

Linear Algebra revision

Eigenvalues/eigenvectors/Singular Value Decomposition

• A : real-valued square matrix

We will assume that A has full rank

$$\text{rank}(A) = p$$

[A has p -linearly independent columns (A is invertible)]

We say that λ_i is an eigenvalue of A with v_i being the corresponding eigenvector if

$$A v_i = \lambda_i v_i$$

The eigenvalues of A are the roots of its characteristic polynomial

$$\det(A - \lambda I) = 0$$

Example

$$A = \begin{pmatrix} 0 & 1 \\ 2 & 1 \end{pmatrix}, \quad A - \lambda I = \begin{bmatrix} -\lambda & 1 \\ 2 & 1-\lambda \end{bmatrix}$$

$$\det(A - \lambda I) = \lambda(\lambda - 1) - 2 = \lambda^2 - \lambda - 2$$

$$\lambda_1 = 2 \quad \lambda_2 = -1$$

$$\lambda_1 = 2$$

$$Av_1 = \lambda_1 v_1 \leadsto (A - \lambda_1 I)v_1 = 0$$

$$\begin{bmatrix} -2 & 1 \\ 2 & -1 \end{bmatrix} \begin{bmatrix} v_1^{(1)} \\ v_1^{(2)} \end{bmatrix} = 0 \quad \dots \quad v_1 = \begin{bmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix}$$

Typically we select the eigenvector that has norm 1

$$\|v_1\|_2 = \sqrt{\left(\frac{1}{\sqrt{5}}\right)^2 + \left(\frac{2}{\sqrt{5}}\right)^2} = 1$$

- Because A is invertible it has no zero eigenvalues ($\det(A - \lambda I) = \det(A) \neq 0$)
- The matrices we will be typically working with they will have no repeated eigenvalues.
- We will also assume that all the corresponding eigenvectors are linearly independent *

So we have

$$Av_1 = \lambda_1 v_1$$

$$Av_2 = \lambda_2 v_2$$

$$\vdots$$

$$Av_p = \lambda_p v_p$$

$$\left[\underbrace{Av_1 \mid Av_2 \mid \dots \mid Av_p}_{p \times p} \right] = \left[\underbrace{\lambda_1 v_1 \mid \lambda_2 v_2 \mid \dots \mid \lambda_p v_p}_{p \times p} \right]$$

$$A \underbrace{\begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_p \\ | & | & \dots & | \end{bmatrix}}_Q = \underbrace{\begin{bmatrix} | & | & \dots & | \\ v_1 & v_2 & \dots & v_p \\ | & | & \dots & | \end{bmatrix}}_Q \underbrace{\begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_p \end{bmatrix}}_\Lambda \quad \textcircled{3}$$

$$\Rightarrow A Q = Q \Lambda$$

$$A Q Q^{-1} = Q \Lambda Q^{-1}$$

$$A = Q \Lambda Q^{-1}$$

Because of $\textcircled{*}$
 Q is invertible!

diagonalisation of A

Now suppose in addition that A is symmetric
 $(A = A^T)$

\Rightarrow then the eigenvectors corresponding to different eigenvalues are orthogonal

$$\left. \begin{aligned} A v_1 &= \lambda_1 v_1 \\ A v_2 &= \lambda_2 v_2 \end{aligned} \right\} \Rightarrow v_1 \perp v_2 \quad \left(\begin{aligned} v_1 \cdot v_2 &= 0 \\ v_1^T \cdot v_2 &= 0 \end{aligned} \right)$$

$$\begin{bmatrix} v_1^T \\ \vdots \end{bmatrix} \begin{bmatrix} v_2 \\ \vdots \end{bmatrix} \quad \begin{matrix} 1 \times p & p \times 1 \\ \hline 1 \times 1 \end{matrix}$$

These eigenvectors can also be chosen to have norm equal to 1

$$\Rightarrow \begin{cases} v_i^T v_j = 0 & i \neq j \\ v_i^T v_i = \|v_i\|_2^2 = 1 & \forall i \end{cases}$$

We call this an orthonormal set of eigenvectors

In matrix form

$$A = Q \Lambda Q^{-1}$$

$$Q = \left[\begin{array}{c|c|c} v_1 & \dots & v_p \end{array} \right], \quad Q^{-1} = \left[\begin{array}{c} \hline v_1^T \\ \vdots \\ v_p^T \\ \hline \end{array} \right]$$

$$Q^{-1} Q = \left[\begin{array}{ccc} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{array} \right] \quad (\text{check } Q^{-1} Q = I)$$

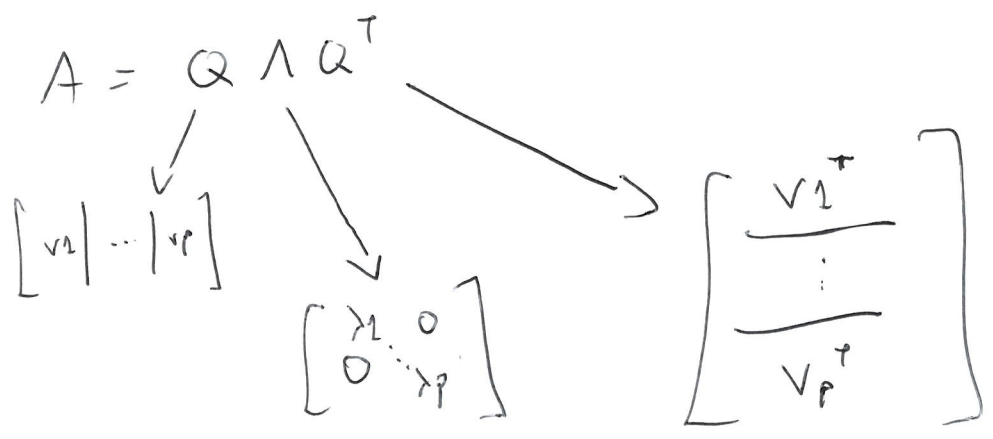
But this means

$$Q^{-1} = Q^T$$

$$\Rightarrow A = Q \Lambda Q^{-1} = Q \Lambda Q^T$$

$$A = Q \Lambda Q^T$$

spectral decomposition
(or Karhunen-Loeve)



$$A = Q \Lambda Q^T = \sum_{i=1}^p \lambda_i v_i v_i^T \quad (\text{check at home})$$



Every $v_i v_i^T$ is a ~~rank-1~~ rank-1 matrix

Singular Value Decomposition

(5)

Consider on, in general, non-square matrix X of size $n \times p$ (real-valued)

- We assume that X has full-rank:

$$\text{rank}(X) = \min(n, p)$$

The singular value decomposition (SVD) is the following factorization:

$$X = U D V^T$$

$n \times p$ \downarrow $p \times \min(n, p)$
 $n \times \min(n, p)$ $\min(n, p) \times \min(n, p)$ (V^T is $\min(n, p) \times p$)

D : is diagonal, contains all the singular values
In this case, these will be the square root
of eigenvalues of $\underbrace{X X^T}$.
symmetric, all the eigenvalues are positive.

U : Satisfies $U^T U = I$
The columns of U are eigenvectors of $X X^T$

V : Satisfies $V^T V = I$
The columns of V are eigenvectors of $X^T X$

Similarly as before we can write X as the sum of rank-1 matrices. (6)

~~TMK~~

Important result

If X is square, symmetric, and positive (all the eigenvalues are real and positive)

then SVD coincides with the spectral decomposition.

Derivatives of linear and Quadratic forms ①

Linear forms

Recall: $f: \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = \alpha x$
 $f'(x) = \frac{df}{dx} = \alpha$

A is a matrix of size $n \times p$

We can naturally define the following

function: $f: \mathbb{R}^p \rightarrow \mathbb{R}^n$ as.

$$f(s) = A s \in \mathbb{R}^n, \quad s \in \mathbb{R}^p$$

$\begin{matrix} \downarrow & \downarrow \\ n \times p & p \times 1 \\ \hline n \times 1 \end{matrix}$

$$s = (s_1, s_2, \dots, s_p)^T$$

$$f(s) = (f_1(s), f_2(s), \dots, f_n(s))^T$$

$f_i: \mathbb{R}^p \rightarrow \mathbb{R}$ for every $i=1, \dots, n$

$$\frac{\partial f}{\partial s}(s_0) = \begin{bmatrix} \frac{\partial f_1(s_0)}{\partial s_1} & \frac{\partial f_1(s_0)}{\partial s_2} & \dots & \frac{\partial f_1(s_0)}{\partial s_p} \\ \frac{\partial f_2(s_0)}{\partial s_1} & \frac{\partial f_2(s_0)}{\partial s_2} & \dots & \frac{\partial f_2(s_0)}{\partial s_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n(s_0)}{\partial s_1} & \frac{\partial f_n(s_0)}{\partial s_2} & \dots & \frac{\partial f_n(s_0)}{\partial s_p} \end{bmatrix} \quad n \times p$$

It turns out that for $f(s) = A s$ (exercise) (2)

$$\frac{\partial f}{\partial s}(s_0) = A \quad \text{for every } s_0 \in \mathbb{R}^p$$

Example

$$A = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{matrix} 1 \times 2 \\ n \times p \end{matrix} \quad \begin{matrix} n=1 \\ p=2 \end{matrix}$$

$$f(s) = A s = A \begin{bmatrix} s_1 \\ s_2 \end{bmatrix} = \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} s_1 \\ s_2 \end{bmatrix} = s_1 + 2 s_2$$

$$\frac{\partial(A s)}{\partial s_1} = 1 \quad \frac{\partial(A s)}{\partial s_2} = 2$$

$$\frac{\partial(A s)}{\partial s} = \begin{bmatrix} \frac{\partial(A s)}{\partial s_1} & \frac{\partial(A s)}{\partial s_2} \end{bmatrix} = \begin{bmatrix} 1 & 2 \end{bmatrix} = A$$

Quadratic forms

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(x) = b x^2$$

$$f'(x) = \frac{d f(x)}{d x} = 2 b x$$

Let B be a symmetric $p \times p$ matrix
we can define $f: \mathbb{R}^p \rightarrow \mathbb{R}$

$$f(s) = s^T B s$$

$\begin{matrix} \downarrow & \downarrow & \downarrow \\ 1 \times p & p \times p & p \times 1 \\ \hline & 1 \times 1 & \end{matrix}$

$$\frac{\partial f}{\partial s} (s_0) = \left[\frac{\partial f}{\partial s_1} (s_0) \quad \frac{\partial f}{\partial s_2} (s_0) \quad \dots \quad \frac{\partial f}{\partial s_p} (s_0) \right] \quad (3)$$

It turns out:

$$\frac{\partial f}{\partial s} (s_0) = 2 \underbrace{s_0^T}_{1 \times p} \underbrace{B}_{p \times p}$$

Note: If B was not symmetric $\frac{\partial f}{\partial s} (s_0) = s_0^T (B + B^T)$

Sometimes for ease of notation we will write this as a column vector:

$$\left[\frac{\partial f}{\partial s} (s_0) \right]^T = (2 s_0^T B)^T = 2 B^T (s_0^T)^T = 2 B^T s_0 = 2 B s_0 \quad (B = B^T)$$

Example

$$B = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}, \quad \text{then } BS = B \begin{bmatrix} s_1 \\ s_2 \end{bmatrix} = \begin{bmatrix} s_1 + 2s_2 \\ 2s_1 + s_2 \end{bmatrix}$$

$$f(s) = s^T B s = \begin{bmatrix} s_1 & s_2 \end{bmatrix} \begin{bmatrix} s_1 + 2s_2 \\ 2s_1 + s_2 \end{bmatrix}$$

$$= s_1^2 + s_2^2 + 4s_1s_2$$

Then according to the definition $\frac{\partial f}{\partial s} = \left(\frac{\partial f}{\partial s_1} \quad \frac{\partial f}{\partial s_2} \right)$

$$= \begin{bmatrix} 2s_1 + 4s_2 & 2s_2 + 4s_1 \end{bmatrix} = 2 \begin{bmatrix} s_1 + 2s_2 & s_2 + 2s_1 \end{bmatrix}$$

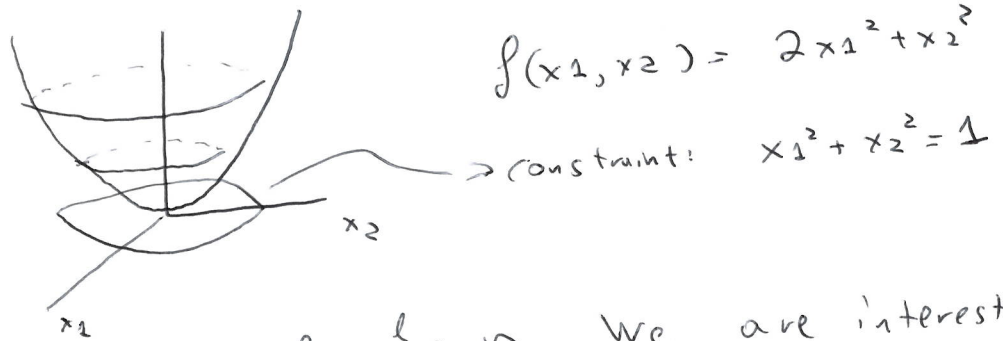
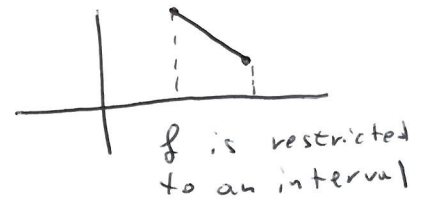
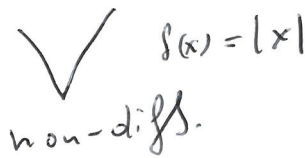
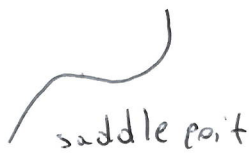
$$= 2 \begin{bmatrix} s_1 & s_2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$$

$$= 2 s^T B \quad \checkmark$$

Computing extreme points of a function ①
 subject to constraints

What do you think about the following phrase?
 "Let f be a function. In order to find its minimum / maximum I just need to solve

$$f'(x) = 0 \quad "$$



We have a function $f: \mathbb{R}^l \rightarrow \mathbb{R}$. We are interested in computing extreme points (maximum or minimum)

subject to constraints: $h_1(x) = 0$
 $h_2(x) = 0$

$$\vdots$$

$$h_m(x) = 0$$

Here f and $h_i, i=1, \dots, m$ are differentiable

$$\begin{cases} \max_x f(x) \\ \text{s.t. } h_1(x) = 0, h_2(x) = 0, \dots, h_m(x) = 0 \end{cases}$$

This means we are looking for $x^* \in \mathbb{R}^l$ such that $h_1(x^*) = h_2(x^*) = \dots = h_m(x^*) = 0$ and if \tilde{x} is another such point then

$$f(x^*) \geq f(\tilde{x})$$

Example

$$\begin{cases} \max & f(x) = x_1 + x_2 \\ \text{s.t.} & h(x) = x_1^2 + x_2^2 - 1 = 0 \end{cases}$$

$$x = (x_1, x_2)$$

(2)

Lagrange multipliers method

a) Form the Lagrangian L

$$L(x) = \underbrace{f(x)}_{\text{target}} - \sum_{j=1}^m \lambda_j \underbrace{h_j(x)}_{\text{Lagrange multipliers}} \rightarrow j\text{-th constraint}$$

In the previous example ($m=1$)

$$L(x) = L(x_1, x_2) = x_1 + x_2 - \lambda (x_1^2 + x_2^2 - 1)$$

Observation: L is a function of both x and λ
($\lambda_1, \dots, \lambda_m$)

$$L(x, \lambda) = L(x_1, \dots, x_l, \lambda_1, \dots, \lambda_m)$$

\leadsto Transformation of the original ~~problem~~ constrained problem into an unconstrained one.

b) Construct a system of $l+m$ -equations.

$$\frac{\partial L}{\partial x} = 0$$

$$\frac{\partial L}{\partial \lambda} = 0$$

→ Derivative of L with respect to every x_i , $i=1, \dots, l$. Set every one of them to 0 (5)

$$\frac{\partial L}{\partial x} = \left(\underbrace{\frac{\partial L}{\partial x_1}, \frac{\partial L}{\partial x_2}, \dots, \frac{\partial L}{\partial x_l}}_{l \text{ equations}} \right) = (0, 0, \dots, 0)$$

→ Derivative of L with respect to every λ_i , $i=1, \dots, m$. Set them to 0.

$$\frac{\partial L}{\partial \lambda} = \left(\underbrace{\frac{\partial L}{\partial \lambda_1}, \frac{\partial L}{\partial \lambda_2}, \dots, \frac{\partial L}{\partial \lambda_m}}_{m \text{ equations}} \right) = (0, 0, \dots, 0)$$

c) Solutions of the above system of $(l+m)$ -equations (x^*, λ^*) are candidates to solve the original constrained optimisation problem

Exercise

Apply b) to the previous example

$$L(x, \lambda) = x_1 + x_2 - \lambda (x_1^2 + x_2^2 - 1)$$

and you find candidate for extremes, and identify the maximum.

