

## Topics:

### 1) Data reduction

→ reduce # of variables (dimensions of data) while preserving information

- Principle Component Analysis
- Multidimensional scaling

### 2) Cluster analysis

→ unsupervised - uses the structure of the data to create labels

- Hierarchical clustering
- K-means

### 3) Classification (Discriminant analysis)

→ supervised - uses existing labels to assign labels to new data

- K-nearest neighbours
- Linear and Quadratic discriminant analysis
- Logistic regression

## Multivariate data

→  $m$ -dimensional data =  $m$  attributes (variables) for every object (data point)

↳ data point:  $\tilde{x} = (x_1, \dots, x_m)^T$

→ attributes

- numerical continuous  $x_i \in \mathbb{R}$
- numerical discrete  $x_i \in \mathbb{N}$
- categorical  $x_i = \text{gender}$

Def: Multivariate data sets are represented by a matrix  $X \in \mathbb{R}^{n \times m}$

- $n = \#$  observations
- $m = \#$  variables

$$X = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} \quad \text{where} \quad \tilde{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{im} \end{pmatrix} = i^{\text{th}} \text{ observation}$$

Def: A random variable is a mapping  $X: \Omega \rightarrow \mathbb{R}$  which assigns real numbers to possible outcomes.

Def: For a random variable  $X$  we define its

- ① probability mass function  $p_X(x) := P(X=x)$
- ② cumulative distribution function  $F_X(x) := P(X \leq x)$

Def: The random variable  $X$  is said to be

- ① discrete  $\equiv \text{Im}(X)$  is countable
- ② continuous  $\equiv \exists f_X: \mathbb{R} \rightarrow \mathbb{R}$  s.t.  $F_X(x) = \int_{-\infty}^x f_X(t) dt$   
 $\hookrightarrow f_X$  is called the probability density function of  $X$

Def: The expected value of a r.v.  $X$  is

- ①  $X$  discrete ...  $\mu = E[X] := \sum_{x \in \text{Im}(X)} x \cdot p_X(x)$
- ②  $X$  continuous ...  $\mu = E[X] := \int_{-\infty}^{\infty} x f_X(x) dx$

Theorem (PNS): The expected value of a function  $g(X)$  is given by

- ①  $E[g(X)] = \sum_x g(x) p_X(x)$
- ②  $E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$

Def: Consider random variables  $X_1, X_2, \dots, X_m$ . Define

① variance  $\text{Var}[X_i] := E[(X_i - E[X_i])^2] = E[X_i^2] - (E[X_i])^2 = \sigma^2$  standard deviation  
↓

② covariance  $\text{Cov}[X_i, X_j] := E[(X_i - E[X_i])(X_j - E[X_j])] = E[X_i X_j] - E[X_i] E[X_j]$

③ correlation  $\text{Cor}[X_i, X_j] := \frac{\text{Cov}[X_i, X_j]}{\sqrt{\text{Var}[X_i] \cdot \text{Var}[X_j]}} = \frac{\text{Cov}[X_i, X_j]}{\sigma_i \cdot \sigma_j}$  ... normalized covariance

④ covariance matrix of  $\tilde{X} = (X_1, \dots, X_m)$

$$\Sigma = \text{Cov}[\tilde{X}] \text{ where } \Sigma_{ij} = \text{Cov}[X_i, X_j] \quad \dots \quad \odot \Sigma_{ii} = \text{Var}[X_i]$$

⑤ correlation matrix of  $\tilde{X} = (X_1, \dots, X_m)$

$$C = \text{Cor}[\tilde{X}] \text{ where } C_{ij} = \text{Cor}[X_i, X_j] \quad \dots \quad \odot C_{ii} = 1$$

Fact: Correlation is scaled s.t.  $-1 \leq \text{Cor}[X_i, X_j] \leq 1$

Def:  $X, Y$  are independent,  $X \perp Y \equiv P[X=x \& Y=y] = P[X=x] \cdot P[Y=y]$

Theorem: If  $X$  and  $Y$  are independent, then

$$E[XY] = E[X]E[Y] \Rightarrow \text{Cov}(X, Y) = \text{Cov}(X, Y) = 0$$

Theorem: Linear combinations of random variables give

$$\textcircled{1} E[aX+b] = aE[X] + b$$

$$\textcircled{2} E[aX+bY] = aE[X] + bE[Y]$$

$$\textcircled{3} \text{Var}[aX+b] = a^2 \text{Var} X$$

$$\textcircled{4} \text{Var}[aX+bY] = a^2 \text{Var} X + b^2 \text{Var} Y + 2ab \underbrace{\text{Cov}[X, Y]}_{0 \text{ if } X \perp Y}$$

$$\textcircled{5} \text{Cov}[aX+b, cY+d] = ac \cdot \text{Cov}[X, Y]$$

$$\textcircled{6} \text{Cov}[aX+bW, cY+dZ] = ac \text{Cov}[X, Y] + ad \text{Cov}[X, Z] + bc \text{Cov}[W, Y] + bd \text{Cov}[W, Z]$$

Matrix representation

• Expected value

$$E[a_1 X_1 + \dots + a_m X_m] = a_1 \mu_1 + \dots + a_m \mu_m$$

$$a = (a_1, \dots, a_m)^T \in \mathbb{R}^m$$

$$\tilde{X} = (X_1, \dots, X_m)^T$$

$$\tilde{\mu} = (\mu_1, \dots, \mu_m)^T$$

$$\Rightarrow E[a^T \tilde{X}] = a^T \tilde{\mu}$$

• Variance

$$\text{Var}[a_1 X_1 + \dots + a_m X_m] = \sum_{i=1}^m a_i^2 \overbrace{\text{Var} X_i}^{\Sigma_{ii}} + \sum_{\substack{i, j=1 \\ i \neq j}}^m a_i a_j \overbrace{\text{Cov}[X_i, X_j]}^{\Sigma_{ij}}$$

$$\Rightarrow \underline{\text{Var}[a^T \tilde{X}] = a^T \Sigma a}$$

• Covariance

$$\text{Cov}[a_1 X_1 + \dots + a_m X_m, b_1 X_1 + \dots + b_m X_m] = \sum_{i, j=1}^m a_i b_j \text{Cov}[X_i, X_j]$$

$$\Rightarrow \underline{\text{Cov}[a^T \tilde{X}, b^T \tilde{X}] = a^T \Sigma b = b^T \Sigma a}$$

# Principle Component Analysis - PCA → eigenvector

Def:  $\lambda \in \mathbb{R}$  is an eigenvalue of  $A \in \mathbb{R}^{m \times m} \equiv \exists v \neq 0$  s.t.  $Av = \lambda v$

Finding eigenvalues:  $Av = \lambda v \Rightarrow Av - \lambda v = 0 \Rightarrow (A - \lambda I)v = 0 \Rightarrow \det(A - \lambda I) = 0$

⊛ if  $\det(A - \lambda I) = 0$ , then  $v = 0$  is the only solution ⊛

Def:  $v$  is an unit eigenvector  $\equiv \|v\| = 1$  ... standard norm  $\|v\| = \sqrt{v^T v}$

Def: vectors  $u, v \in \mathbb{R}^m$  are

- orthogonal  $u \perp v \equiv u^T v = 0$
- orthonormal  $\equiv u \perp v$  &  $\|u\| = \|v\| = 1$

## Properties of the covariance matrix

Theorem:  $\text{Cov}[\tilde{X}] = \Sigma$  is symmetric and positive semi-definite.

Pf: Symmetric because  $\text{Cov}[x_i, x_j] = \text{Cov}[x_j, x_i]$

note:  $\text{Var}[a^T \tilde{X}] = a^T \Sigma a$

↳ variance is always  $\geq 0 \Rightarrow \Sigma$  is positive semi-definite

Corollary: The eigenvalues of  $\Sigma$  are all non-negative

Pf:  $\Sigma v = \lambda v \Rightarrow v^T \Sigma v = v^T \lambda v = \lambda v^T v \Rightarrow \lambda = \frac{v^T \Sigma v}{v^T v} \geq 0$

Theorem:  $\Sigma \in \mathbb{R}^{m \times m}$  has  $m$  orthonormal eigenvectors.

Pf:  $\Sigma$  is symmetric  $\Rightarrow \exists R \in \mathbb{R}^{m \times m}$  s.t.  $R^T \Sigma R$  is diagonal &  $R^T R = I_m$

↳ spectral decomp. of a symmetric matrix

- since  $\Sigma$  is diagonalizable, it has  $m$  lin. ind. eigenvectors
- they are orthonormal because  $R^T R = I_m$  and the columns of  $R =$  eigenvectors of  $\Sigma$

## • PCA

→ idea: our variables might not be conveying the information hidden in the data very efficiently

⇒ it might be possible to express most of the info with just a few carefully chosen linear combinations of the variables

→ if two variables are highly correlated, then we really only need one (kinda)

Goal: Describe the variation in a set of correlated variables  $X_1, \dots, X_m$  using a new set of uncorrelated variables  $Y_1, \dots, Y_k$  and hopefully  $k \ll m$

$Y_1 =$  lin. comb. of  $\tilde{X}$  such that it accounts for the most variation possible

$Y_2 =$  lin. comb. of  $\tilde{X}$  accounting for as much of the remaining variance while subject to the constraint  $\text{Cov}(Y_1, Y_2) = 0 \dots Y_1 \perp Y_2$

Def: We have data:  $n$  observations of  $m$  variables  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_m \in \mathbb{R}^m$

① sample mean of each variable  $\bar{X}_j := \frac{1}{n} \sum_{i=1}^n X_{ij}$

② sample covariance matrix  $Q \in \mathbb{R}^{m \times m}$

$$q_{ij} = \frac{1}{n-1} \sum_{z=1}^n (X_{zi} - \bar{X}_i)(X_{zj} - \bar{X}_j)$$

↳  $z$ th observation of variables  $X_i$  and  $X_j$

Fact: Previously discussed results for  $\Sigma$  also apply for  $Q$

Note: We divided by  $n-1$  rather than  $n$ . This is because by using the sample mean, we have lost a degree of freedom. If we knew  $EX_i$  and  $EX_j$ , then we would divide by  $n$ .

