

Unsupervised Learning

First half: supervised problems

↳ a prediction problem:
use covariates (X_s) to predict a
response (Y)

→ training data to supervise a \hat{f}

finding/summarizing patterns in the data

→ no distinct X_s vs. Y
all we have is X_s

→ refine data in an interpretable way

Ex. ① dimensionality reduction:

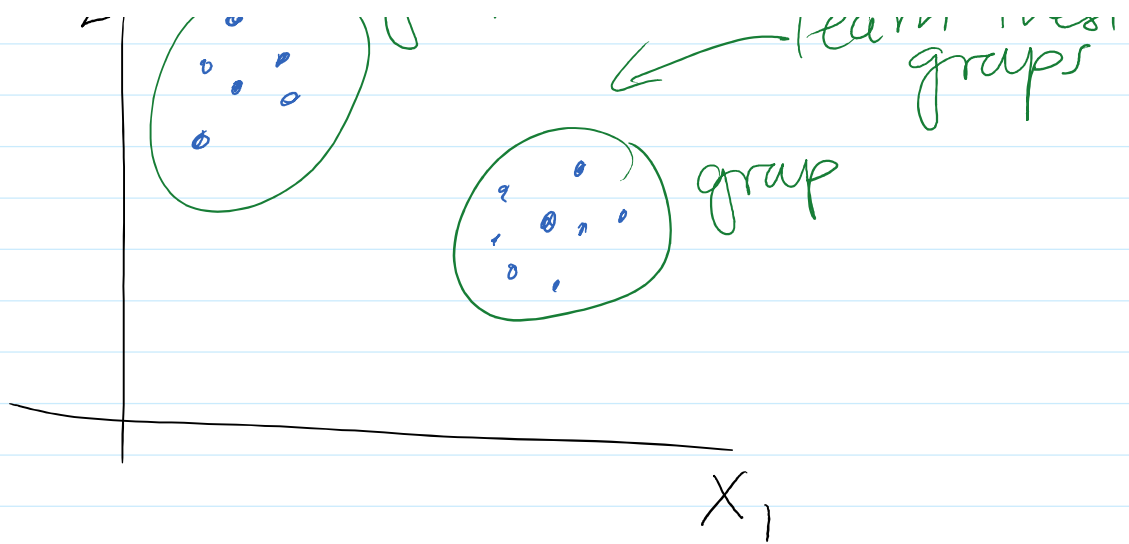
represent the data using a smaller number
of features

P covariates $\xrightarrow[\text{learning}]{\text{unsupervised}}$ $q \ll P$ features
that approx.
descr. the data.

② Clustering: group data into similar clusters

X_2 |  group 1

← learn these
groups

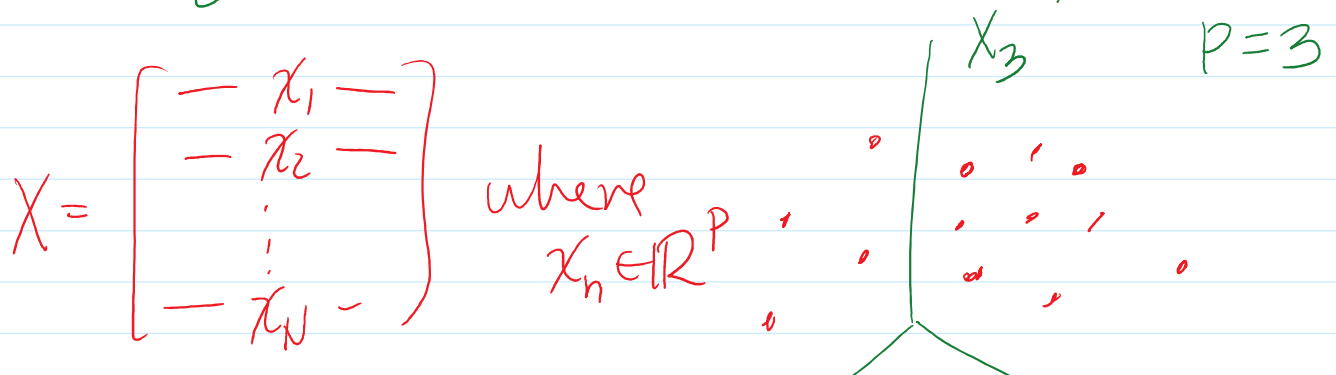


Today: Principal Components Analysis (PCA)
 Technique for dim'l reduction.