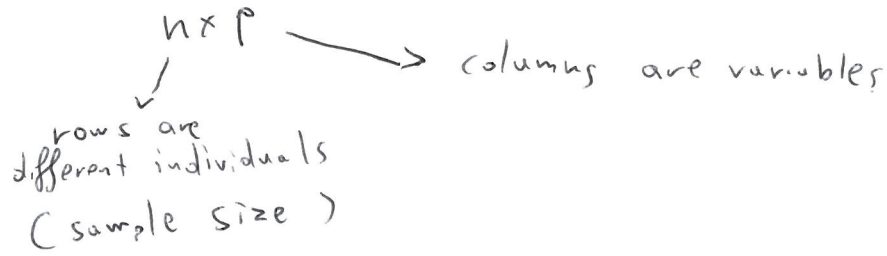
(Later on PCA):

$$\left\{ \begin{array}{l} f(x_1, x_2) = [x_1 \ x_2]^T A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ \text{s.t.} \quad x_1^2 + x_2^2 = 1 \end{array} \right.$$

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$$

Review of some statistical concepts

Consider data in a matrix X of size $n \times p$



Large $n, p \rightarrow$ "Big data"

$n < p \rightarrow$ "Fat data"

Although entries of X could be categorical ("Yes"/"No") here we consider the simple case in which they are numbers.

Normally $n \gg p$ (much larger), i.e. the number of individuals exceeds the number of variables.

Example

X : running times of 10 persons ($n=10$) in 3 different types of ground ($p=3$)

	H: 11	Flat	Tarmac
1	19	17	19
2	29	21	25
⋮	⋮	⋮	⋮
10	26	18	24

} X
 $n \times p$

We have the following concepts

Population

• $X = (X_1, X_2, \dots, X_p)^T \in \mathbb{R}^p$
 random vector

• $E(X) = (E(X_1), E(X_2), \dots, E(X_p))^T$
 vector of expectations
 (population means)

• Covariance between two random variables X_1, X_2
 $Cov(X_1, X_2) = E(X_1 X_2) - E(X_1)E(X_2)$

• Variance-Covariance matrix (symmetric)
 $Cov(X) = E(X X^T) - E(X) \cdot E(X)^T$

$\begin{matrix} p \times 1 & 1 \times p \\ \hline & p \times p \end{matrix}$
 $\begin{matrix} p \times 1 & 1 \times p \\ \hline & p \times p \end{matrix}$

→ Diagonal contains the variances of X_i 's
 $Var(X_i) = E(X_i^2) - E(X_i)^2$

→ off-diagonal entries (i, j) are the pairwise covariances
 $Cov(X_i, X_j)$

Sample

Data X , matrix of size $n \times p$ (typically $n > p$)

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}$$

$\left(\frac{1}{n} \sum_{i=1}^n X_{i1}, \frac{1}{n} \sum_{i=1}^n X_{i2}, \dots, \frac{1}{n} \sum_{i=1}^n X_{ip} \right)$
 \bar{X} ← vector of means.

$\bar{X} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)$

$= \frac{1}{n} [1 \ 1 \ \dots \ 1] X$

$= \frac{1}{n} (\mathbf{1}^T X)$
 $\begin{matrix} 1 \times n & n \times p \end{matrix}$

• Sample Variance-Covariance matrix

$\Sigma = \frac{1}{n-1} \left(X^T X - n \bar{X} \bar{X}^T \right)$

$\begin{matrix} n \times p & p \times n \\ \hline & p \times p \end{matrix}$
 $\begin{matrix} 1 & 1 \\ \hline p \times 1 & 1 \times p \\ \hline & p \times p \end{matrix}$

→ Diagonal contains the sample variances of each column of X

$\frac{1}{n-1} \sum_{i=1}^n X_{i5}^2 - \frac{n}{n-1} \bar{X}_5^2$

$= \frac{1}{n-1} \sum_{i=1}^n (X_{i5} - \bar{X}_5)^2$

The standard deviation is the square root of the sample variance

(6)

\hookrightarrow off-diagonal entries (i, j) are the covariances between i -th, j -th column of X

$$\frac{1}{n-1} \sum_{k=1}^n X_{ki} X_{kj} - \frac{1}{n-1} \bar{X}_i \bar{X}_j$$

$$= \frac{1}{n-1} \sum_{k=1}^n (X_{ki} - \bar{X}_i)(X_{kj} - \bar{X}_j)$$

If C is a matrix such that CX is well-defined

$$E[CX] = C E[X]$$

$$\text{Cov}(CX) = C^T \text{Cov}(X) C$$

Without loss of generality in our sample data analysis we will assume that every column of X is centered around its mean

This is done by replacing

$$X_{ij} \rightsquigarrow X_{ij} - \bar{X}_j$$

In that case

$$\bar{X} = (0, 0, \dots, 0)$$

$$\Sigma = \frac{1}{n-1} X^T X$$

← this is now the "new X "

Note: This "centering" operation does not change the variance-covariance matrix Σ (prove it!)

Example
 $n=3, p=2$

"raw" data
(before centering)

$$\begin{bmatrix} 1 & 2 \\ 2 & 2 \\ 3 & 8 \end{bmatrix}$$

means

$$\begin{bmatrix} 2 & 4 \end{bmatrix}$$

centering



$$\begin{bmatrix} -1 & -2 \\ 0 & -2 \\ 1 & 4 \end{bmatrix}$$

mean

$$\begin{bmatrix} 0 & 0 \end{bmatrix}$$

X
↓

(7)

$$\begin{aligned} \Sigma &= \frac{1}{n-1} X^T X = \frac{1}{3-1} \begin{bmatrix} -1 & 0 & 1 \\ -2 & -2 & 4 \end{bmatrix} \begin{bmatrix} -1 & -2 \\ 0 & -2 \\ 1 & 4 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} 2 & 6 \\ 6 & 24 \end{bmatrix} = \begin{bmatrix} 1 & 3 \\ 3 & 12 \end{bmatrix} \end{aligned}$$

Exercise: Compute Σ for the "raw" data using

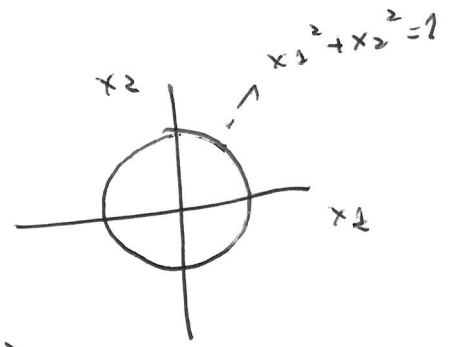
$$\Sigma = \frac{1}{n-1} X^T X - \frac{n}{n-1} \bar{X} \bar{X}^T$$

and check that it coincides with



Previous exercise

$$\begin{cases} \max_x & x_1 + x_2 \\ \text{st} & x_1^2 + x_2^2 = 1 \end{cases} \quad x = (x_1, x_2)$$



$$L(x_1, x_2, \lambda) = x_1 + x_2 - \lambda(x_1^2 + x_2^2 - 1)$$

$$\frac{\partial L}{\partial x_1} = 1 - 2\lambda x_1 \quad \begin{matrix} \frac{\partial L}{\partial x_1} = 0 \\ \Rightarrow \end{matrix} \quad 1 = 2\lambda x_1 \quad \leadsto \quad 1 = 4\lambda^2 x_1^2 \quad (1)$$

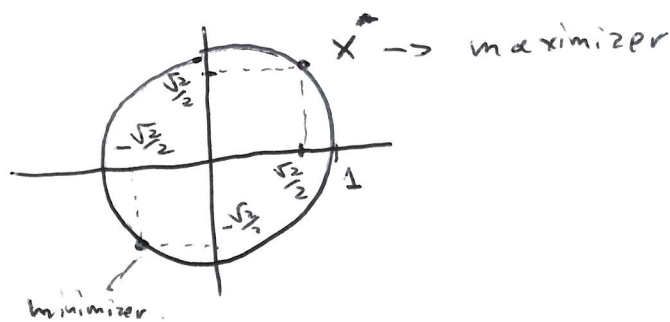
$$\frac{\partial L}{\partial x_2} = 1 - 2\lambda x_2 \quad \begin{matrix} \frac{\partial L}{\partial x_2} = 0 \\ \Rightarrow \end{matrix} \quad 1 = 2\lambda x_2 \quad \leadsto \quad 1 = 4\lambda^2 x_2^2 \quad (2)$$

$$\frac{\partial L}{\partial \lambda} = -(x_1^2 + x_2^2 - 1) \quad \begin{matrix} \frac{\partial L}{\partial \lambda} = 0 \\ \Rightarrow \end{matrix} \quad x_1^2 + x_2^2 = 1 \quad (3)$$

$$(1) + (2) \Rightarrow 2 = 4\lambda^2(x_1^2 + x_2^2) \Rightarrow \lambda = \pm \frac{\sqrt{2}}{2}$$

$$\text{if } \lambda = \frac{\sqrt{2}}{2} \Rightarrow x_1 = x_2 = \frac{\sqrt{2}}{2} \quad \leadsto \quad x_1 + x_2 = \sqrt{2} \quad \checkmark$$

$$\text{if } \lambda = -\frac{\sqrt{2}}{2} \Rightarrow x_1 = x_2 = -\frac{\sqrt{2}}{2} \quad \leadsto \quad x_1 + x_2 = -\sqrt{2}$$



Principal Component Analysis (PCA)

Objective: Study projections of the data X that best reproduce its variability

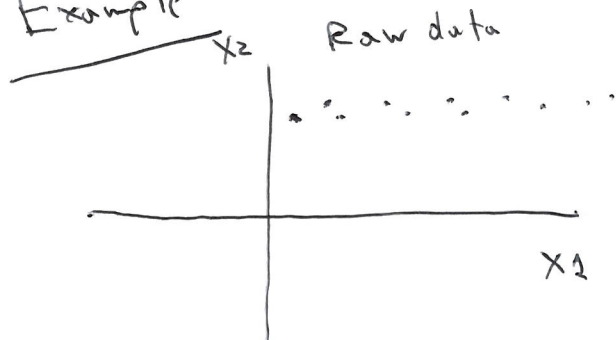
\Rightarrow dimensionality \downarrow , variability: remain as much as original data.

\Rightarrow Useful for: analysis, visualisation, storage

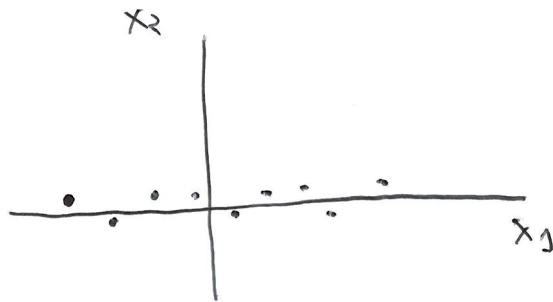
Data X , matrix $n \times p$, and we assume that every column of X has been centered around its mean such that $\bar{X} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_p) = (0, 0, \dots, 0)$

Note: This centering of X does not alter the variance-covariance matrix.

Example 1



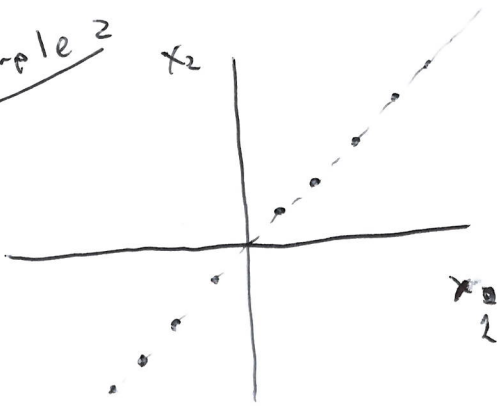
centering \rightarrow



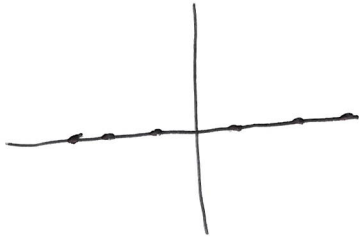
$$X: \begin{bmatrix} X_1 & X_2 \\ 11 & -0.02 \\ 12 & 0.01 \\ 23 & 0.002 \\ -31 & -0.03 \\ \vdots & \vdots \\ 25 & 0.01 \\ -50 & -0.005 \end{bmatrix}_{n \times 2} \quad \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}_{2 \times 1} \quad \begin{matrix} \alpha_1 = 1 \\ \sim \\ \alpha_2 = 0 \end{matrix} \quad \begin{bmatrix} 11 \\ 12 \\ 23 \\ -31 \\ \vdots \\ 25 \\ -50 \end{bmatrix}_{n \times 1}$$

Note: $\alpha_1^2 + \alpha_2^2 = 1$

Example 2



↳ rotate the data



$$X = \begin{matrix} & x_1 & x_2 \\ \begin{bmatrix} 11 \\ 12 \\ 23 \\ -31 \\ \vdots \\ 25 \\ -50 \end{bmatrix} & \begin{bmatrix} 11 \\ 12 \\ 23 \\ -31 \\ \vdots \\ 25 \\ -50 \end{bmatrix} & \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} \end{matrix}$$

It will turn out that

$$d_1 = \frac{\sqrt{2}}{2}$$

$$d_2 = \frac{\sqrt{2}}{2}$$

(we will come back to this)

Aim: Create a linear projection of the data X that has maximal variance

$$z_1 = X d_1$$

vector of projections $n \times 1$ data matrix $n \times p$ vector of constants size $p \times 1$ (unknown at this stage) $\|d_1\| = 1$

Check: z_1 has zero mean

We want to find d_1 such that the variance of z_1 is maximised. The sample variance of z_1 is:

$$\frac{1}{n-1} \sum_{i=1}^n z_{1i}^2 = \frac{1}{n-1} z_1^T z_1 = \frac{1}{n-1} (X d_1)^T X d_1$$

$$= \frac{1}{n-1} d_1^T \underbrace{X^T X}_{\Sigma} d_1 = d_1^T \left(\frac{1}{n-1} X^T X \right) d_1$$

Σ : variance-covariance matrix

$$= \boxed{d_1^T \Sigma d_1}$$

$\downarrow \quad \downarrow \quad \downarrow$
 $1 \times p \quad p \times p \quad p \times 1$
 $\underbrace{\hspace{10em}}_{1 \times 1}$

← Quadratic form

PCA

$$z_1 = X_{n \times p} \cdot d_1$$

$$\text{Var}(z_1) = d_1^T \Sigma d_1$$

maximize $d_1^T \Sigma d_1$ subject to $d_1^T d_1 = 1$
" $\|d_1\|_2^2$ "

"maximize the variance of the projection, using projection vector of unit norm"

$$\sum_{i=1}^p (d_{1,i})^2 = 1$$

We turn to the method of Lagrange multipliers in order to solve this problem.

$$L = \underbrace{d_1^T \Sigma d_1}_{\text{"target" / "objective function" / "sample variance of the projection"}} - \lambda (d_1^T d_1 - 1)$$

"target"
"objective function"
sample variance of the projection

↓
Lagrange multiplier

↘ "constraint"
vector of coefficients has unit norm
 $d_1^T d_1 = 1$ re-written as $d_1^T d_1 - 1 = 0$

Total unknowns: p elements for d_1
1 Lagrange multipliers

Taking derivatives of the Lagrangian

(1) $\frac{\partial L}{\partial d_1} = 2 \Sigma d_1 - 2\lambda d_1 = 2(\Sigma d_1 - \lambda d_1)$

$$\begin{cases} f(s) = s^T B s \\ \frac{\partial f}{\partial s} = 2 B s \\ d_1^T d_1 = d_1^T I d_1 \end{cases}$$

(2) $\frac{\partial L}{\partial \lambda} = -(d_1^T d_1 - 1)$

(2)

• Build the system of equations by setting these derivatives to zero

$$\begin{aligned}
 (1) \quad 2(\Sigma z_1 - \lambda z_1) &= 0 \Rightarrow \Sigma z_1 = \lambda z_1 \\
 (2) \quad -(z_1^T z_1 - 1) &= 0 \Rightarrow z_1^T z_1 = 1
 \end{aligned}$$

This is nothing else than an eigenvalue problem:

λ is an eigenvalue of Σ
 z_1 is the corresponding eigenvector of unit norm

} candidates for solution of the maximization problem

Reinserting the above into

$$\begin{aligned}
 V(z_1) &= z_1^T \Sigma z_1 \\
 &= z_1^T (\Sigma z_1) = z_1^T \lambda z_1 = \lambda \underbrace{z_1^T z_1}_{1} = \lambda
 \end{aligned}$$

\Rightarrow In order to maximize the variance of z_1 , we take λ to be the largest eigenvalue of Σ , with z_1 the corresponding eigenvector.

