## 3.8.1 Principal Component Analysis (PCA)

We begin by considering the problem of representing all of the vectors in a set of $n$ $d$-dimensional samples $\mathbf{x}_1, \ldots, \mathbf{x}_n$ by a single vector $\mathbf{x}_0$. To be more specific, suppose that we want to find a vector $\mathbf{x}_0$ such that the sum of the squared distances between $\mathbf{x}_0$ and the various $\mathbf{x}_k$ is as small as possible. We define the squared-error criterion function $J_0(\mathbf{x}_0)$ by

$$J_0(\mathbf{x}_0) = \sum_{k=1}^{n} ||\mathbf{x}_0 - \mathbf{x}_k||^2, \tag{78}$$

and seek the value of $\mathbf{x}_0$ that minimizes $J_0$. It is simple to show that the solution to this problem is given by $\mathbf{x}_0 = \mathbf{m}$, where $\mathbf{m}$ is the sample mean,

$$\mathbf{m} = \frac{1}{n} \sum_{k=1}^{n} \mathbf{x}_k. \tag{79}$$

This can be easily verified by writing

$$J_0(\mathbf{x}_0) = \sum_{k=1}^{n} ||(\mathbf{x}_0 - \mathbf{m}) - (\mathbf{x}_k - \mathbf{m})||^2$$

$$= \sum_{k=1}^{n} ||\mathbf{x}_0 - \mathbf{m}||^2 - 2 \sum_{k=1}^{n} (\mathbf{x}_0 - \mathbf{m})^t (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2$$

$$= \sum_{k=1}^{n} ||\mathbf{x}_0 - \mathbf{m}||^2 - 2(\mathbf{x}_0 - \mathbf{m})^t \sum_{k=1}^{n} (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2$$

$$= \sum_{k=1}^{n} ||\mathbf{x}_0 - \mathbf{m}||^2 + \underbrace{\sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2}_{independent\ of\ \mathbf{x}_0}. \tag{80}$$

Since the second sum is independent of $\mathbf{x}_0$, this expression is obviously minimized by the choice $\mathbf{x}_0 = \mathbf{m}$.

The sample mean is a zero-dimensional representation of the data set. It is simple, but it does not reveal any of the variability in the data. We can obtain a more interesting, one-dimensional representation by projecting the data onto a line running through the sample mean. Let $\mathbf{e}$ be a unit vector in the direction of the line. Then the equation of the line can be written as

$$\mathbf{x} = \mathbf{m} + a\mathbf{e}, \tag{81}$$

where the scalar $a$ (which takes on any real value) corresponds to the distance of any point $\mathbf{x}$ from the mean $\mathbf{m}$. If we represent $\mathbf{x}_k$ by $\mathbf{m} + a_k\mathbf{e}$, we can find an "optimal" set of coefficients $a_k$ by minimizing the squared-error criterion function

$$J_1(a_1, \ldots, a_n, \mathbf{e}) = \sum_{k=1}^{n} ||(\mathbf{m} + a_k\mathbf{e}) - \mathbf{x}_k||^2 = \sum_{k=1}^{n} ||a_k\mathbf{e} - (\mathbf{x}_k - \mathbf{m})||^2$$

$$= \sum_{k=1}^{n} a_k^2 ||\mathbf{e}||^2 - 2 \sum_{k=1}^{n} a_k \mathbf{e}^t (\mathbf{x}_k - \mathbf{m}) + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2. \tag{82}$$

Recognizing that $||\mathbf{e}|| = 1$, partially differentiating with respect to $a_k$, and setting the derivative to zero, we obtain

$$a_k = \mathbf{e}^t(\mathbf{x}_k - \mathbf{m}).$$

(83)

Geometrically, this result merely says that we obtain a least-squares solution by projecting the vector $\mathbf{x}_k$ onto the line in the direction of $\mathbf{e}$ that passes through the sample mean.

**SCATTER MATRIX**　This brings us to the more interesting problem of finding the *best direction* $\mathbf{e}$ for the line. The solution to this problem involves the so-called *scatter matrix* $\mathbf{S}$ defined by

$$\mathbf{S} = \sum_{k=1}^{n}(\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^t.$$

(84)

The scatter matrix should look familiar—it is merely $n - 1$ times the sample covariance matrix. It arises here when we substitute $a_k$ found in Eq. 83 into Eq. 82 to obtain

$$
\begin{aligned}
J_1(\mathbf{e}) &= \sum_{k=1}^{n} a_k^2 - 2\sum_{k=1}^{n} a_k^2 + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2 \\
&= -\sum_{k=1}^{n}[\mathbf{e}^t(\mathbf{x}_k - \mathbf{m})]^2 + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2 \\
&= -\sum_{k=1}^{n}\mathbf{e}^t(\mathbf{x}_k - \mathbf{m})(\mathbf{x}_k - \mathbf{m})^t\mathbf{e} + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2 \\
&= -\mathbf{e}^t\mathbf{S}\mathbf{e} + \sum_{k=1}^{n} ||\mathbf{x}_k - \mathbf{m}||^2.
\end{aligned}
$$

