

Chapter 9 : Support vector Machines

Maximal Margin Classifiers

separable cases
only

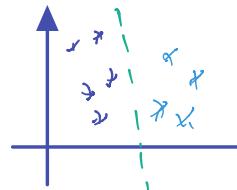
hyperplane: subspace of $p-1$ dim in p dims
flat and affine (doesn't need to go through the origin)

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p = 0$$

divides space in half, if not "on the plane", then above (>0)
or below (<0)

Assume the data $x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{pmatrix}$, x_n is separable by a hyperplane

e.g. $\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p > 0$ if $y_i = 1$
 < 0 if $y_i = -1$

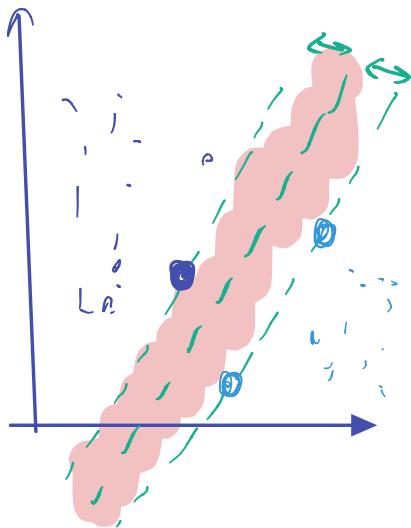


Classify based on sign e.g. $\text{sgn}(f(x^*)) = y^*$

distance e.g. $|f(x^*)|$ tells us about the confidence

The best hyperplane is the one that maximizes
the distance of the observations from the hyperplane
→ maximal margin hyperplane

- This is also the middle of the largest "slab" you can insert between the data
- This only depends on a few points exactly on the margin, not the rest → you could move the rest and the margin wouldn't move



only depends on these three points → the support vectors

Ridgeless

$$\text{minimize } M$$

$$\text{s.t. } \sum \beta_j^2 = 1$$

$$y_i: (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) > M \quad \forall i$$

→ correct side of the hyperplane + padding



this is not really a constraint, but signifies like:

$y_i(\beta_0 + \beta_p x_{-p})$ will be the distance from the hyperplane
to the i^{th} observation

so this means each point is at least a distance M from the hyperplane,
on the correct side

Support Vector Classifiers - Non-separable cases

Separating hyperplanes are very sensitive to individual observations

- 0% error on training
- one observation can move the plane by a lot
- overfitting

We might want

- more robust to individual obs
- better classification to most of the training obs.

e.g. misclassify a few to do better on the rest

"soft margin classifier"

↳ margin can be violated by some obs.

problem: width of margin

maximize M

$$\text{s.t. } \sum \beta_j^2 = 1$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \varepsilon_i)$$

$$\varepsilon_i \geq 0, \sum \varepsilon_i \leq C$$

slack variable

allowing obs. to be on the wrong side

C : budget of how wrong observation can be classified

$\varepsilon_i = 0$ if on correct side of the margin

if C grows, margin grows
if $C=0$, separating hyperplane

$\varepsilon_i > 0$ violates the margin

$\varepsilon > 1$ of many side of hyperplane

C fit w/ $C-V$

C larger: more bias, less variance

C smaller: more variance

observations on the margin and on the way will affect the classifier, the rest do not

↓
support vectors

(e.g., may support vectors)

For non-linear boundaries, you can find only poly

$$\beta_0 + \sum \beta_{j1} x_{ij} + \sum \beta_{j2} x_{ij}^2 + \dots + \sum \beta_{jp} x_{ij}^p \quad \text{etc.}$$

Use SVM instead:

→ more computationally efficient

Support vector Machines

use Kernel for non-linear boundaries

it can be shown that $f(x)$ for support vector classifiers is of the form

$$f(x) = \beta_0 + \sum \alpha_i (x_i \cdot x) \quad \text{where } \alpha_i \neq 0 \text{ iff } x_i \text{ is a support vector}$$

for an SVM, you use a kernel for the inner product

e.g. $f(x) = \beta_0 + \sum \alpha_i k(x, x_i)$

$$k(x_i, x_j) = \sum x_{ij} \alpha_i x_j \rightarrow \text{just the same linear}$$

$$k(x_i, x_{i'}) = (1 + \sum x_{ij} \cdot x_{i'j})^d \rightarrow \text{polynomial}$$

$$k(x_i, x_{i'}) = \exp(-\gamma \sum (x_{ij} - x_{i'j})^2) \quad \text{radial} \quad \gamma > 0$$

More than two classes:

1-vs-1: Train $\binom{k}{2}$ SVM, each copies a pair of classes

for a test observation, we count up how many
each class got in total across the $\binom{k}{2}$, and put it
in the most frequent one

1-vs-all: Train k SVM, each copies one class w/ the
remaining $k-1$

let $\beta_{0k}, \beta_{1k}, \dots, \beta_{pk}$ the params from fitting
an SVM copying the k th class to the rest

