

## Week 9.

This week we shall look at Regression analysis.

Refs. • Wikipedia

- James, Witten, Hastie, Tibshirani "An intro to stat. learning" chap 3.
- Anderson (2003).
- Bai, Jiang, Yao, Zheng (2013) - Testing linear hypotheses.
- Xie, Xiao (2016?) - likelihood ratio test for high-dim linear reg.

Regression is a massive body of literature. I will only focus on a few topics.

"Method of least squares" Legendre (1805), Gauss (1809)

"regression" term coined by Galton to describe biological phenomenon: "regression toward the mean".

Given data  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^2$  and we want to determine model

$$y \approx f(x; \beta)$$

where  $\beta$  is a parameter (vector or scalar). This is the simple univariate case as  $y$  is scalar.

We need to specify model  $f$ .

In linear regression, the model specification is that the dependent variable  $y_i$  is a linear combination of the parameters.

E.g. linear model:  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i=1, \dots, n.$

parabola:  $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i \quad i=1, 2, \dots, n.$

Regression is a supervised learning problem: we are given outputs and inputs and we have to learn the model.

Stats. v. Machine/Statistical learning.

$\Rightarrow$  Stats. learn something about phenomena (Science) and the focus is on supporting and rejecting hypothesis.

$\Rightarrow$  ML: Focus is on optimal out-of-sample prediction. No need or desire to understand the model. "Black box".

More generally in the univariate case, we have observations

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

where  $x_i$  is  $p$ -dimensional. We assume model of the

form

$$y_i \approx \beta' x_i$$

$$\beta := (\beta_0, \beta_1, \dots, \beta_{p-1})'$$

for  $i=1, \dots, n$  and  $n > p$ .

$$x_i := (1, x_{i2}, x_{i3}, \dots, x_{ip})$$

In other words our model is  $f(x; \beta) := \beta' x$ . This gives us a number of decisions  $\mathcal{D} := \{f(x_1; \beta), f(x_2; \beta), \dots, f(x_n; \beta)\}$  and our outcome space  $\mathcal{Y} = \{y_1, y_2, \dots, y_n\}$ .

We score our model by a loss function  $l: \mathcal{D} \times \mathcal{Y} \rightarrow \mathbb{R}$ .

One example is the least squares loss

$$l(\mathcal{D}, \mathcal{Y}) = \sum_{i=1}^n |y_i - \beta' x_i|^2 =: S(\beta) \quad (\text{as only depends on choice of } \beta)$$

We want to minimise our loss by varying  $\beta$ . Then the optimal choice  $\hat{\beta}$  is

$$\hat{\beta} := \arg \min_{\beta} S(\beta).$$

Our problem can be written in matrix notation.

$$X\beta = Y$$

where

$$X = \begin{bmatrix} x_1' \\ x_2' \\ \vdots \\ x_n' \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

So that  $S(\beta) = \sum_{i=1}^n |y_i - \beta' x_i|^2 = \|Y - X\beta\|^2$

and the solution of the minimisation is given by the normal equations

$$X'X\hat{\beta} = X'Y$$

and the solution is  $\hat{\beta} = (X'X)^{-1}X'Y$ .

Defining  $e_i = |y_i - \beta' x_i|$  then  $S(\beta) = e_1^2 + e_2^2 + \dots + e_n^2$   
 which gives the name Residual sum of squares (RSS)  
 or sum of squares error (SSE).

We assume true relationship  $y = \beta'x + \varepsilon$  where  $\varepsilon$  is a mean-zero random error term.

## Multivariate case

observations  $(Y_1, X_1), \dots, (Y_n, X_n)$

$Y_i$ :  $p$ -dimensional     $X_i$ :  $q$ -dimensional.

Case  $p=1$  gives back univariate case.