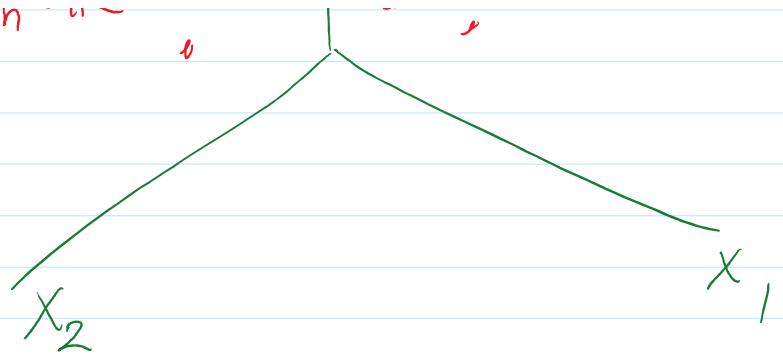
Unsupervised setting: data matrix X

$$X_{N \times P} = \begin{bmatrix} \text{--- obs 1 ---} \\ \text{--- obs 2 ---} \\ \text{--- obs 3 ---} \\ \vdots \\ \text{--- obs N ---} \end{bmatrix} = \begin{bmatrix} | & | & & | \\ \text{Var 1} & \text{Var 2} & \dots & \text{Var p} \\ | & | & & | \end{bmatrix}$$

visualize: Feature space / variable space (\mathbb{R}^P)

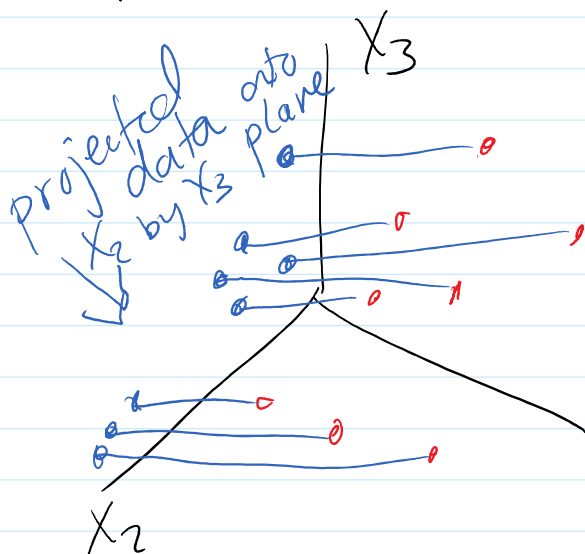


think of data as
living in some
 P -dim'l space.



Dimensionality reduction: can I get away w/
fewer dimensions?

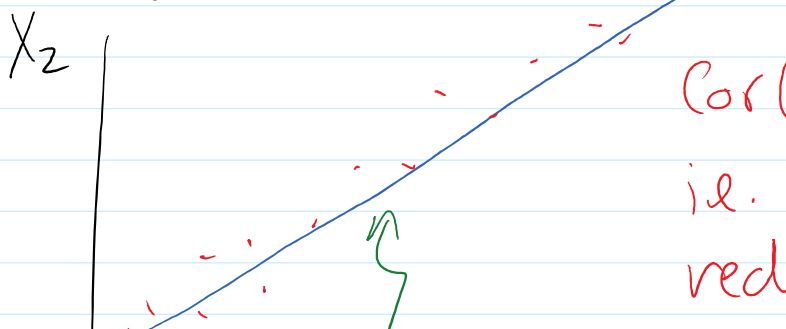
Simple Case: Variable selection



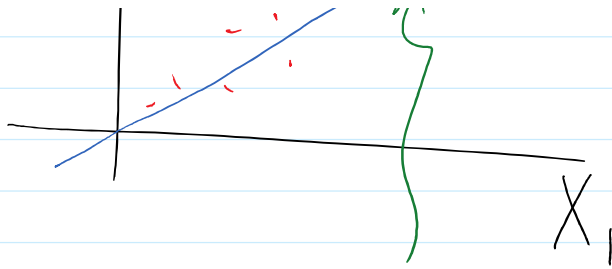
Maybe I don't
lose much if
I remove X_1
from my data

Hope! made my data
simpler w/o losing
much info.