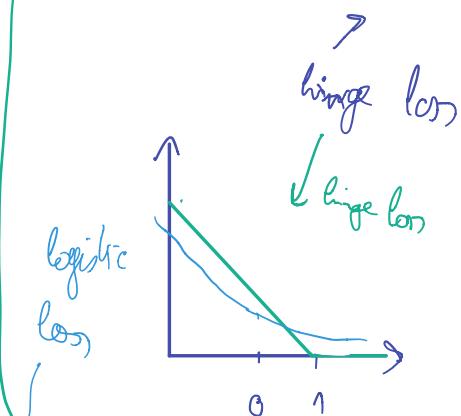
we assign the test observation to the class that has the
highest $\beta_{0k} + \beta_{1k}x_1 + \dots + \beta_{pk}x_p$, as this has best confidence

Relationship w/ logistic regression:

We can rewrite the SVM problem as

$$\text{minimize} \quad \sum_i \max \left[0, 1 - y_i f(x_i) \right] + \lambda \sum_j \beta_j^2$$

$\beta_j^2 / \text{ridge penalty}$



\uparrow hinge loss
 \downarrow logistic loss
 \uparrow hinge loss
 \downarrow logistic loss

↑ tuning param
 small $\lambda \approx$ small C
 when λ small, fewer violations

in this form, margin is at value 1, with width $\sum_j \beta_j^2$
 observations not on the margin are not affected b.c. hinge loss
 gives 0 loss for them no matter how far they are from the
 margin because $y_i (\beta_0 + \dots + \beta_p x_{ip}) \geq 1$

in contrast, this is not true for logistic regression - the loss is not 0,
 but very small

so SVM and logistic regression give very similar results

- when classes are separated, SVM tends to be better.
- when lots of overlap, log. reg. is often better

Chapter 10 - Unsupervised learning

PCA - discussed in chapter 6 for PCR

direction along which data is most variable

$$z_1 = \phi_{11} x_1 + \phi_{21} x_2 + \dots \phi_{p1} x_p \quad \sum \phi_{j1}^2 = 1$$

$\phi_{11}, \dots, \phi_{p1}$ → loadings of 1st PC

how to construct ϕ ?

- center X to have mean 0 and scale it

$$\underset{\phi}{\text{maximize}} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 \right\} \text{ s.t. } \sum \phi_{j1}^2 = 1$$

e.g. maximize variance of z_1 (since z_1 is mean 0)

$$\text{if } z_{i1} = \phi_{11}x_{i1} + \phi_{21}x_{i2} + \dots + \phi_{p1}x_{ip}$$

$$\text{then maximize } \frac{1}{n} \sum_i z_{i1}^2 \quad \text{s.t. } \sum_j \phi_{j1}^2 = 1$$

can be solved by eigen decomposition

$\phi_{11} - \phi_{p1} \Rightarrow$ defines the direction where data varies the most

if we project the n points x_1, \dots, x_n onto this direction, the values are the scores z_{11}, \dots, z_{n1}

For next component z_2 , w/ $z_{i2} = \phi_{12}x_{i1} + \dots + \phi_{p2}x_{ip}$, we need $\phi_2 \perp \phi_1$

We get the components by taking the eigenvectors of $X^T X$, and the variances of the components as the eigenvalues