(85)

Clearly, the vector $\mathbf{e}$ that minimizes $J_1$ also maximizes $\mathbf{e}^t\mathbf{S}\mathbf{e}$. We use the method of Lagrange multipliers (described in Section A.3 of the Appendix) to maximize $\mathbf{e}^t\mathbf{S}\mathbf{e}$ subject to the constraint that $||\mathbf{e}|| = 1$. Letting $\lambda$ be the undetermined multiplier, we differentiate

$$u = \mathbf{e}^t\mathbf{S}\mathbf{e} - \lambda(\mathbf{e}^t\mathbf{e} - 1)$$

(86)

with respect to $\mathbf{e}$ to obtain

$$\frac{\partial u}{\partial \mathbf{e}} = 2\mathbf{S}\mathbf{e} - 2\lambda\mathbf{e}.$$

(87)

Setting this gradient vector equal to zero, we see that $\mathbf{e}$ must be an eigenvector of the scatter matrix:

$$\mathbf{S}\mathbf{e} = \lambda\mathbf{e}.$$

(88)

In particular, because $\mathbf{e}^t\mathbf{S}\mathbf{e} = \lambda\mathbf{e}^t\mathbf{e} = \lambda$, it follows that to maximize $\mathbf{e}^t\mathbf{S}\mathbf{e}$, we want to select the eigenvector corresponding to the largest eigenvalue of the scatter matrix. In other words, to find the best one-dimensional projection of the data (best in the least-
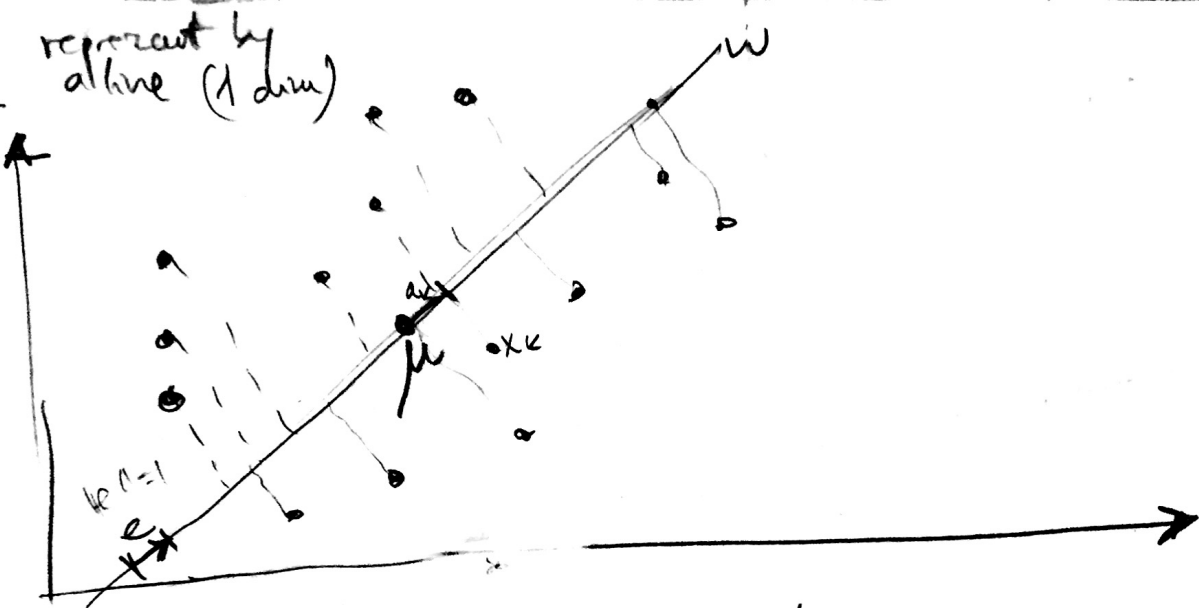
# PCA

① REPREZENT BY A POINT (0 dimansione)

— that would be the mean

$$\mu = \frac{1}{n} \sum_k x_k$$

$$\mu = \underset{x}{\text{argmin}} \sum_k \| x - x_k \|^2$$

② PCA   represent by
         a line (1 dim)



the line will pass through the mean

we need to define and find a good w.

