# K means clustering

1) assign each obs. randomly to each of the K classes, but so that none of the clusters is empty
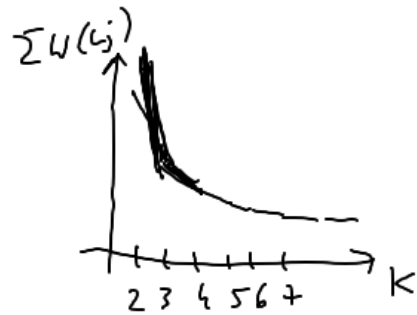
2) $\bar{x}_k$ — centroid $\quad \bar{x}_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i \in \mathbb{R}^p$

one can prove that $\sum_{j=1}^{k} W(C_j)$ decreases in each iteration of step 2)

- K is given in advance

- The above algorithm may be repeated multiple times to find a better solution
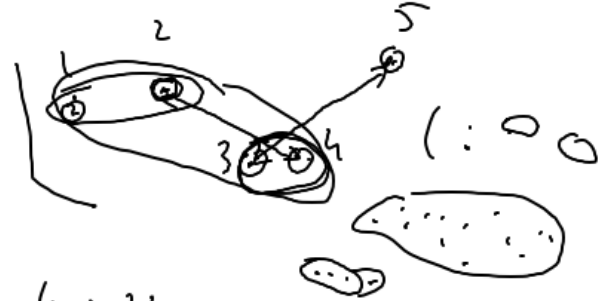
# Hierachical clustering $\quad X_1, \dots, X_n \in \mathbb{R}^p$

1) begin with $n$ observations and a measure (e.g. Euclidean distance) of all $\binom{n}{2}$ pairwise dissimilarities

Treat each observation as its own cluster, so we start with $n$ clusters, each having 1 observation

2) for $i = n, n-1, \dots, 2$ :

(a) Examine all pairwise inter-cluster dissimilarities among $i$ clusters and identify the pair of ~~the most~~ most similar clusters. Fuse these two clusters.

The dissimilarity between these two clusters will give the height at the dendrogram at which fusion should be placed

(b) compute new pairwise inter-cluster dissimilarities between remaining $(i-1)$ clusters
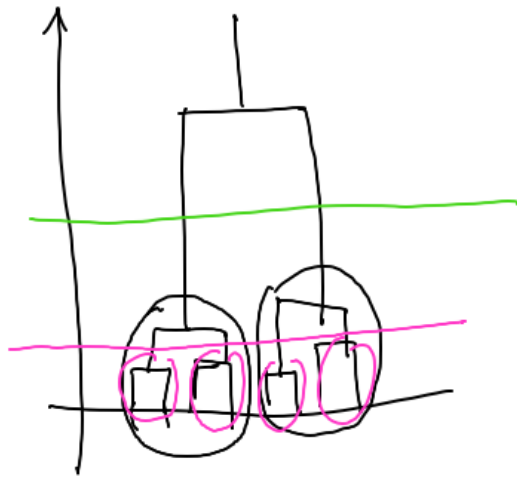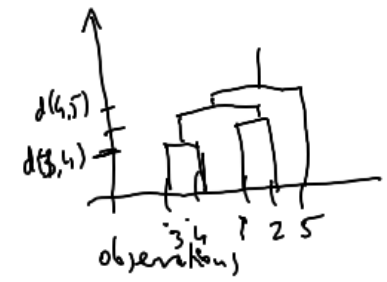
dissimilarities (linkage)

complete $\quad \max_{\substack{x \in K_1 \\ y \in K_2}} d(x,y)$

single $\quad \min_{x \in K_1, y \in K_2} d(x,y)$

average $\quad \dfrac{1}{|K_1||K_2|} \sum_{\substack{x \in K_1 \\ y \in K_2}} d(x,y)$

ward $\quad \sum_{x \in K_1, y \in K_2} d(x,y)$

$d(4,5)$

$d(3,4)$

observations

3 4 1 2 5

This algorithm gives us a sequence of clusterings, each pair of clusterings has the property that the clusters in one of them are subsets of the clusters in the other
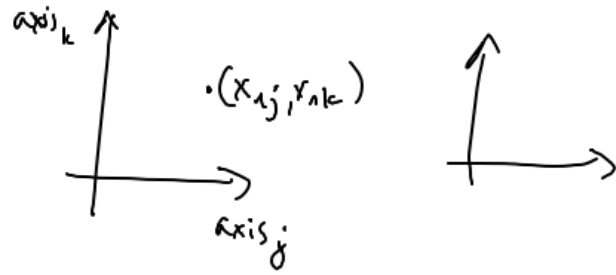
- it is deterministic

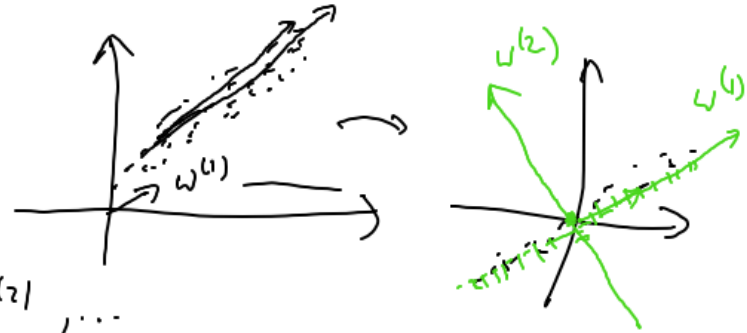# PCA principal component analysis

Used to reduce dimensionality

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ x_{21} & \cdots & x_{2p} \\ & \vdots & \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \in \mathbb{R}^{n \times p}$$

rows — $n$ observations, each in $\mathbb{R}^p$

$x_i = (x_{i1}, x_{i2}, \ldots, x_{ip}) \in \mathbb{R}^p$

$$w^{(1)} = \begin{bmatrix} w_1^{(1)} \\ \vdots \\ w_p^{(1)} \end{bmatrix}$$

orthonormal basis
$$\underbrace{w^{(1)}, w^{(2)}, \ldots, w^{(p)}}$$

We want to find a new coordinate system $w^{(1)}, w^{(2)}, \ldots, w^{(p)}$ such that the data varies the most along $w^{(1)}$, and then $w^{(2)}, \ldots$

