

# Unsupervised Learning

Supervised: have both an input  $X$  and output  $Y$

Want to predict  $Y$  from  $X$   $p(y|x)$

Unsupervised: only have an  $X$

Want to summarize important patterns in  $X$

$p(x)$

Ex. ① dimensionality reduction

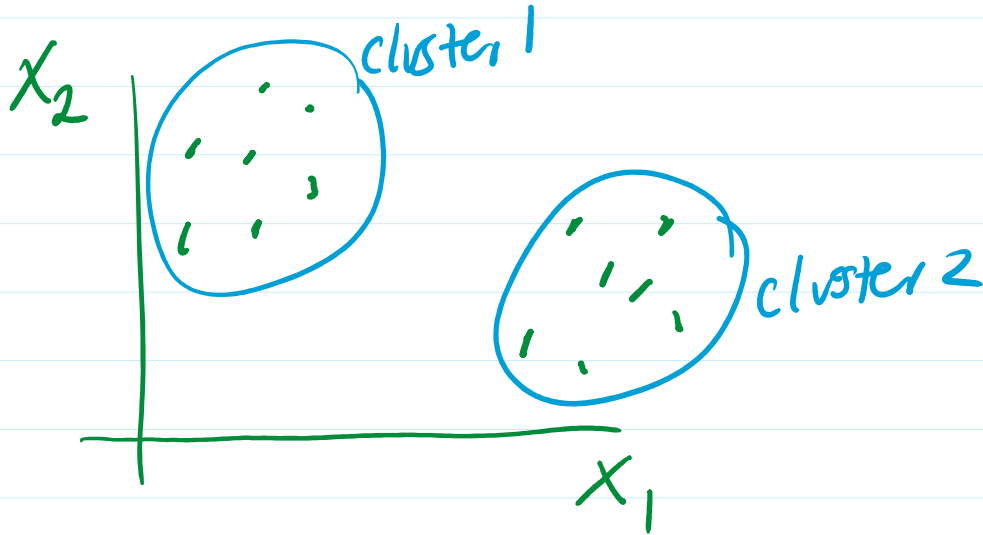
rep data using fewer vars

$P$  covariates  $\rightsquigarrow$   $q$  covariates

$q \ll P.$

$$q \ll P.$$

② Clustering: finding high density subspaces  
group data into similar "clusters"



