

Lecture 15

Can generally fit penalized methods

$$\hat{f} = \underset{f}{\operatorname{argmin}} L(f)$$

↑ loss

penalize:

$$\hat{f} = \underset{f}{\operatorname{argmin}} L(f) + \lambda J(f)$$

↑ penalization strength

↑ measures complexity of f

E.g.

ridge: $\beta \Leftrightarrow f$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta) + \lambda \|\beta\|_2^2$$

↑ complexity

LASSO:

$$L(\beta) + \lambda \|\beta\|_1$$

Can do this w/ other methods too.

E.g. Logistic Regression

$$y_n | \underline{x}_n = \underline{x}_n \stackrel{\text{indep}}{\sim} \text{Bern}(p_{\beta}(\underline{x}_n))$$

↑ logistic fn

$$p_{\beta}(\underline{x}_n) = 1 / (1 + \exp(-\underline{x}_n^T \beta))$$

Find $\hat{\beta}$ as MLE

$$\hat{\beta} = \arg \max_{\beta} P_{\beta}(\underline{y} | \underline{x})$$

$$= \underset{\beta}{\operatorname{argmin}} \underbrace{-\log P_{\beta}(\tilde{Y} | \tilde{X})}_{\text{negative log likelihood (NLL)}}$$

$$= \underset{\beta}{\operatorname{argmin}} L(\beta)$$

Can also penalize this fit,

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta) + \lambda \|\beta\|_2^2$$

$$\alpha + \lambda \|\beta\|_1$$

This has similar effects as in regression case.

Unsupervised Learning

Supervised Problem: have both input X and an output Y

Want to predict Y from X . $P(Y|X)$

Unsupervised: only have X

Want to summarize important patterns in X .

$P(X)$

Ex.

① dimensionality reduction:

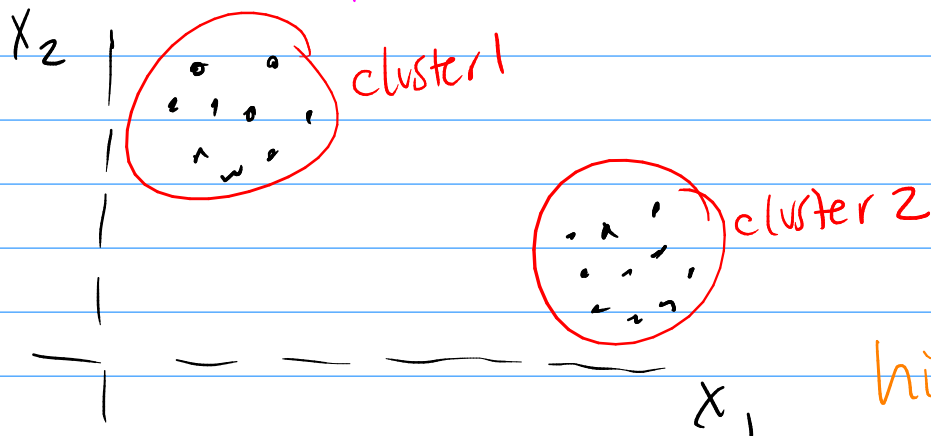
represent data using fewer vars.

P covariates $\rightsquigarrow q$ covariates

$q \ll P$

high density subspaces of $P(X)$

② clustering: group my data into similar "clusters"



high-density regions of $P(X)$

Principal Components Analysis (PCA)

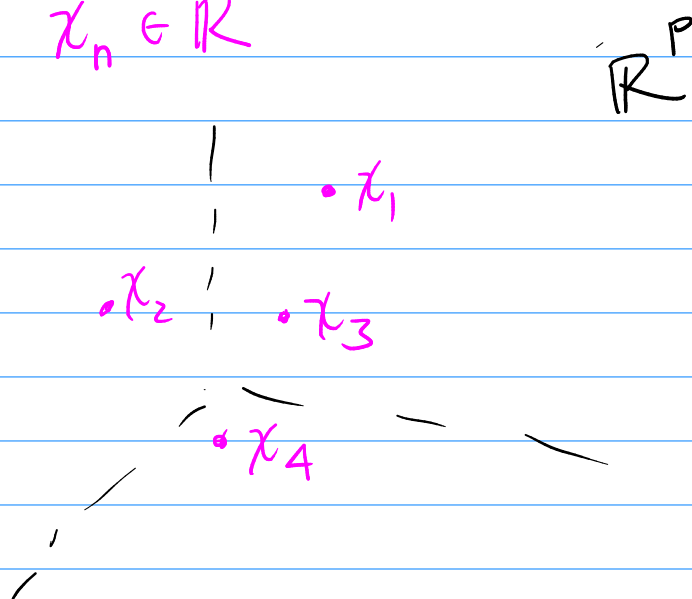
→ dim-reduction approach

Unsupervised!