→ consider a sample of  $X_i$  of size  $n$ . It exists in  $\mathbb{R}^n$  and so has  $n$  degrees of freedom in movement - each sample member can be anything

→ however by fixing the sample mean, we are constraining the sample to have a fixed sum

⇒  $n-1$  members can be anything, but the last one must have the proper value to ensure the sample mean is unchanged

→ we are summing over some values  $n$  times, but it is in truth a sum of only  $n-1$  independent things

## Linear combinations of data

→ let  $Y$  be lin comb of our  $m$  variables  $Y = \sum_{i=1}^m a_i X_i = a^T \tilde{X}$

→ to determine  $\text{Var} Y$  we would need  $\Sigma = \text{Cov}(\tilde{X})$

⇒ estimate  $\Sigma$  with  $Q$

⇒  $\text{Var}(a^T \tilde{X}) \approx a^T Q a$ ,  $\text{Cov}(a^T \tilde{X}, b^T \tilde{X}) \approx a^T Q b = b^T Q a$

Problem: Find  $a \in \mathbb{R}^m$  to maximize  $\text{Var}(a^T \tilde{X}) \approx a^T Q a$  subject to  $a^T a = 1$ .

Solution: Use Lagrange multipliers. We want to

maximize  $f: \mathbb{R}^m \rightarrow \mathbb{R}^+$  subject to  $g(\tilde{a}) := a^T a - 1 = 0$   
 $\tilde{a} \mapsto a^T Q a$

Theorem: A necessary condition for  $\tilde{a}$  being a local maximum of  $f$  subject to  $g(\tilde{a}) = 0$  is existence of  $\lambda \in \mathbb{R}$  s.t.

$$\nabla f(\tilde{a}) = \lambda \nabla g(\tilde{a})$$

⇒ define  $L(\tilde{a}) := f(\tilde{a}) - \lambda g(\tilde{a}) = a^T Q a - \lambda(a^T a - 1)$  ... we want  $\nabla L = 0$

$$L = \sum_{i,j=1}^m a_i a_j q_{ij} - \lambda \sum_{i=1}^m a_i^2 + 1$$

$$\frac{\partial L}{\partial a_k} = 2a_k q_{kk} + \sum_{j \neq k} a_j q_{kj} + \sum_{i \neq k} a_i q_{ik} - 2\lambda a_k$$

$$= 2a_k q_{kk} + 2 \sum_{j \neq k} a_j q_{kj} - 2\lambda a_k \quad \dots Q \text{ is symmetric}$$

$$\Rightarrow \frac{\partial L}{\partial a_k} = 0 \Leftrightarrow \sum_{j=1}^m a_j q_{kj} - \lambda a_k = 0$$

$$\Leftrightarrow (q_{k1}, q_{k2}, \dots, q_{km}) = \lambda \cdot \text{row of } Q$$

$$\Rightarrow \nabla L = 0 \Leftrightarrow Q \tilde{a} = \lambda \tilde{a} \Leftrightarrow \lambda \text{ is an eigenvalue of } Q, \tilde{a} \text{ eigenvector}$$

Conclusion:  $\tilde{a}$  maximizes  $\text{Var}(a^T \tilde{X})$   $\Leftrightarrow \tilde{a}$  is an eigenvector of  $Q$ .

↳ also since we require  $a^T a = 1$ , we want unit eigenvectors

$$\text{Var}(a^T \tilde{X}) = a^T Q a = a^T \lambda a = \lambda a^T a = \lambda \Rightarrow \underline{\text{variance} = \text{eigenvalue}}$$

👁 Consider  $\tilde{a} \perp \tilde{b}$  eigenvectors of  $Q$  with eigenvalues  $\lambda_a, \lambda_b$

$$\text{Cov}(a^T \tilde{X}, b^T \tilde{X}) = a^T Q b = a^T \lambda_b b = \lambda_b a^T b = 0$$

↳ eigenvectors of  $Q$  are orthogonal

⇒  $Y_1 = a^T \tilde{X}$  and  $Y_2 = b^T \tilde{X}$  are independent

## • Principle components

Method: Given  $n$  observations  $\tilde{x}_1, \dots, \tilde{x}_m \in \mathbb{R}^m$  of  $m$  variables  $X_1, \dots, X_m$

1) compute the covariance matrix of  $\tilde{x}_1, \dots, \tilde{x}_m \rightarrow Q \in \mathbb{R}^{m \times m}$

2) calculate the eigenvalues and eigenvectors of  $Q$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \quad \forall i: a_i^T a_i = 1, \quad \forall i \neq j: \tilde{a}_i \perp \tilde{a}_j$$

$$\tilde{a}_1, \tilde{a}_2, \dots, \tilde{a}_m$$

3) the principle components of the data are variables

$$Y_1 = a_1^T \tilde{X}, \quad Y_2 = a_2^T \tilde{X}, \quad \dots \quad Y_m = a_m^T \tilde{X}, \quad \forall i \neq j: Y_i \perp Y_j$$

$$\text{Var } Y_1 = \lambda_1 \geq \text{Var } Y_2 = \lambda_2 \geq \dots \geq \text{Var } Y_m = \lambda_m$$

4) transform the data matrix of  $X_1, \dots, X_m$  to a data matrix of  $Y_1, \dots, Y_m$

$\rightarrow$  suppose  $Y_j = a_j^T \tilde{X}$  and observation  $\tilde{x}_i \in \mathbb{R}^m \Rightarrow \tilde{y}_j = a_j^T \tilde{x}_i$

$$X = \begin{pmatrix} \text{---} & \tilde{x}_1^T & \text{---} \\ \text{---} & \tilde{x}_2^T & \text{---} \\ \vdots & \vdots & \vdots \\ \text{---} & \tilde{x}_m^T & \text{---} \end{pmatrix} \in \mathbb{R}^{m \times m} \rightsquigarrow Y = X \cdot \begin{bmatrix} | & | & \dots & | \\ \tilde{a}_1 & \tilde{a}_2 & \dots & \tilde{a}_m \\ | & | & \dots & | \end{bmatrix} \in \mathbb{R}^{m \times m}$$

Theorem:  $\text{Var } X_1 + \dots + \text{Var } X_m = \text{Var } Y_1 + \dots + \text{Var } Y_m = \sum_i \lambda_i$

Proof: We will show

$$\text{trace}(Q) = \sum \text{Var } X_i = \sum \text{Var } Y_i = \sum \lambda_i$$

$\rightarrow$  consider the char. polynomial of  $Q$

$$p_Q(t) = b_m t^m + b_{m-1} t^{m-1} + \dots + b_0 = \det \begin{bmatrix} q_{11} - t & q_{12} & \dots & q_{1m} \\ q_{21} & q_{22} - t & & \vdots \\ \vdots & & \ddots & \\ q_{m1} & \dots & & q_{mm} - t \end{bmatrix} = \prod_{i=1}^m (a_{ii} - t) + \dots$$

$$\text{eye: } b_{m-1} = (-1)^{m-1} \sum_{i=1}^m a_{ii} = (-1)^{m-1} \text{tr}(Q)$$

$\rightarrow$  since there are  $m$  lin. ind. eigenvectors:  $m = \sum \text{Geom } \lambda \leq \sum \text{Alg } \lambda \leq m$

$$p_Q(t) = (\lambda_1 - t)^{r_1} (\lambda_2 - t)^{r_2} \dots (\lambda_k - t)^{r_k}, \quad r_1 + r_2 + \dots + r_k = m = \sum \text{Alg } \lambda$$

$$\text{eye: } b_{m-1} = (-1)^{m-1} \sum_{i=1}^k r_i \lambda_i$$

$\Rightarrow$  when  $Q$  has eigenvalues (not necessarily distinct)  $\lambda_1, \dots, \lambda_m: \text{tr } Q = \sum_i \lambda_i$   $\blacksquare$

5) using this theorem we can see that

$$\frac{\lambda_k}{\sum_i \lambda_i} = \text{proportion of variance explained by } Y_k$$

## • Scaling the data

Problem: The Iris dataset gives flower measurements in cm

→ would the results of PCA be different if one of the measurements was in mm?

→ yes, because PCA seeks to maximize variance and

$$\text{Var}(aX) = a^2 \text{Var}X \Rightarrow \text{it is very sensitive to data scaling}$$

→ even worse, what if  $X_1$  measured time and  $X_2$  length?