It is convenient to assume that <u>each column of $X$ has mean zero.</u>

The first coordinate of our samples in the coordinate system will have the form:

$$\langle x_i, w^{(1)} \rangle = \sum_{j=1}^{p} x_{ij} \, w_j^{(1)} \qquad , i = 1, \ldots, n$$

We want that

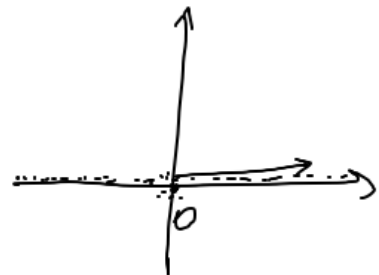$$\sum_i |\langle x_i, w^{(1)} \rangle|^2 = \sum_i \underbrace{\left( \sum_{j=1}^{p} x_{ij} w_j^{(1)} \right)}_{\langle x_i, w^{(1)} \rangle}{}^2 = \sum_i |(X \cdot w^{(1)})_i|^2 = \| X \cdot w^{(1)} \|^2 = \langle X \cdot w^{(1)}, X \cdot w^{(1)} \rangle =$$

$$= \langle w^{(1)}, \underbrace{X^T X \cdot w^{(1)}}_{\text{positive def. symmetric matrix}} \rangle$$

to be the largest possible

So $w^{(1)}$ should be the eigenvector of $X^TX$ corresponding to the largest eigenvalue

$w^{(2)}$ — $\mid\mid$ _____ the 2nd — $\mid\mid$ _____

and so on.

We can rewrite our data in the new coordinate system:

$$\tilde{X} = \begin{bmatrix} \langle x_1, w^{(1)} \rangle & \langle x_1, w^{(2)} \rangle & & \langle x_1, w^{(p)} \rangle \\ \langle x_2, w^{(1)} \rangle & \langle x_2, w^{(2)} \rangle & & \vdots \\ \vdots & \vdots & \cdots & \\ \langle x_n, w^{(1)} \rangle & \langle x_n, w^{(2)} \rangle & & \langle x_n, w^{(p)} \rangle \end{bmatrix}$$
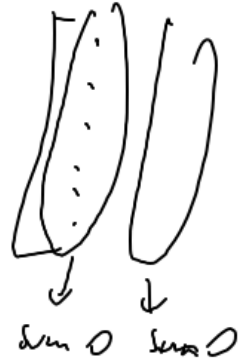
$\in \mathbb{R}^p$

↑
this column carries
the most information
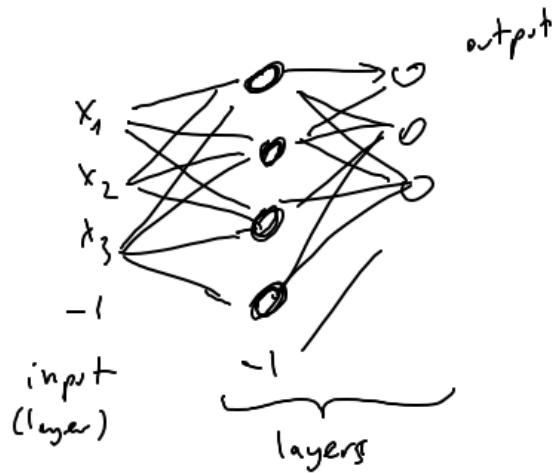about X

$\cdots$ and so on

Above it was important to make the data so that $\sum_i x_{ij} = 0$, $j = 1, 2, \dots, p$

Sometimes it is also desirable to normalise the data, that is rescale each column so

that $\sum_i x_{ij}^2 = 1$.

# Neural nets



input (layer)

layers

output

$x_1$

$x_2$

$x_3$

$-1$

$-1$

For image/pattern recognition dense layers are not very good for several reasons:



28

28

$\rightarrow$

100

dense layer treats such an image 784 as a vector of $28^2$ numbers
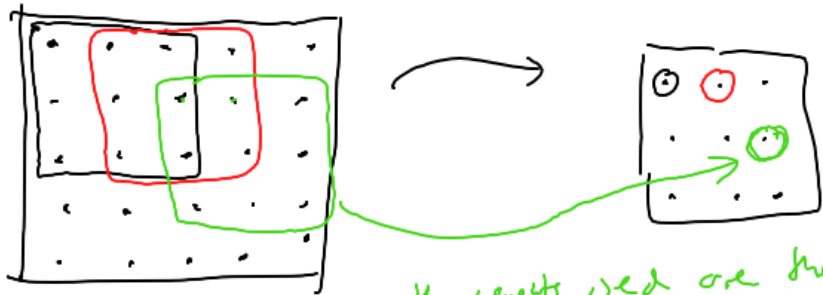
1) the geometrical structure is lost

2) the image after a slight shift looks completely different for the dense layer

3) often one has a huge number of weights   e.g. $28^2$ - input $\Big\}$ $\Rightarrow$ 78'500 weights in the first layer

100 - neurons

$200 \times 200 = 40k$

# Convolutional neural networks

they exploit the geometrical structure



the weights used are the same
for each output
just 10 weights here (9 inputs + 1 bias)

25 inputs $\longrightarrow$ 9 outputs

25 $\longrightarrow$ 9

9·26 weights for dense