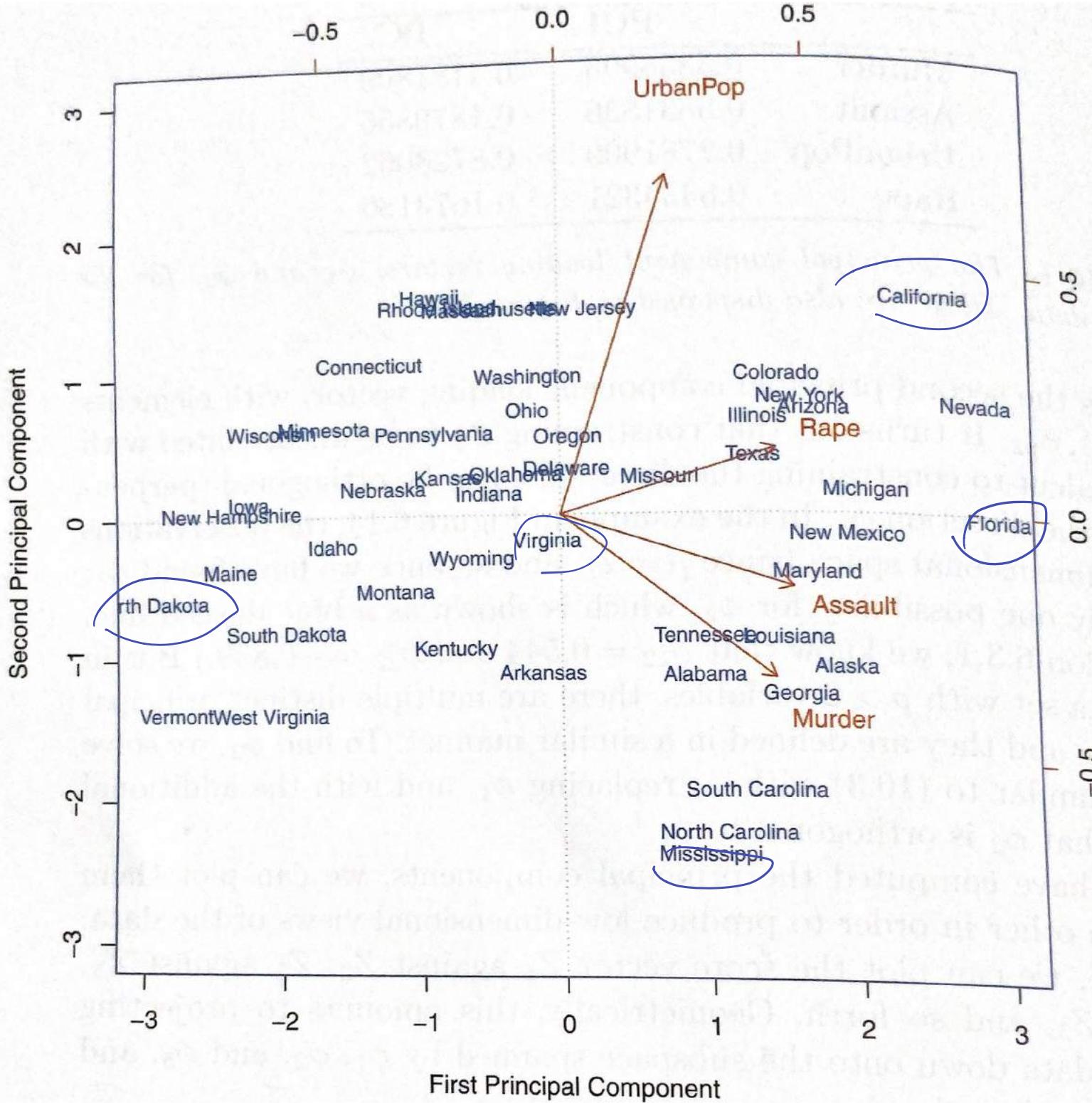
There are at most $\min(n-1, p)$ PCs

$z_{11}, z_{21}, \dots, z_{n1} \rightarrow$ score of PC

with $\frac{1}{n} \sum z_{i1} = 0$

(s.t. $\sum_i z_{ij} = 0$)

↓
to first maximize sample variance of scores



Biplot: 4 variables have projected onto PC component axes
variables close to each other are correlated

e.g. Rape, Assault and Murder

Florida: High crime rate

North Dakota: low crime rate

California: High crime rate + urbanization

Mississippi: low urbanization

Virginia: avg levels of both

Alternative interpretation: low dimensional linear surface that we closest
to the observations (\rightarrow projection)

1st PC: line that is closest to observations

1st + 2nd PC: plane that is closest

1st + 2nd + 3rd PC: 3rd hyperplane that is closest

we can get back x_{ij} for the PCs:

$$x_{ij} \approx \sum_{m=1}^M z_{im} \phi_{jm} \quad \text{for } M \text{ PCs}$$

where if $M = \min(n, p)$, this is exact: $x_{ij} = \sum_{m=1}^M z_{im} \phi_{jm}$
 (I think this is b.c. ϕ are orthogonal to each other.)

Scaling

- Important to scale variables to have std 1
- Or measure variables in same units if you can
 All of the method relies on size of variance in a direction,
 so if a variable is inflated, that will change everything

Proportion of variance explained (PVE)

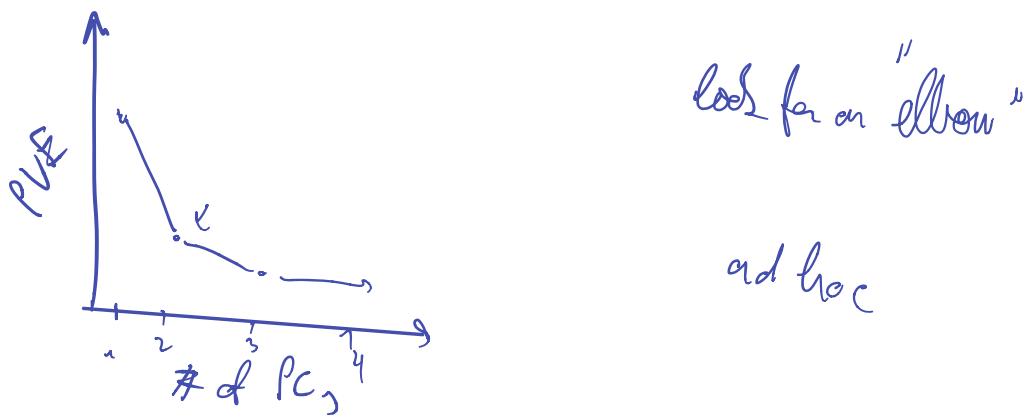
$$\text{Total variance in data} = \sum_{j=1}^p \frac{1}{n} \sum_{i=1}^n x_{ij}^2$$

$$\text{variance explained by } m\text{th PC} = \frac{1}{n} \sum_{i=1}^n \hat{z}_{im}^2 = \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{jm} x_{ij} \right)^2$$

$$\text{PVE} : \frac{\sum_{i=1}^n \left(\sum_{j=1}^p \phi_{jm} x_{ij} \right)^2}{\sum_{j=1}^p \sum_{i=1}^n x_{ij}^2}$$

But how many PC's to use?