Assume true relationship between  $Y$  and  $X$  is given by

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_q x_{iq} + \varepsilon_i, \quad i=1, 2, \dots, n.$$

$$Y_i' = (Y_{i1}, Y_{i2}, \dots, Y_{ip})', \quad X_i' = (x_{i1}, x_{i2}, \dots, x_{iq})'$$

where  $\beta_0, \beta_1, \dots, \beta_q$  are  $p$ -dimensional parameter vectors to be determined. The error terms  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  are  $p$ -dimensional random vectors. We assume they are mutually uncorrelated, and that

$$E[\varepsilon_i] = 0 \quad \text{Var}[\varepsilon_i] = \Sigma \quad i=1, 2, \dots, n.$$

In matrix form

$$Y = \begin{bmatrix} Y_1' \\ \vdots \\ Y_n' \end{bmatrix}, \quad X = \begin{bmatrix} 1 & X_1' \\ \vdots & \vdots \\ 1 & X_n' \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1' \\ \vdots \\ \varepsilon_n' \end{bmatrix}$$

giving

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}.$$

$$\mathbf{B} = (\beta_0, \beta_1, \dots, \beta_k)'$$

$$E[\mathbf{Y}] = \mathbf{X}\mathbf{B}, \quad \text{Var}[\mathbf{Y}_i] = \Sigma \quad \text{for all } i=1, \dots, n.$$

$$\text{Cov}(\mathbf{Y}_i, \mathbf{Y}_j) = 0 \quad \text{for all } i \neq j.$$

$$\text{The RSE is } \sum_{i=1}^n \sum_{j=1}^p \varepsilon_{ij}^2 = \text{tr}(\mathbf{Y} - \mathbf{X}\mathbf{B})'(\mathbf{Y} - \mathbf{X}\mathbf{B}) = S(\mathbf{B})$$

Similar to the univariate case, we seek  $\hat{\mathbf{B}} = \underset{\mathbf{B}}{\text{argmin}} S(\mathbf{B})$  and the solution is given by

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}.$$

Assume  $\varepsilon_i$  are normally distributed, then  $\hat{\mathbf{B}}$  is the MLE of  $\mathbf{B}$  and the MLE of  $\Sigma$  is

$$\hat{\Sigma} = \mathbf{S} = \frac{1}{n}(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})'(\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}}).$$

Thm: For the estimators  $\hat{\beta}$  and  $\hat{\Sigma}$  it holds that:

$$(1) E[\hat{\beta}] = \beta \quad E[\hat{\Sigma}] = \frac{1}{n(n-k+1)} \Sigma \quad \text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = v_{ij} \Sigma$$

$$\text{where } V = (X'X)^{-1} = (v_{ij})$$

(e) If  $\mathbb{Y}$  is normally distributed then  $\hat{\beta} \sim N_{(k+1) \times p}(\beta, \Sigma \otimes V)$   
 $n\hat{\Sigma} \sim W_p(n-k+1, \Sigma)$  and  $\hat{\beta}$  and  $n\hat{\Sigma}$  are independent

Here,  $\otimes$  is Kroneker product [ $A = (a_{ij})$  and  $B$  then  $A \otimes B = (a_{ij}B)$ ].

#

Note:  $\hat{\Sigma}$  is not unbiased. A common unbiased

estimator is

$$\tilde{\Sigma} = \frac{1}{n-k} (\mathbb{Y} - X\hat{\beta})' (\mathbb{Y} - X\hat{\beta})$$

Following Anderson (2003), under normality assumption, we can discuss the equivalent problem:  $X_1, X_2, \dots, X_n$  set of  $n$  independent observations

$$X_k \sim N_p(\beta Z_k, \Sigma)$$

where the vectors  $Z_k$  ( $q$ -dimensional) are called design vectors.

Caution: In other words, my  $y_i$ 's become  $x_i$ 's and my independent variables become  $z_i$ 's (instead of  $x_i$ ).

Inference on the parameters  $B$ .

Assume  $n \geq p+q$  and rank of  $Z = (z_1, z_2, \dots, z_n)$  is  $q$ .

Aim: We may want to test that a subset of the inputs  $\{z_i\}_{i=1, \dots, n}$  play no role in predicting the  $x_i$ .

Partition the parameters  $B = (B_1, B_2)$  so that  $B_1$  has  $q_1$  columns and  $B_2$  has  $q_2$  columns.