Solution: We want comparable units

⇒ make each variable have variance = 1

⇒ divide each variable by its std. dev:  $\text{Var}(\frac{1}{\sigma} X) = \frac{1}{\sigma^2} \text{Var}X = 1$

⇒ do PCA on the Cor. matrix of the transformed data

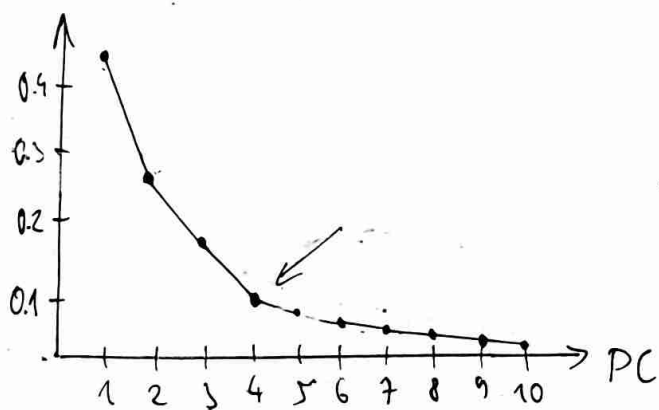
↳ this is the same as doing PCA with the Corr. matrix from the start

→ since  $\forall i: \text{Var} X_i = 1$  we have  $\sum_i \lambda_i = \sum_i \text{Var} X_i = m$

## • How many PCs to choose?

1) keep adding until a fixed proportion of variance (90%) is included

2) look at the Scree plot



$y$  = proportion of variance explained

→ look for an elbow

⇒ here we probably want to go with the first 4 components

## • When can PCA be used?

→ continuous variables

→ or numeric variables which could be interpreted as continuous

→ can't be applied for categorical variables (doesn't make sense)



# Hierarchical Clustering

Idea: combine data to clusters based on how similar they are to each other  
 ↳ start with every data point in its own cluster  
 → gradually combine them  
 ⇒ we need a dissimilarity measure between clusters

Def: The function  $d(\tilde{x}, \tilde{y})$  is a metric ≡

- i)  $d(\tilde{x}, \tilde{y}) \geq 0$  &  $d(\tilde{x}, \tilde{y}) = 0 \iff \tilde{x} = \tilde{y}$
- ii)  $d(\tilde{x}, \tilde{y}) = d(\tilde{y}, \tilde{x})$
- iii)  $d(\tilde{x}, \tilde{y}) \leq d(\tilde{x}, \tilde{z}) + d(\tilde{z}, \tilde{y})$  ... in our case might be ignored

## Examples

- Euclidean =  $\left(\sum_i (x_i - y_i)^2\right)^{\frac{1}{2}}$
- Manhattan =  $\sum_i |x_i - y_i|$
- Maxim =  $\max_i |x_i - y_i|$
- Minkowski =  $\left(\sum_i |x_i - y_i|^p\right)^{\frac{1}{p}}$ ,  $p \geq 1$

## Metrics for binary data

→ look at a cross tabulation of  $\tilde{x}$  and  $\tilde{y}$  and considering how much they agree / disagree

		Point $\tilde{y}$		
		1	0	
Point $\tilde{x}$	1	a	b	a+b
	0	c	d	c+d
		a+c	b+d	a+b+c+d

x \ y	1	0		
1	1	2	3 = # 1 in x	
0	3	1	4 = # 0 in x	
		4	3	7

• Hamming =  $1 - \frac{a+d}{a+b+c+d}$

• Jaccard =  $1 - \frac{a}{a+b+c}$  ... ignore double absence, as may be a redundant variable

• Kulczynski =  $1 - \frac{1}{2} \left( \frac{a}{a+b} + \frac{a}{a+c} \right)$  ... average of ratios of agreement from two samples

• Czekanowski =  $1 - \frac{2a}{2a+b+c}$  ... more emphasis on double presence

## Categorical data

- we can simply do  $d(\bar{x}, \bar{y}) := \# \text{ categories where } \bar{x} \text{ and } \bar{y} \text{ don't agree}$

## Mixed data

- calculate dissimilarities for continuous, binary and categorical variables separately  
⇒ then do a weighted combination

Def: The dissimilarity of two groups  $A = \{x_1, \dots, x_k\}$  and  $B = \{y_1, \dots, y_l\}$  is a linkage

• Single linkage =  $\min_{i,j} d(x_i, y_j)$



... doesn't hold  $\Delta$ -ineq

• Complete linkage =  $\max_{i,j} d(x_i, y_j)$



• Average linkage =  $\frac{1}{|A| \cdot |B|} \sum_{i,j} d(x_i, y_j)$



Example: Use Manhattan metric and complete linkage to cluster the following data

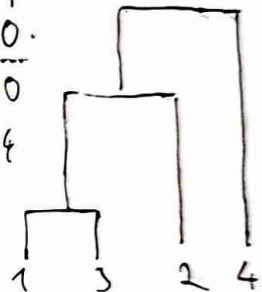
Observation	A	B	C	D
1	11	-6	-4	8
2	15	6	6	9
3	13	-5	-8	10
4	-12	5	-7	6

← variables dissimilarity matrix

	2	3	4
1	27	9	39
2		28	44
3			40

→ first we unify 1 and 3 to  $\{1,3\}$   
→ second 2 and  $\{1,3\}$  into  $\{1,2,3\}$   
→ lastly  $\{1,2,3\}$  and 4

$\{1,3\}$	28	40
$\{1,2,3\}$		44



## Dendrogram

→ the tree-like graph used to visualize hierarchical clustering is called a dendrogram  
→ joining clusters ~ U-link  
↳ length of the two legs of the U-link = distance between clusters

## Scaling the data

→ we should scale the data before constructing the dissimilarity matrix  
→ if the variables are not scaled, then the variable with the greatest variance will dominate the distances and figure most prominently in the clustering solution

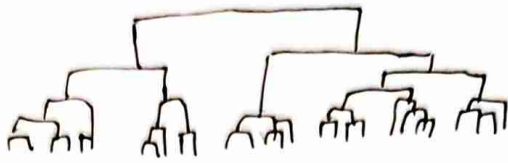
⇒ we divide the variable by its std. dev. ⇒  $\text{Var} = 1$

⇒ then scaling is not good

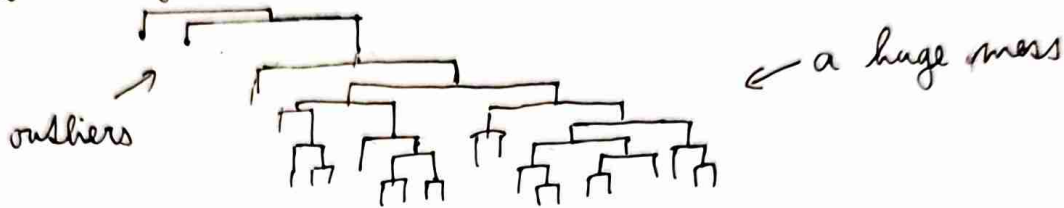
! sometimes we might want to give less weight to a variable which carries less info (has small variance)

## Linkage effects

- complete linkage joins the final clusters at a larger measure of dissimilarity
- complete and average linkage result in "spherical" clusters with good internal similarity



- single linkage displays outliers, which are often hidden in complete linkage

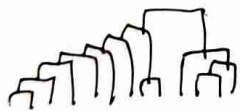


- complete and single linkage are invariant under monotonic transformations of the dissimilarity matrix entries, while average linkage is not
- complete linkage is likely to suggest a smaller number of large clusters

! we should first give careful thought to which linkage makes sense  
↳ trying them all and deciding based on which one looks the best can easily turn an objective solution to a subjective one

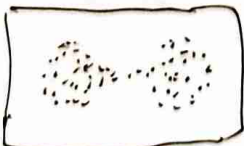
## Chaining

- phenomenon which occurs while using single linkage
- single linkage has the tendency to repeatedly add a single observation to the same group that continues to get larger and larger

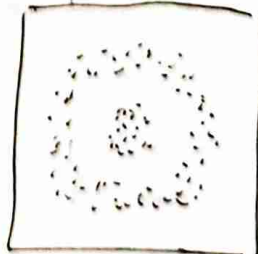


→ occurs because a unit joins a group based on similarity with just one member of that group

- this results in elongated clusters that may include quite dissimilar points
- however, chaining is not always bad



- average works ✓
- single chains all together



- average/complete tries to create spherical clusters and fails (0)
- single manages to make 2 clear clusters thanks to chaining (0)

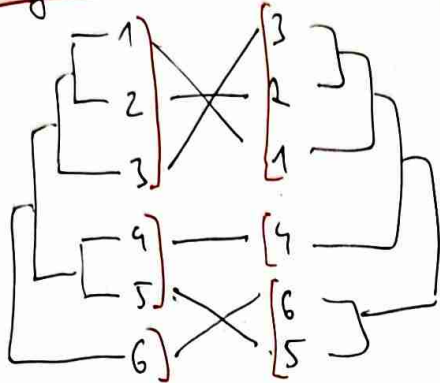
# Cluster Agreement

→ two different clustering methods were applied to the same data

- single × complete linkage hierarchical clustering
- expert uses his knowledge × statistician uses data

→ can we quantify the level of agreement between the two approaches?