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

visualize X row-wise

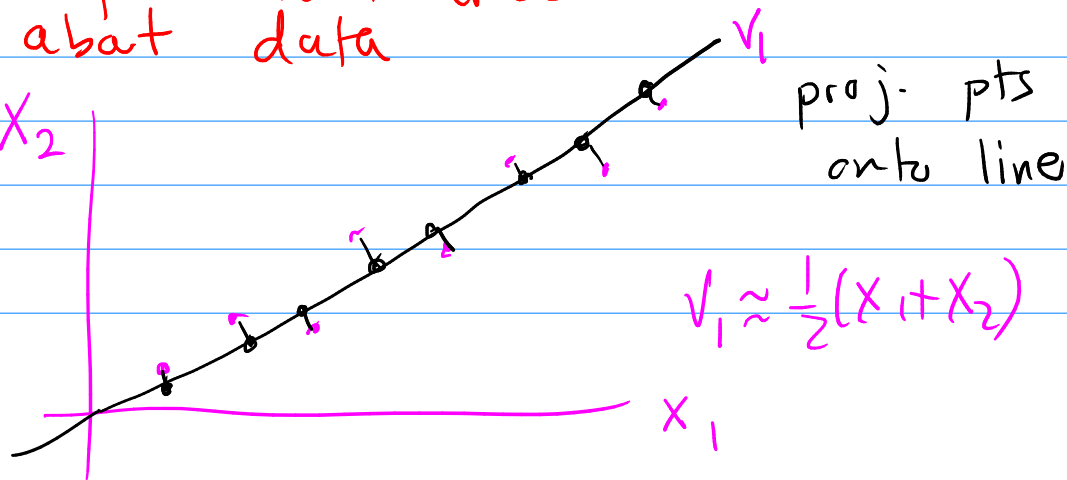
n^{th} row $x_n \in \mathbb{R}^P$



dim-reduction tries to find a lower dim'l subspace that doesn't lose too much info about data

$P=2$

X_2



Goals of PCA:

- ① reduce the number of vars

$$X_1, X_2, \dots, X_p \xrightarrow{\text{reduce}} Z_1, Z_2, \dots, Z_g$$

principal components

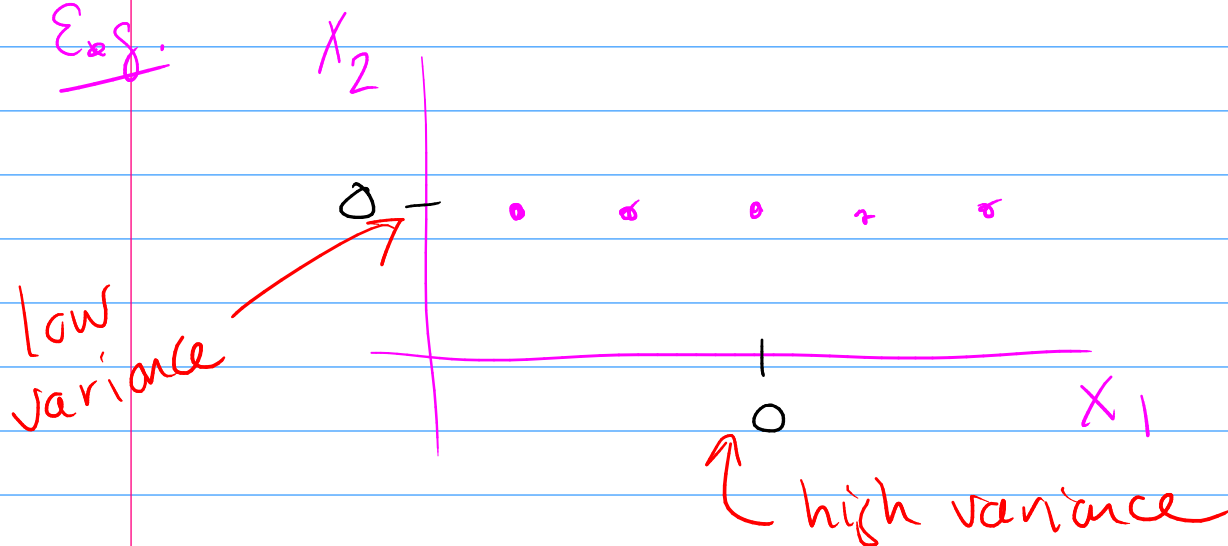
where $g \ll p$.