(what about $-z_1$? this will also work, so some of the relations we will describe next are always "up to change of sign")

Summary: The eigenvector z_1 of Σ gives the coefficients, to build the linear projection $X z_1$ of the data, and the corresponding eigenvalue is the sample variance of the projection.

$z_1 \rightarrow$ First principal component
Second and following components

After we have picked z_1 , now we define

$$z_2 = X z_2$$

③

- We require again:
- 1) $V(z_2)$ is maximized as a function of d_2
 - 2) d_2 has unit norm, $d_2^T d_2 = 1$
 - 3) d_2 has to be orthogonal to d_1

Ex. Show that the Lagrangian for the second principal component is

$$L = d_2^T \Sigma d_2 - \lambda_1 (d_2^T d_2 - 1) - \lambda_2 (d_2^T d_1 - 0)$$

We end up:

$$\begin{cases} \Sigma d_2 = \lambda_2 d_2 \\ d_2^T d_2 = 1 \\ d_2^T d_1 = 0 \end{cases}$$

In order to maximize $V(z_2) = d_2^T \Sigma d_2$, I take the eigenvector d_2 of unit norm, that corresponds to the second largest eigenvalue of Σ

Recall that $\Sigma = \frac{1}{n-1} X^T X$ is symmetric and if $\lambda_1 > \lambda_2$ are eigenvalues of Σ , with d_1, d_2 are the corresponding eigenvectors $\Rightarrow d_2^T d_1 = 0$

The pattern is clear:

For maximizing the variance of the projection in a sequential manner:
 projection vectors d : eigenvectors of variance-covariance matrix Σ

Variance of the projection: $d^T \Sigma d$ eigenvalues λ of Σ

$$\textcircled{4} \quad z_i = X \alpha_i \quad i=1, \dots, p$$

Σ is the sample variance-covariance matrix of data X (centered and possibly scaled) and let Λ and A be the matrices of the Spectral decomposition of Σ (recall $\Sigma = \frac{1}{n-1} X^T X$, symmetric)

$$\Sigma = A \Lambda A^T$$

Here: $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ with λ_i being the eigenvalues of Σ in decreasing order $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$

$A = \begin{bmatrix} \alpha_1 & \alpha_2 & \dots & \alpha_p \end{bmatrix}$ is the matrix with the corresponding eigenvectors α_i as columns.

(In general (in practice) we deal with "well-behaved" matrices, i.e. all eigenvalues are different and non-zero)

PCA terminology

The eigenvectors α_i are also known as "loadings" (PC loadings)

The transformed variables $z_1 = X \alpha_1, z_2 = X \alpha_2, \dots$ are known as "scores" (PC scores)

⑤ We can collect all the PC scores (z_1, z_2, \dots, z_p) as columns in a matrix Z

$$Z = \begin{bmatrix} | & | & & | \\ z_1 & z_2 & \dots & z_p \\ | & | & & | \end{bmatrix}_{n \times p} = \begin{bmatrix} | & | & & | \\ X_{d1} & X_{d2} & \dots & X_{dp} \\ | & | & & | \end{bmatrix}$$

$$= \underbrace{X}_{n \times p} \underbrace{A}_{p \times p}$$

We can now compute the sample variance-covariance matrix of Z

(recall that every column of Z has zero mean (why?))

$$\frac{1}{n-1} Z^T Z = \frac{1}{n-1} (XA)^T (XA)$$

$$= \frac{1}{n-1} A^T X^T X A$$

$$= A^T \left[\frac{1}{n-1} X^T X \right] A$$

$$= A^T \Sigma A$$

$$= \underbrace{A^T A}_{I} \Lambda \underbrace{A^T A}_{I}$$

$$= \Lambda$$

$$(\Sigma = A \Lambda A^T)$$

(~~diag~~ diagonal matrix)
 $\begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_p & \\ 0 & & \lambda_p \end{bmatrix}$

From the above we observe:

- The sample variance of z_i is the eigenvalue λ_i
- The z_i 's (scores) are uncorrelated!
- The total variance of X equals to the total
variance of Z

⑥ Indeed for the third bullet point:

$$\sum_{i=1}^p \text{Var}(X_i) = \text{tr}(\Sigma) = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \text{Var}(z_i)$$

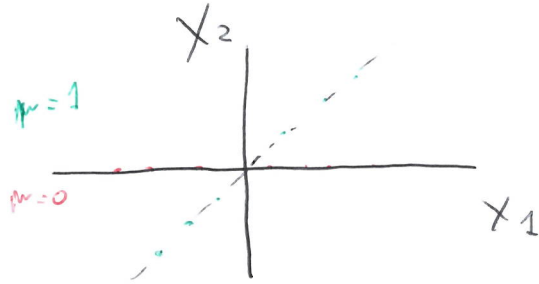
sum of diagonal entries of Σ
(trace of Σ)

General result:

trace of a matrix
= sum of its eigenvalues

Exercise

$$X = \begin{bmatrix} X_1 & \mu X_1 \end{bmatrix}$$



\Rightarrow Do the PCA: find Z, A, Λ

$$A = \begin{bmatrix} d_1 & d_2 \end{bmatrix}$$

$$\text{if } \mu=0 \quad X = \begin{bmatrix} X_1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} X_1 \end{bmatrix}$$

\downarrow
 d_1

What happens: $\mu=1$
 $\mu=-1$?

① PCA continued: Some comments on previous exercises:

Second Principal Component

$$\begin{aligned} \max \quad & d_2^T \Sigma d_2 \\ \text{s.t.} \quad & d_2^T d_2 = 1 \\ & d_1^T d_2 = 0 \end{aligned}$$

$$L = d_2^T \Sigma d_2 - \lambda (d_2^T d_2 - 1) - \mu d_1^T d_2$$

$$\frac{\partial L}{\partial d_2} = 2 \Sigma d_2 - 2\lambda d_2 - \mu d_1 \stackrel{=0}{\Rightarrow} 2 \Sigma d_2 - 2\lambda d_2 = \mu d_1$$

$$\frac{\partial L}{\partial \lambda} = -(d_2^T d_2 - 1) \stackrel{=0}{\Rightarrow} d_2^T d_2 = 1$$

$$\frac{\partial L}{\partial \mu} = -d_1^T d_2 \stackrel{=0}{\Rightarrow} d_1^T d_2 = 0$$

I claim $\mu=0$

Multiply the first equation with d_1^T from the left:

$$2 \underbrace{d_1^T \Sigma d_2}_{=} - 2\lambda \underbrace{d_1^T d_2}_{=0} = \mu \underbrace{d_1^T d_1}_{=1}$$

$$2 \underbrace{d_2^T \Sigma d_1}_{=}$$

$$2\lambda \underbrace{d_2^T d_1}_{=0}$$

$$\boxed{\mu=0}$$

In summary:

$$\boxed{\begin{aligned} \Sigma d_2 &= \lambda d_2 \\ d_2^T d_2 &= 1 \\ d_1^T d_2 &= 0 \end{aligned}}$$

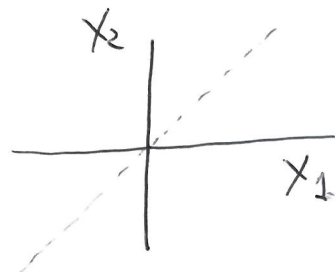
Example

Suppose data matrix X of the following form:

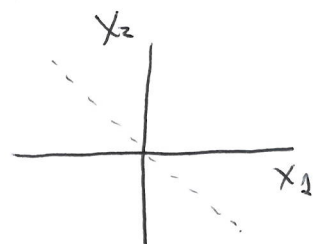
if $\rho = 1$ $X = \begin{bmatrix} 1 & 1 \\ -1 & -1 \end{bmatrix}$
 X_1, X_2 pos. correlated

$$X = \begin{bmatrix} 1 & \rho \\ -1 & -\rho \end{bmatrix}$$

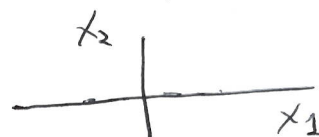
X_1 $X_2 = \rho X_1$



if $\rho = -1$ $X = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$
 X_1, X_2 neg. correlated



if $\rho = 0$ $X = \begin{bmatrix} 1 & 0 \\ -1 & 0 \end{bmatrix}$



$$\Sigma = \frac{1}{n-1} X^T X = \begin{bmatrix} 1 & -1 \\ \rho & -\rho \end{bmatrix} \begin{bmatrix} 1 & \rho \\ -1 & -\rho \end{bmatrix}$$

$$= \begin{bmatrix} 2 & 2\rho \\ 2\rho & 2\rho^2 \end{bmatrix} = 2 \begin{bmatrix} 1 & \rho \\ \rho & \rho^2 \end{bmatrix}$$

Total variance of X : $2(1 + \rho^2)$

Look at eigenvalues of $\begin{bmatrix} 1 & \rho \\ \rho & \rho^2 \end{bmatrix}$

$$\begin{vmatrix} 1-\lambda & \rho \\ \rho & \rho^2-\lambda \end{vmatrix} = (1-\lambda)(\rho^2-\lambda) - \rho^2 = \rho^2 - \lambda - \lambda\rho^2 + \lambda^2 - \rho^2$$

$$= \lambda(\lambda - 1 - \rho^2)$$

The only positive eigenvalue of Σ is $2(1 + \rho^2)$

Eigenvectors: $\begin{bmatrix} 2 - 2(1 + \rho^2) & 2\rho \\ 2\rho & 2\rho^2 - 2(1 + \rho^2) \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = 0$

$$\begin{bmatrix} -2\rho^2 & 2\rho \\ 2\rho & -2 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = 0$$

$$\boxed{d_2 = \rho d_1}$$

$$\begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = d_1 \begin{bmatrix} 1 \\ \rho \end{bmatrix}$$

$$\textcircled{3} \quad \mu=1 \quad \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}$$

$$\left(\frac{\sqrt{2}}{2}\right)^2 + \left(\frac{\sqrt{2}}{2}\right)^2 = 1$$

$$\mu=-1 \quad \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{bmatrix}$$

$$\left(\begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix} \text{ work as well} \right)$$

$$\mu=0 \quad \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

①

PCA and SVD

So far we used the spectral decomposition of $\Sigma = A \Lambda A^T$. All the results can be recovered using the SVD of X :

$$\Sigma = A \Lambda A^T$$

$\begin{matrix} p \times p & p \times p & p \times p & p \times p \end{matrix}$

$$X = U D V^T$$

\downarrow
 $\min(n, p) \times \min(n, p)$

Let us build Σ using the SVD of X

$$\begin{aligned} \Sigma = \frac{1}{n-1} X^T X &= \frac{1}{n-1} V D V^T U D V^T \\ &= \frac{1}{n-1} V D^2 V^T \\ &= V \left(\frac{1}{n-1} D^2 \right) V^T \end{aligned}$$

We conclude:

a) Eigenvectors of Σ are the columns of V is SVD of X
 PC loadings [up to a sign]

$$A = V$$

b) The eigenvalues of Σ are related to D (singular values)

$$\lambda_i = \frac{1}{n-1} d_i^2 \quad D = \begin{bmatrix} d_1 & d_2 & & 0 \\ & & \dots & \\ 0 & & & d_p \end{bmatrix}$$

Let us finish by computing also the PC scores z :

$Z = X A$ and now using the SVD

$$Z = X A = U D V^T A \underset{A=V}{=} U D \underbrace{A^T A}_I = U D$$

$$Z = X A = U D \quad (\text{up to a sign})$$

②

Interpreting PCA output

The main objective of PCA is to reduce dimensionality of the data and select only a few principal components, that best represent the variability of the data.

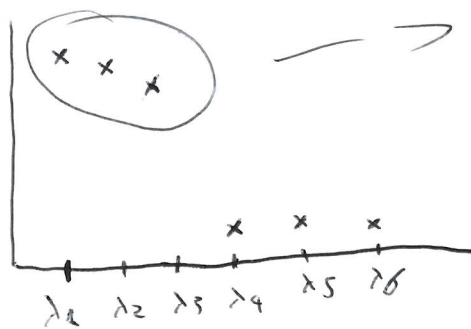
How many components to choose?

Strategies to select the number of components:

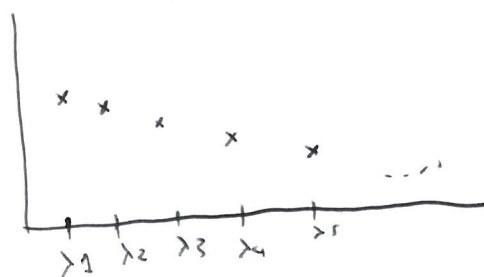
Screen plot

Plot eigenvalues (sample variances of PC scores) in a decreasing order of magnitude, look where the plot has a sharp decrease/becomes flat.
=> Select number of components according to this.

Select 3 components



(easy case)



(hard case)

③ Tabulate eigenvalues and cumulative sum > scale it to be a percentage.

	PC 1	PC 2	PC 3
Variance	52.2	15.11	1.98
Proportion of variance	0.75	0.21	0.02
Cumulative Proportion	0.75	0.96	1.00

→ Set a threshold: of how much variability you want to explain and then take the minimum number of components that explain this variability

Example: if threshold 80% = 0.8
 ⇒ take PC1, PC2
 if threshold 99% = 0.99
 ⇒ take PC1, PC2, PC3

Summarize

The proportion of variance explained by the first k -components

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p}$$

→ variance of first k -components

→ total variance of all components.

For a fixed threshold pick the smallest value of k such that the above proportion exceeds this threshold.

④

Interpreting PCA as weighted averages and contrasts

Look at PC loadings (eigenvectors) and we inspect their signs and magnitudes

These are used to interpret PCA as weighted averages and/or as contrasts

	PC1	PC2	} PC2 is essentially a weighted average of "Hill" and "Tarmac" flat
Hill	0.89	-0.17	
Flat	-0.09	-0.98	
Tarmac	0.75	0.02 \rightarrow too small	

\downarrow
weighted average of "Hill" and "Tarmac" (same signs)
compared against "Flat" (different signs)

Remember: These interpretations should be the same when all signs are reversed

$$\left[\begin{array}{l} \max \alpha^T \Sigma \alpha \\ \|\alpha\|_2 = 1 \end{array} \quad \alpha \rightarrow -\alpha \right]$$

Should we scale the data or not?