## Tanglegrams



# clusters = 3

→ the two methods agree on the first cluster

→ but they disagree on the other two

! they might disagree with the same # clusters but agree on different



## Cross tabulation

		Method B		
		C <sub>1</sub>	C <sub>2</sub>	
Method A	C <sub>1</sub>	50	10	40
	C <sub>2</sub>	15	60	75
		45	70	115

agree on  $\frac{90}{115}$

	C <sub>1</sub>	C <sub>2</sub>	
C <sub>1</sub>	10	30	40
C <sub>2</sub>	60	15	75
	70	45	115

same as the first table  
→ just rename C<sub>1</sub> and C<sub>2</sub>

	C <sub>1</sub>	C <sub>2</sub>	
C <sub>1</sub>	29	16	40
C <sub>2</sub>	46	29	75
	70	45	115

→ don't really agree on anything  
 $\frac{53}{115}$  and  $\frac{62}{115}$  are both meh

	C <sub>1</sub>	C <sub>2</sub>	
C <sub>1</sub>	10	30	40
C <sub>2</sub>	20	5	25
C <sub>3</sub>	40	10	50
	70	45	115

→

	C <sub>1</sub>	C <sub>2</sub>	
C <sub>1</sub>	10	30	40
C <sub>2+C<sub>3</sub></sub>	60	15	75
	70	45	115

→ this is in fact the same table as in example 2

→ in more complex cases it is very difficult to describe agreements just by looking

⇒ in 1941 has Rand proposed an index for measuring the agreement as a number between 0 and 1

# The Rand Index

A \ B	C1	C2	...	C <sub>l</sub>	
C1	M <sub>11</sub>	M <sub>12</sub>	...	M <sub>1l</sub>	M <sub>1*</sub>
C2	M <sub>21</sub>	M <sub>22</sub>	...	M <sub>2l</sub>	M <sub>2*</sub>
⋮	⋮	⋮	⋮	⋮	⋮
C <sub>l</sub>	M <sub>l1</sub>	M <sub>l2</sub>		M <sub>ll</sub>	M <sub>l*</sub>
	M <sub>*1</sub>	M <sub>*2</sub>	...	M <sub>*l</sub>	M

→ n observations  $\tilde{X}_1, \dots, \tilde{X}_n$

$\alpha := \#(\tilde{X}_i, \tilde{X}_j)$ : method A: same cluster  
method B: same cluster

$\beta := \#(\tilde{X}_i, \tilde{X}_j)$ : method A: different clusters  
method B: different clusters

$\gamma := \#(\tilde{X}_i, \tilde{X}_j)$ : A: same  
B: different

$\delta := \#(\tilde{X}_i, \tilde{X}_j)$ : A: different  
B: same

👁️ # all possible pairs of data points  $\binom{n}{2} = \alpha + \beta + \gamma + \delta$

⇒ Rand index  $R = \frac{\alpha + \beta}{\alpha + \beta + \gamma + \delta} = \frac{\binom{n}{2} - \gamma - \delta}{\binom{n}{2}}$

$\gamma = \sum_{i=1}^k \binom{M_{i*}}{2} - \sum_{i,j} \binom{M_{ij}}{2}$  } • same = same - same  
 $\delta = \sum_{j=1}^l \binom{M_{*j}}{2} - \sum_{i,j} \binom{M_{ij}}{2}$  } • diff = \* - same  
 • same = \* - same

→ the rand index of the tables on previous page is

1: 0.654      2: 0.654      3: 0.999      4: 0.588

⚠️ the value for table 3 seems quite large, but there wasn't much agreement

Example: Calculate the rand index of

	C1	C2	
C1	5	25	30
C2	20	5	25
C3	40	5	45
	65	35	100

$\alpha = \sum_{i,j} \binom{M_{ij}}{2} = 3 \cdot \binom{5}{2} + \binom{25}{2} + \binom{20}{2} + \binom{40}{2} = 1300$

$\gamma = \sum_i \binom{M_{i*}}{2} - \alpha = \binom{30}{2} + \binom{25}{2} + \binom{45}{2} - 1300 = 425$

$\delta = \sum_j \binom{M_{*j}}{2} - \alpha = \binom{65}{2} + \binom{35}{2} - 1300 = 1375$

$\binom{n}{2} = \binom{100}{2} = 4950 \Rightarrow R = \frac{4950 - 425 - 1375}{4950} \approx 0.636$

Problem: The Rand index tends to give large values even when the clustering methods are in substantial disagreement

## Adjusted Rand index

$\alpha = \sum_{i,j} \binom{M_{ij}}{2}$

$\beta_A = \sum_i \binom{M_{i*}}{2}$

$\beta_B = \sum_j \binom{M_{*j}}{2}$

$ARI = \frac{\alpha \binom{n}{2} - \beta_A \beta_B}{\frac{1}{2}(\beta_A + \beta_B) \cdot \binom{n}{2} - \beta_A \beta_B} = \frac{RI - E[RI]}{\text{Max RI} - E[RI]}$

↳ use hypergeometric distribution

→ ARI can be negative, but it can't be greater than 1

→ for tables 1 to 4: 1: 0.311, 2: 0.311, 3: -0.006, 4: 0.185

↳ now we have a small value for table 3

→  $ARI \leq 0 \Rightarrow$  no agreement

## • Classification with kNN = k-nearest neighbours

- we have data with labels and want to classify new data points

- kNN is non-parametric and doesn't make any assumptions on the spread of the data  
↳ the distribution of the data

⇒ there is no measurement of uncertainty when assigning labels

- kNN simply looks at the k closest points and assigns the new point to the group which has the majority

- Two things to consider

1) do we scale the data? results may vary

2) how will we calculate the distance?

3) how big should k be? ← important

- choosing k

→ split the data into

• training set - will be used to classify "unlabeled" data

• test set - treated as unlabeled and is used to find the best k

• validation set - treated as unlabeled, used to estimate the classification error of the best k

→ plot the proportion of incorrectly classified data from the test set



→  $k=8$  has the best rate ( $8 < 10$ )

⇒ now test the correct classification rate on the validation set to estimate the error

→ this is a general technique

• training data - data used to fit several models

• test data - data used to compare these models and choose the best one

• validation data - used to assess the performance of the chosen model

→ typical split is 50% 25% 25%

## • Cross-Validation

- another general performance assessment technique we can use to pick  $k$

→ do this for every model:

1) split the data into  $k$  subsets  $X_1, \dots, X_k$

2) For  $i=1 \dots k$ :

- fit the model using  $\{X_1, \dots, X_k\} \setminus \{X_i\}$

- then count how many points from  $X_i$  would be correctly classified

3) calculate the total classification rate

⇒ now pick the model with the best classification rate

• leave-one-out crossvalidation = # subsets = # datapoints

•  $k$ -fold crossvalidation = # subsets =  $k$

## Discriminant Analysis - classification

- supervised statistical techniques where we assume some info about the classes and use that to classify new data

- used when we know that there are  $k$  groups within the data and that there is a subset of the data which is labeled

- Linear DA (LDA) and Quadratic DA (QDA) both assume a distribution over the data

⇒ now we can use probability theory to calculate the probability of a point belonging to a group under the assumptions we have made

Def: Let  $\tilde{X} = (X_1, \dots, X_m)^T$  be a vector of random variables.

We say that  $\tilde{X}$  follows a Multi-Variate Normal (MVN) distribution

$\tilde{X} \sim \text{MVN}(\tilde{\mu}, \Sigma)$ , where  $\tilde{\mu} \in \mathbb{R}^m$  and  $\Sigma \in \mathbb{R}^{m \times m}$  is positive semi-definite

≡ the probability-density function of  $\tilde{X}$  is

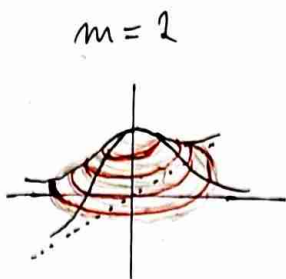
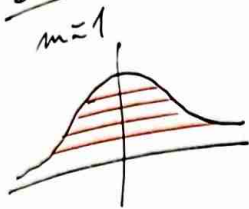
$$f(\tilde{X} | \tilde{\mu}, \Sigma) = \left( (2\pi)^m |\Sigma| \right)^{-\frac{1}{2}} \cdot \exp \left( -\frac{1}{2} (\tilde{X} - \tilde{\mu})^T \Sigma^{-1} (\tilde{X} - \tilde{\mu}) \right)$$

↳  $\det(\Sigma)$

Theorem: The covariance matrix of  $\tilde{X}$  is  $\text{cov}(\tilde{X}) = \Sigma$ .