Simpler yet: $X_1 \approx X_2$



$\text{Cor}(X_1, X_2)$ high,
ie. X_1 and X_2 are
redundant.



redundant.

$$X_1 \text{ and } X_2 \rightsquigarrow \frac{1}{2}(X_1 + X_2)$$

Goals of PCA:

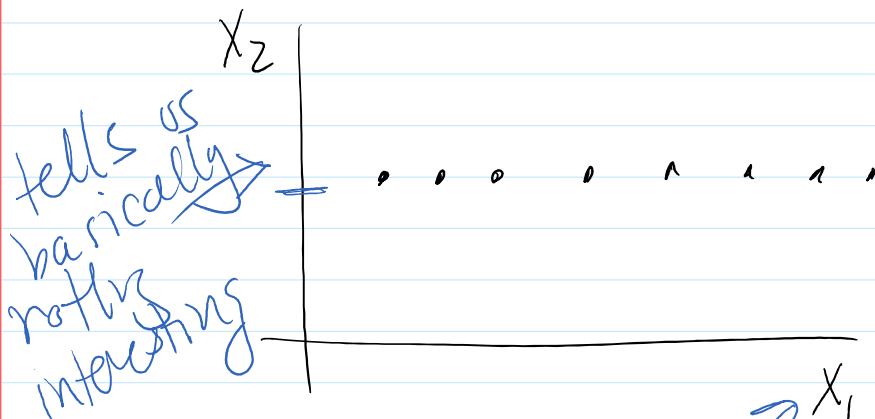
- ① reduce the data (# of vars)
dimension
 $X_1, \dots, X_p \xrightarrow{\text{PCA}} Z_1, \dots, Z_g \text{ where } g \ll p.$
- ② not lose too much info.
equiv.: maximize the amount of info retained

min
loss



max
retain

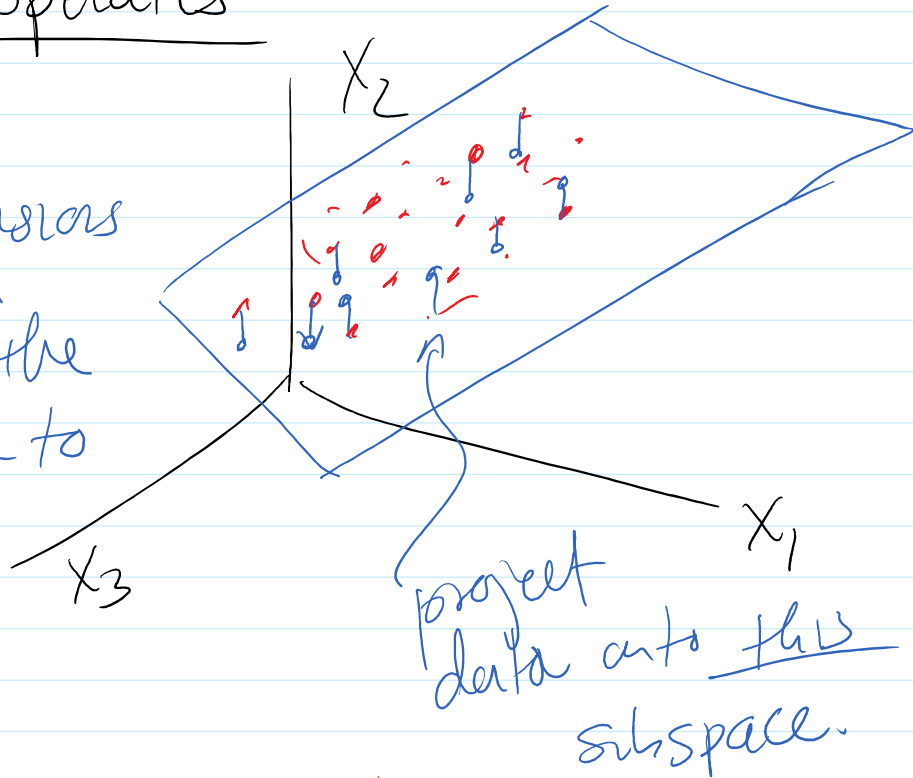
central dogma: Variance = interestingness
or info