A crucial feature of PCA is that analyses and all conclusions differ between data that have been scaled and data have not been scaled.

③ Common practice is to analyse scaled variables (centered and scaled) only if the units of measurements are the same for all the variables. Note in that case, the proportion of variability explained by the first k components is $\frac{\sum_{i=1}^k \lambda_i}{p}$

Sample variance-covariance matrix for the scaled data

$$R = S^{-1/2} \Sigma S^{-1/2} \rightarrow S: \text{diagonal matrix whose diagonal is equal to } \Sigma$$

$$R = S^{-1/2} \Sigma S^{-1/2} = S^{-1/2} A \Lambda A^T S^{-1/2} \Rightarrow \text{multiply } S^{1/2} A$$

$$R S^{1/2} A = S^{-1/2} A \Lambda$$

Compare $\Sigma A = A \Lambda$

$\Rightarrow R$ & Σ have different eigenvalues / eigenvectors.

Biplots

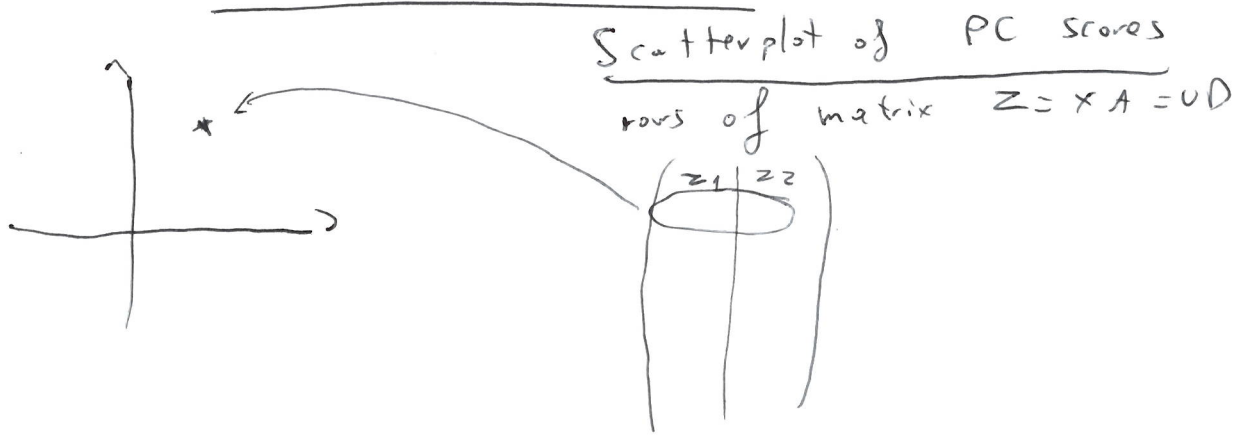
A biplot is a plot that simultaneously shows the PC scores and the PC loadings.

PC scores: we ~~see~~ usually plot the first two in a scatterplot.

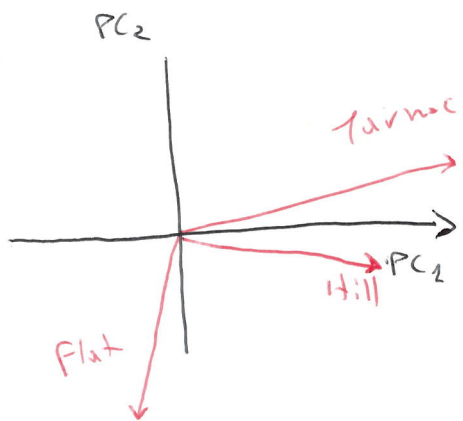
PC loadings: arrows, typically as many as variables.

⑥ It is used to associate observations with variables in the first two principal components. We also interpret the loadings in a graphical way.

Elements of a biplot



Scatterplot of PC loadings



	PC ₁	PC ₂
Hill	0.09	-0.17
Flat	-0.09	-0.98
Tarmac	0.75	-0.02

Typically in scaled data, there is no clear dominance of any variable \rightarrow arrows have approximately the same length.

PCA is often the first step in reducing complexity of a dataset, by keeping linear combinations of data that maximize the variance.

This method is typically used in conjunction with other methods of machine learning:

Example: \rightarrow do PCA \rightarrow clustering
 \searrow do PCA \rightarrow regression

(1400) 3 - Clustering

The aim of clustering is to detect groups "clusters" in the individuals that compose the data.

Data matrix $X = \begin{pmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{pmatrix}$ ← rows are individuals
 ← entries are numbers
 ← "a vector of p numbers"
 $n \times p$
 ← columns are variables

Cluster - group of individuals

We will look at $\begin{cases} \text{agglomerative clustering} \\ \text{K-means (medoids) method} \end{cases}$

3.1 Distances and distance matrix

Assume that we have data on $n=5$ individuals and $p=2$ variables

$$X = \begin{matrix} & \begin{matrix} \text{individual labels (index)} \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 0 & 4 \\ 3 & 6 \\ 6 & 2 \\ 0 & 5 \\ 1 & 1 \end{pmatrix} \end{matrix}$$

We will consider distances between points/individuals that are a metric.

Metric { The distance d_{ij} between individuals i and j satisfies the following properties: non-negative, symmetric, satisfies the triangle inequality and takes value zero between an individual and itself.

→ Non-negative $d_{ij} \geq 0$ For all i, j

→ Symmetric $d_{ij} = d_{ji}$ for all i, j

→ Triangle inequality $d_{ij} + d_{jk} \geq d_{ik}$ for all i, j, k

→ Zero distance with itself $d_{ii} = 0$ for all i

Let us consider the 'Manhattan distance between points' in the example

$$d_{12} = |0-3| + |4-6| = 5$$

$$d_{13} = |0-6| + |4-2| = 8$$

$$d_{14} = |0-0| + |4-5| = 1$$

$$d_{15} = |0-1| + |4-1| = 4$$

Exercise. Check/compute $d_{23}, d_{24}, d_{25}, d_{34}, d_{35}, d_{45}$ as well as d_{11}, \dots
 d_{21}, d_{31}, \dots

Distance matrix: collection of distances d_{ij} in a $n \times n$ matrix

Notation $\underline{D} = (d_{ij})$

$$\underline{D} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 0 & 5 & 8 & 1 & 4 \\ 5 & 0 & 7 & 4 & 7 \\ 8 & 7 & 0 & 9 & 6 \\ 1 & 4 & 9 & 0 & 5 \\ 4 & 7 & 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

can also write

$$\underline{D} = \begin{pmatrix} - & 5 & 8 & 1 & 4 \\ 5 & - & 7 & 4 & 7 \\ 8 & 7 & - & 9 & 6 \\ 1 & 4 & 9 & - & 5 \\ 4 & 7 & 6 & 5 & - \end{pmatrix}$$

We consider the following distances between individuals/points $\underline{x}_i = (x_{i1} \ x_{i2} \ \dots \ x_{ip})$ and

$$\underline{x}_j = (x_{j1} \ x_{j2} \ \dots \ x_{jp})$$

→ Manhattan $d_{ij} = \sum_{k=1}^p |x_{ik} - x_{jk}| = \|\underline{x}_i - \underline{x}_j\|_1$ L_1 norm

→ Euclidean $d_{ij} = \sqrt{\sum_{k=1}^p |x_{ik} - x_{jk}|^2} = \|\underline{x}_i - \underline{x}_j\|_2$ L_2 norm

→ Minkowski: $d_{ij} = \left(\sum_{k=1}^p |x_{ik} - x_{jk}|^m \right)^{1/m} = \|\underline{x}_i - \underline{x}_j\|_m$ L_m norm
(L_p norm in books)

→ Mahalanobis $d_{ij} = \sqrt{(\underline{x}_i - \underline{x}_j) \Sigma^{-1} (\underline{x}_i - \underline{x}_j)^T}$
 $p \times p$ covariance matrix between columns

→ Absolute correlation $d_{ij} = \sqrt{1 - |\rho_{ij}|}$ ← correlation between points

3.2 Agglomerative clustering

This is a sequential procedure to cluster 'n' individuals.

- a) Initially, each individual is its own cluster
- b) At every step, an updated table of dissimilarities between clusters is considered and minimized to merge an individual to a cluster.
- c) At the end of the procedure, a single cluster is achieved.

①

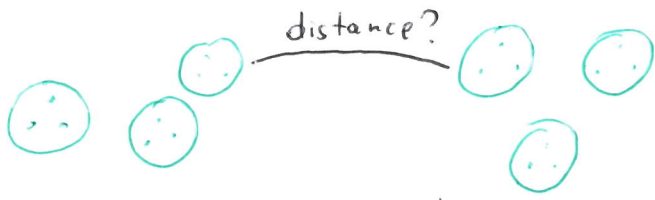
Agglomerative clustering

Two ingredients to be specified:

distance
↓
dissimilarity between points.

linkage
↓
dissimilarity between sets of points (clusters)

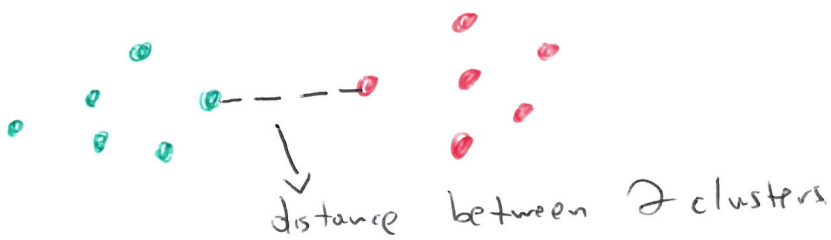
Dissimilarities: "distance" between clusters that contain more than one point



The way these dissimilarities are computed is dictated by the linkage

Single linkage

The dissimilarity between clusters is the distance between the closest 2 points in each cluster (this is also known as the "nearest neighbour")



Average linkage

The dissimilarity between clusters is the average of distances between points in each cluster ("average neighbour")



② Complete linkage

The dissimilarity between clusters is the distance between the two points in each cluster that are farthest apart ("furthest neighbour")



To perform agglomerative clustering we select a

- distance (e.g. Euclidean, Manhattan, ...)
- linkage (single, average, complete)

Example Using Manhattan
single linkage

$$X = \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} \begin{pmatrix} 0 & 4 \\ 3 & 6 \\ 6 & 2 \\ 0 & 5 \\ 1 & 1 \end{pmatrix}$$

$$D = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} - & 5 & 8 & 1 & 4 \\ 5 & - & 7 & 4 & 7 \\ 8 & 7 & - & 9 & 6 \\ 1 & 4 & 9 & - & 5 \\ 4 & 7 & 6 & 5 & - \end{pmatrix} \end{matrix}$$

a) Initially, every individual is its own cluster.
At this stage the distance matrix is the dissimilarity matrix

$$\{1, 2, 3, 4, 5\}$$

b) Merge clusters 1 and 4, at distance 1

$$\Rightarrow \{14, 2, 3, 5\}$$

Update dissimilarity matrix

$$\begin{matrix} & \begin{matrix} 14 & 2 & 3 & 5 \end{matrix} \\ \begin{matrix} 14 \\ 2 \\ 3 \\ 5 \end{matrix} & \begin{pmatrix} - & 4 & 8 & 4 \\ & - & 7 & 7 \\ & & - & 6 \\ & & & - \end{pmatrix} \end{matrix}$$

③ Compute dissimilarity between 14 & 2, $d_{14,2}$

We have $d_{1,2} = 5$, $d_{4,2} = 4$

"Single" $d_{14,2} = \min(d_{1,2}, d_{4,2}) = \min(5, 4) = 4$

"Average" $d_{14,2} = \frac{d_{1,2} + d_{4,2}}{2} = 4.5$

"complete" $d_{14,2} = \max(d_{1,2}, d_{4,2}) = \max(5, 4) = 5$

\leadsto Clusters 14 & 2 are merged at distance 4

$\Rightarrow \{124, 3, 5\}$

Update dissimilarity matrix:

$$\begin{matrix} & \begin{matrix} 124 & 3 & 5 \end{matrix} \\ \begin{matrix} 124 \\ 3 \\ 5 \end{matrix} & \begin{pmatrix} - & 7 & 4 \\ & - & 6 \\ & & - \end{pmatrix} \end{matrix}$$

$$d_{124,3} = \min(8, 7, 9) = 7$$

$$d_{124,5} = 4$$

\Rightarrow Clusters 124 & 5 are merged at distance 4

$\Rightarrow \{1245, 3\}$

Update dissimilarity matrix

$$\begin{matrix} & \begin{matrix} 1245 & 3 \end{matrix} \\ \begin{matrix} 1245 \\ 3 \end{matrix} & \begin{pmatrix} - & 6 \\ & - \end{pmatrix} \end{matrix}$$

$$d_{1245,3} = \min(8, 7, 9, 6) = 6$$

c) All clusters are merged at distance 6

$\{12345\}$

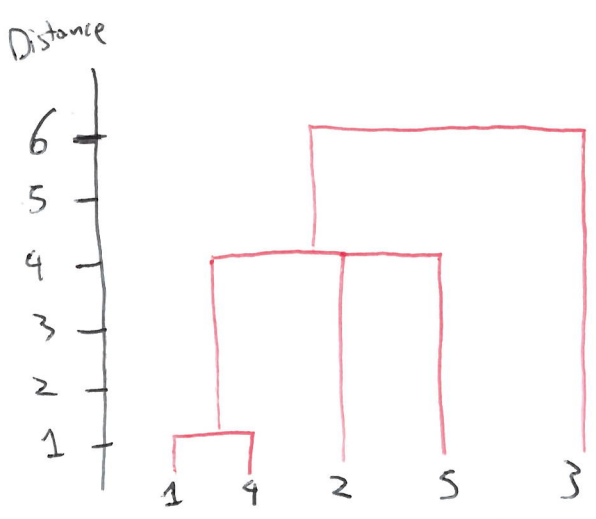
④

Dendrograms

"Dendro" ~ "tree" in Greek

Example (Single, Manhattan)

1	{1, 2, 3, 5}
4	{1, 2, 4, 3, 5}
4	{1, 2, 4, 5, 3}
6	{1, 2, 3, 4, 5}



The clustering process is summarized in a diagram called dendrogram. In this diagram, individual points are represented as "leaves".

As we move upwards the diagram, ~~leaves~~ clusters are formed and merged in a gradual manner joining the main part (~~the~~ "trunk") of the tree.

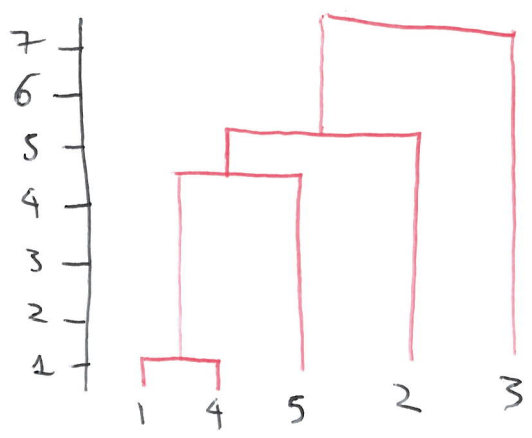
This joining happens at heights determined by the dissimilarity at every step.

Ordering of leaves in the dendrogram, ~~may~~ may not correspond to data order.

Another example (try to replicate!)

(Average, Manhattan)

Distance	Clusters formed
1	{1, 4, 2, 3, 5}
4.5	{1, 4, 5, 2, 3}
5.3	{1, 2, 4, 5, 3}
7.5	{1, 2, 3, 4, 5}



$$* \begin{cases} d_{1,5} = \frac{d_{1,4} + d_{4,5}}{2} \\ = 4.5 \end{cases}$$

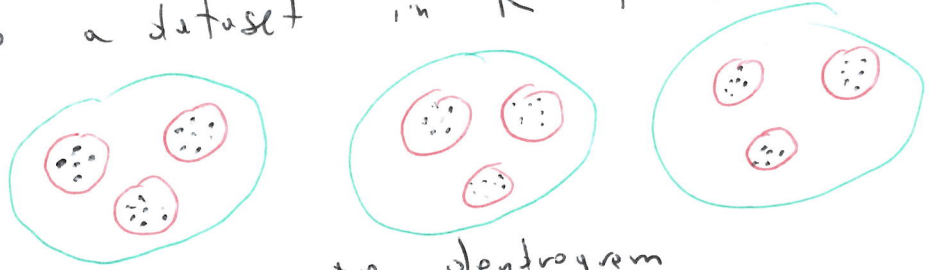
5) Height of Jumps

The height of jumps in the dendrogram suggests where to "cut" to obtain "meaningful" groups of individuals.