finding high density regions

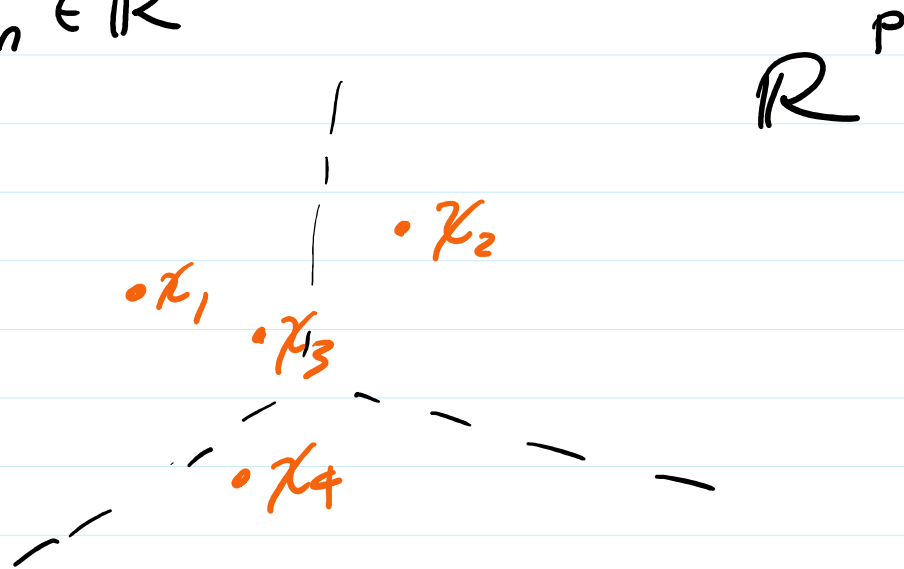
---

## Principal Components Analysis (PCA)

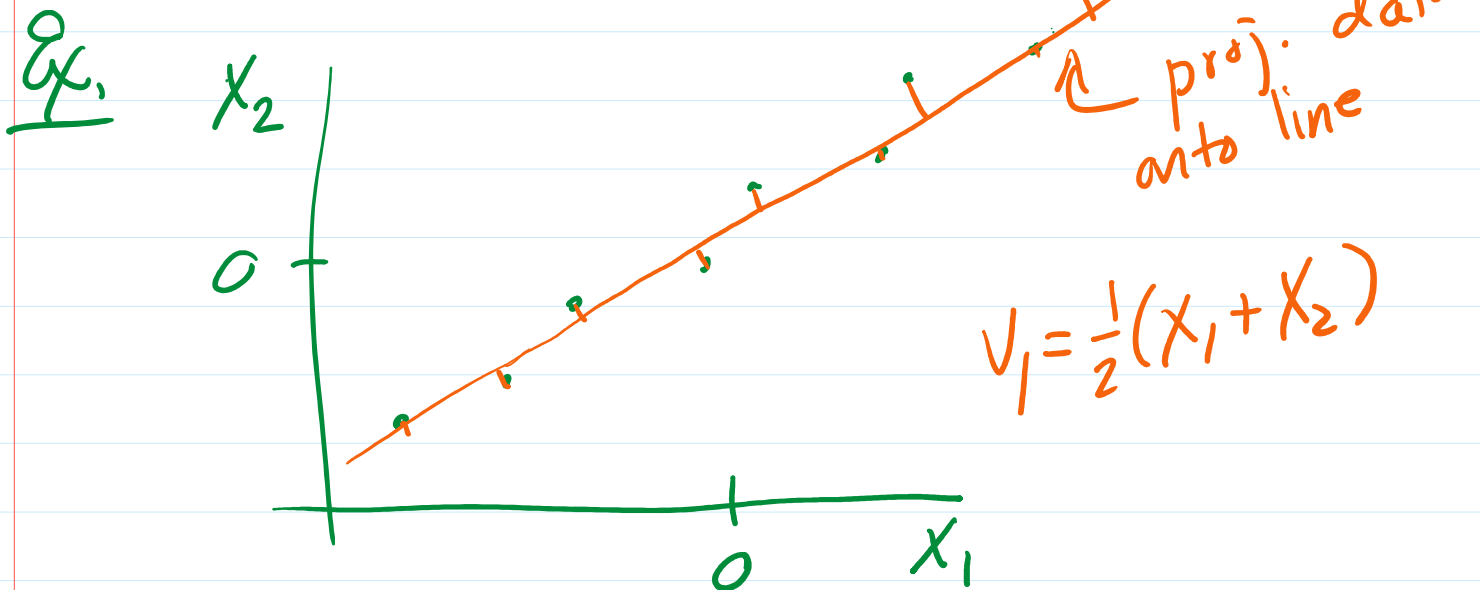
Dim'l reduction technique.

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

visualize  $X$  row-wise so that  $n^{\text{th}}$  row  
 $x_n \in \mathbb{R}^p$



Goal: find a lower dim'l subspace of  
 $\mathbb{R}^p$  that doesn't lose too much info  
about my data