Screen plot:



Clustering methods

cluster obs. on the basis of features to find groups in the observations

or reverse: cluster features on the basis of observations to group features

K-means clustering:

k clusters C_1, \dots, C_k

minimize variations within clusters

$$\min_{C_1, \dots, C_k} \left\{ \sum_{g=1}^k \frac{1}{|C_g|} \sum_{i,j \in C_g} \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})^2 \right\}$$

Alg:

1, Randomly assign a label 1...k to each observation

2, Iterate until assignments change:

- compute cluster centroid
- Assign each observation to the closest cluster centroid

This is a local optimum, because

$$\frac{1}{|C_g|} \sum_{i,j \in C_g} \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})^2 = 2 \sum_{i \in C_g} \sum_{j=1}^p (x_{ij} - \bar{x}_{gj})^2$$

and step 2 decreases this every step

Important to rerun from different starting points, leave local optimum

Hierarchical clustering

no need to specify K before the algo

look by eye where to cut the dendrogram

works better on hierarchical data

- Algo:
1. Begin w/ n obs + a matrix for all $\binom{n}{2}$ dissimilarities
 2. For $i=1, n-1, \dots, 2$:
 - (a) Examine all pairwise cluster dissimilarities and find least dissimilar clusters. Fuse them
 - (b) Compute the new dissimilarities

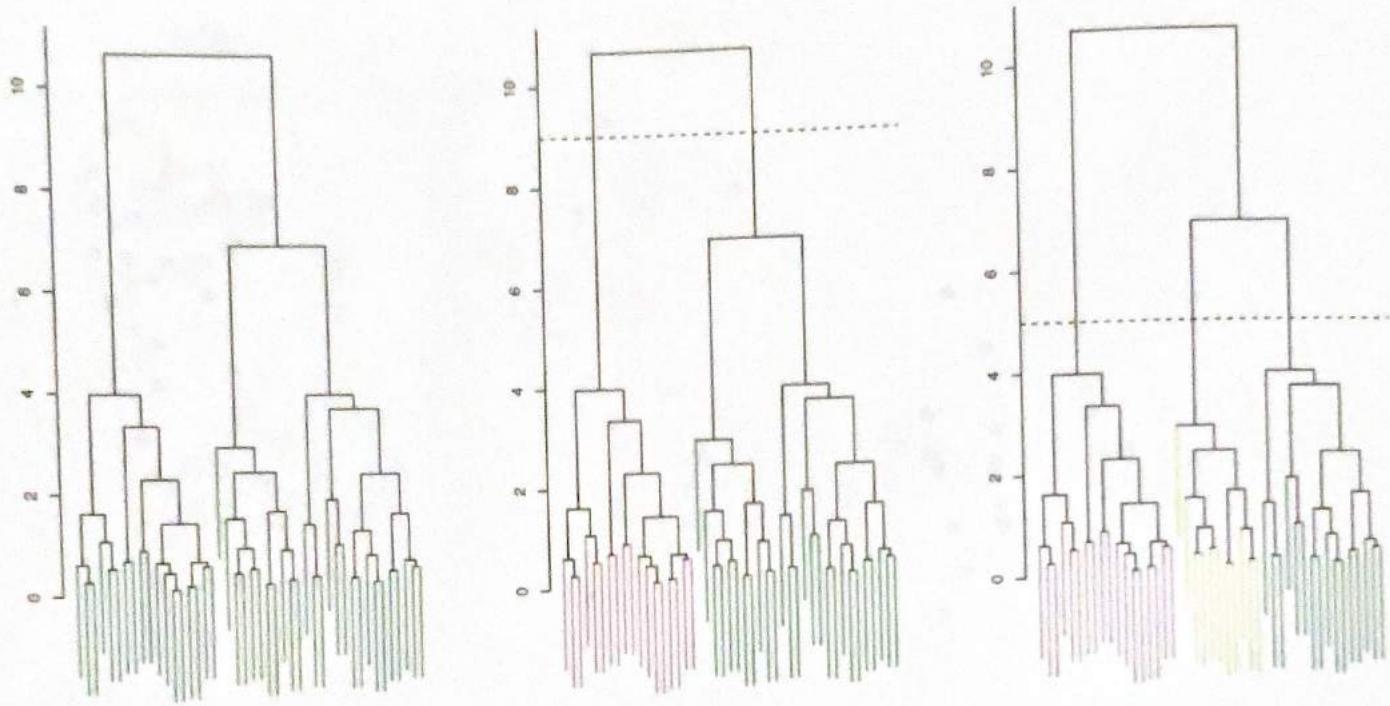


FIGURE 10.9. Left: *dendrogram obtained from hierarchically clustering the data from Figure 10.8 with complete linkage and Euclidean distance.* Center: *the dendrogram from the left-hand panel, cut at a height of nine (indicated by the dashed line).* This cut results in two distinct clusters, shown in different colors. Right: *the dendrogram from the left-hand panel, now cut at a height of five.* This cut results in three distinct clusters, shown in different colors. Note that the colors were not used in clustering, but are simply used for display purposes in this figure.