Big jumps in the dendrogram indicates clusters that are separate (far away from each other) and should not be merged.

Do the following for practice

Create a dataset in R that looks like:



⇒ Examine the dendrogram

• Results depend on the choice of metric and linkage. Different choices should be explored, to examine the clustering from different angles.

• Recall that clustering is a problem of unsupervised learning and thus the true number of clusters is not available.

R-demo

ac: agglomerative coefficient

Measure the dissimilarity of a point to the first cluster it joins, divided by the dissimilarity of the final merging → average across points

In R: $1 - \curvearrowright$

→ High values (close to 1) indicates well-formed clusters

→ Low values (close to 0) indicates less-well-formed clusters.

①

Method of k-means

This clustering methodology works on a different premise than agglomerative clustering. In a sense it is more automatic as the user only has to specify the number of clusters K .

An algorithm seeks to minimize a measure of closedness to clusters, by changing allocation of individuals to clusters.

Example. $K=2$

We would like to define a map

$$\left\{ \begin{array}{l} \text{all possible } K\text{-cluster} \\ \text{configurations} \end{array} \right\} \longrightarrow \mathbb{R}^T$$

such that



→ this type of cluster has small value



→ this type of cluster has high value.

Once we have defined such a map, then we choose the clustering configuration that minimises this map

② The function to minimise is the following:

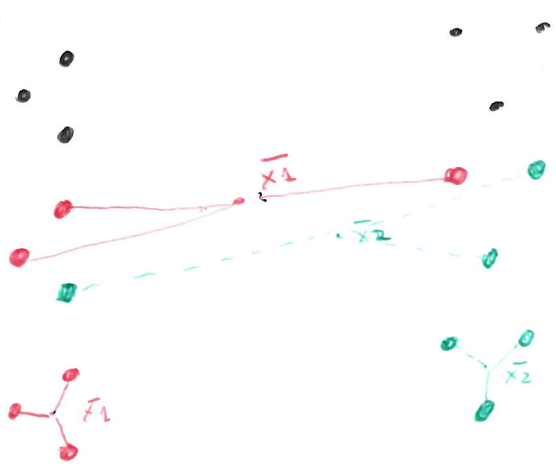
$$ESS = \sum_{k=1}^K \sum_{\substack{i \text{ st} \\ c(i)=k}} \|x_i - \bar{x}_k\|_2^2$$

\downarrow error sum of squares
 \downarrow sum over all clusters
 \downarrow sum over all points in cluster k

\downarrow i -th row of matrix X (belongs to the cluster k because of $c(i)=k$)

\bar{x}_k is a centroid for the cluster k . For instance, it can be the mean (vector of means per column/variable)

Example



\rightarrow these are my data X

\rightarrow ESS is large! obviously this not a "good" clustering

\rightarrow ESS is small! obviously this a "good" clustering

ESS essentially sums up all the distances between all the points and the centres of the clusters they belong to. (means)

Apart from the mean, it could be the case that \bar{x}_k is one of the points in cluster k such that the inner sum in ESS is minimised. In that case \bar{x}_k is called the medoid and the method is called k -medoids.

k -means \leadsto \mathbb{R} k -means
 k -medoids \leadsto \mathbb{R} pam

③ Example

Recall synthetic data:

$$X = \begin{pmatrix} 0 & 4 \\ 3 & 6 \\ 6 & 2 \\ 0 & 5 \\ 1 & 1 \end{pmatrix}$$

$n = 5$ individuals

Output K-means, $K = 2$

Cluster 1

Cluster 2

{ 1, 2, 4, 5 }

{ 3 }

Members
Center

(1, 4)

(6, 2)



①

K-means

In summary, K-means (K-medoids) is a procedure that splits the dataset into K clusters. This is done by minimising ESS and if we were to do this "by hand" we would have to use an exhaustive approach

Example

$$X = \begin{pmatrix} 0 & 4 \\ 3 & 6 \\ 6 & 2 \\ 0 & 5 \\ 1 & 1 \end{pmatrix}$$

n = 5 observations, K = 2 clusters.

- Number of ways of splitting these 5 observations into 2 clusters of 3 and 2 individuals

$$\binom{5}{2} = \binom{5}{3} = 10 \text{ different ways}$$

- Number of ways of doing this splitting into 2 clusters of 4 and 1 individuals.

$$\binom{5}{1} = \binom{5}{4} = 5 \text{ different ways}$$

Clusters
 123, 45
 124, 35
 ⋮
 12, 345
 1, 2345
 1245, 3
 1235, 4

ESS
 K-medoid
 55
 37
 ⋮
 56
 69
 24
 48

K-means
 39.5
 21
 ⋮
 38
 34.75
 20
 35.75

② Example: ESS value for cluster $\{1, 2, 4, 5, 3\}$ k -means, $k=2$

x_1 (0, 4) → cluster 1 with center $\bar{x}_1 = \left(\frac{0+3+0+1}{4}, \frac{4+5+5+1}{4}\right) = (1, 4)$
 x_2 (3, 6)
 x_3 (6, 2) → cluster 2 with center $\bar{x}_2 = (6, 2)$
 x_4 (0, 9)
 x_5 (1, 1)

$$ESS = \sum_{i=1,3,4,5} \|x_i - \bar{x}_1\|_2^2 + \sum_{i=3} \|x_i - \bar{x}_2\|_2^2$$

sum over elements in 1st cluster sum over elements in 2nd cluster
 = 0 since $\|x_3 - \bar{x}_2\|_2^2 = 0$

$\|x_1 - \bar{x}_1\|_2^2 = (0-1)^2 + (4-4)^2 = 1$
 $\|x_2 - \bar{x}_1\|_2^2 = 8$, $\|x_4 - \bar{x}_1\|_2^2 = 2$, $\|x_5 - \bar{x}_1\|_2^2 = 9$ sum ≈ 20

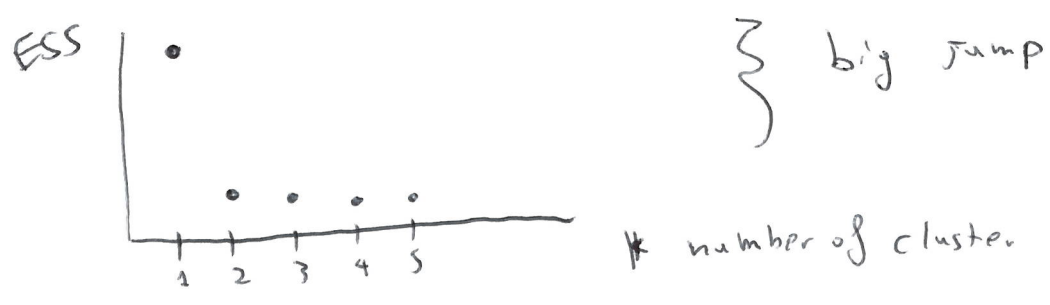
This exhaustive approach is too expensive and the number of possible clusters grows exponentially with n . In practice ESS is minimised numerically using

kmeans, pam

How many k ?

The quantity ESS is used as a tool to determine a "meaningful" number of clusters.

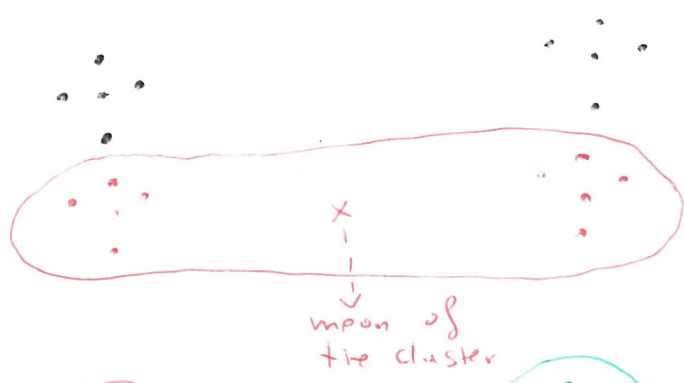
This is a subjective tool and there is no 100% guarantee that the numbers of clusters that we get from the following procedure is the "correct".



In this easy scenario, where we have a big jump from $k=1$ to $k=2$, we choose $k=2$ as the best option.

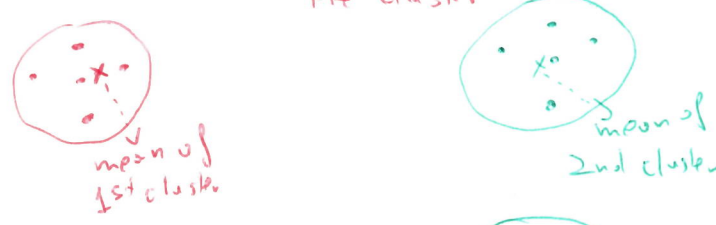
③ Example

$k=1$



- ESS value will be large

$k=2$



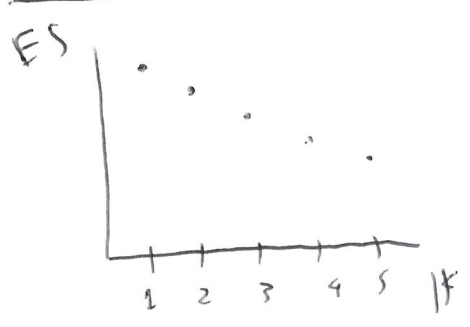
- ESS value will be much smaller! It will make a big jump

$k=3$



- ESS will be maybe a little bit smaller but it will not make a big jump

More difficult scenario



Difficult to determine a good value for k

That could happen when we have a dataset which is not so well clustered



Summary:

This kind of plot can be used to select potential number of clusters by looking at the decrease of ESS, while keeping k as small as possible

[Recall: $k=n \Rightarrow ESS=0$]

④

Silhouette plot

This is a tool that calculates and shows the silhouette value $s(i)$ for every individual i in the data.

This quantity is usually between 0 and 1 but sometimes it can also take negative values up to -1

Values of $s(i)$

close to 1

close to 0

close to -1 / negative

Interpretation

→ individual i is well-clustered

→ individual i is probably between 2 or more clusters.

→ individual i is most likely in the wrong cluster.

This $s(i)$ is defined as follows:

$$s(i) = \frac{b(i) - a(i)}{\max(b(i), a(i))}$$

$a(i)$: average dissimilarity (distance) between i and all other points of the cluster to which i belongs
(if i is alone, in its own cluster: $a(i) = 0$)



$a(i)$ is large



$a(i)$ is small

⑤ $d(i, C)$ = average dissimilarity between i and its ~~nearest~~ all observations in cluster C
 ↓
 cluster

$$b(i) = \min_C d(i, C)$$

dissimilarity between i and its nearest cluster

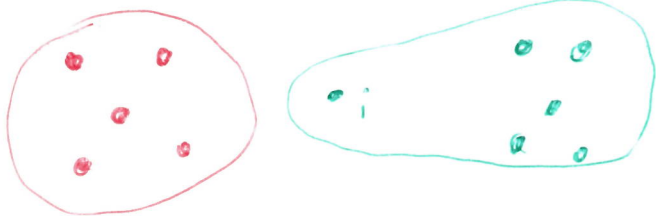
Examples



$\alpha(i)$ is large

$b(i)$ is small (close to zero)

$$s(i) = \frac{b(i) - \alpha(i)}{\max(\alpha(i), b(i))} \approx \frac{-\alpha(i)}{\alpha(i)} = -1$$



Here $\alpha(i) \approx b(i)$

$$s(i) \approx 0$$



$\alpha(i)$ is small

$b(i)$ is large

$$s(i) = \frac{b(i) - \alpha(i)}{\max(\alpha(i), b(i))} \approx \frac{b(i)}{b(i)} = 1$$

⑥

Clustering combined with PCA

We can use PCA as a pre-processing step to reduce dimensionality of the data and then ~~per~~ perform clustering to the PC scores.

That is instead of using the data matrix X we use the first columns of $Z = XA$ of PCA

Advantages: reduce dimensionality and thus computational time before clustering. Better visualisation.

Disadvantages: interpretation of PC scores may not be the same as the one for the raw data. Moreover, there is not always guarantee that we will obtain a better clustering like that.

Supervised Learning

classification

regression

In supervised learning we have the variable X (predictor features) but also a response variable Y
target ~~output~~ output

Essentially for ~~each~~ every $x = (x_1, \dots, x_p)$ there is a value y associated to it

$x \mapsto y$

The target is to "learn" this map given some data

Data looks like this:

$\left[\begin{array}{c|c} X & Y \end{array} \right]$
input variables $n \times p$ matrix output variable matrix of size $n \times 1$

$n \times (p+1)$

In classification, Y takes values from a discrete set (labels). For our purposes, we will consider binary classification: 0 or 1

In regression, Y takes values from a potential infinite set (real values)

On the other hand, the variables X (especially) are real numbers.

Aim of supervised learning

Give some data (set of input-output variables)

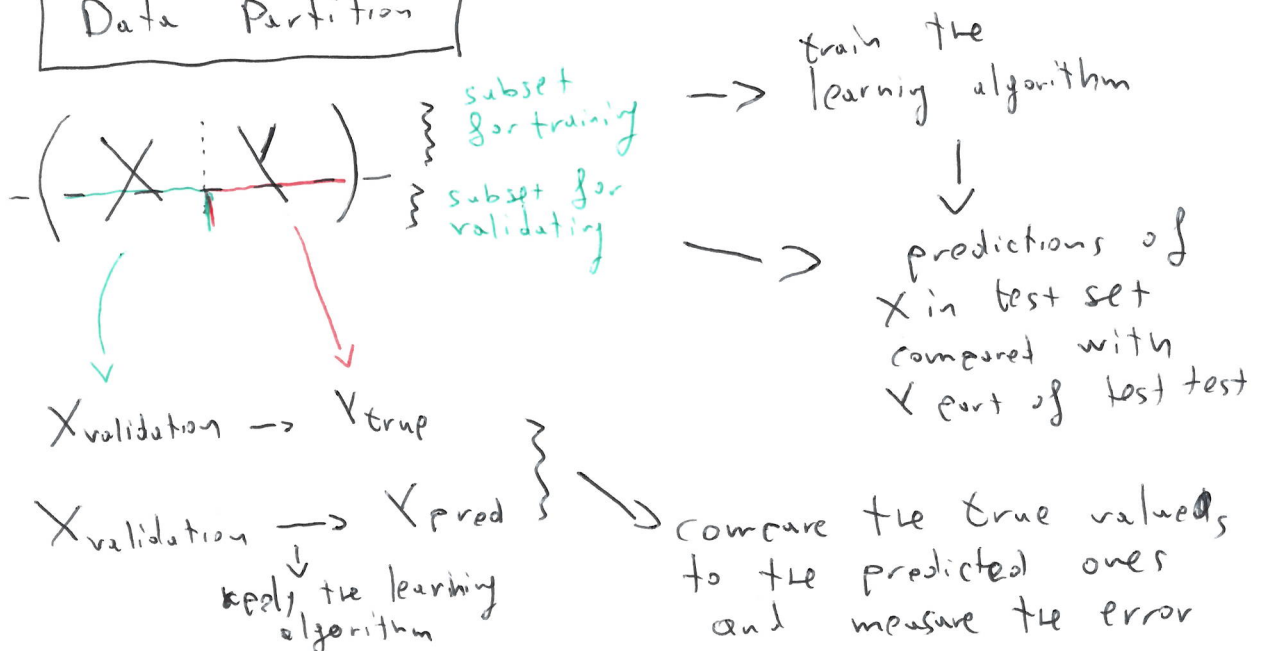
the aim is to create a model/algorithm such that it can predict correctly the output values for new values of X

Usually, we determine the predictive accuracy of an algorithm by measuring its accuracy for a data set that has not been used.