# Goals of PCA :

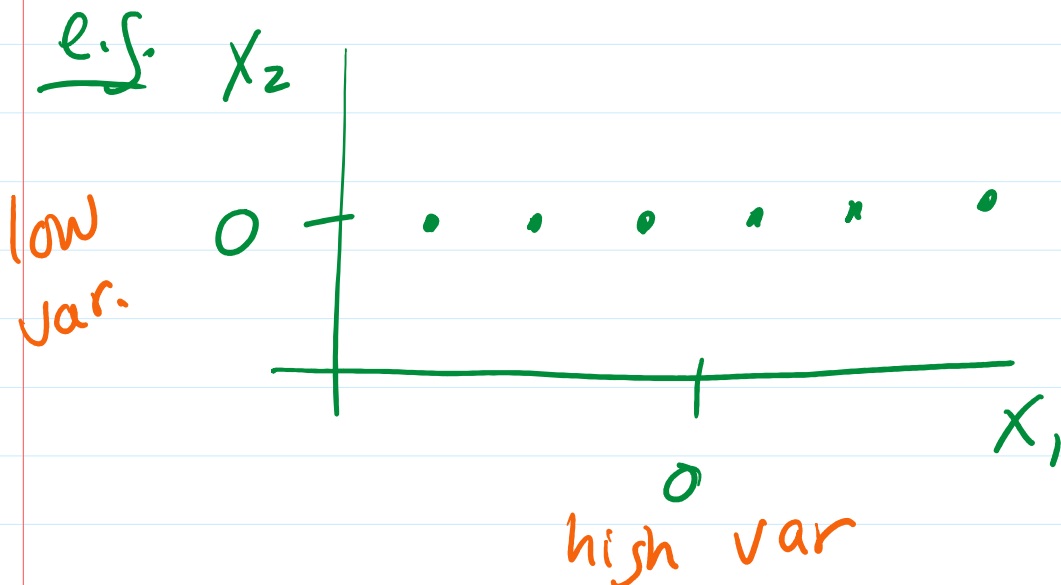
① reduce the number of vars

$$X_1, \dots, X_p \xrightarrow{\text{reduce}} Z_1, \dots, Z_q$$

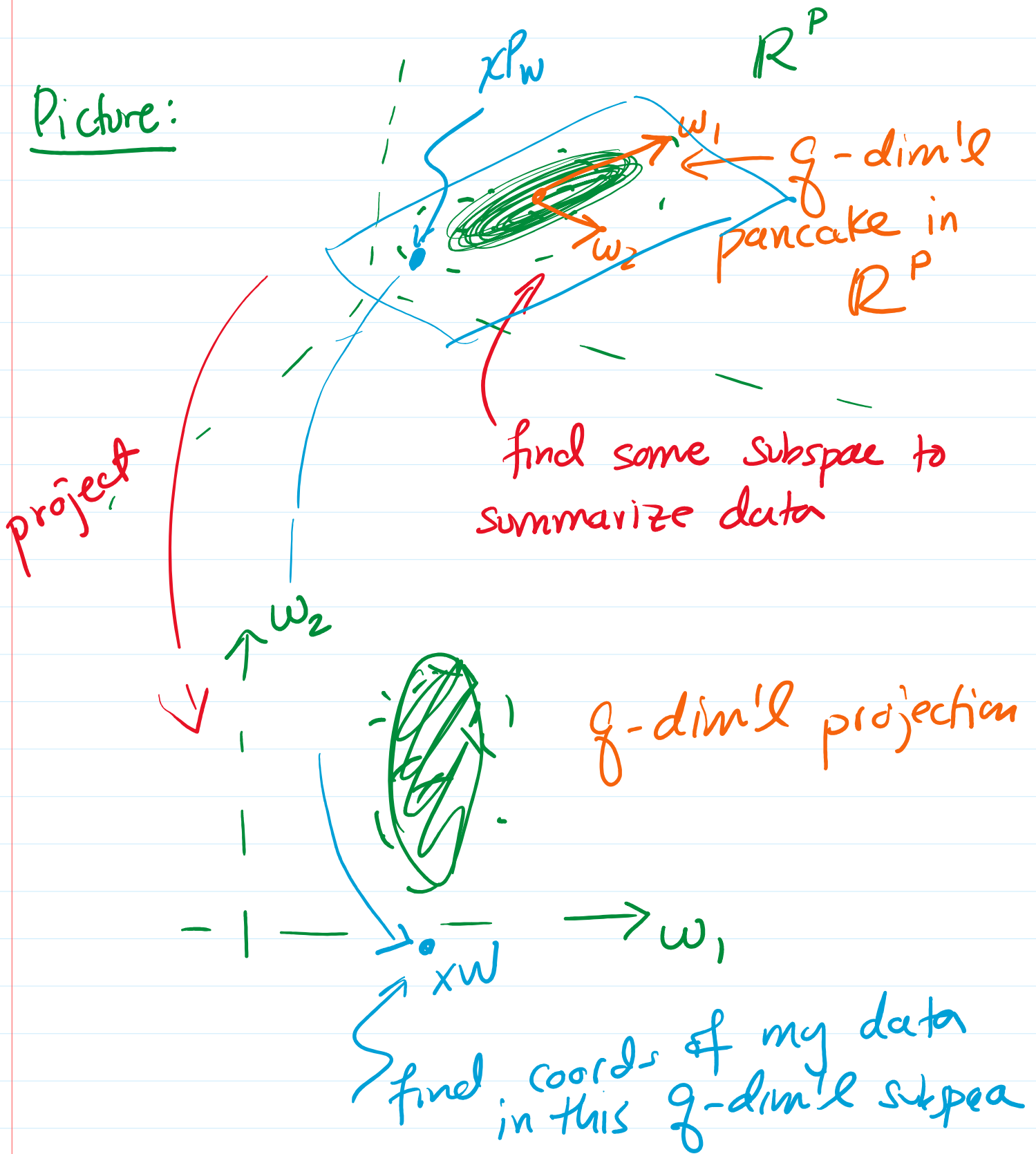
where  $q \ll p$ .

