf{y}) \geq L(\theta, \mathbf{y}) \quad \text{for all } \theta \in \Theta \quad (\text{parameter space}). \quad (1.3.4)$$

Equivalently, $\hat{\theta}$ is the value that maximizes the log-likelihood function $l(\theta, \mathbf{y}) = \ln L(\theta, \mathbf{y})$.

1.3 Fundamental definitions

- **Cont. Maximum likelihood estimation (MLE)**

- **Method**

1. Differentiate $l(\theta, \mathbf{y})$ with respect to each element θ_j of θ
2. Solving the simultaneous equations

$$\frac{\partial l(\theta, \mathbf{y})}{\partial \theta_j} = 0 \quad j = 1, \dots, p \quad (1.3.6)$$

3. If the matrix of second derivatives

$$\frac{\partial^2 l(\theta, \mathbf{y})}{\partial \theta_j \partial \theta_k} \quad (1.3.7)$$

evaluated at $\theta = \hat{\theta}$ is negative definite, then $\hat{\theta}$ maximizes $l(\theta, \mathbf{y})$ in the **interior** of Θ .

Note: Check if there are any values of θ at the **edges** of Θ that give local maxima of $l(\theta, \mathbf{y})$. When all local maxima have been identified, the value of $\hat{\theta}$ corresponding to the largest one is the MLE.

1.3 Fundamental definitions

- **Cont. Maximum likelihood estimation:**

- **Example** (Poisson distribution)

Let Y_1, Y_2, \dots, Y_n be independent random variables with Poisson distribution

$$f(y_i, \theta) = \frac{\theta^{y_i} e^{-\theta}}{y_i!},$$

$y_i = 0, 1, \dots$, with the same parameter θ . Find the MLE of θ .

1.3 Fundamental definitions

- **Property of MLE**

- **Invariance**: If $g(\theta)$ is any function of θ , then the MLE of $g(\theta)$ is $g(\hat{\theta})$.
- **Consistency**: If $\hat{\theta}$ is the MLE of θ , then $\hat{\theta}$ is consistent for θ ; i.e.,

$$\Pr\left(|\hat{\theta}_n - \theta| \geq \varepsilon\right) \rightarrow 0, \text{ as } n \rightarrow \infty \quad \equiv \quad \hat{\theta}_n \xrightarrow{P} \theta$$

- **Asymptotic normality**: If $\hat{\theta}$ is the MLE of θ , then

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} N\left(0, \frac{1}{I(\theta)}\right), \quad I(\theta) : \text{Fisher Information of } \theta$$

- **Asymptotic efficiency**: Generally, a CAN estimator $\hat{\theta}_n$ of θ is said to be asymptotically efficient if the asymptotic variance of $\hat{\theta}_n$ equals $(I(\theta))^{-1}$.
- **Sufficiency**: If a sufficient statistic exists for θ , then the MLE can be expressed as a function of this statistic.

1.3 Fundamental definitions

- **Least squares estimation**

Y_1, \dots, Y_n : independent r.v. with $\mathbb{E}(Y_i) = \mu_i(\beta)$, where $\beta = [\beta_1, \dots, \beta_p]^\top$, $p < n$, is the parameter vector.

We want to **estimate** β .

Least squares estimator: finding the estimator $\hat{\beta}$ that **minimizes the sum of squares (SS)** of the differences between Y_i 's and their expected values

$$SS = \sum_{i=1}^n [Y_i - \mu_i(\beta)]^2. \quad (1.3.8)$$

and $\hat{\beta}$ is obtained by differentiating SS with respect to the elements β_j :

$$\frac{dSS}{d\beta_j} = 0 \quad j = 1, \dots, p \quad (1.3.9)$$

1.3 Fundamental definitions

- **Weighted least squares**

If Y_i 's have variances σ_i^2 (not necessarily equal), **minimize the weighted sum of squared (WSS) differences**

$$WSS = \sum_{i=1}^n w_i [Y_i - \mu_i(\beta)]^2. \quad (1.3.10)$$

where $w_i = (\sigma_i^2)^{-1}$.

Idea: less reliable observations have less influence on the estimates.

General form: $\mathbf{y} = [Y_1, \dots, Y_n]^\top$ is a random vector with mean vector $\boldsymbol{\mu} = [\mu_1, \dots, \mu_n]^\top$ and variance-covariance matrix \mathbf{V} , then

$$WSS = (\mathbf{y} - \boldsymbol{\mu})^\top \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu}). \quad (1.3.11)$$

1.3 Fundamental definitions

- **Comments**

- Least squares estimator: No assumptions on the distribution of the response variables (In MLE, the distribution is needed).
- For many situations, MLE and least squares estimates are identical.
- In many cases, numerical methods are used for parameter estimation.

1.3 Fundamental definitions

- **Model fitting**

- Step 1. **Model specification**: an equation linking the response and the explanatory variables and a probability distribution for the response variable
- Step 2. **Estimation of the parameter of the model**
- Step 3. **Checking the adequacy of the model**
- Step 4. **Inference**: calculating confidence interval, testing hypotheses

1.3 Fundamental definitions

- **Some Examples**

- **Example**: Australian longitudinal study on women's health, Lee et al. (2005)
- **Example**: Relating income to years of education

R Code

1.3 Fundamental definitions

- **What are the questions to be answered when analysing a set of data?**

1. What is the scale of measurement?
2. What is a reasonable distribution to model the data?
3. What is the relationship with other variables?

$$\mathbb{E}[Y] = \alpha + \beta x$$

$$\log[\mathbb{E}(Y)] = \alpha + \beta \sin(\gamma x)$$

4. What is the best parameter estimation process? MLE, Least Squares, Bayesian methods
5. Why choosing a restrictive (parametric method) instead of a very flexible (nonparametric) approach?
6. Model checking: residual checking, plots, checking of the assumptions

1.3 Fundamental definitions

- Don't under-evaluate exploratory statistics!