We are going to look at the likelihood ratio criterion for testing the hypothesis

$$H_0 : B_1 = B_{01}$$

Where  $B_{01}$  is a given matrix (eg.  $B_{01} = \mathbf{0}$  matrix)



The maximum likelihood  $L = f_{B, \Sigma}(\mathbf{x})$  of a sample  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  is

$$\max_{B, \Sigma} L = (2\pi)^{-\frac{1}{2}pn} |\hat{\Sigma}_R|^{-\frac{1}{2}n} e^{-\frac{1}{2}pn}$$

where  $\hat{\Sigma}_R = \hat{\Sigma}$ .

We need to restrict the MLE for the parameters to the subspace  $\omega$  induced by the null hypothesis.

Set  $\mathbf{y}_k = \mathbf{x}_k - B_0 \mathbf{z}_k \quad k=1, \dots, n.$

where  $\mathbf{z}_k = \begin{pmatrix} \mathbf{z}_{1k} \\ \mathbf{z}_{2k} \end{pmatrix}$  for  $k=1, 2, \dots, n$ , is partitioned in same way as  $B$ .

Then  $\mathbf{y}_k$  can be considered as an observation from  $N(B_0 \mathbf{z}_k, \Sigma)$ .

Using our new notation. (Anderson-style)

$$\hat{B} = CA^{-1}$$

$$C = \sum_{k=1}^n \mathbf{x}_k \mathbf{z}_k' \quad A = \sum_{k=1}^n \mathbf{z}_k \mathbf{z}_k'$$

old notation:

$$\hat{B} = \underbrace{(\mathbf{X}'\mathbf{X})^{-1}}_{A^{-1}} \underbrace{\mathbf{X}'\mathbf{Y}}_C$$

By partitioning,

$$\begin{aligned}\hat{B}_{2\omega} &= \sum_{k=1}^n y_k z_{2k}' A_{22}^{-1} = \sum (x_k - B_{01} z_{1k}) z_{2k}' A_{22}^{-1} \\ &= (C - B_{01} A_{12}) A_{22}^{-1}\end{aligned}$$

with  $C = (C_1, C_2)$  and  $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$

The estimator of  $\Sigma$  is

$$\begin{aligned}n \hat{\Sigma}_{\omega} &= \sum_{k=1}^n (y_k - \hat{B}_{2\omega} z_{2k})(y_k - \hat{B}_{2\omega} z_{2k})' \\ &= \sum_{k=1}^n y_k y_k' - \hat{B}_{2\omega} A_{22} \hat{B}_{2\omega}' \\ &= \sum_{k=1}^n (x_k - B_{01} z_{1k})(x_k - B_{01} z_{1k})' - \hat{B}_{2\omega} A_{22} \hat{B}_{2\omega}'\end{aligned}$$

This gives the MLE over  $\omega$ :

$$\max_{B_2, \Sigma} L = (2\pi)^{-\frac{1}{2}pn} |\hat{\Sigma}_{\omega}|^{-\frac{1}{2}n} e^{-\frac{1}{2}pn}.$$

The likelihood ratio criterion for testing  $H_0$  is

$$\lambda = \frac{|\hat{\Sigma}_{\Omega}|^{\frac{1}{2}n}}{|\hat{\Sigma}_{\omega}|^{\frac{1}{2}n}} \quad \text{reject for } \lambda < \lambda_0.$$

## Distribution of $\lambda$ when $H_0$ is true

Likelihood ratio criterion  $\lambda$  can be written in terms of

$$u = \lambda^{2/n} = \frac{|\hat{\Sigma}_\Omega|}{|\hat{\Sigma}_\omega|}$$

and (p296, Anderson 2003):

$$n\hat{\Sigma}_\omega = n\hat{\Sigma}_\Omega + (\hat{B}_{12} - B_{01})A_{11 \cdot 2}(\hat{B}_{12} - B_{01})'$$

where  $A_{11 \cdot 2} := A_{11} - A_{12}A_{22}^{-1}A_{21}$ .