In practice:

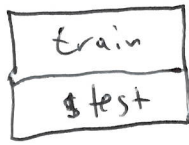
- 1) We **train (build)** the algorithm/predictor using some part of the data (**training set**)
- 2) We **validate (evaluate)** the algorithm on some different part of the data (**validation/test set**)

Data Partition



There are three main strategies for splitting the data

1) Divide the data into a training set and a test set



Partitions could be 80/20, 70/30, 50/50
 ↙ ↘
 train test

2) K-fold cross validation.

Split the data into k parts (k folds), and use all folds but one to train and the remaining to test. Exhaust for all combinations.

Example 5-fold CV

Data
1
2
3
4
5

Train
1, 2, 3, 4
1, 2, 3, 5
1, 3, 4, 5
1, 3, 4, 5
2, 3, 4, 5

Test
5
4
3
2
1

Error
error 1
error 2
error 3
error 4
error 5

↓
 average ~~error~~ over all errors.

3) n-fold cross validation.

Also known as leave-one-out

In each of n instances, one observation is left out of the training set and it is used to evaluate the predicts every time.

The measure of accuracy is usually the average error over all instances.

9 The partition of data is usually done at random
but to make results reproducible, we set
a value for the random seed
R command: `set.seed()`

Machine learning vs Statistics

In machine learning, the focus is in
comparing models with respect to their
prediction capabilities.

This is contrast to statistics, where models
are fitted to data and then significance tests
for the terms of the models (predictive
capability is not the main goal)

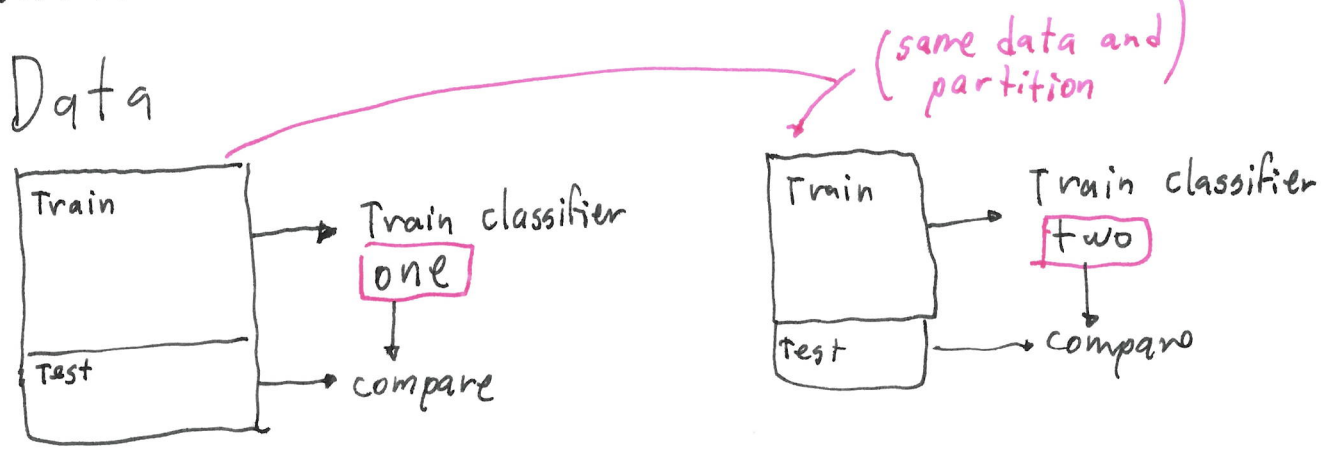
4-Classification

In classification, outputs take values from a discrete set. This is a set of labels and we are interested in predicting labels. We will concentrate on the simplest case with only two labels.

Labels <	0	"-"	healthy individual	normal transaction	person has not defaulted
	1	"+"	ill individual	Fraudulent transaction	person has defaulted
			detecting disease	detecting fraud in debit card payment	assessing default on loans

4.1 Performance of methods

Using available data, the machine learning task of classification involves comparison of classifiers.



Example with two classifiers.

As an example of comparisons, consider the following data. In automatic fraud detection

Y_{true}	$Y_{Model 1}$	
1	0	"Algorithm 'Model 1' concluded no fraud yet it was." False negative
0	1	"Algorithm flagged a fraud yet there was none" False positive
1	1	True positive
1	1	Misclassification

Define the following concepts

True positive (TP) is a positive classified as a positive.

False positive (FP) is a negative classified as a positive.

False negative (FN) is a positive classified as a negative.

True negative (TN) is a negative classified as a negative.

We use the acronyms to denote counts of data
 TP, FP, FN, TN.

$P = TP + FN$ ← number of positives
 $N = TN + FP$ ← number of negatives

(1500) Confusion matrix - a table that summarizes the different cases (correct classifications, misclassifications)

Preferred

		Predicted class		N
		0	1	
True class	0	TN	FP	P
	1	FN	TP	

A different way of showing

		Predicted		P
		1	0	
True	1	TP	FN	N
	0	FP	TN	

same table same information

Example with $P=3, N=2$
 ($P+N=5$ observations)

		Pred.		2
		0	1	
True	0	0	2	3
	1	1	2	

$TP=2$ ← correct
 $FP=2$
 $FN=1$ } Misclassified
 $TN=0$ ← correct

		Pred		3
		1	0	
True	1	2	1	2
	0	2	0	

$TPR = \frac{2}{3} = 0.66$

$FPR = \frac{2}{2} = 1$

Accuracy = $\frac{2}{5} = 0.4$

Error rate = $\frac{3}{5} = 0.6$

Some measures of classifiers' performance: * Important 3

* True positive rate (TPR) $TPR = TP/P$ (aka Hit rate, recall or sensitivity)
proportion of positives correctly classified as positives.

False positive rate (FPR) $FPR = FP/N$ (equal to $1 - \text{specificity}$)
proportion of negatives wrongly classified as positives

False negative rate (FNR) $FNR = \frac{FN}{P} = 1 - TPR$
proportion of positives wrongly classified as negatives.

True negative rate (TNR) $TNR = \frac{TN}{N} = 1 - FPR$ (aka specificity)
proportion of true negatives correctly classified as negatives

Precision = $\frac{TP}{TP+FP}$
proportion of cases classified as positive that are really positive (aka positive predictive value)

Accuracy = $\frac{TP+TN}{P+N}$
proportion of instances correctly classified.

Error rate = $\frac{FP+FN}{P+N}$
proportion of instances classified incorrectly

F1 score = $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
Combination of precision and recall

Example. We have $P+N=8$ observations available for testing two models (classifiers) that have been trained already.

Y_{true}	Y_1	Y_2
1	1	1
0	1	1
1	1	0
0	0	1
0	0	1
1	0	0
0	0	1
0	1	1

M1
P. 1

	0	1	
T. 0	3	2	5
1	1	2	3

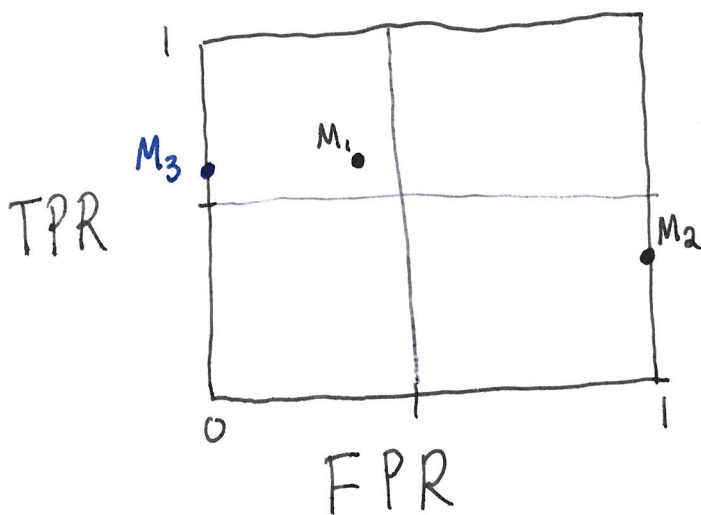
TPR = 0.67
FPR = 0.4

M2
P. 1

	0	1	
T. 0	0	5	5
1	2	1	3

TPR = 0.33
FPR = 1

We compare classifiers by plotting TPR vs. FPR. This is known as the Receiver Operating Characteristic (ROC) graph.



Model	TPR	FPR
1	0.67	0.4
2	0.33	1

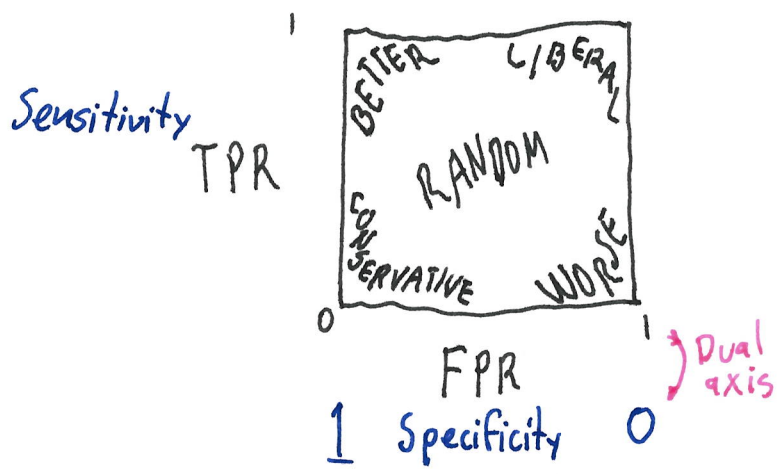
"Model 1 is a better classifier than Model 2"

140324 Model 3 reverses predictions Y_2

Model	TPR	FPR
3	0.67	0

140324

(1400) The ROC graph has regions with clear interpretation. The adjectives apply to classifiers.



Note that flipping predicted values $0 \rightarrow 1$ will improve a bad classifier (worsen a good classifier)

Confusion matrices for some of these extreme cases

	Ideal	Worst	Liberal	Conservative	Random
0	$\begin{matrix} 0 & 1 \\ N & O \end{matrix}$	$\begin{matrix} 0 & 1 \\ O & N \end{matrix}$	$\begin{matrix} 0 & 1 \\ O & N \end{matrix}$	$\begin{matrix} 0 & 1 \\ N & O \end{matrix}$	$\begin{matrix} 0 & 1 \\ N/2 & N/2 \end{matrix}$
1	$\begin{matrix} O & P \end{matrix}$	$\begin{matrix} P & O \end{matrix}$	$\begin{matrix} O & P \end{matrix}$	$\begin{matrix} P & O \end{matrix}$	$\begin{matrix} P/2 & P/2 \end{matrix}$
TPR =	1	0	1	0	1/2
FPR =	0	1	1	0	1/2

4.2 The ROC curve and AUC

Classifiers are points in ROC graph and we compare them according to the position in the graph

- closeness to top left corner $FPR=0, TPR=1$
- relative positions and importance of 'positive/negative'.

There are cases where we can tune the classifier through some parameter. Then the plot of points in the ROC graph depends on such parameter and we have the ROC curve.

Example (logistic model) with linear predictor $\beta_0 + \beta_1 x$

Training data

x	y
1	0
-7	0
-3	0
5	1
13	1
9	0

and we have $\hat{\beta}_0 = -2.0948$
 $\hat{\beta}_1 = 0.2862$

so we can predict $\Pr(Y=1|x) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x}} = p(x)$

Test data "set a threshold and produce labels"

x	y	p(x)	0	0.05	0.2	0.3	1
3	0	0.2251	1 ← FP	1	1	0	0
15	1	0.9	1 ✓	1	1	1	0
-1	1	0.0846	1 ✓	1	0	0	0
-5	0	0.0286	1 ← FP	0	0	0	0
11	1	0.7414	1 ✓	1	1	1	0
7	1	0.4771	1 ✓	1	1	1	0

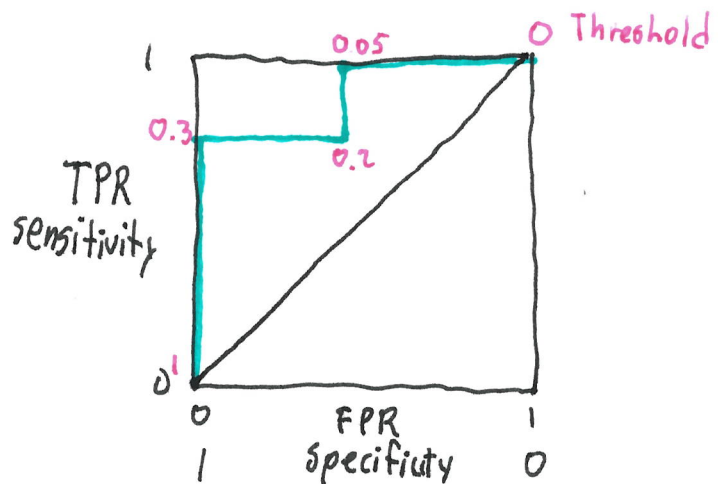
0	1
1	2
1	4

0	1
1	1
1	3

0	1
1	2
1	4

Threshold	TPR	FPR
0	1	1
0.05	1	0.5
0.2	0.75	0.5
0.3	0.75	0
1	0	0

Using thresholds and values TPR, FPR, we build the ROC curve.



AUC = 0.875

190324 (1400) The area under the (ROC) curve AUC is a useful quantity to compare the performance of classifiers. The closer AUC gets to one, the better the performance. AUC is also known as the c-statistic.

Models for classification
We'll do a survey of models for this task. The emphasis is on the comparison between classifiers.

4.3 Linear classifier

The simplest model for the vector of labels \underline{y} treats it as if it was continuous. This is a regression with expectation

$$E(Y) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

intercept terms for the p explanatory variables

Function lm in R

Using fitted coefficients $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ and test data, we obtain predictions \hat{y} . A simple rule is to allocate label "1" if $\hat{y} > 0.5$ and label "0" otherwise.

Advantages - Simple and intuitive
- Well known and understood

Disadvantages - Estimates \hat{y} may lie outside $[0, 1]$ and be far away from labels "0", "1" (*)
- Difficult ~~also~~ to use with more than two labels.

4.4 Logistic classifier

logistic regression is an established methodology which is part of generalized linear models (GLM). For example techniques for $\left\{ \begin{array}{l} \text{standard linear regression} \\ \text{count data regression} \\ \text{binary data regression} \end{array} \right\}$ are all GLMs.

The logistic classifier avoids the inconvenience(*) modelling the probability $p(\underline{x}) = \Pr(Y=1|\underline{x})$ in such a way that $p(\underline{x}) \in [0, 1]$.

We use a distribution (Bernoulli) to model data and a special parameterization (logistic link) to model the probabilities. The Bernoulli probabilities $p_i = P(x_i)$ obey

$$p_i = P(x_i) = \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}$$

← data is $Y_i \sim \text{Ber}(p_i)$

Inverting the logistic puts the log-odds as a linear function of covariates

multiple logistic regression

$$\log \frac{p(x_i)}{1 - p(x_i)} = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

To train the logistic model, we use maximum likelihood estimation (outside the scope of MTH6101).