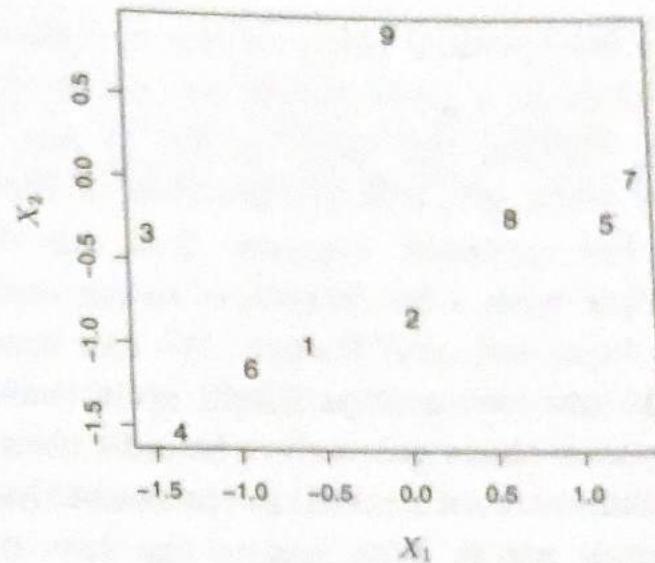
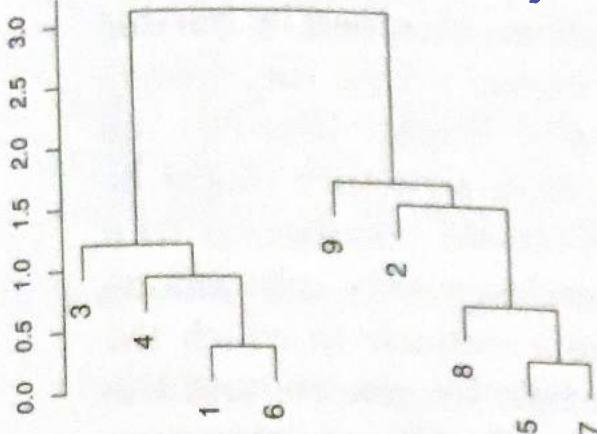


FIGURE 10.10. An illustration of how to properly interpret a dendrogram with nine observations in two-dimensional space. Left: a dendrogram generated using Euclidean distance and complete linkage. Observations 5 and 7 are quite similar to each other, as are observations 1 and 6. However, observation 9 is no more similar to observation 2 than it is to observations 8, 5, and 7, even though observations 9 and 2 are close together in terms of horizontal distance. This is because observations 2, 8, 5, and 7 all fuse with observation 9 at the same height, approximately 1.8. Right: the raw data used to generate the dendrogram can be used to confirm that indeed, observation 9 is no more similar to observation 2 than it is to observations 8, 5, and 7.

Linkage : dissimilarity between two group

- Can be euclidean distance of the observations (Table below)
 - Or correlation between the features for two observations
e.g. useful for retail shopping carts :
 - data is 1 if bought product, 0 o/w
 - euclidean would put everyone buying ↑ product close to each other
- Right need scaling too
to $\text{std} = 1$
- socks are bought more often than computers

Euclidean Measures:

Linkage	Description
Complete	Maximal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>largest</i> of these dissimilarities.
Single	Minimal intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>smallest</i> of these dissimilarities. Single linkage can result in extended, trailing clusters in which single observations are fused one-at-a-time.
Average	Mean intercluster dissimilarity. Compute all pairwise dissimilarities between the observations in cluster A and the observations in cluster B, and record the <i>average</i> of these dissimilarities.
Centroid	Dissimilarity between the centroid for cluster A (a mean vector of length p) and the centroid for cluster B. Centroid linkage can result in undesirable <i>inversions</i> .