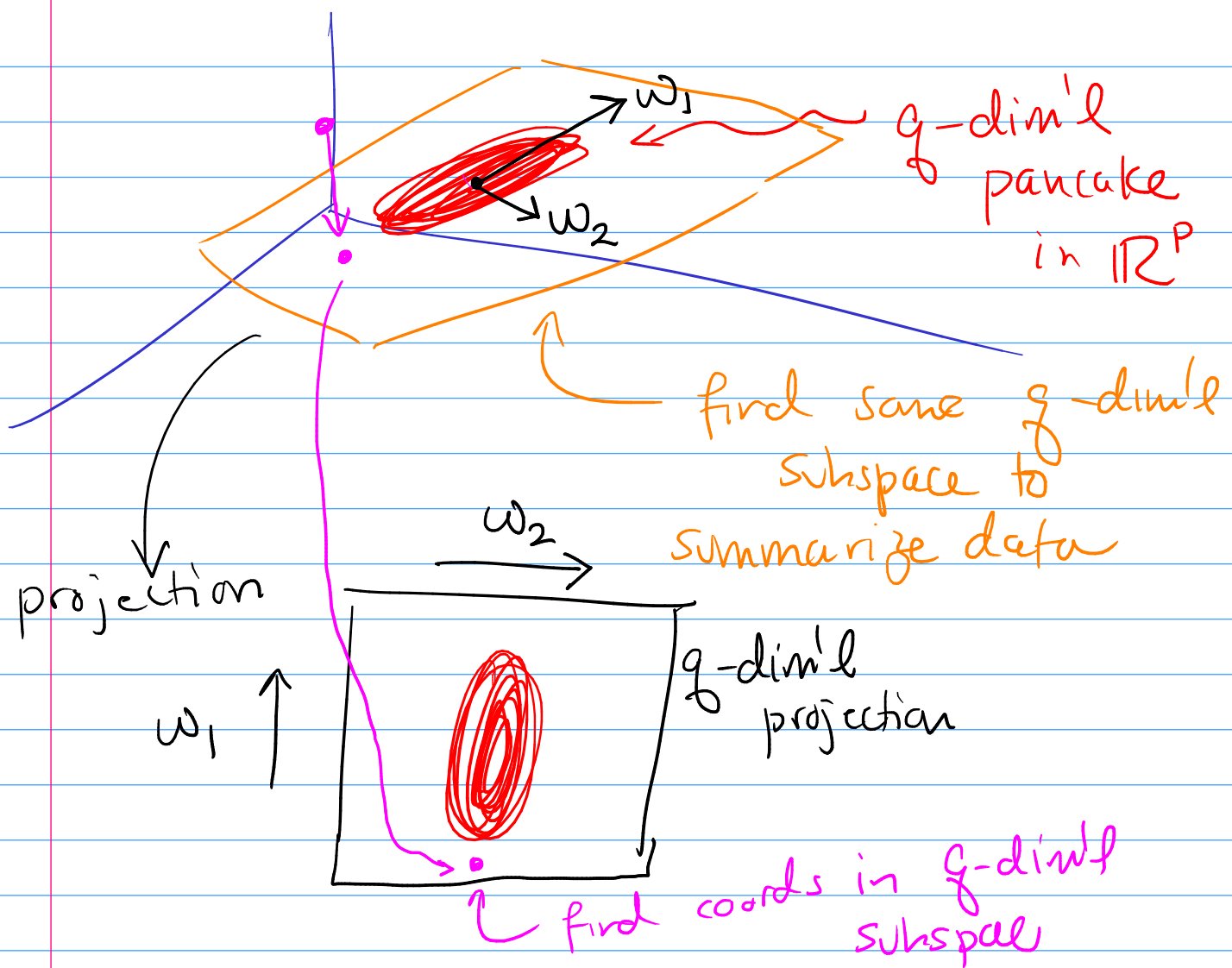
- ② don't want to lose too much info about my data



Central dogma: variance = information

Ex.



\mathbb{R}^p 

Proj. matrices

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

The proj. mtrix onto $\text{col}(W)$ is

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

so that if $x \in \mathbb{R}^p$ then $x P_W$ is the coord (in \mathbb{R}^p) of proj. pt.

If W has orthonormal cols then

$$W^T W = I \quad \text{so}$$

$$P_W = W W^T \quad (P \times P)$$

so

$$\underbrace{X}_{P} \underbrace{P_W}_{g} = \underbrace{X}_{\text{coords in org. basis}} \underbrace{W}_{\text{coordinates of proj. in } W \text{ basis}} \underbrace{W^T}_{P \times P}$$

coords in org. basis

coordinates of proj. in W basis

PCA wants to calc positions of pts in the lower dim'l subspace.

If $X_{N \times P}$ is data mtx, W is $P \times g$ basis mtx then

$$Z = X W$$

\nearrow $N \times P$ $P \times g$

$N \times g$ mtx of data in lower-dim'l coords

If z_i is the i^{th} col of Z
 x_j " j^{th} col of X

$$z_i = X W_i = x_1 w_{i1} + x_2 w_{i2} + \dots + x_p w_{ip}$$

linear combn of x_j 's w/ weights specified by W

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 to be unit vectors