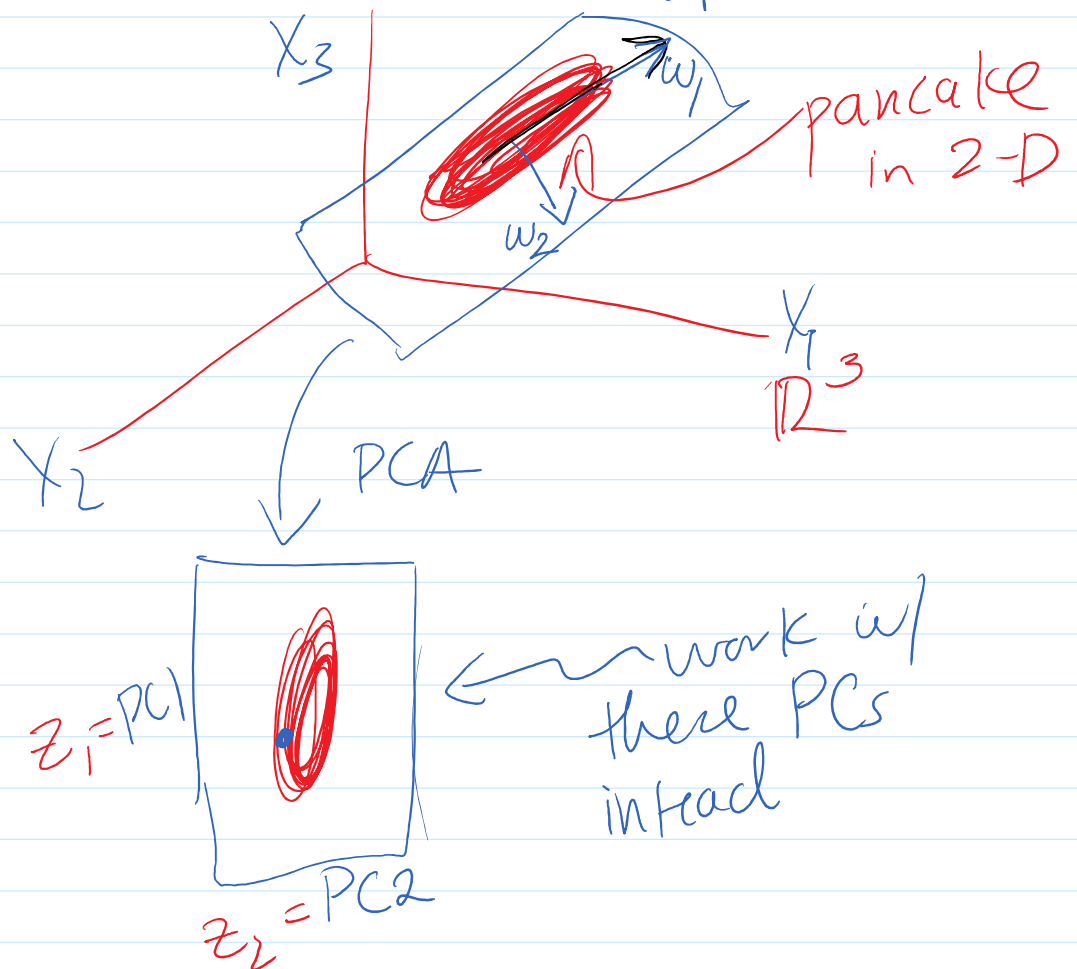
X_1 tells us something about the data

How PCA operates

Instead of
removing dimensions
of our data
aligned w/ the
axes, reduce to
only good
subspace



Idea!



PCA:

PCA:

Take orig. X_1, \dots, X_p $\xrightarrow{\text{create}}$ z_1, \dots, z_g

If w_i is the i^{th} basis vector for this subspace then

$$Z = XW$$

$N \times g$ $N \times p$

where

$$W = \begin{bmatrix} | & | & & | \\ w_1 & w_2 & \dots & w_g \\ | & | & & | \end{bmatrix}$$

← change of basis mtx

alt.

← i^{th} PC

$$z_i = Xw_i = w_{i1}X_1 + w_{i2}X_2 + \dots + w_{ip}X_p$$

Alternative Formulation:

Let $z_i = \text{LCs of } X_1, \dots, X_p$

where z_1 maximizes variance

z_2 max. variance st. $\text{cor}(z_1, z_2) = 0$

z_3 max. var. st. it is " "

Uncorr w/ z_1 and z_2

⋮

z_g max. var. s.t. uncorr w/
rest.

Mathematically

$X_{N \times p}$ we want to find $W = \begin{bmatrix} w_1 & \dots & w_g \\ 1 & & 1 \end{bmatrix}$

so that we maximize PC_i

$$\text{Var}(z_i) = \text{Var}(Xw_i)$$

subject to the constraint that $\text{Cor}(z_i, z_j) = 0$
for all $j < i$.

Constraint: W is orthogonal.
equiv. w_i 's are unit vectors.

Review: $x = (x_1, \dots, x_N)$ (same variable/col of X)

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n = \frac{1}{N} \mathbf{1}^T x = \frac{1}{N} x^T \mathbf{1}$$

where $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^N$

WLOG: assume $\bar{X} = 0$ (if not we'll center our data)

PCA assume mean-centered data.

$$\text{Var}(x) = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2 = \frac{1}{N} \sum_n x_n^2 = \frac{1}{N} \|x\|^2 \\ = \frac{1}{N} x^T x$$