So we can write

$$|\hat{\Sigma}_\Omega|$$

$$u = \frac{|\hat{\Sigma}_\Omega|}{|\hat{\Sigma}_\Omega + (\hat{B}_{12} - B_{01})A_{11 \cdot 2}(\hat{B}_{12} - B_{01})'|}$$

Lemma: Set  $A := n\hat{\Sigma}_\Omega$   $H := (\hat{B}_{12} - B_{01})A_{11 \cdot 2}(\hat{B}_{12} - B_{01})'$

Then  $A \sim W(n-q, \Sigma)$ ,  $H \sim W(q, \Sigma)$ ,  
and they are independent.

#

The distribution of  $u$  can be characterised in terms of a product of beta variables.

$$u = V_1 V_2 \dots V_p.$$

where  $V_1 = g_{11}/(g_{11} + h_{11})$  and  $V_i = \frac{|G_i|}{|G_{i-1}|} \bigg/ \frac{|G_i + H_i|}{|G_{i-1} + H_{i-1}|}$

and  $G_i$  and  $H_i$  are submatrices of  $G$  and  $H$ , respectively, of the first  $i$  rows and columns.

Theorem: When  $H_0$  true,  $u = \prod_{i=1}^p V_i$  where  $V_1, V_2, \dots, V_p$  independent and  $V_i$  has the beta density

$$\begin{aligned} & \beta\left[v; \frac{1}{2}(n-q+1-i), \frac{1}{2}q_i\right] \\ &= \frac{\Gamma(\frac{1}{2}(n-q+q_i+1-i))}{\Gamma(\frac{1}{2}(n-q+1-i))\Gamma(\frac{1}{2}q_i)} v^{\frac{1}{2}(n-q+1-i)-1} \\ & \quad \times (1-v)^{\frac{1}{2}q_i-1}. \end{aligned}$$

~~##~~

Typically this is too difficult to work with.

One can develop an asymptotic approximation.

Let  $U =: U_{p,q_1, n-q}$  where  $q = q_1 + q_2$  and let

$u_{p,q_1, n-q}(\alpha)$  be the significance point for  $U_{p,q_1, n-q}$

that is,

$$P(U_{p,q_1, n-q} \leq u_{p,q_1, n-q}(\alpha) \mid H_0) = \alpha.$$

It can be shown (Section 8.5, Andersen 2003) that

$$-(n-q - \tfrac{1}{2}(p-q_1+1)) \log[U_{p,q_1, n-q}]$$

has a limiting  $\chi^2$ -distribution with  $pq_1$  degrees of freedom (There is also a Normal and F distribution approx).

Let  $\chi_{pq_1}^2(\alpha)$  denote the  $\alpha$  significance point of  $\chi_{pq_1}^2$  and

$$\text{let } C_{p,q_1, n-q-p+1}(\alpha) := \frac{-(n-q - \tfrac{1}{2}(p-q_1+1)) \log(u_{p,q_1, n-q}(\alpha))}{\chi_{pq_1}^2(\alpha)}.$$

Reject hypothesis if

$$-(n-q - \tfrac{1}{2}(p-q_1+1)) \log U_{p,q_1, n-q} > C_{p,q_1, n-q-p+1}(\alpha) \chi_{pq_1}^2(\alpha)$$

The values of  $C_{pq_1, n-q-p+1}(\alpha)$  are usually tabulated somewhere (eg. Anderson, Appendix B, Table 1), and serve as a correction factor wrt. asymptotic  $\chi^2$  quantile.

This  $\chi^2$  approximation is not very good and a number of researchers have tried to correct it.