The means of  $\tilde{X}$  are given by  $\tilde{\mu}$  ...  $E[X_i] = \mu_i$

Intuition:



→ the pdf has a bell shape

→  $\{\tilde{x} | f(\tilde{x}) = c\}$  makes circles / ellipses

→  $\{\tilde{x} | f(\tilde{x}) \geq c\}$  makes filled in ellipsoids

→ in general:

$$f(\tilde{x}) \geq c \Leftrightarrow \frac{1}{((2\pi)^m |\Sigma|)^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\tilde{x}-\tilde{\mu})^T \Sigma^{-1}(\tilde{x}-\tilde{\mu})\right) \geq c$$

$$\Leftrightarrow \exp(\dots) \geq c(\dots)^{\frac{1}{2}}$$

$$\Leftrightarrow -\frac{1}{2}(\tilde{x}-\tilde{\mu})^T \Sigma^{-1}(\tilde{x}-\tilde{\mu}) \geq \ln(\dots)$$

$$\Leftrightarrow (\tilde{x}-\tilde{\mu})^T \Sigma^{-1}(\tilde{x}-\tilde{\mu}) \leq -2 \ln(c \cdot \sqrt{(2\pi)^m |\Sigma|})$$

$$\Rightarrow \{\tilde{x} | f(\tilde{x}) < c\} = \{\tilde{x} | (\tilde{x}-\tilde{\mu})^T \Sigma^{-1}(\tilde{x}-\tilde{\mu}) \leq \text{some number}\}$$

↳ which for positive-semidefinite  $\Sigma$  is an ellipsoid centered at  $\tilde{\mu}$

→ if we assume that the data within group  $k$  follows  $MVN(\tilde{\mu}_k, \Sigma_k)$ , then the scatter of the data should be roughly elliptical

↳  $\tilde{\mu}_k \sim$  location of the ellipsoid &  $\Sigma_k \sim$  shape of the ellipsoid

Distance approach

Def: The Mahalanobis distance of  $\tilde{x} \in \mathbb{R}^m$  from the center  $\tilde{\mu}$  is  $D$ , where

$$D^2 = (\tilde{x}-\tilde{\mu})^T \Sigma^{-1}(\tilde{x}-\tilde{\mu})$$

☞ Two points  $\tilde{x}_1$  and  $\tilde{x}_2$  are on the same ellipsoid shell  $\Leftrightarrow D_1 = D_2$



→ we can use the Mahalanobis distance to which cluster should  $\tilde{x}$  belong

↳ want to find cluster  $k$  s.t.  $(\tilde{x}-\tilde{\mu}_k)^T \Sigma_k^{-1}(\tilde{x}-\tilde{\mu}_k)$  is minimized

$\Rightarrow \tilde{x}$  is closer to cluster 1 than to cluster 2  $\Leftrightarrow$

$$(\tilde{x}-\tilde{\mu}_1)^T \Sigma_1^{-1}(\tilde{x}-\tilde{\mu}_1) < (\tilde{x}-\tilde{\mu}_2)^T \Sigma_2^{-1}(\tilde{x}-\tilde{\mu}_2) \quad \leftarrow \text{Quadratic expression in } \tilde{x}$$

→ if  $\Sigma$  is the same for all clusters, this simplifies:

$$(\tilde{x}-\tilde{\mu}_1)^T \Sigma^{-1}(\tilde{x}-\tilde{\mu}_1) < (\tilde{x}-\tilde{\mu}_2)^T \Sigma^{-1}(\tilde{x}-\tilde{\mu}_2)$$

$$\underbrace{x^T \Sigma^{-1} x}_{\text{same}} - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 < \underbrace{x^T \Sigma^{-1} x}_{\text{same}} - x^T \Sigma^{-1} \mu_2 - \mu_2^T \Sigma^{-1} x + \mu_2^T \Sigma^{-1} \mu_2$$

$$-2x^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mu_1 < -2x^T \Sigma^{-1} \mu_2 + \mu_2^T \Sigma^{-1} \mu_2$$

$$\Rightarrow \underline{x^T \Sigma^{-1} (\mu_1 - \mu_2) > \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2)} \quad \leftarrow \text{Linear expression in } \tilde{x}$$

$\Sigma$  is symmetric  
 $\Rightarrow \Sigma^{-1}$  is as well



## • Probabilistic approach

→ we want to classify a new observation  $\tilde{x}$  into one of the  $K$  clusters

⇒ let  $X: \mathbb{R}^m \rightarrow [K]$  be a random variable assigning points to clusters

Define  $\pi_k := P[X=k]$  = proportion of population objects belonging to cluster  $k$

→ Bayes theorem states

$$P[X=k|\tilde{x}] = \frac{P[X=k] \cdot P[\tilde{x}|X=k]}{P[\tilde{x}]} \rightsquigarrow P[X=k|\tilde{x}] \propto \pi_k P[\tilde{x}|X=k]$$

→ we are however dealing with continuous r.v., so we need to use the pdf,

$$P[X=k|\tilde{x}] \propto \pi_k \cdot f(\tilde{x}|X=k) = \pi_k \cdot f_k(\tilde{x}) \leftarrow$$

→ we know the pdf for cluster  $k$  is from  $MVN(\tilde{\mu}_k, \Sigma_k)$

→ we can calculate these probabilities by estimating  $\tilde{\mu}_k$  and  $\Sigma_k$

→ then we will assign  $\tilde{x}$  to the cluster with the largest probability

$$P[X=k|\tilde{x}] > P[X=l|\tilde{x}] \Leftrightarrow \pi_k \cdot f_k(\tilde{x}) > \pi_l \cdot f_l(\tilde{x})$$

$$\Leftrightarrow \ln \pi_k + \ln f_k(\tilde{x}) > \ln \pi_l + \ln f_l(\tilde{x})$$

Recall:  $f_k(\tilde{x}) = \left( (2\pi)^m |\Sigma_k| \right)^{-\frac{1}{2}} \cdot \exp\left(-\frac{1}{2}(\tilde{x}-\tilde{\mu}_k)^T \Sigma_k^{-1}(\tilde{x}-\tilde{\mu}_k)\right)$

$$\Leftrightarrow \ln \pi_k - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2}(\tilde{x}-\tilde{\mu}_k)^T \Sigma_k^{-1}(\tilde{x}-\tilde{\mu}_k) > \ln \pi_l - \frac{1}{2} \ln |\Sigma_l| - \frac{1}{2}(\tilde{x}-\tilde{\mu}_l)^T \Sigma_l^{-1}(\tilde{x}-\tilde{\mu}_l)$$

① Linear DA: assumes all  $\Sigma$  for all clusters

$$\ln \pi_k - \frac{1}{2} \left[ x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu_k + \mu_k^T \Sigma^{-1} \mu_k \right] > \ln \pi_l - \frac{1}{2} \left[ x^T \Sigma^{-1} x - 2x^T \Sigma^{-1} \mu_l + \mu_l^T \Sigma^{-1} \mu_l \right]$$

$$\ln \pi_k + x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k > \ln \pi_l + x^T \Sigma^{-1} \mu_l - \frac{1}{2} \mu_l^T \Sigma^{-1} \mu_l$$

$$\ln \frac{\pi_k}{\pi_l} + x^T \Sigma^{-1} (\mu_k - \mu_l) > \frac{1}{2} (\mu_k^T \Sigma^{-1} \mu_k - \mu_l^T \Sigma^{-1} \mu_l) \Leftrightarrow P[\tilde{x} \in k] > P[\tilde{x} \in l]$$

↳ if we assume  $\forall k: \pi_k = \frac{1}{K}$ , then  $\ln \frac{\pi_k}{\pi_l} = 0$

and we get the formula we have gotten from the distance approach

② Quadratic DA: different clusters have different  $\Sigma_k$

→ no simplification arises, so  $P[\tilde{x} \in k] > P[\tilde{x} \in l] \Leftrightarrow$

$$\ln \pi_k - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2}(\tilde{x}-\tilde{\mu}_k)^T \Sigma_k^{-1}(\tilde{x}-\tilde{\mu}_k) > \ln \pi_l - \frac{1}{2} \ln |\Sigma_l| - \frac{1}{2}(\tilde{x}-\tilde{\mu}_l)^T \Sigma_l^{-1}(\tilde{x}-\tilde{\mu}_l)$$

## • Estimating Covariances

→ labeled data classified into  $K$  groups of sizes  $n_1, n_2, \dots, n_K$

1. For  $\forall$  group  $k$ : calculate its sample covariance matrix  $Q_k \in \mathbb{R}^{m \times m}$  and its sample mean  $\bar{x}_k \in \mathbb{R}^m$

• QDA: we are finished

• LDA: we assume that  $\forall k, l: \Sigma_k = \Sigma_l \Rightarrow$  need to pool the matrices

→  $Q_k$  has  $(n_k - 1)$  degrees of freedom ... fixes  $\tilde{\mu}_k$

→ the pooled matrix  $Q$  will have  $N - K$  degrees of freedom ...  $n_1, \dots, n_K$

$$\Rightarrow Q = \frac{1}{N - K} \sum_{k=1}^K (n_k - 1) \cdot Q_k$$

→  $N$  data points in total

## 2. Perform classification

→ define  $\pi_k$ :

•  $\pi_k := \frac{1}{K}$  or  $\frac{n_k}{N}$  or something else

$$\bullet f_k(\tilde{x}) := \left( (2\pi)^m |Q_k| \right)^{-\frac{1}{2}} \cdot \exp\left(-\frac{1}{2}(\tilde{x} - \bar{x}_k)^T Q_k^{-1} (\tilde{x} - \bar{x}_k)\right)$$

→ the estimated probability that  $\tilde{x} \in$  group  $k$  is proportional to

$$\underline{P[k|\tilde{x}] \propto \pi_k f_k(\tilde{x})} \quad \leftarrow \text{posterior probabilities}$$

→ assign  $\tilde{x}$  to the group with the largest probability

→ the decision boundary between class  $k$  and class  $l$  is given by

$$\frac{P[k|\tilde{x}]}{P[l|\tilde{x}]} = \frac{\pi_k f_k(\tilde{x})}{\pi_l f_l(\tilde{x})} = 1 \quad \Rightarrow \quad \underline{\log \frac{P[k|\tilde{x}]}{P[l|\tilde{x}]} = \log \frac{\pi_k}{\pi_l} + \log \frac{f_k(\tilde{x})}{f_l(\tilde{x})} = 0}$$

## Summary:

- LDA and QDA are model based parametric classifiers where the data of each group is assumed to follow a MVN distribution

- model based  $\Rightarrow$  we can estimate the probability of correct assignment

- MVN  $\Rightarrow$  the groups are assumed to have an elliptical shape

↳ LDA: all groups have the same covariance matrix

↳ QDA: different covariance matrices between groups

## • k-means Clustering

- simple algorithm for dividing data into  $k$  groups

1. initialise the cluster means  $\tilde{\mu}_1, \dots, \tilde{\mu}_k$

2. while not happy:

3. assign every point to the nearest cluster mean

4. recalculate the cluster means based on those assignments

→ how can we tell that it's converging?