② don't lose too much info

Central dogma: Variance = info



Picture:



Proj. matrices

## Proj. matrices

Let  $W$  be the  $P \times q$  matrix of basis elements.

The proj. matrix onto  $\text{Col}(W)$  is

$$P_W = W(W^T W)^{-1} W^T$$

so that if  $x \in \mathbb{R}^P$  (row-vector) then

$x P_W$  is the

coord (in  $\mathbb{R}^q$ ) of the proj. pt.

If  $W$  is "orthogonal" (col orthonormal)

then

$$P_W = W W^T$$

and so

$$\underbrace{x P_W}_{\text{coords of proj. pt in new basis}} = \underbrace{x W W^T}_{\substack{1 \times P \quad P \times q \quad q \times P \\ \text{bases}}}$$

coords of proj.  
in orig. space  
( $\mathbb{R}^p$ )

$1 \times p$   $p \times q$   $q \times p$  — bases

PCA wants to find pos. of data pts  
in the lower dim'l space.

If  $X_{N \times p}$  is my data mtr then

$$Z = XW$$

$N \times p$   $p \times q$

(expressed)  
is the data mtr embedded in this  
lower dim'l coords.

If  $Z_i$  is the  $i^{\text{th}}$  col of  $Z$  ( $i^{\text{th}}$  PC)  
 $X_j$  is the  $j^{\text{th}}$  col of  $X$  (PCs)

$w_i$  is  $i^{\text{th}}$  col  $w$  LC of cols of  $X$

then

$$\boxed{z_i = X w_i} = X_1 w_{i1} + X_2 w_{i2} + \dots + X_p w_{ip}$$

for  $i=1, \dots, q$

---

PCA: Find LCs of  $X_j$ s to

(1) max var. of resulting  $z_i$ s

Subject to

(2)  $z_i$ s are uncorrelated  
(no redundancy)

(3)  $w_i$  are unit vectors

---

Aside: let  $x \in \mathbb{R}^N$

assume  $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n = 0$ .



$$\begin{aligned} \text{Then } \hat{\text{var}}(x) &= \frac{1}{N-1} \sum_n (x_n - \bar{x})^2 \\ &= \frac{1}{N-1} \sum_n x_n^2 \\ &\propto x^T x \end{aligned}$$

Similarly if  $y \in \mathbb{R}^N$  is some other var then

$$\hat{\text{cor}}(x, y) \propto \hat{\text{cov}}(x, y) \propto x^T y$$

So uncorrelated  $\approx$  orthogonal

Similarly if  $X$  is an  $N \times P$  data matrix (cols are mean-centered)

$\hat{\text{Cov}}(X)$  is a  $P \times P$  matrix

$$\text{where } \hat{\text{Cov}}(X)_{ij} = \hat{\text{cov}}(x_i, x_j)$$

where  $\text{Cov}(X)_{ij} = \text{Cov}(x_i, x_j)$

$$\widehat{\text{Cov}}(X)_{ii} = \widehat{\text{Var}}(x_i)$$

then

$$\widehat{\text{Cov}}(X) \propto \underbrace{X^T X}_{P \times P}$$

$P \times N$     $N \times P$

PCA:  $Z = XW$

Want: (1) max diag elements of  $\widehat{\text{Cov}}(Z)$   
(max sum of diag)

(2) off diag elements of  $\widehat{\text{Cov}}(Z)$  are zero

(3) Cols of  $W$  to be unit vectors.

⊆ sets of -- vectors.

---