$x_k = \mu + a_k e$      $\|e\|=1$ is the direction of the line
                        $a_k = $ scalar

$J(a_1, a_2, \ldots, e) = \sum_n \| \mu + a_k e - x_k \|^2 = \sum_k \| a_k e - (x_k - \mu) \|^2 = $

$= \sum_k a_k^2 \|e\|^2 - 2 \sum_k a_k e^T (x_k - \mu) + \sum \| x_k - \mu \|^2$

$\dfrac{\partial J}{\partial a_i} = 0 \Rightarrow a_k = e^T (x_k - \mu)$

$E[\mu + a_k e] = \mu + E[a_k e] = \mu + e^T E[x_k - \mu] e =$
$= \mu$

MOST IMPORTANT: what $\boxed{e}$ is a good direction?

- minimize $J$
- maximize the variance of the projections on $e$.

---

variance of projections

$$E\left[\left(\mu + a_k e - E[\mu + a_k e]\right)^2\right] = E\left[\mu + a_k e - \mu\right] = E\left[(a_k e)^2\right]$$

$$= E\left[e^T(x_k - \mu) \cdot e^T(x_k - \mu)\right] = E\left[e^T(x_k - \mu)(x_k - \mu)^T e\right] =$$

$$= e^T \Sigma e$$

where $\Sigma = \sum_k (x_k - \mu)(x_k - \mu)^T = $ ①

$$\begin{bmatrix} \sum_k (x_k^1 - \mu^1)(x_k^1 - \mu^1) & & \textcircled{j} & & \sum_k (x_k^1 - \mu^1)(x_k^d - \mu^d) \\ & & \sum_k (x_k^i - \mu^i)(x_k^i - \mu^i) & & \\ \sum_k (x_k^d - \mu^d)(x_k^1 - \mu^1) & & & & \sum_k (x_k^d - \mu^d)(x_k^d - \mu^d) \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_{11} & & \sigma_{1d} \\ & \sigma_{ij} & \\ \sigma_{d1} & & \sigma_{dd} \end{bmatrix} = \text{covariance matrix}$$

---

$$J(e) = \sum_k a_k^2 - 2\sum_k a_k^2 + \sum_k \|x_k - \mu\|^2 =$$

$$= -\sum_k \left(e^T(x_k - \mu)\right)^2 + \sum_k \|x_k - \mu\|^2 =$$

$$= -\sum_k e^T(x_k - \mu)(x_k - \mu)^T e + \sum_k \|x_k - \mu\|^2$$

$$\underbrace{\qquad\qquad} - e^T \Sigma e$$

So minimizing the error $\iff$ maximize the variance of projections.

maximize $e^T \Sigma e$

subject to $\|e\| = 1 \implies e^T e = 1$

Lagrangian $\quad L = \max\left[ e^T \Sigma e - \alpha (e^T e - 1) \right]$

$\dfrac{\partial L}{\partial e} = 0 \iff 2\Sigma e - 2\alpha e = 0 \implies \Sigma e = \alpha e$

$\Downarrow$

$\boxed{\begin{array}{l} e = \text{eigen vector of } \Sigma \\ \alpha = \text{eigen value of } \Sigma \end{array}}$

$e^T \Sigma e = e^T \alpha e = \alpha$

$\Uparrow$

we need to choose the (eigen vector, eigen value) pair with the biggest eigen value.

---

Say we want a second "biggest" dimension.

- constrained that is orthogonal on the first dim $e_1 = e$ so that measures a diff variance component

Lagrangian: $\max\left[ L = e_2^T \Sigma e_2 - \alpha (e_2^T e_2 - 1) - \beta (e_2^T e_1 - 0) \right]$

$\dfrac{\partial L}{\partial e} = 0 \implies 2\Sigma e_2 - 2\alpha e_2 - \beta e_1 = 0$

$\Downarrow$

$2 e_1^T \Sigma e_2 - 2 e_1^T \alpha e_2 - e_1^T \beta e_1 = 0 \implies \beta = 0$

$\underbrace{e_1^T e_2 = 0}_{} \qquad \underbrace{e_1^T e_2 = 0}_{}$

$\Downarrow$

$\Sigma e_2 = \alpha e_2 \implies e_2 \text{ eigenvector}$

$\alpha = \lambda_2 = \text{second eigenvalue}$

Spectral decompozition     $\Sigma$ symmetric, por def? $\Rightarrow$ ei orthogonrual.

$$C = \begin{pmatrix} | & | & & | \\ e_1 & e_2 & \cdots & e_d \\ | & | & & | \end{pmatrix} \qquad e_i = \text{eigenvectrs of } \Sigma, \text{ normalized}$$