SS :=  $\sum_i (x_i - \mu(x_i))^2$  = sum of squared distances of each point to its centroid

⇒ keep iterating until

- 1, the assignments stop changing
- 2, the improvement in SS  $\approx 0$

→ picking initial cluster centroids

→ the algorithm might converge to a local minima

1. initialise randomly and do multiple runs
2. select them based on prior knowledge of the data
3. perform hierarchical clustering first as a basis

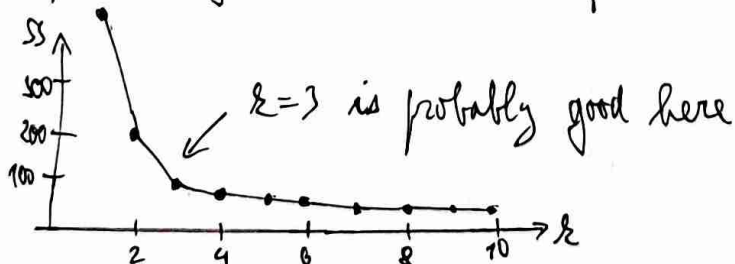
→ choosing  $k$

- what if we don't know how many clusters there should be?

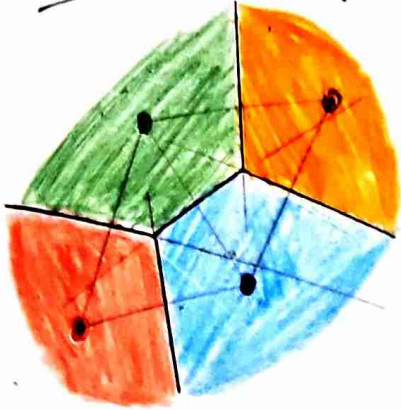
→ run for  $k=1, 2, \dots, 10, \dots$  and try to minimize SS

! SS will naturally decrease as  $k$  will go up

⇒ plot  $k$  against SS and look for an elbow in the graph



→ when k-means fails



→ k-means classifies data by dividing the plane by lines

→ it cannot deal with data that does not have compact spherical groups



### • Silhouette width

→ a general technique for evaluating the performance of a clustering solution

→ for each observation  $\tilde{x}_i$  from  $\tilde{x}_1, \dots, \tilde{x}_n$  compute

- $a_i :=$  average distance from  $\tilde{x}_i$  to the other points in its cluster
- $b_i :=$  average distance from  $\tilde{x}_i$  to the nearest cluster it is not in

↳ calculate the dist for every cluster and take the minimum

$$\bullet \Delta_i := \frac{b_i - a_i}{\max(a_i, b_i)} \quad \leftarrow \text{silhouette width of } \tilde{x}_i$$

👁  $-1 \leq \Delta_i \leq 1$

- $\Delta_i \approx 1 \Rightarrow b_i \gg a_i \Rightarrow$  good cluster separation
- $\Delta_i \approx 0 \Rightarrow b_i \approx a_i \Rightarrow$  clusters are poorly separated / overlapping
- $\Delta_i < 0 \Rightarrow \tilde{x}_i$  has probably been assigned to the wrong cluster

→ we can look at the average silhouette width in each cluster and for the data overall

# • Multidimensional Scaling - MDS

→ we have data  $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^m$

↳ distances  $d_{ij} := d(\tilde{x}_i, \tilde{x}_j)$

⇒ we want to find  $\tilde{y}_1, \dots, \tilde{y}_n \in \mathbb{R}^d$ ,  $d < m$

↳ distances  $\delta_{ij} := d(\tilde{y}_i, \tilde{y}_j)$

such that  $\delta_{ij}$  are as close to  $d_{ij}$  as possible for all  $i, j$

⇒ object of MDS = provide an optimal configuration of observations in  $\mathbb{R}^d$

→ we generally want a mapping  $\delta_{ij} = f(d_{ij})$

• metric MDS:  $f$  continuous and monotonic

• non-metric MDS:  $f$  monotonic

$$d_{ij} < d_{k\ell} \Rightarrow \delta_{ij} < \delta_{k\ell}$$

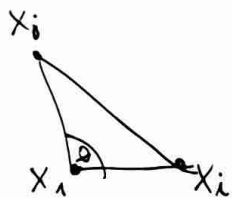
## ① Classical Metric Scaling

→ assume that  $f$  is the identity function

→ lets just say that we want to find an equivalent representation in  $\mathbb{R}^m$  which is centered in the origin

→ we have calculated the distances  $d_{ij}$  and want to reconstruct  $\tilde{x}_1, \dots, \tilde{x}_n$

→ let  $\tilde{x}_1$  be at the origin and consider  $\tilde{x}_j$  and  $\tilde{x}_k$



$$d_{ij}^2 = d_{1i}^2 + d_{1j}^2 - 2d_{1i}d_{1j} \cos(\theta) \Rightarrow -\frac{1}{2}(d_{ij}^2 - d_{1i}^2 - d_{1j}^2) = d_{1i}d_{1j} \cos \theta$$

$$x_i^T x_j = \|x_i\| \cdot \|x_j\| \cdot \cos \theta = d_{1i} d_{1j} \cos \theta, \text{ where } \|\cdot\| = \text{Euclid norm}$$

$$\Rightarrow \underbrace{x_i^T x_j}_{\text{want}} = -\frac{1}{2} \underbrace{(d_{ij}^2 - d_{1i}^2 - d_{1j}^2)}_{\text{know}}$$

⇒ we construct  $B \in \mathbb{R}^{m \times m}$  s.t.  $b_{ij} = -\frac{1}{2}(d_{ij}^2 - d_{1i}^2 - d_{1j}^2) = x_i^T x_j$

⇒  $B = X^T X$  for some  $X \rightarrow k^{\text{th}}$  column of  $X$  is  $\tilde{x}_k$

⇒  $B$  is symmetric ⇒ as in PCA can be decomposed as

$$B = R D R^T = R \sqrt{D} \sqrt{D} R^T = R \sqrt{D} \sqrt{D}^T R^T = (R \sqrt{D}) \cdot (R \sqrt{D})^T$$

⇒ hence  $X^T = R \sqrt{D}$ , where  $R$  contains unit eigenvectors of  $B$  and  $\sqrt{D}$  contains the square roots of the eigenvalues

⇒ the vector  $\tilde{x}_k$  can be expressed as

$$X^T = RVD \Rightarrow \begin{pmatrix} \text{---} x_1^T \text{---} \\ \text{---} x_2^T \text{---} \\ \vdots \\ \text{---} x_m^T \text{---} \end{pmatrix} = \begin{pmatrix} | & | & & | \\ \sqrt{\lambda_1} v_1 & \sqrt{\lambda_2} v_2 & \dots & \sqrt{\lambda_m} v_m \\ | & | & & | \end{pmatrix}$$

- $v_1, \dots, v_m$   
↳ unit eigenvectors
- $\lambda_1, \dots, \lambda_m$   
↳ eigenvalues

$$\Rightarrow \tilde{x}_{ki} = \sqrt{\lambda_i} v_{ki}$$

→ the squared euclidean distance between  $\tilde{x}_k$  and  $\tilde{x}_j$  is

$$d_{jk}^2 = \sum_{i=1}^m (\tilde{x}_{ji} - \tilde{x}_{ki})^2 = \sum_{i=1}^m (\sqrt{\lambda_i} v_{ij} - \sqrt{\lambda_i} v_{ik})^2 = \sum_{i=1}^m \lambda_i (v_{ij} - v_{ik})^2$$

⇒ lets say we want to shrink the data to  $\tilde{y}_1, \dots, \tilde{y}_m \in \mathbb{R}^d$ ,  $d < m$

1. order the eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$

2. define  $\tilde{y}_k \in \mathbb{R}^d$  as  $\tilde{y}_{ki} := \sqrt{\lambda_i} v_{ki}$

3. then the distance between  $\tilde{y}_k$  and  $\tilde{y}_j$  is

$$\tilde{\sigma}_{jk}^2 = \sum_{i=1}^d (\tilde{y}_{ji} - \tilde{y}_{ki})^2 = \sum_{i=1}^d \lambda_i (v_{ij} - v_{ik})^2$$

4. to choose appropriate  $d$  we can consider

$$\left( \sum_{i=1}^d \lambda_i \right) / \left( \sum_{i=1}^m \lambda_i \right) = \text{proportion of variance explained by } d \text{ dimensions}$$

→ using PCA knowledge

### - performance evaluation

→ we want the stress of the MDS to be as little as possible

$$\text{Stress} := \sum_{i=2}^m \sum_{j < i} (d_{ij} - \tilde{\sigma}_{ij})^2 \quad \dots \text{ basically } \begin{matrix} 1 & 2 & 3 & 4 \\ 2 & \cdot & & \\ 3 & \cdot & \cdot & \\ 4 & \cdot & \cdot & \cdot \\ 5 & \cdot & \cdot & \cdot \end{matrix}$$

👁 it is better to consider the relative error

$$\frac{(d_{ij} - \tilde{\sigma}_{ij})^2}{d_{ij}} \Rightarrow \text{small } d_{ij} \text{ want better accuracy of } \tilde{\sigma}_{ij} \quad d=5, \tilde{\sigma}=5$$

! not the same →  $d=50, \tilde{\sigma}=48$

$$\text{Samon Stress} := \left( \sum_{i \neq j} \frac{(d_{ij} - \tilde{\sigma}_{ij})^2}{d_{ij}} \right) / \left( \sum_{i \neq j} d_{ij} \right)$$

## ② Metric least squares scaling

→ finds a configuration  $\tilde{y}_1, \dots, \tilde{y}_m \in \mathbb{R}^d$  which minimizes a loss function  $S$  e.g. stress or Samon-stress.