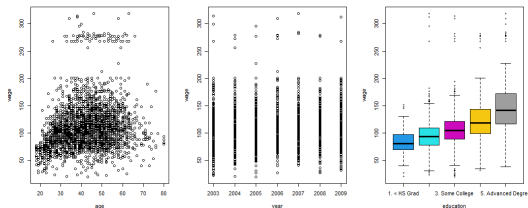


Figure: Representations of wage as a function of age, year and education

- Many statistical learning methods are relevant and useful in a wide range of disciplines.
- Statistical learning should not be viewed as a series of black boxes.
- While it is important to know what job is performed by each tool, it is not necessary to have the skills to construct the machine inside the box.

1.3 Fundamental definitions

- Measuring the quality of a fit

Goal: How well do the model predictions match the data??

- In a regression setting, use the **mean squared error (MSE)**

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad (1.3.15)$$

- select the model which **minimises the MSE**.
- Low *MSE* can hide problems of **overfitting**.

1.3 Fundamental definitions

- Measuring the quality of a fit

Goal (Corrected): What is the accuracy of the predictions when we apply the method on unseen/new data??

Step 1. Select (training) observations $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$.

Step 2. Estimate $f(\mathbf{x})$.

Step 3. Consider some new set of observations (x_0, y_0) , called testing observations.

Step 4. Use \hat{f} , the obtained estimation in Step 2, to find the **test MSE**

$$\text{Ave}(\hat{f}(x_0) - y_0)^2$$

Then:

- select the model with **minimum test MSE**, if test observations are available
- select the model with minimum training MSE, if test observations are not available
- use estimation method for the test MSE, like cross-validation

1.3 Fundamental definitions

- The bias-variance trade-off

The expected test MSE, for x_0 , can be decomposed into **three quantities**:

$$MSE(x_0) = \mathbb{E} \left[y_0 - \hat{f}(x_0) \right]^2 = Bias^2 \left(\hat{f}(x_0) \right) + Var \left(\hat{f}(x_0) \right) + Var[\varepsilon]$$

The **expected test MSE** can **never** be **lower than** $Var[\varepsilon]$

- Two competing properties of statistical learning methods
 - **Bias**: the error that is introduced by approximating a potentially complicated relationship between Y and X with a simpler model:

$$Bias \left(\hat{f}(x_0) \right) = \left[f(x_0) - \mathbb{E}(\hat{f}(x_0)) \right]$$

- **Variance**: the amount by which \hat{f} would change when changing the training dataset:

$$Var \left(\hat{f}(x_0) \right) = \mathbb{E} \left[\hat{f}(x_0) - \mathbb{E}(\hat{f}(x_0)) \right]^2$$

1.3 Fundamental definitions

- The bias-variance trade-off
 - Comments
 - **Variance:** flexible methods have larger variance
 - **Bias:** restrictive methods have a larger bias

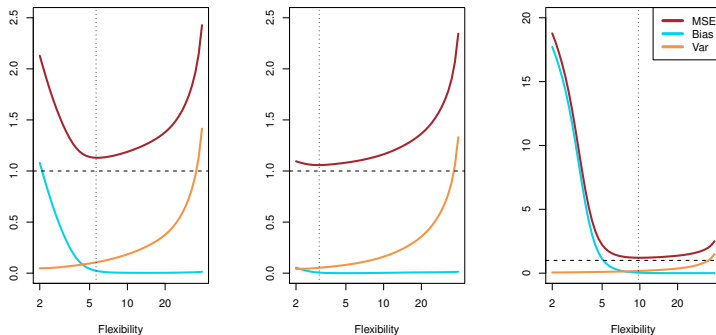


Figure: MSE decomposition to Bias and Variance based on ISL

1.3 Fundamental definitions

- **The classification setting**

Goal: How to measure model accuracy for **categorical outputs**??

The model accuracy is measured by the (training) **error rate**

$$\text{ER} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i \neq \hat{y}_i) \quad (1.3.16)$$

where \hat{y}_i is the predicted label and

$$\mathbb{I}(y_i \neq \hat{y}_i) = \begin{cases} 0 & \text{if } y_i = \hat{y}_i \text{ (correct classification)} \\ 1 & \text{if } y_i \neq \hat{y}_i \text{ (miss-classification)} \end{cases}$$

- **Test Error rate:**

$$\text{Ave}(\mathbb{I}(y_0 \neq \hat{y}_0))$$

where \hat{y}_0 is the prediction obtained by applying the classifier on input x_0 .

1.3 Fundamental definitions

- **The Bayes classifier** The test error rate is minimised, on average, by the **Bayes classifier**.
 - **Rule:** Assigns each observation to the most likely class given its predictor values, i.e.,

$$\hat{y}_0 = j \quad \text{if } \Pr(Y = j|X = x_0) > \Pr(Y = i|X = x_0) \text{ for all categories } i \neq j$$

- **Example** If there are only 2 classes, the Bayes classifier is

$$\hat{y}_0 = \begin{cases} 0 & \text{if } \Pr(Y = 0|X = x_0) > 0.5 \\ 1 & \text{if } \Pr(Y = 1|X = x_0) > 0.5 \end{cases}$$

1.3 Fundamental definitions

- The **expected prediction error** is

$$\text{EPE} = \mathbb{E} [\mathbb{I}(y_0 \neq \hat{y}_0)] = \mathbb{E}_X \sum_{j=0}^1 [\mathbb{I}(y_0 \neq \hat{y}_0)] \Pr(y_0 = j | X = x)$$

- To minimise the expected error, you need to minimize the probability of being wrong.
- Remember that

$$\hat{y}_0 = 1 \quad \text{if} \quad \Pr(y_0 = 1 | X = x_0) = \max_{j \in \{0,1\}} \Pr(y_0 = j | X = x_0)$$

- The **expected Bayes error rate** is then

$$1 - \mathbb{E}_X \left[\max_{j \in \{0,1\}} \Pr(Y_0 = j | X) \right]$$

Note: In practice, we do not know the conditional distribution of Y given X , so the Bayes classifier is an **unattainable gold standard**.

