

INF8953 CE - Fall 2020

Machine Learning

- Sarath Chandar

6. Dimensionality Reduction.



6. Dimensionality Reduction

Let $X \in \mathbb{R}^d$ be a random vector. We wish to find $k < d$ directions that captures as much as possible of the variance of X (It contains the most important info in the dataset)

Applications: ① feature extraction

② data visualization

③ lossy data compression.



1-D Case: Find the single direction that captures as much as possible of the variance of X .

Formally: We want $p \in \mathbb{R}^d$ (the direction) such that $\|p\|=1$, so as to maximize $\text{var}(p^T X)$.

why $\|p\|=1$?

Else you ~~can~~ maximize the variance by letting $\|p\| \rightarrow \infty$.

Let $\mu = E[x]$

$$S = \text{Cov}(x) = E[(x-\mu)(x-\mu)^T]$$

For any $P \in \mathbb{R}^d$, the projection $P^T x$ has

$$\text{mean: } E[P^T x] = P^T \mu$$

$$\text{variance: } \text{Var}(P^T x) = E[(P^T x - P^T \mu)^2]$$

$$= E[P^T(x-\mu)(x-\mu)^T P]$$

$$= P^T S P$$

Goal:

$$\max_{\|P\|=1} P^T S P = \max_{P \neq 0} \frac{P^T S P}{P^T P}$$

$$= \max_{P \neq 0} \frac{P^T Q \Lambda Q^T P}{P^T Q Q^T P}$$

matrix with
orthonormal row
vectors
diagonal matrix
with eigen values
↓
eigen
vectors.

$[Q \Lambda Q^T$ is the spectral decomposition of $S]$
 $[\because Q Q^T = I]$ since it is orthonormal

$$= \max_{y \neq 0} \frac{y^T \Lambda y}{y^T y} \quad [\text{writing } y = Q^T P]$$

$$= \max_{y \neq 0} \frac{\lambda_1 y_1^2 + \dots + \lambda_d y_d^2}{y_1^2 + \dots + y_d^2}$$

where
 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$
 are the eigen values
 and corresponding
 eigen vectors u_1, \dots, u_d .

$$= \max_{y \neq 0} \lambda_1 \frac{y_1^2}{\sum_i y_i^2} + \dots + \lambda_d \frac{y_d^2}{\sum_i y_i^2}$$

$$\leq \lambda_1$$

where equality is attained in the last step when

$$y = e_1, \quad \text{i.e. } Z = Qe_1 = u_1,$$

where u_1 is the first eigen vector of S .

We call this the principal component.

First principle Component \rightarrow first eigen vector of $\text{Cov}(x)$.

Similarly, k -dimensional subspace that captures as much as possible of the variance of X is simply the subspace spanned by the top- k eigen vectors of $\text{Cov}(x)$: u_1, \dots, u_k .

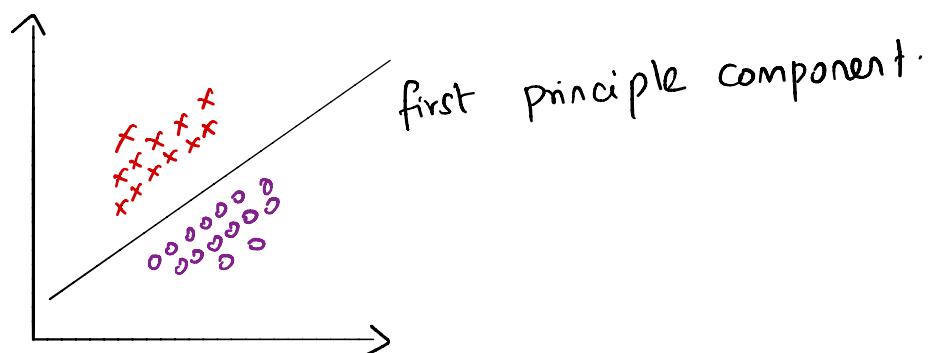
This is known as Principal Component Analysis (PCA).

Procedure:-

(1) Compute the mean μ and the Covariance matrix S of the data X .

- ② Compute the top-k eigen vectors $u_1 \dots u_k$ of S .
- ③ Project $X \rightarrow P^T X$, where P^T is the $k \times d$ matrix whose rows are $u_1 \dots u_k$.
-

Can we use PCA dimensions for classification?



Projecting data points to the first principal component
collapses the classes!

Better Solution:- Find the projection that maximizes
the class separation.

Consider a 2-class problem in which
there are N_1 points of class C_1 and N_2
points of class C_2 .

Mean vectors of C_1, C_2 are:

$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x^{(n)}$$

$$m_2 = \frac{1}{N_2} \sum_{n \in C_2} x^{(n)}$$

Choose w so as to maximize

$$w^T m_2 - w^T m_1$$

However this expression can be made arbitrarily large by increasing the magnitude of w .