If $\bar{x} = 0$ then $\|\cdot\|^2 \approx \text{Var}$

If $y \in \mathbb{R}^N$ and $\bar{y} = 0$ then

$$\text{Cov}(x, y) = \frac{1}{N} \sum_n x_n y_n = \frac{1}{N} x^T y$$

If $\bar{x} = \bar{y} = 0$ then inner product \approx Covariance

so if $x^T y = 0 \Leftrightarrow \text{Cov}(x, y) = 0 \Leftrightarrow \text{Cor}(x, y) = 0$

Simplify PCA:

$$\text{Var}(Xw_i) = \text{Var}(z_i)$$

$$(*) \quad w_1 = \underset{\|w\|=1}{\text{argmax}} \|Xw\|^2$$

$$w_2 = \underset{\|w\|=1}{\text{argmax}} \|Xw\|^2$$

$$w^T w_1 = 0$$

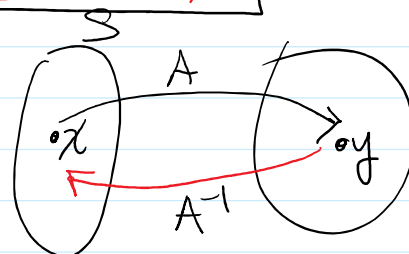
$$w_3 = \underset{\substack{\|w\|=1 \\ w^T w_1 = 0 \\ w^T w_2 = 0}}{\operatorname{argmax}} \|Xw\|^2$$

$$\vdots$$

$$w_g = \dots$$

Aside: if A is an invertible matrix and

$$x^* = \underset{x \in S}{\operatorname{argmax}} f(x)$$



$$y^* = \underset{y \in A^{-1}S}{\operatorname{argmax}} f(Ay)$$

find $x^* = Ay^*$

Ex. U orthogonal matrix $\Leftrightarrow U^{-1} = U^T$

$$\underset{\|z\|=1}{\operatorname{argmax}} f(z) = U^T \underset{\|U^T z\|=1}{\operatorname{argmax}} f(Uz)$$

Ex. $\underset{z^T a = 0}{\operatorname{argmax}} f(z) = U^T \underset{z^T U^T a = 0}{\operatorname{argmax}} f(Uz)$

PC1 $w_1 = \underset{\|w\|=1}{\operatorname{argmax}} \|Xw\|^2$

$$\text{PC1 } w_1 = \underset{\|w\|=1}{\operatorname{argmax}} \|Xw\|^2$$

$$\left(\text{let } X = UDV^T \right.$$

(shown prev.)
 $\|u\| = \|v\|$

$$= \underset{\|w\|=1}{\operatorname{argmax}} \|UDV^T w\|^2$$

$$= \underset{\|w\|=1}{\operatorname{argmax}} \|DV^T w\|^2$$

$$= V^T \underset{\|V^T w\|=1}{\operatorname{argmax}} \|DV^T V w\|^2$$

$$= V^T \underset{\|w\|=1}{\operatorname{argmax}} \|Dw\|^2$$

$$D = \left[\begin{array}{ccc|c} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_r \\ \hline & & & 0 \end{array} \right]$$

$$\left(\frac{1}{\sigma_1} w_1 \sigma_1 + \cancel{w_2 \sigma_2} + \dots + \cancel{w_r \sigma_r} \right)^2$$

$$\text{set } w_1 = 1$$

$$\text{others} = 0$$

$$\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$$

$$= V^T \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = V_1 \leftarrow \text{first col of } V.$$

$$\text{PC2 } w_2 = \underset{\substack{\|w\|=1 \\ V_1^T w = 0}}{\operatorname{argmax}} \|Xw\|^2$$

$$= \underset{\|w\|=1}{\operatorname{argmax}} \|UDV^T w\|^2$$

$$\begin{aligned}
 &= \underset{\substack{\|w\|=1 \\ V_1^T w = 0}}{\operatorname{argmax}} \|\cancel{U} D V^T w\|^2 \\
 &= V^T \underset{\substack{\|w\|=1 \\ \underbrace{V_1^T V^T w = 0}_{\text{simply to const. that first comp. of } w = 0}}}{\operatorname{argmax}} \|D w\|^2
 \end{aligned}$$

$V_1^T w = 0$
 \downarrow
 $V_1^T V^T w$
 $\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}^T w = 0$

$$= V^T \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = V_2 \leftarrow \text{col 2 of } V.$$

Keep playing this game.

$$W_{n \times q} = V_{1:q} \leftarrow \text{first } q \text{ cols of } V$$

when $X = U D V^T$.

Punchline: PCA procedure:

- ① mean center cols of X (call it X_c)
- ② let $X_c = U D V^T$

③ $W = \text{first } g \text{ cols of } V = V_{1:g}$

④ $Z_{N \times g} = XW = XV_{1:g}$