1.3 Fundamental definitions

- **K-nearest neighbours**

- The **K-nearest neighbours (KNN) classifier** **estimates** the conditional probability of Y given X and classifies a given observation to the class with the **highest estimated probability**.
 - simplicity
 - often close to the Bayes classifier.
- Let K be a positive integer and x_0 be a test input observation, then:
 - Identifies the K points in the training dataset closest to x_0 - a neighbourhood of x_0 , N_0
 - Computes $\Pr(Y = j|X = x_0) = \frac{1}{K} \sum_{i \in N_0} \mathbb{I}(y_i = j)$
 - Applies the Bayes rule $\Pr(X = x_0|Y = j) = \frac{\Pr(Y=j|X=x_0)}{\Pr(Y=j)}$
 - Classifies the test observation x_0 to the class with the largest $\Pr(X = x_0|Y = j)$

What is the best choice of K ??

1.3 Fundamental definitions

- **K-nearest neighbours**

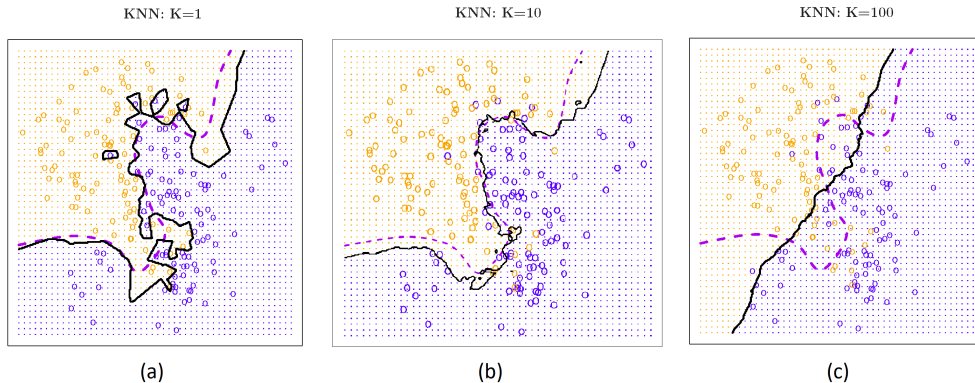


Figure: KNN with $K = 1$, $K = 10$ and $K = 100$. KNN decision boundary compared with Bayes decision boundary.

1.4 Estimation procedure

- **Introduction**

This section is about obtaining

- **Point estimates**
- **Interval estimates**

no analytical solutions \Rightarrow numerical methods (e.g., Newton-Raphson algorithm).

1.4 Estimation procedure

- **The Weibull distribution [Weibull(θ, λ)]**

- commonly used model for times to failure (or survival times) with pdf

$$f(y|\lambda, \theta) = \frac{\lambda y^{\lambda-1}}{\theta^\lambda} \exp \left[- \left(\frac{y}{\theta} \right)^\lambda \right], \quad y > 0, \quad (1.4.1)$$

where λ is the **shape parameter**, θ is the **scale parameter**.

Equation (1.4.1) can be rewritten as

$$f(y|\lambda, \theta) = \exp \left[\log \lambda + (\lambda - 1) \log y - \lambda \log \theta - \left(\frac{y}{\theta} \right)^\lambda \right], \quad y > 0, \quad (1.4.2)$$

1.4 Estimation procedure

- **Likelihood function for Weibull(θ, λ)**

- Let $Y_1, \dots, Y_N \stackrel{\text{iid}}{\sim} \text{Weibull}(\theta, \lambda)$. The likelihood function is

$$f(y_1, \dots, y_N | \lambda, \theta) = \prod_{i=1}^N \frac{\lambda y_i^{\lambda-1}}{\theta^\lambda} \exp \left[- \left(\frac{y_i}{\theta} \right)^\lambda \right]$$

and the log-likelihood function is

$$\begin{aligned} \ell(\theta, \lambda; y_1, \dots, y_N) &= \sum_{i=1}^N \left[(\lambda - 1) \log y_i + \log \lambda - \lambda \log \theta - \left(\frac{y_i}{\theta} \right)^\lambda \right] \\ &= (\lambda - 1) \sum_{i=1}^N \log y_i + N \log \lambda - N \lambda \log \theta - \sum_{i=1}^N \left(\frac{y_i}{\theta} \right)^\lambda \end{aligned} \quad (1.4.3)$$

1.4 Estimation procedure

- **Estimation of the Weibull parameters**

- First, focus on θ and consider λ as known (Exponential Family).
- To find MLE of θ , called $\hat{\theta}$,
 - derive the derivative of the log-likelihood with respect to θ (**score function**):

$$\frac{d\ell}{d\theta} = U = \sum_{i=1}^N \left[-\frac{\lambda}{\theta} + \frac{\lambda y_i^\lambda}{\theta^{\lambda+1}} \right] \quad (1.4.4)$$

- set $U = 0$

$$\sum_{i=1}^N \left[-\frac{\lambda}{\theta} + \frac{\lambda y_i^\lambda}{\theta^{\lambda+1}} \right] = 0 \iff -\frac{N\lambda}{\theta} + \frac{\lambda}{\theta^{\lambda+1}} \sum_{i=1}^N y_i^\lambda = 0 \iff \frac{-N\lambda\theta^\lambda + \lambda \sum_{i=1}^N y_i^\lambda}{\theta^{\lambda+1}} = 0$$

Consequently,

$$\hat{\theta} = \left(\frac{\sum_{i=1}^N y_i^\lambda}{N} \right)^{1/\lambda}$$

1.4 Estimation procedure

Example

Consider a dataset including lifetimes (times to failure in hours) of Kevlar epoxy strand pressure vessels at 70% stress level, (Andrews and Herzberg, 1985). If this data set follows the Weibull distribution, find the MLE of θ . R Code

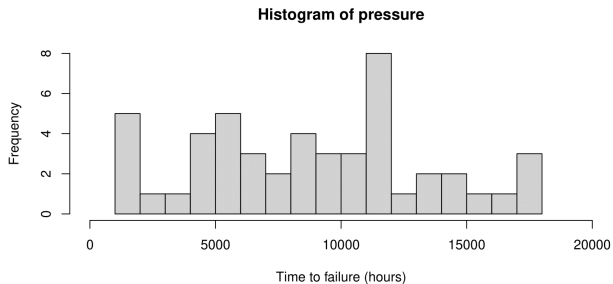


Figure: Histogram of times to failure (hours) for Kevlar epoxy strand pressure vessels.