→ iterative numeric approach

👁️ classical MDS uses the Euclidian distance model and minimizes the stress value

## ③ Non-metric multidimensional scaling

$$\text{Kruskal Stress} := \left( \sum_{i \neq j} (f(d_{ij}) - \delta_{ij})^2 \right) / \left( \sum_{i \neq j} \delta_{ij}^2 \right)$$

→ it isn't trying to achieve  $\frac{d_{ij}}{\delta_{ij}} \approx 1$

→ it only seeks to preserve the rank order of the distances

$$d_{ij} < d_{kl} \Rightarrow \delta_{ij} < \delta_{kl} \quad \forall i, j, k, l$$

→ iterative numeric algorithm again

→ good for non-metric  $\approx$  non-geometric data → ordinal data

→ to choose  $d$  we can plot the stress vs.  $d$  and look for an elbow

## • Procrustes Analysis

👁️ if we take a MDS configuration and rotate/translate/reflect it, all of the distances remain unchanged  $\Rightarrow$  it's an equivalent solution

$\Rightarrow$  Say two MDS methods have been applied to a set of  $m$  points  $\in \mathbb{R}^m$  resulting in coordinate matrices  $X \in \mathbb{R}^{m \times d}$  and  $Y \in \mathbb{R}^{m \times d}$

↳ we want to know if they are equivalent

$\Rightarrow$  we want to match the  $i^{\text{th}}$  point in  $X$  to the  $i^{\text{th}}$  point in  $Y$

$$\Rightarrow \text{minimize } R^2 := \sum_{i=1}^m (\tilde{y}_i - \tilde{x}_i)^T (\tilde{y}_i - \tilde{x}_i) = \sum_{i=1}^m \sum_{j=1}^d (y_{ij} - x_{ij})^2$$

$\Rightarrow$  we will keep  $Y$  fixed (reference configuration) and transform  $X$  by rotating, translating and reflecting it to minimize  $R^2$

→ we will also allow uniform scaling of the points in  $X$  - preserves distance ratios

→ the point  $\tilde{x}_i$  will be transformed to

$$\tilde{x}_i \mapsto \tilde{x}'_i = S A^T \tilde{x}_i + b$$

where

- $S \in \mathbb{R}$  ... scaling factor note:  $S \cdot \text{Id}$  is called the dilation matrix
- $A \in \mathbb{R}^{d \times d}$  ... orthogonal matrix causing rotation and reflection
- $b \in \mathbb{R}^d$  ... translation factor

→ new sum of squared distances is

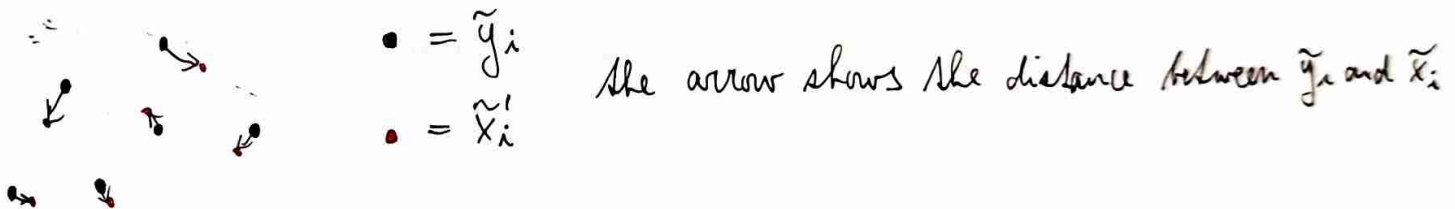
$$R^2 = \sum_{i=1}^m (\tilde{y}_i - S A^T \tilde{x}_i - b)^T (\tilde{y}_i - S A^T \tilde{x}_i - b)$$

→ by seeking the minimal  $R^2$  we can estimate the optimal  $S$ ,  $A$  and  $b$

⇒ Procrustes sum of squares := minimal  $R^2$

↳ measure of "match" between  $X$  and  $Y$

⇒ we can plot point-wise residuals between the reference configuration and the final transformed configuration



• t-SNE = Stochastic Neighborhood Embedding

- recent and popular approach for dimension reduction
- idea: distances can be converted into probabilities

→ similarity between  $x_i$  and  $x_j$  := probability  $p_{ji|i}$  that  $x_i$  would pick  $x_j$  as its neighbour if neighbours were picked in proportion to their probability density under a Gaussian centered at  $x_i$

$$p_{ji|i} := \frac{\exp\left(-\frac{1}{2} \left(\frac{\|x_i - x_j\|}{\sigma}\right)^2\right)}{\sum_{k \neq i} \exp\left(-\frac{1}{2} \left(\frac{\|x_i - x_k\|}{\sigma}\right)^2\right)}$$

Note:  $p_{ji|i} \neq p_{ij}$  ⇒ we use  $p_{ij} = \frac{1}{2}(p_{ji|i} + p_{ij})$



- we want to represent high-dim  $X$  with low-dim  $Y$
- we will represent the similarity between  $y_j$  and  $y_i$  using prob. density  $q_{ij}$
- assume that the similarities in the low-dim space are governed by a Student-t distribution with one degree of freedom, resulting in

$$q_{ij} = (1 + \|y_i - y_j\|)^{-2} / \sum_{z \neq i} (1 + \|y_i - y_z\|)^{-2} \quad \rightarrow \text{to find optimal } Y$$

- we want to minimize a loss function based on  $f_{ij}$  and  $q_{ij}$
- ⇒ the one which is used is

$$S = \sum_{i \neq j} f_{ij} \log\left(\frac{f_{ij}}{q_{ij}}\right) \quad \dots \quad \frac{f}{q} =: r \Rightarrow \begin{cases} r = 1 & \Rightarrow +0 \\ r < 0 & \Rightarrow -\epsilon \\ r > 0 & \Rightarrow +\epsilon \end{cases}$$

- $S$  can be minimized using an adaptive learning algorithm

! we should also specify the variance of the Gaussian  $\sigma$

- in practice, perplexity is specified instead

↳ it is a function of  $\sigma$  and is interpreted as a smooth measure of the effective number of neighbors considered

↳ it is typically set to be somewhere between 5 and 50

↳ results are sensitive to the choice and should be checked

# Logistic Regression

→ binary classification

→ we have some data and the answer is Yes/No

↳ predicting the presence/absence of a health condition

↳ assessing the likelihood of treatment success

↳ determining the likelihood of a customer to purchase a product

↳ spam detection

→ logistic regression gives us the probabilities of Yes and No

→ similar to LDA and QDA it is a parametric technique making distributional assumptions over the data

→ motivational example:

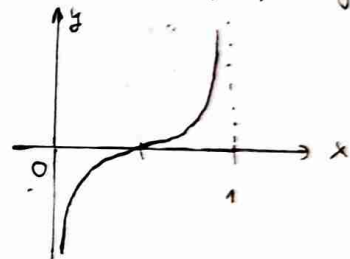
→ 92 subjects recorded: Resting Pulse (Low/High), Smokes (Yes/No), Weight (lb)

⇒ we want to classify whether the person has Low or High resting pulse

↳ we want  $P(\text{Low}) = f(\text{Smokes}, \text{Weight})$  where  $\text{Smokes} \in \{0, 1\}$ ,  $\text{Weight} \in \mathbb{R}^+$

Def: We define the logit function as

$$\text{logit}: [0, 1] \rightarrow \mathbb{R}, \quad x \mapsto \ln\left(\frac{x}{1-x}\right)$$



→ Consider the following model

$$\text{logit}(P(\text{Low})) = \alpha + \beta_S \cdot \text{Smokes} + \beta_W \cdot \text{Weight}$$

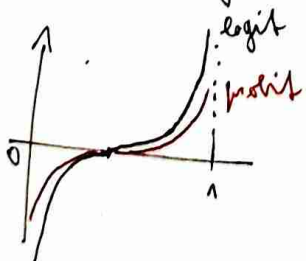
$$\Rightarrow \frac{P}{1-P} = \exp(\alpha + \beta_S \cdot S + \beta_W \cdot W) \Rightarrow P(\text{Low}) = \frac{\exp(\alpha + \beta_S \cdot S + \beta_W \cdot W)}{1 + \exp(\alpha + \beta_S \cdot S + \beta_W \cdot W)}$$