$e_i^T e_j = 0$ if $i \neq j$ ;   $\|e_i\| = 1$. Then.

$$C \cdot C^T = I_d$$

$$\Sigma = \Sigma C C^T = \Sigma \begin{pmatrix} | & | & & | \\ e_1 & e_2 & \cdots & e_d \\ | & | & & | \end{pmatrix} C^T = \left(\Sigma e_1, \Sigma e_2 \cdots \Sigma e_d\right) C^T$$

$$= \left(\lambda_1 e_1, \lambda_2 e_2, \cdots \lambda_d e_d\right) C^T = \lambda_1 e_1 e_1^T + \lambda_2 e_2 e_2^T \cdots + \lambda_d e_d e_d^T$$

$$= C D C^T \quad \text{where} \quad D = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \circ \\ & & \ddots & \\ \circ & & & \lambda_d \end{pmatrix}$$

$$\Sigma = \begin{bmatrix} | & | & & \\ e_1 & e_2 & & \\ | & | & & \end{bmatrix} \cdot \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & \lambda_d \end{bmatrix} \begin{bmatrix} \text{---} e_1 \text{---} \\ \text{---} e_2 \text{---} \\ \\ \text{---} e_d \text{---} \end{bmatrix}$$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \geq 0$$

– if some $\lambda_{t+1} = \lambda_{t+2} = \lambda_d = 0$ then    we    can only keep

$e_1, e_2 \cdots e_t$ ; $\lambda_1, \lambda_2 \cdots \lambda_t$    so    we    reduced    to $t$ dimens.

– if underlined{approximate} to $t$ dimensions    eliminate $t+1 \rightarrow d$
                                            eigenv.  even if    they    are underlined{not} 0

$\Sigma$ symmetric $\Rightarrow$ $e_i$ orthogonal :
$$(\lambda_1 e_1)^T e_2 = (A e_1)^T e_2 = e_1^T A^T e_2 = e_1^T (A e_2) = e_1^T \lambda_2 e_2.$$
But $\lambda_1 \neq \lambda_2 \Rightarrow e_1^T e_2 = 0.$

eigenvalue

eigenvector

variance explained

eigenvector

# FISHER LINEAR DISCRIMINANT

Find dimensions that help classify data. (PCA might loose these)



PCA (first dim)
$W$

class better dim.

$m_1$ points $\in C_1$
$m_2$ points $\in C_2$

$\mu_1 = \frac{1}{m_1} \sum_{x \in C_1} x$

$\mu_2 = \frac{1}{m_2} \sum_{x \in C_2} x$

$\|w\| = 1$  $w$ = the line ; projected points $z = wx$

projected means: $\bar{\mu}_1 = \frac{1}{m_1} \sum_{C_1} w^T x = w^T \mu_1$

distance between projected means is $|\bar{\mu}_1 - \bar{\mu}_2| = |w^T (\mu_1 - \mu_2)|$

$\Sigma_i = \sum_{C_i} (x - \mu_i)(x - \mu_i)^T$     $\Sigma_w = \Sigma_1 + \Sigma_2$

$\sigma_i^2 = \underset{C_i}{\text{Var}} [w^T x] = w^T \Sigma_i w$     $\sigma_1^2 + \sigma_2^2 = w^T \Sigma_1 w + w^T \Sigma_2 w = w^T \Sigma_w w$

$\Sigma_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$

$(\bar{\mu}_1 - \bar{\mu}_2)^2 = (w^T \mu_1 - w^T \mu_2)^2 = w^T \Sigma_B w$

Fihers linear disc antena

$$\text{(maximize)} \quad J(w) = \frac{|\bar{\mu}_1 - \bar{\mu}_2|^2}{\sigma_1^2 + \sigma_2^2} =$$

$$= \frac{w^T \Sigma_B w}{w^T \Sigma_W w} \quad \left(\begin{array}{l}\text{sometimes called} \\ \text{"Rayleigh" quotient}\end{array}\right)$$

Solution must satisfy $\Sigma_B w = \lambda \Sigma_W w$
(max J)

Ⓐ if $\Sigma_W$ nonsingular $\Rightarrow \Sigma_W^{-1} \Sigma_B w = \lambda w \Rightarrow$

$(\lambda, w)$ eigens of $\Sigma_W^{-1} \Sigma_B$

Ⓓ for us, we do not need $\Sigma_W^{-1} \Sigma_B$ because $\Sigma_B w$ is in the direction of $(\mu_1 - \mu_2)$, so $w = \Sigma_W^{-1} (\mu_1 - \mu_2)$

Ⓒ All it is left is to find a threshold on the projections.