1.4 Estimation procedure

- **Newton-Raphson algorithm in 1D**

- Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable and attains the value 0 at x_0 .
- The **Newton-Raphson algorithm** finds x_0 iteratively.
 - The slope of f at a value $x^{(m-1)}$ is given by

$$\left[\frac{df}{dx} \right]_{x=x^{(m-1)}} = f'(x^{(m-1)}) \approx \frac{f(x^{(m)}) - f(x^{(m-1)})}{x^{(m)} - x^{(m-1)}}$$

if $x^{(m)} - x^{(m-1)}$ is small (mean value theorem).

- If $x^{(m)} = x_0$ is the solution so that $f(x^m) = 0$ holds, then

$$x^{(m)} = x^{(m-1)} - \frac{f(x^{(m-1)})}{f'(x^{(m-1)})}$$

- Iterating this algorithm creates a sequence $x^{(m)}$ which approaches x_0 .

Visualisation on Wikipedia

1.4 Estimation procedure

- **Cont. Newton-Raphson algorithm in 1D**

For the Weibull distribution with $\lambda = 2$ we have

$$U = -\frac{2N}{\theta} + \frac{2 \sum_{i=1}^N y_i^2}{\theta^3}$$

and the derivative

$$\frac{dU}{d\theta} = U' = \frac{2N}{\theta^2} - \frac{2 \times 3 \times \sum_{i=1}^N y_i^2}{\theta^4} \quad (1.4.5)$$

1.4 Estimation procedure

- **Newton-Raphson algorithm in pD**

- Let $\boldsymbol{\theta} = [\theta_1, \dots, \theta_p]^T$ be the vector of parameters.
- Define the **score function with respect to** θ_j as $U_j := \partial \ell / \partial \theta_j$.
- Set $\mathbf{u} := (U_1, \dots, U_p)$.
- Construct a vector sequence $\boldsymbol{\theta}^{(m)}$ converging to $\hat{\boldsymbol{\theta}}$, where $\mathbf{u}(\hat{\boldsymbol{\theta}}) = 0$.
 - Write $\mathbf{u}^{(m)} := \mathbf{u}(\boldsymbol{\theta}^{(m)})$.
 - By the mean value theorem, we have

$$\mathbf{u}^{(m)} - \mathbf{u}^{(m-1)} \approx H_\ell^{(m-1)} \left(\boldsymbol{\theta}^{(m)} - \boldsymbol{\theta}^{(m-1)} \right)$$

where the Hessian H_ℓ is the $p \times p$ -matrix with entry $\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}$ at position (j, k) .

- Aiming for $\mathbf{u}^{(m)} = 0$, we arrive at a Newton-Raphson step

$$\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)} - (H_\ell^{(m-1)})^{-1} \mathbf{u}^{(m-1)}. \quad (1.4.6)$$

1.4 Estimation procedure

- **Method of scoring**

- It has been shown empirically that no significant loss of accuracy is suffered when $H_{\ell}^{(m)}$ is replaced by minus the information matrix $-\mathcal{I} := \mathbb{E}[\mathbf{U}']$.
- So, you may iterate:

$$\boldsymbol{\theta}^{(m)} = \boldsymbol{\theta}^{(m-1)} + (\mathcal{I}^{(m-1)})^{-1} \mathbf{u}^{(m-1)}; \quad (1.4.7)$$

This is called the **method of scoring**.

1.4 Estimation procedure

- Likelihood maximisation

- The curvature of the function ℓ in the neighborhood of the maximum gives information about how reliable the MLE is.
- The curvature of ℓ is defined by the rate of change of U , i.e. U' or $\mathbb{E}[U']$:
 - if U' is small, then ℓ is flat and U is small for a wide interval of values for $\theta \Rightarrow \hat{\theta}$ is not well determined and its s.e. is large.
 - if U' is large, then ℓ is concentrated around $\hat{\theta}$
- The variance of $\hat{\theta}$ is inversely related to $\mathcal{I} = \mathbb{E}[-U']$, i.e. $\text{s.e.}(\hat{\theta}) = \sqrt{\frac{1}{\mathcal{I}}}$

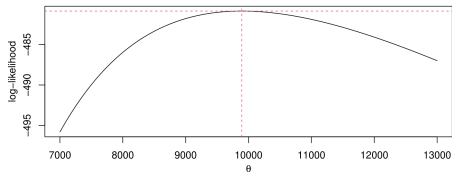


Figure: Log-likelihood function of pressure data as a function of θ .

1.5 Exponential family of distributions

Definition

Let Y be a random variable with values in \mathbb{R} (or \mathbb{N}_0 or \mathbb{Z}), with pdf $f(y; \theta)$. Then, the distribution belongs to the **exponential family** if

$$f(y; \theta) = \underbrace{s(y)}_{a(y)} \underbrace{t(\theta)}_{b(\theta)} e^{a(y)b(\theta)}, \quad \text{for some functions } a, b, s, t$$

or equivalently

$$s(y) = \exp\left(\underbrace{\ln(s(y))}_{a(y)}\right) \quad t(\theta) = \exp\left(\underbrace{\ln(t(\theta))}_{b(\theta)}\right)$$

$$f(y; \theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)] \quad (1.5.1)$$

where $s(y) = e^{d(y)}$ and $t(\theta) = e^{c(\theta)}$.