In R use function `glm` with option `family = "binomial"`

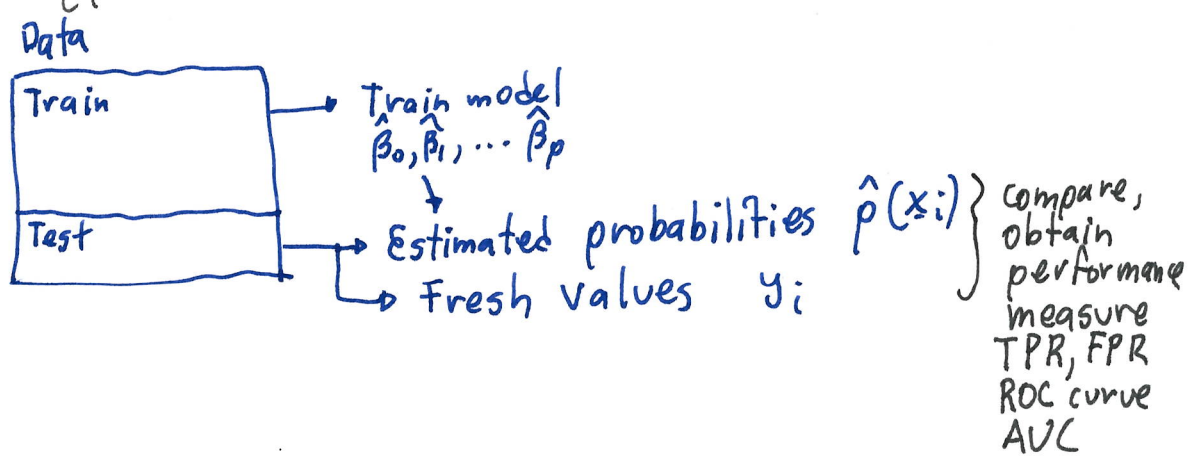
(By default, the link to "binomial" is logistic)

After training, we have estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ and compute probabilities with fresh data

$$\hat{p}(x_i) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}$$

↑ fresh data point

Together with (fresh) response values "true observations" we use these probabilities to describe the performance of the classifier (and compare).



Example - Default data, logistic classifier. We have three classifiers for response 'default'!

1: 'default' $\rightarrow \beta_0 + \beta_1$ 'balance'

2: 'default' $\rightarrow \beta_0 + \beta_1$ 'income'

3: 'default' $\rightarrow \beta_0 + \beta_1$ 'student'

Data has $n=10,000$ observations, split 80/20 (5 folds)
4 folds to train, 1 fold to validate/test
8000 obs. 2000 obs.

Classifier	AUC
1	0.961 ← Good classifier
2	0.5232
3	0.4895

} Like random classifiers

4.5 K-nearest neighbors classifier (KNN classifier)

A new observation is classified (labeled) according to the labels of the K nearest neighbors to it.

Labels are unknown
this is part of testing data

Labels are known
this is part of training data

R function
Knn

For this classifier, the training and validation steps are performed in a single run; there are no parameter estimates produced (no ROC curve).

The value of K chosen is usually odd so that there will be no ties among nearest neighbors. In the case of K even (or more than two labels), the label is allocated at random.

Example KNN, default data, same preparation 80/20 as earlier. We consider $k=1, 2, 3, 4, 5$.
 Analysis performed with standardized variables.

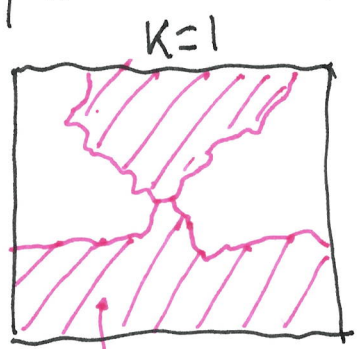
210324
(1400)

k	1	2	3	4	5
TPR	0.3	0.316	0.416	0.3	0.383
FPR	0.024	0.019	0.009	0.009	0.004

← as high as possible
 ← as low as possible

↓ suggested classifier

The KNN classifier is simple to explain and can produce complex regions for classification. It requires the covariates to be continuous, though.



← see page 83 notes



classification done at random

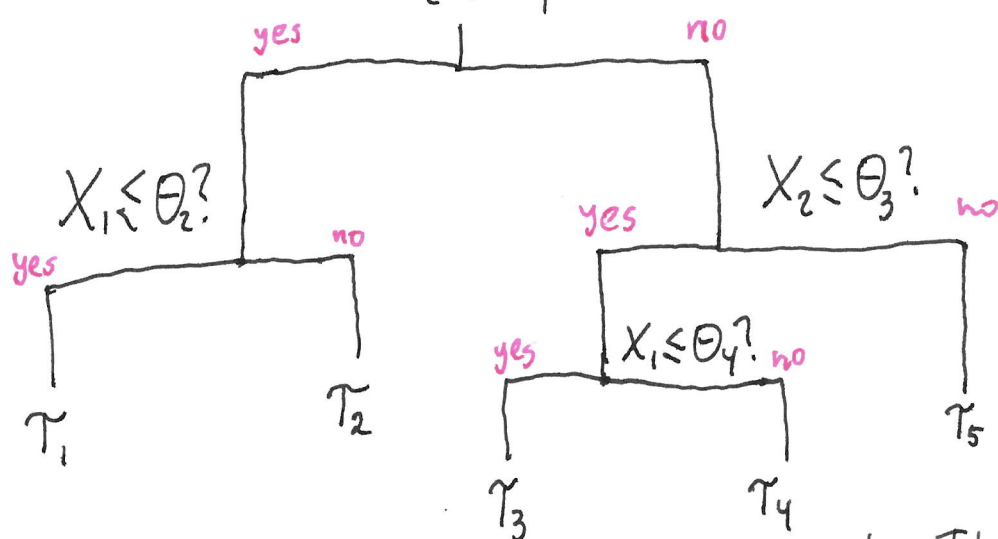
4.6 Classification tree

A classification tree is a structure in which 'leaves' represent class labels and 'branches' represent conjunctions of variables that lead to those class labels.

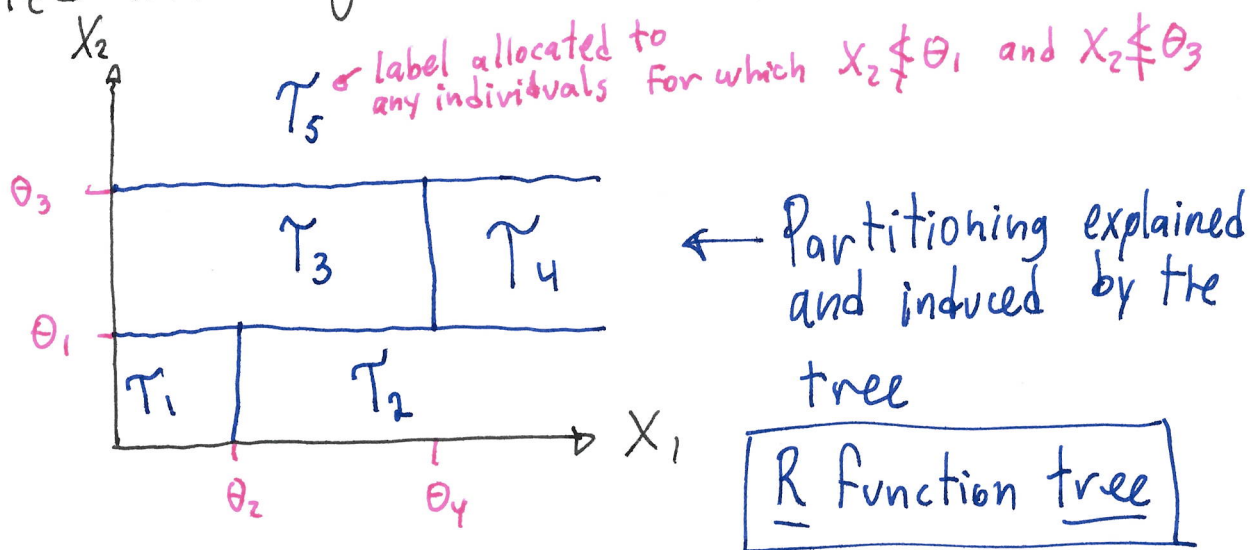
Example of tree with variables X_1, X_2

$$X_2 \leq \theta_1?$$

Here τ_1, τ_2, \dots are terminal nodes "leaves" i.e. labels



For us, leaves will be our 0/1 labels. They are associated with regions in the space of covariates.



Advantages

- Simple to understand and interpret
- Able to handle categorical and numerical data
- Simple to analyze, requiring basic data preparation
- Can be validated with data
- No distributional assumptions, simple in terms of its statistics
- Mirrors human decision-making more than other approaches

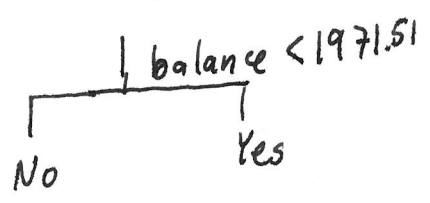
Limitations

- The problem of training a tree is known to be hard (computationally costly). Algorithms use greedy search and there is no guarantee to reach the optimal tree for a ~~given~~ given data set.
- The method of fitting trees to data can be very non-robust. A small change in data may have huge implications in the resulting tree and predictions.
- It may overfit the data and be overly complicated.

260324
(1400)

Example 'Default' data set. Same partition of data as earlier, For validation $P=60, N=1940$.

Tree (after pruning and removing redundancies)



Confusion matrix

	Pr. 0	Pr. 1
Tr. 0	1934	6
Tr. 1	43	17

FPR = $\frac{6}{1940} = 0.003$

TPR = $\frac{17}{60} = 0.283$

4.7 Linear discriminant analysis (LDA)

We use a linear combination of variables to separate two (or more) classes of objects.

LDA assumes $\left[\begin{matrix} p(x|y=0) \\ p(x|y=1) \end{matrix} \right]$ are both normally distributed

conditional distⁿs of x given y

$(\underline{\mu}_0, \underline{\Sigma}_0)$
vector of means

$(\underline{\mu}_1, \underline{\Sigma}_1)$
variance covariance matrix

The Bayes optimal solution predicts points as being from the second class ($y=1$, positive) if the log-likelihood ratios are bigger than a threshold T

$$\underbrace{(\underline{x} - \underline{\mu}_0)^T \underline{\Sigma}_0^{-1} (\underline{x} - \underline{\mu}_0) + \log |\underline{\Sigma}_0|}_{\text{log of multivariate normal density } (\underline{\mu}_0, \underline{\Sigma}_0)} - \underbrace{(\underline{x} - \underline{\mu}_1)^T \underline{\Sigma}_1^{-1} (\underline{x} - \underline{\mu}_1) - \log |\underline{\Sigma}_1|}_{\text{log of multivariate normal } (\underline{\mu}_1, \underline{\Sigma}_1)} > T$$

If matrices $\underline{\Sigma}_0, \underline{\Sigma}_1$ are different, this is quadratic discriminant analysis.

In LDA, we assume that the covariance matrices are identical $\underline{\Sigma} = \underline{\Sigma}_0 = \underline{\Sigma}_1$, and that $\underline{\Sigma}$ is full rank. After cancelling terms, the decision criterion above is a threshold on the dot product

R function lda

(1500)

$$\underline{w} \cdot \underline{x} > c$$

with $\underline{w} = \underline{\Sigma}^{-1} (\underline{\mu}_1 - \underline{\mu}_0)$ and $c = \frac{1}{2} (T - \underline{\mu}_0^T \underline{\Sigma}^{-1} \underline{\mu}_0 + \underline{\mu}_1^T \underline{\Sigma}^{-1} \underline{\mu}_1)$.

In other words, observation \underline{x} belongs to a certain class depending on its location regarding a hyperplane perpendicular to \underline{w} . The location of the hyperplane is defined by the constant c .

Example 'Default' data, same partition 4:1 as before (80/20)
P=60 and N=1940

Confusion matrix

		Pr.	
		0	1
Tr.	0	1937	3
	1	44	16

FPR = $\frac{3}{1940} = 0.0015$
 TPR = $\frac{16}{60} = 0.266$

AUC = 0.9638

Other topics in classification: quadratic discriminant, neural networks' classifier, Bayesian "naive" classifier; more than two labels...

5- Lasso and regularization (penalized likelihood techniques)

The aim of penalization techniques is to improve over the performance of traditional estimators (regression).

Brief review of regression.

The model is $\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$; usually \underline{Y} has been centered around its mean and the columns of \underline{X} been centered as well (possibly scaled), so that the model has no intercept.

((In what follows, we will use the singular value decomposition $\underline{X} = \underline{U}\underline{D}\underline{V}^T$.)

To estimate parameters, we consider the usual quadratic criterion

$$SS_{\varepsilon} = \left\| \underset{\substack{\uparrow \\ \text{observed}}}{\underline{Y}} - \underset{\substack{\uparrow \\ \text{predicted}}}{\underline{X}\underline{\beta}} \right\|_2^2.$$

The standard estimate $\hat{\underline{\beta}}$ is the least squares estimate

$$\begin{aligned} \hat{\underline{\beta}} &= (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y} \\ &= \left(\underbrace{\underline{V}\underline{D}\underline{U}^T}_{\underline{X}^T} \underbrace{\underline{U}\underline{D}\underline{V}^T}_{\underline{X}} \right)^{-1} \underbrace{\underline{V}\underline{D}\underline{U}^T}_{\underline{X}^T} \underline{Y} = \underbrace{\underline{V}\underline{D}^{-2}\underline{V}^T}_{\underline{X}^T \underline{X}} \underbrace{\underline{U}\underline{D}\underline{U}^T}_{\underline{Y}} \\ &= \underline{V}\underline{D}^{-1}\underline{U}^T \underline{Y}. \end{aligned}$$

Consider the predicted values $\hat{Y} = X\hat{\beta}$ that we now develop

$$\hat{Y} = X\hat{\beta} = \underbrace{U D U^T}_{\hat{X}} \underbrace{V D^{-1} U^T}_{\hat{\beta}} Y = \underbrace{U U^T}_{*} Y$$

$$H = X(X^T X)^{-1} X^T = U U^T$$

$$H^2 = U U^T U U^T = U U^T$$

columns of U

The predictions \hat{Y} can be written as $\hat{Y} = \sum_{j=1}^p u_j u_j^T Y$ and we note that $* U^T Y$ are the coordinates of Y with respect to orthonormal basis U .

290324
(1400)

Regularization

- β functions
- ridge (own)
- lars
- glmnet

$$\|Y - X\beta\|_2^2 \leftarrow \text{OLS}$$

$$\|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \leftarrow \text{Ridge (L}_2 \text{ penalization)}$$

$$\frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \leftarrow \text{Lasso (L}_1 \text{ penalization)}$$

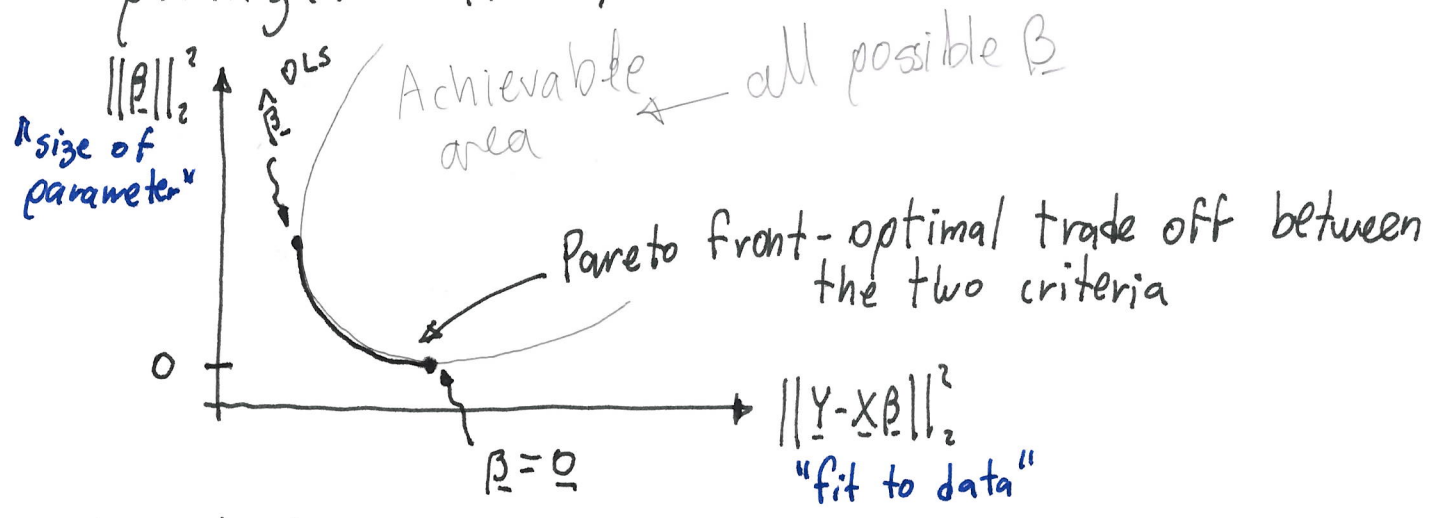
$$\|Y - X\beta\|_2^2 + \lambda [\alpha \|\beta\|_1 + (1-\alpha) \|\beta\|_2^2] \leftarrow \text{Elastic net (Mixture of L}_1, \text{L}_2)$$

5.1 Ridge regression

We want to find β that achieves a small value of $\|Y - X\beta\|_2^2$ in such a way that β

is not too large. We measure the size of β using square euclidean norm $\|\beta\|_2^2$

We have a dual objective problem in optimization theory.