↳ Constrain $\|w\| = 1$

$$\text{maximize } w^T m_2 - w^T m_1$$

$$\text{subject to } w^T w = 1$$

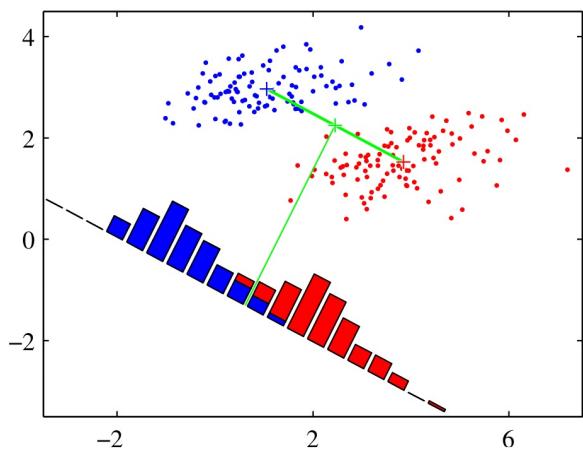
Lagrangian formulation:

$$L = w^T (m_2 - m_1) + \lambda (w^T w - 1)$$

$$\frac{\partial L}{\partial w} = (m_2 - m_1) + 2\lambda w = 0$$

$$\omega = -\frac{1}{2\lambda} (m_2 - m_1)$$

$$\omega \propto (m_2 - m_1)$$



Two classes are well separated
in the original Space.
But have considerable overlap
in the projected space.
→ because of the strongly
nondiagonal covariances of
class distributions.

Solution proposed by Fisher:-

Maximize a function that will give a large separation between the projected class means while also give a small variance within each class, thereby minimizing class overlap. within-class variance of the transformed data for class C_k is

$$S_k^2 = \sum_{n \in C_k} (\omega^\top x^{(n)} - \omega^\top m_k)^2$$

$$\text{Total within-class variance} = S_1^2 + S_2^2$$

$$\text{Fisher criterion: } J(\omega) = \frac{(\omega^\top m_2 - \omega^\top m_1)^2}{S_1^2 + S_2^2}$$

$$\begin{aligned} (\omega^\top m_2 - \omega^\top m_1)^2 &= (\omega^\top (m_2 - m_1))^2 \\ &= \omega^\top (m_2 - m_1) (m_2 - m_1)^\top \omega \\ &= \omega^\top S_B \omega \quad \text{--- (1)} \end{aligned}$$

\hookrightarrow between class covariance matrix.

$$\begin{aligned} S_1^2 + S_2^2 &= \sum_{n \in C_1} (\omega^\top (x^{(n)} - m_1))^2 + \sum_{k \in C_2} (\omega^\top (x^{(k)} - m_2))^2 \\ &= \sum_{n \in C_1} \omega^\top (x^{(n)} - m_1) (x^{(n)} - m_1)^\top \omega \\ &\quad + \sum_{k \in C_2} \omega^\top (x^{(k)} - m_2) (x^{(k)} - m_2)^\top \omega \\ &= \omega^\top S_\omega \omega \quad \text{--- (2)} \end{aligned}$$

$$\begin{aligned} \text{where } S_\omega &= \sum_{n \in C_1} (x^{(n)} - m_1) (x^{(n)} - m_1)^\top + \\ &\quad \sum_{k \in C_2} (x^{(k)} - m_2) (x^{(k)} - m_2)^\top \end{aligned}$$

Using ① and ②,

$$J(\omega) = \frac{\omega^T S_B \omega}{\omega^T S_W \omega}$$

diff. w.r.t. ω and set it to 0.

$$(\omega^T S_B \omega) S_W \omega - (\omega^T S_W \omega) S_B \omega = 0$$

$$(\omega^T S_B \omega) S_W \omega = (\omega^T S_W \omega) S_B \omega \quad \text{--- } \textcircled{A}$$

We do not care about the magnitude of ω , only its direction. So drop the scalars $(\omega^T S_B \omega)$, $(\omega^T S_W \omega)$ is \textcircled{A} .

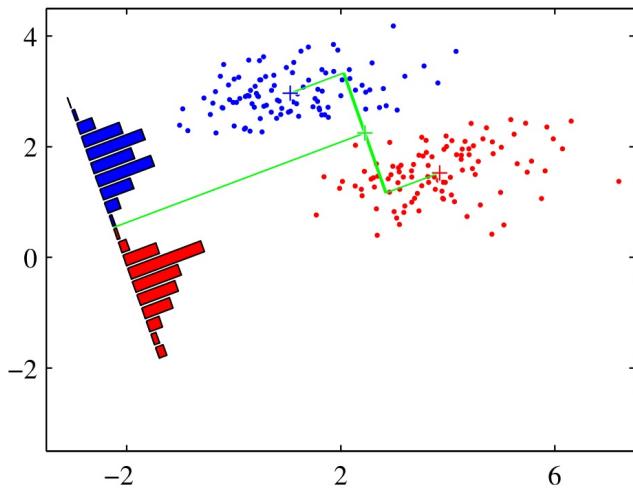
$$S_W \omega = S_B \omega.$$

$$\text{We know that } S_B = (m_2 - m_1)(m_2 - m_1)^T.$$

So $S_B \omega$ is always in the direction of $(m_2 - m_1)$.

$$S_W \omega \propto (m_2 - m_1)$$

$$\omega \propto S_W^{-1} (m_2 - m_1)$$



Note: If the within-class covariance is isotropic,
so that S_w is proportional to the unit matrix,
 w is proportional to the difference of the class
means.

This model is known as Fisher's linear discriminant.

However it is not a discriminant. It is a specific
choice of direction for projection of data.

The projected data can be used to construct
a discriminant - by choosing a threshold θ .

$$w^T x \geq \theta \Rightarrow c_1$$

$$w^T x < \theta \Rightarrow c_2$$

θ can be found by minimizing the
training error.

You should know!

1. Dimensionality reduction.
 2. Applications of dimensionality reduction.
 3. Principal Component Analysis (PCA).
 4. Fisher's linear discriminant analysis.
-