- $a(y) = y \Rightarrow$ the distribution of Y is called **canonical**.
- The distribution is in its **canonical form**, $\Rightarrow b(\theta)$ is the **natural parameter**.
- If there are other parameters, they are called **nuisance parameters**

1.5 Exponential family of distributions

- **Binomial distribution**

- Suppose $Y \sim \text{Bin}(n, p)$. Y represents ...
- The probability mass function of Y is:

$$f(y|p) = \binom{n}{y} p^y (1-p)^{n-y}, \quad y \in \{0, 1, \dots, n\}.$$

- Let n be known and p be the parameter of interest, which indicates the probability of success in a single experiment.
- **Does the Binomial belong to the Exponential family of distributions?**
 - Rewritten $f(y|p)$ as

$$f(y|p) = \exp \left[y \log p - y \log(1-p) + n \log(1-p) + \log \binom{n}{y} \right]$$

- Compare with (1.5.1) :

$$a(y) = y, \quad b(p) = \log \frac{p}{1-p}, \quad c(p) = n \log(1-p), \quad d(y) = \log \binom{n}{y}. \quad (1.5.2)$$

1.5 Exponential family of distributions

- **Cont. Binomial distribution**

- The binomial distribution is used as a model for observations of a process with binary outcomes:
 - The number of candidates who pass a test (Result of the test: pass/fail)
 - The number of patients with some disease who are alive at a specified time since diagnosis (Possible outcomes: survival/death)

Code

1.5 Exponential family of distributions

- **Normal distribution**

- Suppose $Y \sim N(\mu, \sigma^2)$ with the pdf:

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

- **Does the normal belong to the Exponential family of distributions?**

- Rewritten $f(y|\mu)$ as

$$f(y|\mu) = \exp \left[\underbrace{-\frac{y^2}{2\sigma^2}}_{d(y)} + \underbrace{\frac{\boxed{y\mu}}{\sigma^2}}_{c(\mu)} - \underbrace{\frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2)}_{c(\mu)} \right]$$

- Compare with (1.5.1) :

$$\underbrace{a(y)} = y, \quad \underbrace{b(\mu)} = \mu/\sigma^2, \quad c(\mu) = -\frac{\mu^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2), \quad d(y) = -\frac{y^2}{2\sigma^2} \quad (1.5.3)$$

σ^2 is
the variance
parameter

1.5 Exponential family of distributions

- **Cont. Normal distribution**

- **Why is normal distribution important?**

- Many phenomena are well described by the normal distribution, ex. height or blood pressure of people;
 - Even if data are not normal, the average or total of a random sample of values will be approximately normally distributed (Central Limit Theorem);
 - Statistical theory is developed in a large extent for the Normal distribution.



Code

1.5 Exponential family of distributions

- **Poisson distribution**

- Suppose $Y \sim \text{Pois}(\lambda)$, with the pdf:

$$f(y|\lambda) = \frac{\lambda^y}{e^\lambda y!}, \quad y \in \{0, 1, 2, \dots\}, \lambda > 0$$

$$E(Y) = \text{Var}(Y) = \lambda$$

- **Does the Poisson belong to the Exponential family of distributions?**

- Rewritten $f(y|\lambda)$ as

$$f(y|\lambda) = \exp(y \log \lambda - \lambda - \log y!)$$

- Compare with (1.5.1) :

arg. y → canonical form

→ natural parameter

$$a(y) = y, \quad b(\lambda) = \log \lambda, \quad c(\lambda) = -\lambda, \quad d(y) = -\log y!. \quad (1.5.4)$$

- The Poisson distribution expresses the probability of a given number of events occurring in a **fixed interval of time and/or space** if these events occur with a **known average rate** and **independently** of the time since the last event.

1.5 Exponential family of distributions

- **Properties of distributions in the exponential family**
 - The expected value and variance of $a(Y)$ are given by

$$\mathbb{E}[a(Y)] = -\frac{c'(\theta)}{b'(\theta)} \quad (1.5.5)$$

and

$$\mathbb{V}ar[a(Y)] = \frac{b''(\theta)c'(\theta) - c''(\theta)b'(\theta)}{b'(\theta)^3} \quad (1.5.6)$$

Proof

1.5 Exponential family of distributions

$$f(y; \theta) = \exp(a(y)b(\theta) + c(\theta) + d(y))$$

- **Score and information**

- The **log-likelihood function** for the exponential family is:

$$\ell(\theta; y) = a(y)b(\theta) + c(\theta) + d(y). \quad \checkmark$$

- The **score statistic** is

$$U(\theta; y) = \frac{d\ell(\theta; y)}{d\theta} = a(y)b'(\theta) + c'(\theta)$$

which depends on y and can be interpreted as a random variable:

$$U := U(\theta; Y) = a(Y)b'(\theta) + c'(\theta). \quad E(a|Y) = -\frac{c'(\theta)}{b'(\theta)}$$

- By (1.5.5), we have

$$E(U) = b'(\theta)E[a(Y)] + c'(\theta) = 0 \quad (1.5.7)$$

- The variance of U or the information, \mathcal{I} , is:

$$\mathcal{I} = \text{Var}(U) = b'(\theta)^2 \text{Var}[a(Y)] = \frac{b''(\theta)c'(\theta)}{b'(\theta)} - c''(\theta). \quad (1.5.8)$$

1.5 Exponential family of distributions

- Cont. Score and information

- Another property of the score function U is

$$\mathbb{E}(U^2) = \text{Var}(U) = -\mathbb{E}\left(\frac{dU}{d\theta}\right).$$

Handwritten note: $\mathbb{E}(U - \underbrace{\mathbb{E}(U)}_0)^2$ with an arrow pointing from $\text{Var}(U)$ to the expression.