Theorem (Box and Bartlett). With  $k = n - \frac{1}{2}(p - q_1 + 1)$  the CDF of  $-k \log(U_{pq_1, n-q})$  has the following expansion

$$P(-k \log(U_{pq_1, n-q}) \leq z) = \varphi_{pq_1}(z) + \frac{\gamma_2}{k^2} \{ \varphi_{pq_1+4}(z) - \varphi_{pq_1}(z) \} \\ + \frac{1}{k^4} \left[ \gamma_4 \{ \varphi_{pq_1+8}(z) - \varphi_{pq_1}(z) \} - \gamma_2^2 \{ \varphi_{pq_1+4}(z) - \varphi_{pq_1}(z) \}^2 \right] + R_n$$

where  $\varphi_m(z) := P(\chi_m^2 \leq z)$  and

$$\gamma_2 = \frac{pq_1(p^2 + q_1^2 - 5)}{48}$$

$$\gamma_4 = \frac{\gamma_2^2}{2} + \frac{pq_1}{1920} [3p^4 + 3q_1^4 + 10p^2q_1^2 - 50(p^2 + q_1^2) + 159]$$

$R_n$  order  $O(n^{-6})$ .

We have considered some techniques for testing the general linear hypothesis:

- Distribution of  $U$  as product of beta variables; See page 12. (too difficult to implement?)
- $\chi^2$  approximation (poor performance  $p > 2$ ).
- Box and Bartlett correction (still poor for large  $p$ ).

We now look at the high-dimensional regime using RMT, following Bai et al. (2013).

$$\text{Recall } U \stackrel{d}{=} \frac{|A|}{|H|} = |AH^{-1}| \quad \begin{array}{l} H \sim W(q, \Sigma) \\ A \sim W(n-q, \Sigma) \end{array}$$

$$\text{So } T_n = -n \log U = n \sum_{j=1}^p \log(1 + \ell_j)$$

where  $\ell_j$  are eigenvalues of  $HA^{-1}$ .

We can assume without loss of generality

$$\Sigma = I_p.$$

For fixed  $p$  and  $q$  (# params fixed) the eigenvalues behave like  $\ell_j \xrightarrow{p} 0$  at a rate  $1/n$ . (as  $n \rightarrow \infty$ )

This gives the approximation

$$T_n = n \sum_{j=1}^n \ell_j + O_p(1/n).$$

and  $T_n$  is distributed according to  $\chi^2_{pq,1}$ .

When  $p, q, n$  become large together this is not the expected behaviour of the eigenvalues

Assume  $\frac{p}{q_1} \rightarrow y_1 \quad \frac{p}{n-q} \rightarrow y_2 \in (0, 1)$

Define  $F = \frac{n-q}{q_1} H \Lambda^{-1}$



If populations are normal distributed then we have seen that  $F$  is distributed to a random Fisher matrix with  $(q, n-q)$  degrees of freedom.

Let  $\ell_j$  be eigenvalues of  $F$  and define the finite horizon proxies

$$\frac{p}{q} =: y_{n_1}, \quad \frac{p}{n-q} =: y_{n_2}.$$

The statistic can be rewritten

$$T_n = n \sum_{j=1}^p [\log(y_{n_1} + y_{n_2} \ell_j) - \log y_{n_1}] =: n \sum_{j=1}^p f(\ell_j)$$

with function  $f(x) = \log(1 + \frac{y_{n_2}}{y_{n_1}} x)$ .

A CLT can be derived; see Bai et al (2013).

Theorem: Under  $H_0$ , true,

$$\frac{1}{n} T_n - \mu_n \xrightarrow{\mathcal{D}} N(\eta, \sigma^2).$$

Where

18

$$\begin{aligned} \mu_n = & -(n-q-p) \log c_n - (q_1-p) \frac{y_{n_1}-1}{y_{n_1}} \log(c_n - d_n h_n) \\ & + (n-q+q_1) \log\left(\frac{c_n h_n - d_n y_{n_2}}{h_n}\right) \end{aligned}$$

$$\eta = \frac{1}{2} \log(y_1 + y_2 - y_1 y_2)$$

$$\sigma^2 = 2 \log\left(\frac{y_1 + y_2 - y_1 y_2}{(y_1 + y_2)(1 - y_2)}\right)$$

$$c_n := \frac{h_n}{\sqrt{y_{n_1}} (1 - y_{n_2})}$$

$$d_n := \frac{y_{n_2}}{\sqrt{y_{n_1}} (1 - y_{n_2})}$$

$$h_n := y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}.$$

~~///~~