Note: Economists often use the probit model:

$$\Phi^{-1}(z) = w \quad \equiv \quad x \sim N(0,1), \text{ then } P(x < w) = z \quad \dots \quad \Phi = \text{cumulative dist. function of } N(0,1)$$

$$\rightarrow \Phi^{-1}(P(\text{Low})) = \alpha + \beta_S \cdot S + \beta_W \cdot W$$

! Probit grows slower than logit



⇒ if we fix  $\alpha + \beta_S \cdot S + \beta_W \cdot W =: Q$  and consider

$\text{logit}(P(\text{Low})) = \text{probit}(P(\text{Low})) = Q$ , then clearly

$P_{\text{logit}} < P_{\text{probit}} \Rightarrow$  Logit is more strict

↳ requires more evidence for high P

⇒ Consider  $n$  independent observations  $\tilde{X}_1, \dots, \tilde{X}_n$  where

$$\tilde{X}_i = (\text{High/Low}, \text{Smokes}_i, \text{Weight}_i)$$

↳  $y_i := 1$  if Low, 0 if High

⇒ let's say that the resting pulses are L, L, H, L, H, H, H, ... L)

↳ we can look at the probability

$$P(\tilde{X}_1, \dots, \tilde{X}_n) = P(L, L, H, L, H, H, H, \dots, L) = P(L_1)P(L_2)P(H_3) \dots P(L_n) \quad , \text{note } P(H_i) = 1 - P(L_i)$$

$$= \prod_{i=1}^n \left( \frac{\exp(d + \beta_s S_i + \beta_w W_i)}{1 + \exp(d + \beta_s S_i + \beta_w W_i)} \right)^{y_i} \cdot \prod_{i=1}^n \left( \frac{1}{1 + \exp(d + \beta_s S_i + \beta_w W_i)} \right)^{1 - y_i}$$

→ we seek  $d, \beta_s, \beta_w$  that maximize this probability

- best  $d, \beta_s, \beta_w$  = maximum likelihood estimates
- best  $P(\tilde{X}_1, \dots, \tilde{X}_n)$  = likelihood

→ if we ask R to do this for us, we get

	Estimate	Std. dev.	z-value	Pr(> z )	
$d$	-1.99	1.68	-1.18	0.24	
$\beta_s$	-1.19	0.55	-2.16	0.03	*
$\beta_w$	0.025	0.012	2.04	0.04	*

} statistically significant

→ we also need to consider the std. dev. of the estimates

$$\underline{z\text{-value}} = \frac{\text{estimate}}{\text{std. dev.}} = \text{relative std. dev.} \Rightarrow \text{larger } |z| = \text{more significant}$$

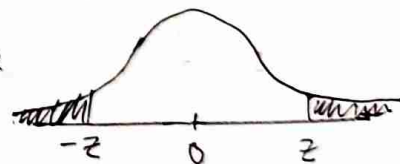
→ so say if the estimate has statistical value, we consider the test

$$H_0: \text{real value} = 0 \quad \rightarrow \text{under } H_0: \frac{\text{estimate} - \mu}{\hat{\sigma}} \approx \frac{\text{estimate}}{\hat{\sigma}} = z\text{-value} \sim N(0,1)$$

$H_1: \text{real value} \neq 0$

⇒ we consider the Pr of  $z$  being this large

$$\Rightarrow \underline{\text{Pr}(>|z|)} = \text{shaded area} + \text{shaded area} = p\text{-value of the experiment}$$



↳ for  $\beta_s$  we have  $p = 0.03 < 0.05 \Rightarrow$  statistically significant

↳ for  $d$  we have  $p = 0.24 \Rightarrow$  we aren't very sure about the true  $d$

Interpretation:  $\text{logit}(P(\text{low})) = d + \beta_s \cdot \text{Smokes} + \beta_w \cdot \text{Weight}$



- Smokes = True  $\Rightarrow P(\text{low})$  smaller  $\Rightarrow$  higher pulse rate
- ↑ weight  $\Rightarrow P(\text{low})$  larger  $\Rightarrow$  lower pulse rate

→ the standard errors tend to decrease as sample size increases  
↳ we can use them to construct 95% confidence intervals for the estimates

$$95\% \text{ CI} = (\text{Estimate}) \pm 2(\text{Std err})$$

$$\Rightarrow \beta_s \in (-2.30, -0.09) \quad \text{and} \quad \beta_w \in (0.001, 0.05)$$

$$\Rightarrow \exp(\beta_s) \in (0.1, 0.9) \quad \text{and} \quad \exp(\beta_w) \in (1.00, 1.05)$$

### • Interactions

→ if the effect that weight has upon resting pulse would differ depending on whether or not the individual smoked, then we might consider the model

$$\text{logit}(P(\text{Low})) = \alpha + \beta_s \cdot \text{Smokes} + \beta_w \cdot \text{Weight} + \beta_{sw} \cdot \text{Smokes} \cdot \text{Weight} \quad \text{interaction}$$

→ when appropriate, interactions can greatly increase the model's performance

### • Akaike's information criterion

→ choosing a model ~ balancing two opposite goals

- model fit ... how good the  $P(\tilde{x}_1, \dots, \tilde{x}_n) = \text{likelihood}$  is
- model complexity ...  $p := \# \text{ parameters}$

$$\Rightarrow \text{AIC} := -2 \log(\text{Likelihood}) + 2p$$

↳ the model with minimal AIC is the best compromise

### • Logistic regression × LDA

→ logistic regression can be used to classify data to two clusters  $C_1$  and  $C_2$

$$\text{logit}(P(\tilde{x} \in C_1 | \tilde{x})) = \frac{P(\tilde{x} \in C_1 | \tilde{x})}{P(\tilde{x} \in C_2 | \tilde{x})} = \alpha + \beta^T \tilde{x}$$

→ compare this to LDA

$$\frac{P(\tilde{x} \in C_1 | \tilde{x})}{P(\tilde{x} \in C_2 | \tilde{x})} = \ln \frac{\pi_1}{\pi_2} + \ln \frac{f_1(\tilde{x})}{f_2(\tilde{x})} = 0$$

→ the models have the same form

## • Deviance

- another measure for determining a model's quality

→ consider this: data points  $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^d$

$$\tilde{x}_j = (x_1, x_2, \dots, x_{d-1}, y)$$

covariate vector  $\hookrightarrow$  classification  $\in \{0, 1\}$

→ let  $M := \#$  distinct covariate vectors

$m_i := \#$  data points with covariate vector  $i$

$r_i := \#$  data points with covariate vector  $i$  & assigned to class 1

→ we want to create a full (saturated) model with  $M$  parameters

• if  $M = m$ , then we can simply let  $j^{\text{th}}$  parameter  $I_j$  to be an indicator for  $\tilde{x}_j$

• if  $M < m$ , let  $\pi_i := P[\text{class 1} \mid \text{covariate vector } i]$

$\hookrightarrow$  based on observed data  $\pi_i = \frac{r_i}{m_i}$

$\Rightarrow$  then the likelihood of observed data is

$$L[\text{Data} \mid \pi_1, \dots, \pi_M] \propto \prod_{i=1}^M \pi_i^{r_i} (1 - \pi_i)^{m_i - r_i}$$

$$l[\text{Data} \mid \pi_1, \dots, \pi_M] = \text{const} + r_i \log(\pi_i) + (m_i - r_i) \log(1 - \pi_i)$$

log likelihood

→  $l$  is maximized when  $\pi_i = \frac{r_i}{m_i}$ , (assign  $0 \cdot \log 0 = 0$  if  $r_i = 0$ )

$\Rightarrow$  let  $L_{\max}$  be the max. value of  $L$

$\hookrightarrow$  this represents the best likelihood under any model

→ assume we prefer a simpler model

$\Rightarrow$  let  $L_{\text{mod}}$  be the max. likelihood for this model

Def: The deviance for this model is  $Dev = 2[\log(L_{\max}) - \log(L_{\text{mod}})]$

$\Rightarrow$  the smaller the deviance, the better the model

→ it is also sometimes called the residual deviance

Def: Now consider the simplest model with just 1 parameter based on the number of data points assigned to class 1 ...  $\pi = \frac{\# \text{ in } 1}{n}$

$\hookrightarrow$  the deviance for this model is the null deviance  $= 2[\log(L_{\max}) - \log(L_{\text{null}})]$

Note:  $L_{\text{null}} \leq L_{\text{mod}} \leq L_{\max}$

→ let's assume  $H_0$ : The proposed model is true

↳ it can be shown that under  $H_0$ , the deviance is approximately distributed as a  $\chi^2_{M-k-1}$  distribution where  $k = \#$  of parameters of the model

⇒ if the deviance is greater than the 95% point of the appropriate chi-squared distribution, then  $H_0$  probably isn't true and the model is false  
⇒ we can not assign a prob. to  $H_0$  holding true, but we can question it

How to do this in R?

calculate:  $1 - pchisq(\text{residual-deviance}, \text{deg-of-freedom})$

↳ this gives us the p-value of the experiment

→ when this value is very small ( $\leq 5\%$ ) then there is statistical evidence that there is a significant diff. between the model of interest and the saturated model

→ we might also want to check against the null deviance to see if the model is better than nothing

→ the smaller the  $M$  and larger  $n$ , the more trustworthy this test is