- The first follows from $\text{Var}(X) = \mathbb{E}(X^2) - [\mathbb{E}(X)]^2$, and $\mathbb{E}(U) = 0$.
- For the second equality, note that

$$\begin{aligned}\mathbb{E}\left(\frac{dU}{d\theta}\right) &= \mathbb{E}(a(Y)b''(\theta) + c''(\theta)) = b''(\theta)\mathbb{E}[a(Y)] + c''(\theta) \\ &= b''(\theta) \underbrace{\left[-\frac{c'(\theta)}{b'(\theta)}\right]}_{\text{Handwritten arrow}} + c''(\theta) = -\underbrace{\text{Var}(U)}_{\text{Handwritten underline}} = \underbrace{-\mathcal{I}}_{\text{Handwritten underline}}\end{aligned}$$

1.5 Exponential family of distributions

- **Weibull distribution**

- Suppose $Y \sim \text{Weibull}(\theta, \lambda)$, with the pdf:

$$f(y; \theta, \lambda) = \frac{\lambda y^{\lambda-1}}{\theta^\lambda} \exp \left[- \left(\frac{y}{\theta} \right)^\lambda \right]$$

- **Does the Weibull belong to the Exponential family of distributions?**

- Rewritten $f(y|\theta, \lambda)$ as

$$f(y; \theta) = \exp \left[\log \lambda + (\lambda - 1) \log y - \lambda \log \theta - \left(\frac{y}{\theta} \right)^\lambda \right]$$

- Compare with (1.5.1) :

$$a(y) = y^\lambda, \quad b(\theta) = -\theta^{-\lambda}, \quad c(\theta) = \log \lambda - \lambda \log \theta, \quad d(y) = (\lambda - 1) \log y$$

λ is a nuisance parameter.

1.6 Inference

The two main tools to make statistical conclusions are

- **confidence intervals**: the width of a confidence interval provides a measure of the precision of the **point estimates**.
- **hypothesis testing**: compares how well two related models fit the data.

The logic for hypothesis testing:

- specify a model M_0 corresponding to H_0 and a more general model M_1 corresponding to H_1
- fit model M_0 and compute measure of goodness of fit, G_0 ; repeat for M_1 to obtain G_1
- calculate the improvement in fit
- test the null hypothesis $G_0 = G_1$
- if $G_0 = G_1$ is not rejected, then H_0 is not rejected and M_0 is the preferred model

1.6 Inference

For **confidence intervals** and **hypothesis testing**, **sampling distributions** are required.

- **normally distributed** r.v. \Rightarrow **exact**.
- **other distributions** \Rightarrow large-sample asymptotic results based on the **CLT**.
Under appropriate conditions (such as i.i.d), the **statistic S** is approximately

$$\frac{S - \mathbb{E}(S)}{\sqrt{\text{Var}(S)}} \sim \mathcal{N}(0, 1) \quad \checkmark \quad (1.6.1)$$

or equivalently

$$\frac{[S - \mathbb{E}(S)]^2}{\text{Var}(S)} \sim \chi_1^2 \quad \checkmark \quad (1.6.2)$$

and, in case of **p -multivariate statistics**

$$[\mathbf{S} - \mathbb{E}(\mathbf{S})]^\top \mathbf{V}^{-1} [\mathbf{S} - \mathbb{E}(\mathbf{S})] \sim \chi_p^2 \quad (1.6.3)$$

Note: \sim symbol means "approximately distributed as".

1.6 Inference

- **Sampling distribution for score statistics**

Suppose Y_1, \dots, Y_N are independent random variables with a distribution from the exponential family, with parameter $\theta = (\theta_1, \dots, \theta_p)$.

- The **score statistics** are such that

$$E[U_j] = 0, \quad \text{for all } j = 1, \dots, p. \quad (1.6.4)$$

- ~~The~~ **The variance-covariance matrix of the score statistics** is the information matrix \mathcal{I} , with elements

$$\mathcal{I}_{jk} = \mathbb{E}[U_j U_k] \quad (1.6.5)$$

- If $p = 1$,

$$\frac{U}{\sqrt{\mathcal{I}}} \sim \mathcal{N}(0, 1), \quad \equiv \quad \frac{U^2}{\mathcal{I}} \sim \chi_1^2$$

- If $p > 1$,

$$\mathbf{u} \sim \text{MVN}(\mathbf{0}, \mathcal{I}) \quad \equiv \quad \mathbf{u}^\top \mathcal{I}^{-1} \mathbf{u} \sim \chi_p^2$$

1.6 Inference

- Example: Binomial distribution** Let $Y \sim \text{Bin}(n, p)$, the log-likelihood function and the score statistic are, respectively,

$$\ell(p) = y \log p + (n-y) \log(1-p) + \log \binom{n}{y}, \text{ and } U = \frac{d\ell}{dp} = \frac{Y}{p} - \frac{n-Y}{1-p} = \frac{Y-np}{p(1-p)}$$

Therefore,

$$\mathbb{E}(Y) = np \Rightarrow \mathbb{E}(U) = 0$$

and

$$\text{Var}(Y) = np(1-p) \Rightarrow \mathcal{I} = \text{Var}(U) = \frac{1}{p^2(1-p)^2} \text{Var}(Y) = \frac{n}{p(1-p)}$$

and, hence,

$$\frac{U}{\sqrt{\mathcal{I}}} = \frac{Y-np}{\sqrt{np(1-p)}} \dot{\sim} \mathcal{N}(0, 1)$$

This is known as the normal approximation to the binomial distribution.

1.6 Inference

- **Taylor approximation**

Taylor approximations for generic functions f in a neighbourhood of t :

$$f(x) = f(t) + (x - t) \left[\frac{df}{dx} \right]_{x=t} + \frac{1}{2}(x - t)^2 \left[\frac{d^2f}{dx^2} \right]_{x=t} + \dots$$

The first three terms of the Taylor approximation for log-likelihood

- with $p = 1$:

$$\ell(\theta) = \ell(\tilde{\theta}) + (\theta - \tilde{\theta})U(\tilde{\theta}) - \frac{1}{2}(\theta - \tilde{\theta})^2 \mathcal{I}(\tilde{\theta}) \quad (\tilde{\theta} \text{ is an estimate of } \theta), \quad (1.6.6)$$

Handwritten notes above the equation: $\frac{\partial \ell(\theta)}{\partial \theta} = U$ and $\frac{\partial^2 \ell(\theta)}{\partial \theta^2} = -\frac{\partial U}{\partial \theta}$

and $U' = \frac{d^2 \ell}{d\theta^2}$ is approximated by $E(U') = -\mathcal{I}$.