Correlation-based measure

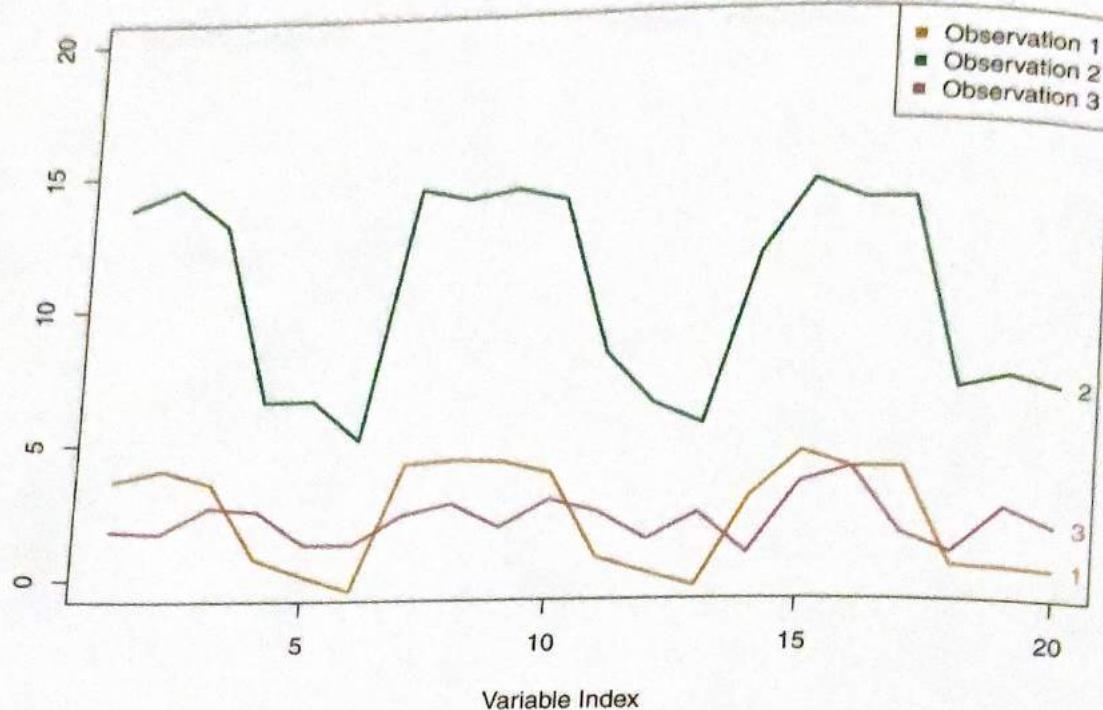


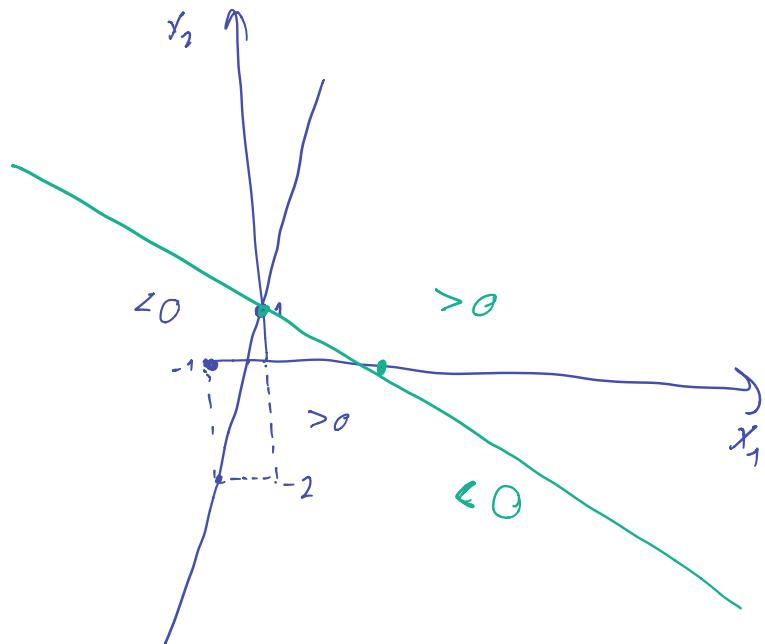
FIGURE 10.13. Three observations with measurements on 20 variables are shown. Observations 1 and 3 have similar values for each variable and so there is a small Euclidean distance between them. But they are very weakly correlated, so they have a large correlation-based distance. On the other hand, observations 1 and 2 have quite different values for each variable, and so there is a large Euclidean distance between them. But they are highly correlated, so there is a small correlation-based distance between them.

Practical issues for clustering:

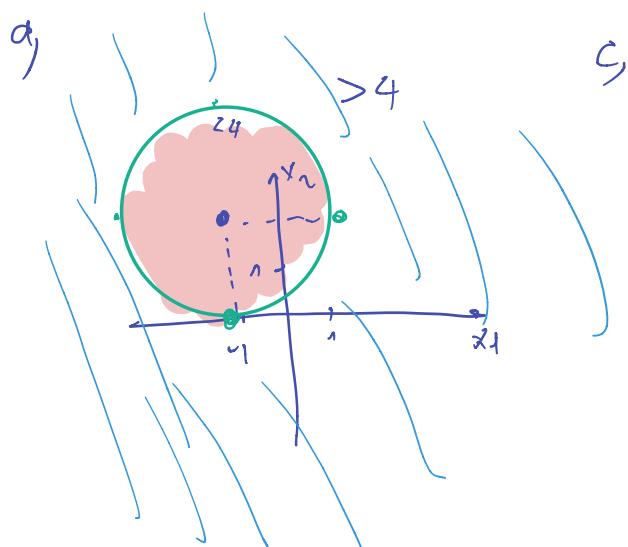
- Should I center and scale the variables?
 - For hierarchical, which dissimilarity measure?
 - Which type of linkage?
 - Where should we cut dendrogram?
 - For k-means, what k?
-
- ④ No consensus on how to validate clusters
 - ④ Not robust to perturbations