This dual problem can be scalarized to have a single objective function. After simplification, the objective function is

$$R = \underbrace{\|Y - X\beta\|_2^2}_{SS_E} + \lambda \underbrace{\|\beta\|_2^2}_{\substack{\text{parameter} \\ \text{sq. Eucl.} \\ \text{norm of} \\ \sum_{i=1}^p \beta_i^2}}$$

to be minimized for (fixed) β positive λ . This is a regularized (L_2 regularized) least squares problem.

Minimizing R .

$$\begin{aligned} R &= (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta \\ &= Y^T Y - 2\beta^T X^T Y + \underbrace{\beta^T X^T X \beta + \lambda \beta^T \beta}_{\beta^T (X^T X + \lambda I) \beta} \end{aligned}$$

$$\frac{\partial R}{\partial \beta} = -2X^T Y + 2X^T X \beta + 2\lambda I \beta$$

The estimate that minimizes R is the following closed form:

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

020424
(1400)

$$= \underbrace{V(D^2 + \lambda I)^{-1} D}_{\text{Diagonal matrix with entries } \frac{d_1}{d_1^2 + \lambda}, \dots, \frac{d_p}{d_p^2 + \lambda}} U^T Y \quad \leftarrow \text{using s.v.d. of } X = U D V^T$$

For $\lambda \rightarrow 0$, $\hat{\beta}^R \rightarrow \hat{\beta}$ ordinary least squares

$\lambda \rightarrow \infty$, $\hat{\beta}^R \rightarrow 0$ in the limit, ridge estimate shrinks to zero

When we compute predictions, we have

$$\hat{Y}^R = X \hat{\beta}^R = \underbrace{U D (D^2 + \lambda I)^{-1} D}_{\text{diagonal}} U^T Y = \sum_{j=1}^p \underbrace{u_j}_{\text{columns of } U} \underbrace{\frac{d_j^2}{d_j^2 + \lambda}}_{\text{Factor of shrinkage of ridge regression compared with regression predictions}} u_j^T Y$$

Advantages of the ridge estimator $\hat{\beta}^R$

As we increase λ , moving along the ridge, shrinkage of coefficients leads to a reduction in the variance of predictions, at the expense of increasing bias.

Ridge is simpler than subset selection and may predict better yet it is not a variable selection technique.

A key problem in ridge regression is the selection of λ for the final model. This can be done by cross-validation.

The ridge trace is the plot of $\hat{\beta}^R = \hat{\beta}^R(\lambda)$ against λ . This is done for a range of values of λ , and each entry of $\hat{\beta}^R$ is plotted. Alternatively, the horizontal axis can be percentage of shrinkage

$$\frac{\|\hat{\beta}^R(\lambda)\|_2^2}{\|\hat{\beta}\|_2^2}$$

sq. Eucl. norm of ridge estimate $\hat{\beta}^R$ / P.S.E.v. of OLS. $\hat{\beta}$

and another way of plotting uses the "effective degrees of freedom" $\text{tr}(H_\lambda)$ in the horizontal axis.

We have $\text{tr}(H_\lambda) = \text{tr}(\underline{U} \underline{D} (\underline{D}^2 + \lambda \underline{I})^{-1} \underline{D} \underline{U}^T)$ ← $\text{tr}(AB) = \text{tr}(BA)$
"circular property of the trace"

(1500) $= \text{tr}(\underline{U}^T \underline{U} \underline{D} (\underline{D}^2 + \lambda \underline{I})^{-1} \underline{D})$

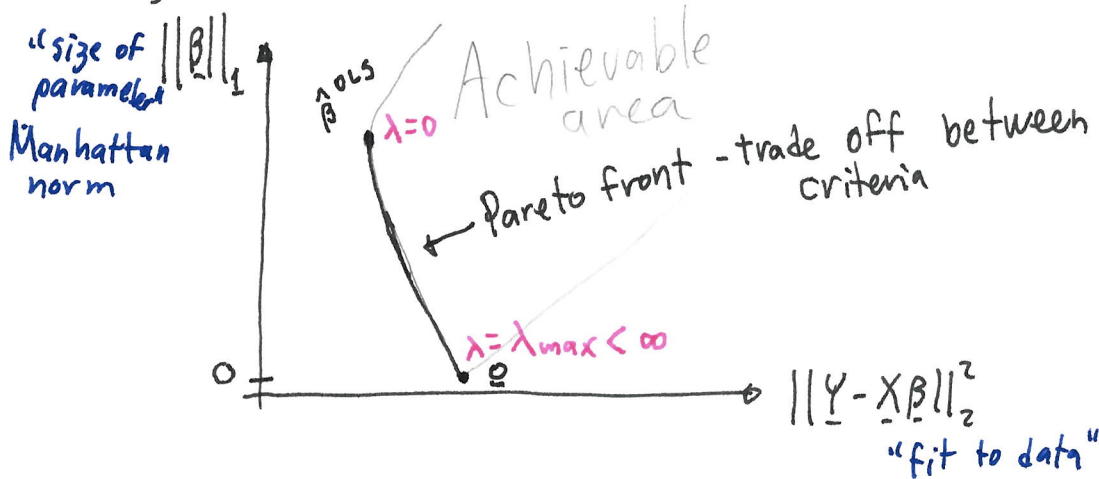
$= \text{tr}(\underline{D} (\underline{D}^2 + \lambda \underline{I})^{-1} \underline{D}) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$

$0 < \text{tr}(H_\lambda) \leq p$

5.2 Lasso

A drawback of $\hat{\beta}^R$ is ~~that~~ that it generally involves all variables in the sense that its entries are not zero for finite λ . Lasso (least absolute shrinkage and selection operator) overcomes this disadvantage using L_1 penalization.

Using a similar dual-objective diagram...



The objective function for the lasso is an L_1 penalization of the sum of squares of the error

$$L = \underbrace{\frac{1}{2} ||Y - X\beta||_2^2}_{\substack{\text{sum of squares} \\ \text{of the error}}} + \underbrace{\lambda}_{\substack{\text{parameter} \\ \lambda \geq 0}} \underbrace{||\beta||_1}_{\substack{\text{Manhattan} \\ \text{norm}}} = \sum_{i=1}^p |\beta_i|$$

need for technical reason
parameter $\lambda \geq 0$
Manhattan norm

Example with data (one explanatory variable, $n=7$)

$$\begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \end{pmatrix} \leftarrow \text{before centering}$$

X Y

$$\frac{1}{7} \begin{pmatrix} -4 & -1 \\ -4 & -1 \\ -4 & -1 \\ 3 & -1 \\ 3 & -1 \\ 3 & 6 \\ 3 & -1 \end{pmatrix} \leftarrow \text{after centering}$$

$$X^T X = \frac{12}{7}$$

$$X^T Y = \frac{3}{7}$$

$$Y^T Y = \frac{6}{7}$$

OLS $\hat{\beta} = (X^T X)^{-1} X^T Y = \frac{7}{12} \cdot \frac{3}{7} = \frac{3}{12} = \frac{1}{4} = 0.25$

Ridge $\hat{\beta}^R = (X^T X + \lambda I)^{-1} X^T Y = \left(\frac{12}{7} + \lambda\right)^{-1} \cdot \frac{3}{7} = \frac{3}{12} \left(\frac{12}{12 + 7\lambda}\right)$

Lasso $\hat{\beta}^L = \frac{3}{12} \left(1 - \frac{7}{3} \lambda\right)$
when $\lambda = \frac{3}{7}$ we have $\hat{\beta}^L = 0$

Let us compute the lasso estimate

$L = \frac{1}{2} Y^T Y - \beta^T X^T Y + \frac{1}{2} \beta^T X^T X \beta + \lambda \|\beta\|_1$ which for this single parameter case becomes

$$L = \frac{6}{14} - \frac{3}{7}\beta + \frac{12}{14}\beta^2 + \lambda|\beta|$$

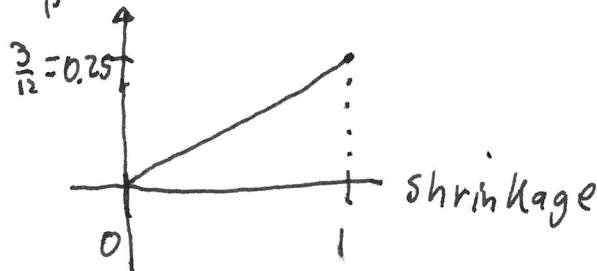
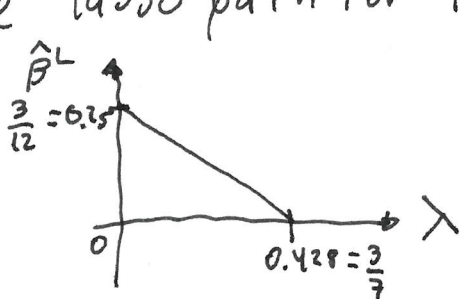
← β equals $|\beta|$ in the case $\beta > 0$. We start here as $\hat{\beta} = 0.25 > 0$.

$$\frac{\partial L}{\partial \beta} = -\frac{3}{7} + 2 \cdot \frac{12}{14} \beta + \lambda$$

$$= -\frac{3}{7} + \frac{12}{7} \beta + \lambda$$

(equate to zero and solve for β)

The lasso path for the small example



040224
(1400)

Contrary to ridge, there is no general closed formula for the lasso estimate $\hat{\beta}^L = \hat{\beta}^L(\lambda)$ that minimizes criterion L .

Lasso path is the collection of $\hat{\beta}^L$ as λ takes values in $[0, \infty)$. However $\hat{\beta}^L$ shrinks to zero for a finite λ . It consists of piecewise linear trajectories, which change direction at breakpoints.

Lasso has a dual function as estimation procedure and as variable selection method. Models from lasso are simpler to ~~it~~ interpret than those from ridge and are sparser (have less non-zero coefficients).

Example (notes) Credit card data. Six explanatory

R function lars

Variables, data split 75/25.

As is common practice, data is centered (and scaled) so no intercept in the analysis.

\$beta from lars

shrinkage

λ	Income	Limit	Rating	Cards	Age	Education	$\ \hat{\beta}^\lambda\ _1 / \max \ \hat{\beta}\ _1$
259.79	-	-	-	-	-	-	0 / 2.01 = 0
85.35	-	-	0.59	-	-	-	0.59 / 2.01 = 0.29
37.40	-	0.11	0.64	-	-	-	0.75 / 2.01 = 0.37
14.09	-0.37	0.35	0.76	-	-	-	1.5 / 2.01 = 0.74
7.86	-0.47	0.42	0.80	-	-0.02	-	1.71 / 2.01 = 0.85
3.12	-0.54	0.58	0.71	0.02	-0.03	-	1.98 / 2.01 = 0.93
0	-0.59	0.70	0.63	0.03	-0.04	-0.01	2.01 / 2.01 = 1

\$lambda

Rows correspond to breakpoints, indexed by either λ or by shrinkage

Each column has estimates at every breakpoint

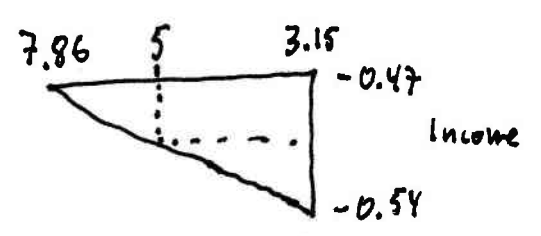
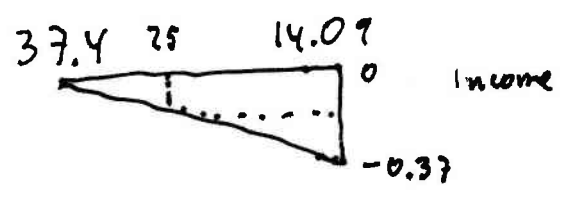
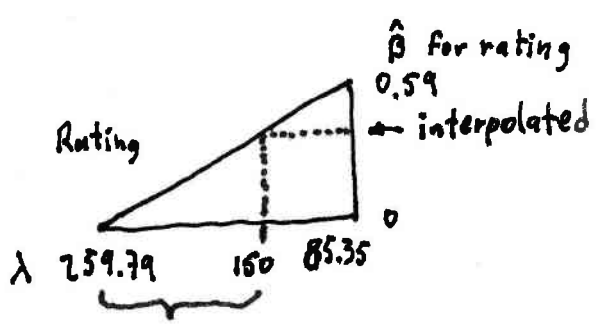
Top row \leftarrow all coefficients shrink to zero

Bottom row \leftarrow O.L.S.

090424 (1400)

Let us interpolate and compute $\hat{\beta}^\lambda(150)$

$$\hat{\beta}^\lambda(150) = (0, 0, 0.3711, 0, 0, 0)^T$$



$$\frac{|150 - 259.79|}{|85.35 - 259.79|} = \frac{109.79}{174.439} = 0.629$$

Elastic net

An extension of the L_1 penalization of lasso and the squared L_2 penalization of ridge regression are elastic nets. The penalization is a combination of L_1 and squared L_2 , controlled with an additional parameter α :

$$E = \underbrace{\|Y - X\beta\|_2^2}_{SS_E} + \lambda \left[\underbrace{\alpha}_{\text{parameter}} \underbrace{\|\beta\|_1}_{L_1} + \underbrace{(1-\alpha)}_{\text{controls the mixture}} \underbrace{\|\beta\|_2^2}_{L_2} \right]$$

R function
glmnet

$\alpha=1 \Rightarrow E=L$ lasso

$\alpha=0 \Rightarrow E=R$ ridge

Exercise. For the small example of 020424, show that minimizing E leads to the path

$$\hat{\beta}^E(\lambda, \alpha) = \frac{6/7 - \lambda \alpha}{24/7 + \lambda(1-\alpha)}$$

The elastic net path is not anymore a series of piecewise linear trajectories as with the lasso. The trajectories are continuous, ridge-like but with the benefit that some coefficients shrink to zero, and that the shrinkage of all is achieved for finite λ . (Although the glmnet does not show) the elastic net path has a series of breakpoints similar to lasso, where the path changes direction.

5.3 Penalized likelihood

So far we have considered the standard

least squares criterion $\|Y - X\beta\|_2^2$ on top of which we added a penalty based on the size of the parameter vector

minimization

L_1 lasso
 L_2 ridge
 L_1, L_2 elastic net.

We could start from a different point, that is, ~~is~~ to consider likelihood as the basis for estimation, and on top of the likelihood, add a penalty of the type elastic net.

If the likelihood is denoted as $l(\beta)$, log-likelihood
parameter vector

the penalized version is

maximized

$$l_p(\beta) = \underbrace{l(\beta)}_{\text{likelihood}} + \lambda \underbrace{g(\beta, \alpha)}_{\text{penalty by the size of } \beta}$$

R Function
glmnet

Where we use $g(\beta, \alpha) = \alpha \|\beta\|_1 + (1-\alpha) \|\beta\|_2^2$ Elastic net penalty.

Example. (as in the lab) Classification of glass, of 7 types of glass we are interested in type=7. There are 9 variables available. We use glmnet with "binomial" distribution and $\alpha = 1/10$ (close to ridge penalization). K=3 folds are used.