- with p -dimensional vector θ :

$$\ell(\theta) = \ell(\tilde{\theta}) + (\theta - \tilde{\theta})^\top \mathbf{u}(\tilde{\theta}) - \frac{1}{2}(\theta - \tilde{\theta})^\top \mathcal{I}(\tilde{\theta})(\theta - \tilde{\theta}) \quad (1.6.7)$$

1.6 Inference

- **Taylor approximation**

The first two terms of the Taylor approximation for the score function

- of a one-dimensional parameter θ :

$$U(\theta) \approx U(\tilde{\theta}) + (\theta - \tilde{\theta})U'(\tilde{\theta}) = U(\tilde{\theta}) - (\theta - \tilde{\theta})\mathcal{I}(\tilde{\theta}) \quad (1.6.8)$$

- of a p -dimensional parameter θ :

$$\mathbf{u}(\theta) \approx \mathbf{u}(\tilde{\theta}) - \mathcal{I}(\tilde{\theta})(\theta - \tilde{\theta}) \quad (1.6.9)$$

$$\Rightarrow u(\theta) \approx \underbrace{u(\hat{\theta})}_{=0} - \mathcal{I}(\hat{\theta})(\theta - \hat{\theta})$$
$$\Rightarrow u(\theta) = -\mathcal{I}(\hat{\theta})(\theta - \hat{\theta})$$

1.6 Inference

- Sampling distribution of $\hat{\theta}$ (MLE of θ)

- The MLE is the estimator which maximises $\ell(\hat{\theta})$, i.e. $\mathbf{u}(\hat{\theta}) = \mathbf{0}$. Therefore,

$$\mathbf{u}(\theta) \approx -\mathbf{I}(\hat{\theta})(\theta - \hat{\theta}) \quad \checkmark$$

$$-\mathbf{I}(\hat{\theta})^{-1} \mathbf{u}(\theta) \approx (\theta - \hat{\theta})$$

- Properties

- consistency ✓

$$\mathbb{E}(\hat{\theta}) = \theta$$

$$\begin{aligned} E(\mathbf{u}(\theta)) &= \mathbf{0} \\ \rightarrow E(-\mathbf{I}(\hat{\theta})\mathbf{u}(\theta)) &= \mathbf{0} \\ \rightarrow E(\theta - \hat{\theta}) &= \mathbf{0} \rightarrow E(\hat{\theta}) = \theta \end{aligned}$$

- variance-covariance matrix

$$\begin{aligned} \mathbb{E} \left[\underbrace{(\hat{\theta} - \theta)(\hat{\theta} - \theta)^{\top}}_{\mathbf{I}^{-1} \mathbf{u}(\theta)} \right] &= \mathbb{E}[\underbrace{\mathbf{I}^{-1} \mathbf{u} \mathbf{u}^{\top} \mathbf{I}^{-1}}] \\ &= \mathbf{I}^{-1} \underbrace{\mathbb{E}[\mathbf{u} \mathbf{u}^{\top}]}_{\mathbf{I}} \mathbf{I}^{-1} = \mathbf{I}^{-1} \end{aligned}$$

- asymptotic sampling distribution (using (1.6.3)) $\hookrightarrow \mathbf{I}$

$$(\hat{\theta} - \theta)^{\top} \mathbf{I}(\hat{\theta})(\hat{\theta} - \theta) \sim \chi^2(p)$$

1.6 Inference

$$\underline{(\hat{\theta} - \theta)^T \mathcal{I}(\theta) (\hat{\theta} - \theta)}.$$

- **Remarks**

- $(\hat{\theta} - \theta)^T \mathcal{I}(\hat{\theta})(\hat{\theta} - \theta)$ is also known as **Wald statistics**.
- for $p = 1$,

$$\hat{\theta} \sim \mathcal{N}(\theta, \mathcal{I}^{-1}).$$

- if the response variable is normally distributed, the results are exact; for other GLM, the results are asymptotic

1.7 Deviance

To assessing the adequacy of a model, compare it with a more general model with the **maximum number of parameters** that can be estimated, called **saturated (maximal/full) model**.

- For observations $Y_i, i = 1, \dots, N$, a saturated model can be specified with N parameters.
 - In general, the maximum number of parameters m that can be estimated is smaller than N , e.g., if observations are repeated.
- Let θ_{\max} be the parameter vector of the saturated model with the MLE $\hat{\theta}_{\max}$.
- The likelihood for the saturated model $L(\hat{\theta}_{\max}; \mathbf{y})$ will be larger than any other likelihood function for these observations, because it provides the most complete description of the data.

1.7 Deviance

$L(\hat{\theta}; \mathbf{y})$: maximum value of the likelihood function for the model of interest.

- To assess the goodness of fit for the model, use the likelihood ratio

$$\lambda = \frac{L(\hat{\theta}_{max}; \mathbf{y})}{L(\hat{\theta}; \mathbf{y})}$$

- In practice, use

$$\log \lambda = \ell(\hat{\theta}_{max}; \mathbf{y}) - \ell(\hat{\theta}; \mathbf{y})$$

- Large values of $\log \lambda$ suggest poor performance of the model of interest relative to the saturated model.