Exercises - chapter 9

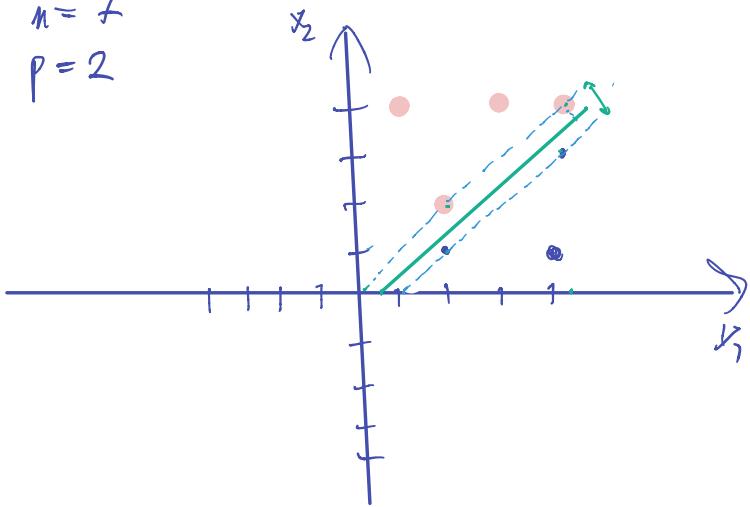
$$1) \quad 1 + 3x_1 - x_2 = 0 \quad -2 + x_1 + 2x_2 = 0$$



$$3) \quad (1+x_1)^2 + (2-x_2)^2 = 4$$



$$3) \quad n=7 \\ p=2$$



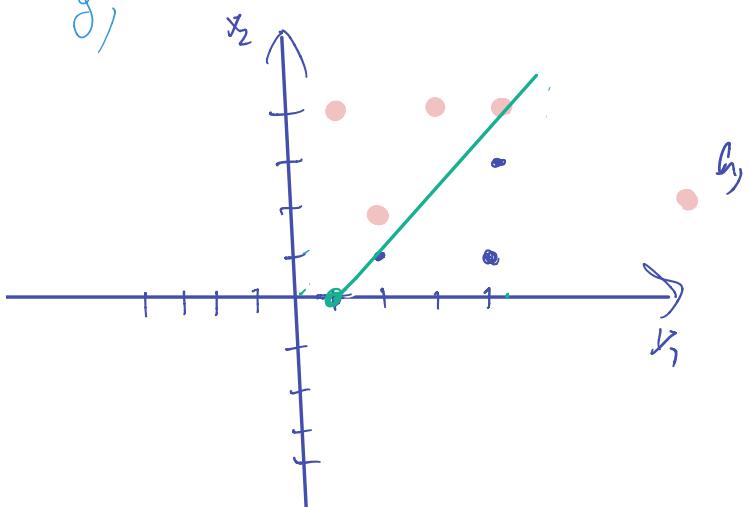
$$x_1 - x_2 = b_2$$

$$o > \bullet$$

$$b_2 + x_2 - x_1 = 0$$

$$o < \bullet$$

Q)



Equation of the line $(3\frac{5}{4})$ and $(4, \frac{15}{4})$

$$x_2 - \frac{3}{4}x_1 + \frac{5}{4} = 0$$

Exercises chapter 10

$$a) \frac{1}{|C_k|} \sum_{i,j \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})^2$$

$$LHS = \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^p (x_{ij}^2 + \bar{x}_{ij}^2 - 2x_{ij}\bar{x}_{ij})$$

$$= \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^p \left(2|C_k|\bar{x}_{ij} - 2x_{ij} \sum_{i' \in C_k} x_{i'j} \right)$$

$$= \frac{2 \sum_{i \in C_k} \sum_j (x_{ij}^2 - x_{ij} \sum_{i' \in C_k} \bar{x}_{i'j})}{|C_k|}$$

$$= \frac{2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij}^2 - x_{ij} \bar{x}_{kj})}{|C_k|}$$

$$= \frac{2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij}^2 + \bar{x}_{kj}^2 - 2 \underbrace{x_{ij}\bar{x}_{kj}}_{|C_k|})}{|C_k|}$$

$$= \frac{2 \sum_i \sum_j (x_{ij} - \bar{x}_{ij})^2}{|C_k|}$$

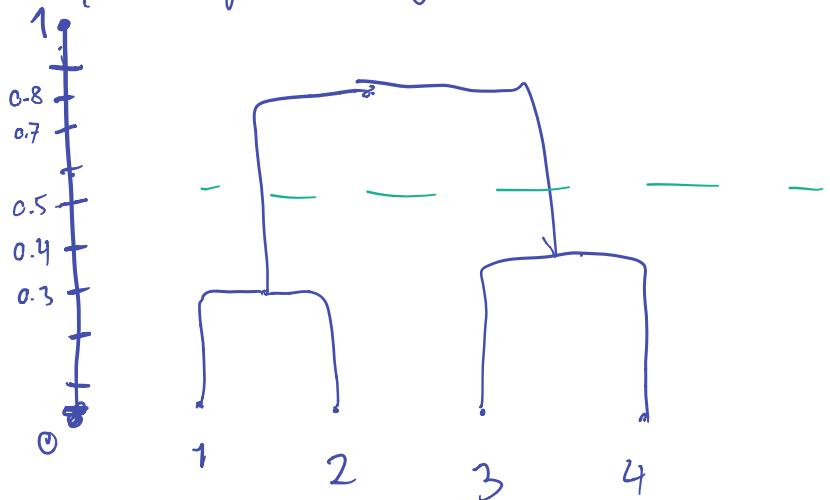
$$\frac{\sum_i x_{ij} \bar{x}_{ij}}{|C_k|} = \overline{x_{ij}}^2$$

