

1.8.1

PCA

- consider 1-D data, n points
 - I imagine we want to characterize the data by a single value
- could take the mean & minimize distance of data set to representative point (would be better to take the median)

Now, imagine d dimension → we can take the mean vector as the representative $\frac{1}{n} \sum_{k=1}^n x_k$

note: the mean vector minimizes the sum of squared distances between a single point and the points in the data set.

note: mean vector is a 0-dimensional representation of the data set

Q: what would a 1-D representation be?

A: a line projected through the mean vector

Let e be a unit vector in the direction of our line

$x = m + ae$ is the equation of the line

distance of x from m

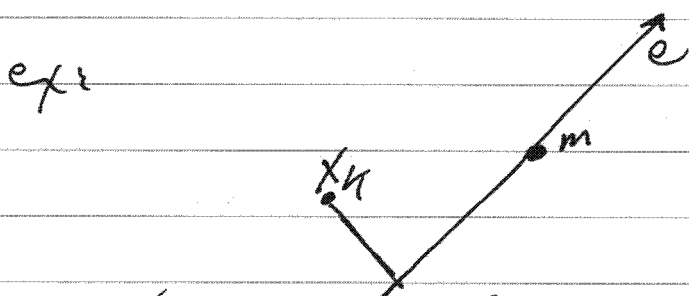


$$x_k = m + a_k e$$

note: optimal coefficients for our data set can be derived by minimizing the squared error criteria function

$$J_1(a_1, \dots, a_n, e) = \sum_{k=1}^n \| (m + a_k e) - x_k \|^2$$

In eqs 82 & 83 text show $a_k = e^t (x_k - m)$
 i.e. the ~~best~~ least squares soln is a projection of x_k onto the line in the direction of e that passes thru m .



→ the interesting problem is to find the best direction e for the data set.

Ex: consider 3-D data, a cloud of points in the shape of a football.

Q: What's the best direction e for our football cloud of points?

A: a line through to 2 pointy ends

obs: this is the eigenvector of the largest eigenvalue of the scatter matrix

$$S = \sum_{k=1}^n (x_k - m)(x_k - m)^t$$

We can extend the idea of a linear projection to d' dimension projection

recall: for linear projection we have $x = m + ae$

\Rightarrow for d' projection $x = m + \sum_{i=1}^{d'} a_i e_i$ $d' \leq d$

obs: criteria fun $J_{d'} = \sum_{k=1}^n \left\| \left(m + \sum_{i=1}^{d'} a_{ki} e_i \right) - x_k \right\|^2$

is minimized when e_i are the d' ^{largest} eigenvectors of S
_↑
 minimized

obs: these eigenvectors are orthogonal

\rightarrow these form a natural basis set for the points x

Q: So what are the "principal components"?

A: the coefficients a_i in $x = m + \sum_{i=1}^{d'} a_i e_i$

practical issue: PCA searches for directions in which the scatter cloud S is greatest
 \rightarrow it is accounting for variance

!! noise features can be a major problem

1/31/06

①

8.2 Fisher Linear Discriminant

PCA allows us to project onto fewer dimensions
 → these dimensions best represent the data

Q: What if we wanted to project onto dimensions that best discriminate between classes?

→ ~~the~~ In the extreme case project onto 1 dimension

goal: find the best line^o that classes are separate

$$D = \{x_1, \dots, x_n\} \quad \begin{array}{l} \text{labelled } w_1 \\ n_1 \text{ in } D_1 \end{array} \quad \begin{array}{l} \text{labelled } w_2 \\ n_2 \text{ in } D_2 \end{array}$$

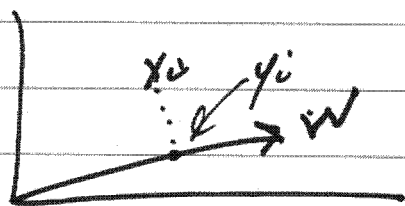
we get $y = w^T x$

\uparrow scalar \uparrow dot product

resulting in y_1, \dots, y_n in subsets y_1 & y_2

if $\|w\| = 1$, i.e., a unit vector

→ y_i is the projection of x_i onto a line in the direction w



goal: we want a w st. x_i in ω_1 cluster separately from x_i in ω_2

note: such a w may not exist if the distributions ^{are} ~~multimodal~~ overlap

Q: How do we find the best w ?

~~This~~ this is the Fisher linear discriminant problem

First, find the sample mean $m_i = \frac{1}{n_i} \sum_{x \in D_i} x$

find the ~~sample mean of projected points~~ ^{sample mean of projected points} $\tilde{m}_i = \frac{1}{n_i} \sum_{y \in \tilde{D}_i} y$

the $\tilde{m}_i = \frac{1}{n_i} \sum_{x \in D_i} w^t x = w^t m_i$

i.e. the "sample mean of projected points" is the same as the "projection of the mean vector"

Q: What is the distance between the projected means?

A: $|\tilde{m}_1 - \tilde{m}_2| = |w^t(m_1 - m_2)|$

note: we want to make this large relative to the scatter of each class

scatter: define scatter of projected points $S_i^2 = \sum_{y \in \tilde{D}_i} (y - \tilde{m}_i)^2$

(3)

def $\tilde{S}_1^2 + \tilde{S}_2^2$ is the within-class scatter

goal: find w that maximizes $J(w) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$

First def. scatter matrices

$$S_i = \sum_{x \in D_i} (x - m_i)(x - m_i)^t$$

$$S_w = S_1 + S_2$$

called "within-class" scatter matrix

→ \tilde{S}_i^2 can be expressed in terms of S_i

recall $y = w^t x$, $\tilde{m}_i = w^t m_i$, $\tilde{S}_i^2 = \sum_{x \in D_i} (y - \tilde{m}_i)^2$

$$\Rightarrow \tilde{S}_i^2 = \sum_{x \in D_i} (w^t x - w^t m_i)^2$$

$$= \sum w^t (x - m_i)(x - m_i)^t w$$

$$= w^t S_i w$$

hence $\tilde{S}_1^2 + \tilde{S}_2^2 = w^t S_w w$

called "between-class" scatter matrix

$$S_B = (m_1 - m_2)(m_1 - m_2)^t$$

similarly $(\tilde{m}_1 - \tilde{m}_2)^2 = w^t S_B w$ (from eq 101)

$$\Rightarrow J(w) = \frac{w^t S_B w}{w^t S_w w}$$

④

11 11

So now we know the line we want to project onto

A: choose point w_0 where posteriors are equal

i.e. $p(w_1 | \#) = p(w_2 | \#)$