1.7 Deviance

The **Deviance** or **log likelihood ratio statistic** is defined as

$$D = 2[\ell(\hat{\theta}_{max}; \mathbf{y}) - \ell(\hat{\theta}; \mathbf{y})].$$

We know that

$$\ell(\theta; \mathbf{y}) - \ell(\hat{\theta}; \mathbf{y}) \stackrel{\text{Taylor approximation}}{=} -\frac{1}{2} \underbrace{(\theta - \hat{\theta})^\top \mathcal{I}(\hat{\theta})(\theta - \hat{\theta})}_{\text{wald statistic}} \equiv 2[\ell(\hat{\theta}; \mathbf{y}) - \ell(\theta; \mathbf{y})] = \underbrace{(\theta - \hat{\theta})^\top \mathcal{I}(\hat{\theta})(\theta - \hat{\theta})}_{\text{wald statistic}}$$

and consequently,

$$2[\ell(\hat{\theta}; \mathbf{y}) - \ell(\theta; \mathbf{y})] \sim \chi_p^2$$

1.7 Deviance

We now have:

$$\begin{aligned} D &= 2[\ell(\hat{\theta}_{max}; \mathbf{y}) - \ell(\hat{\theta}; \mathbf{y})] \\ &= 2[\cancel{\ell(\hat{\theta}_{max}; \mathbf{y})} \pm \cancel{\ell(\theta_{max}; \mathbf{y})} \pm \ell(\theta; \mathbf{y}) - \ell(\hat{\theta}; \mathbf{y})] \\ &= 2[\boxed{\ell(\hat{\theta}_{max}; \mathbf{y}) - \ell(\theta_{max}; \mathbf{y})}] - 2[\boxed{\ell(\hat{\theta}; \mathbf{y}) - \ell(\theta; \mathbf{y})}] \\ &\quad + 2[\underbrace{\ell(\theta_{max}; \mathbf{y}) - \ell(\theta; \mathbf{y})}_{\geq 0}] \end{aligned}$$

χ_m^2 χ_p^2

- $2[\ell(\hat{\theta}_{max}; \mathbf{y}) - \ell(\theta_{max}; \mathbf{y})] \sim \chi_m^2$, where m is the number of parameters in the saturated model
- $2[\ell(\hat{\theta}; \mathbf{y}) - \ell(\theta; \mathbf{y})] \sim \chi_p^2$, where p is the number of parameters in the model of interest
- $2[\ell(\theta_{max}; \mathbf{y}) - \ell(\theta; \mathbf{y})] \geq 0$, which is zero if the model of interest has a fit which is as good as the saturated model

1.7 Deviance

Therefore, the sampling distribution of the deviance is

$$D \sim \chi^2_{m-p}(\nu)$$

under $H_0, \nu = 0$

(1.7.1)

where $\nu = 2[\ell(\theta_{\max}; \mathbf{y}) - \ell(\theta; \mathbf{y})]$ is a non-centrality parameter.

- **Remark**

- The distribution is exact if the response variable is normally distributed
- For some other distributions, D can be calculated and used directly as a goodness of fit statistic

1.7 Deviance

- Example: Binomial distribution**

For the independent response variables Y_1, \dots, Y_N with $Y_i \sim \text{Bin}(n_i, p_i)$,

$$\ell(\mathbf{p}; \mathbf{y}) = \sum_{i=1}^N \left[Y_i \log p_i - Y_i \log(1 - p_i) + n_i \log(1 - p_i) + \log \binom{n_i}{Y_i} \right]$$

- For a saturated model, the p_i 's are all different with $\hat{p}_i = \frac{Y_i}{n_i}$. So,

$$\ell(\hat{\mathbf{p}}_{\max}; \mathbf{y}) = \sum_{i=1}^N \left[Y_i \log \left(\frac{Y_i}{n_i} \right) - Y_i \log \left(\frac{n_i - Y_i}{n_i} \right) + n_i \log \left(\frac{n_i - Y_i}{n_i} \right) + \log \binom{n_i}{Y_i} \right]$$

- For any other model, $p < N$; let's $\hat{\mathbf{p}}^*$ be the MLE for a non-saturated model and $\hat{Y}_i = n_i \hat{p}_i^*$ the fitted values; then

$$\ell(\hat{\mathbf{p}}^*; \mathbf{y}) = \sum \left[Y_i \log \left(\frac{\hat{Y}_i}{n_i} \right) - Y_i \log \left(\frac{n_i - \hat{Y}_i}{n_i} \right) + n_i \log \left(\frac{n_i - \hat{Y}_i}{n_i} \right) + \log \binom{n_i}{Y_i} \right]$$

1.7 Deviance

- **Cont. Example: Binomial distribution**

Therefore, the deviance for binomial distribution is:

$$D = 2 \sum_{i=1}^N \left[Y_i \log \left(\frac{Y_i}{\hat{Y}_i} \right) + (n_i - Y_i) \log \left(\frac{n_i - Y_i}{n_i - \hat{Y}_i} \right) \right] \checkmark$$

1.7 Deviance

- **Nested model**

We say that model M_0 is nested in model M_1 if M_0 results as a special case of M_1 .

Example

If we partition θ as

$$\theta^\top = (\theta^{(1)\top}, \theta^{(2)\top})$$


where θ has length p and $\theta^{(1)}$ has length q , then model M_1 could assume unrestricted θ , whereas M_0 restricts, e.g, $\theta^{(2)} = \mathbf{0}$.

1.7 Deviance

- Cont. Nested model

Let M_0 (with q parameters) be nested in M_1 (with p Parameters) [$q < p < N$].
Then

$$\begin{aligned}\Delta D = D_0 - D_1 &= 2[\ell(\hat{\theta}_{\max}; \mathbf{y}) - \ell(\hat{\theta}_0; \mathbf{y})] - 2[\ell(\hat{\theta}_{\max}; \mathbf{y}) - \ell(\hat{\theta}_1; \mathbf{y})] \\ &= 2[\ell(\hat{\theta}_1; \mathbf{y}) - \ell(\hat{\theta}_0; \mathbf{y})]\end{aligned}$$

Since $D_0 \sim \chi^2(N - q)$ and $D_1 \sim \chi^2(N - p)$, it can be concluded that, for large N ,

$$\Delta D \sim \chi^2(p - q)$$

If the values of ΔD is in the critical region, then model M_1 provides a significantly better description of the data.