2

-	0.3	0.4	0.7
0.3	-	0.5	0.8
0.4	0.5	-	0.45
0.7	0.8	0.45	-

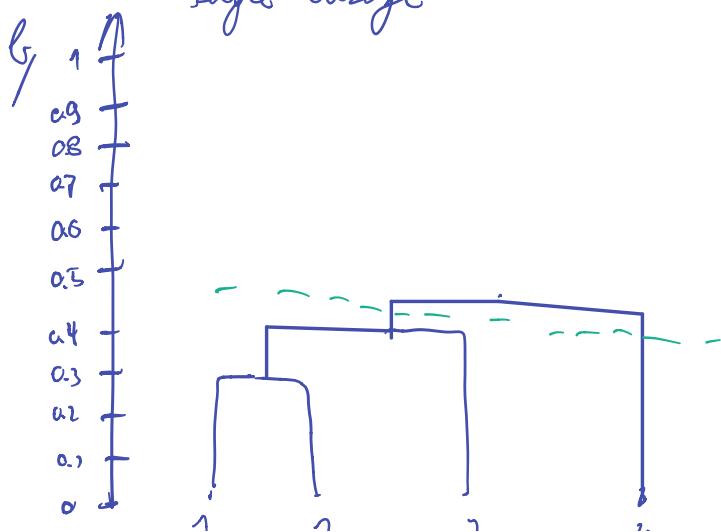
Dendogram:

a) complete linkage

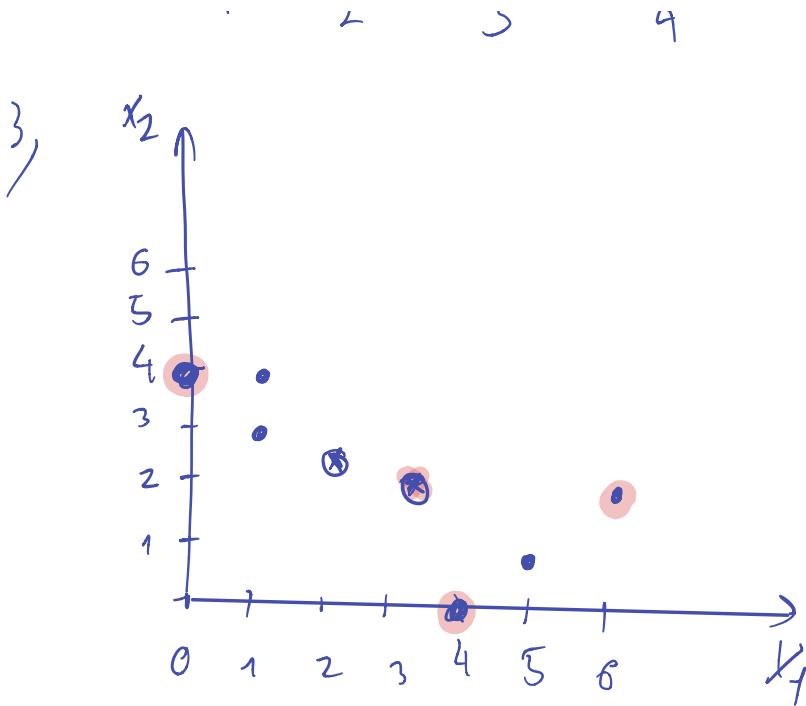


↪ (1,2), (3,4)

b) single linkage



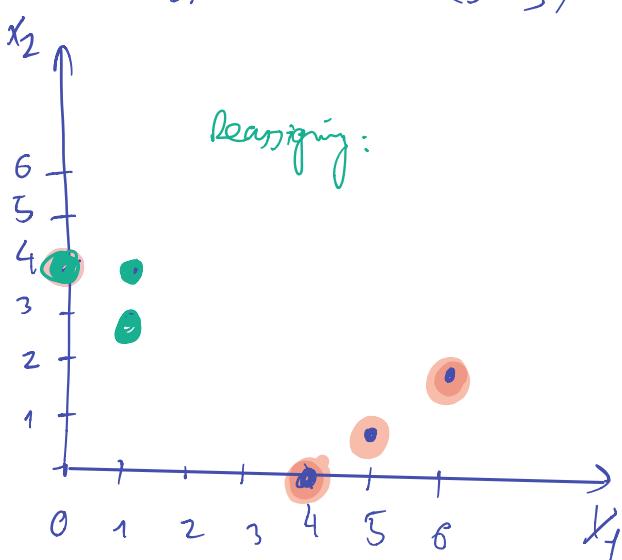
↪ (1,2,3) (4)



Red: $(0, 4)$
 $(4, 0)$
 $(6, 2)$

Blue: $(1, 4)$
 $(1, 3)$
 $(5, 1)$

centroid = $(\frac{10}{3}, 2)$ $(\frac{7}{3}, \frac{8}{3})$



4) a) single layer

b) same height

5)