

# Lecture october 29

## PCA Theorem

$$X \in \mathbb{R}^{n \times p}$$

$$\begin{aligned}\Sigma_x &= \frac{1}{n} X X^T \in \mathbb{R}^{p \times p} \\ &= E[X X^T]\end{aligned}$$

$$X \in \mathbb{R}^{p \times n}$$

$$\Sigma_x = E[X X^T]$$

Define an orthogonal transformation

$$S = [s_0 s_1 s_2 \dots s_{p-1}]$$

$$S \in \mathbb{R}^{p \times p}$$

$$S S^T = S^T S = \mathbf{1}$$

$$S X = y$$

$$\text{var}(S X) =$$

$$= E[(S X)^2] = E[S X X^T S^T]$$

$X$  is a scaled matrix with subtracted mean value

$$\begin{aligned}
 X &\in \mathbb{R}^{p \times n} \\
 &= S \Sigma_X S^T = \Sigma_Y \\
 &= \text{diag}(\lambda_0, \lambda_1, \dots, \lambda_{p-1})
 \end{aligned}$$

$$\begin{aligned}
 \Sigma_X \cdot S_0 &= \lambda_0 S_0 = \text{var}(y_0) S_0 \\
 &\quad \quad \quad \parallel \\
 &\quad \quad \quad \text{var}(y_0)
 \end{aligned}$$

PCA - theorem ;  
want  $d < p$

$$\begin{aligned}
 y_i' &= S_i^T X \quad i = 0, 1, \dots, p-1 \\
 \lambda_0 &\geq \lambda_1 \geq \lambda_2 \dots \geq \lambda_p \\
 \text{var}(y_0) &= \lambda_0 \geq \text{var}(y_1) \geq \dots \text{var}(y_{p-1})
 \end{aligned}$$

$$\text{var}(y_0) = \lambda_0 \quad \text{with eigenvector } S_0$$

$$\begin{aligned}
 \max_{S_0 \in \mathbb{R}^p} \quad & S_0^T \Sigma_X S_0 \\
 \text{subject to} \quad & S_0^T S_0 = \underline{1}
 \end{aligned}$$

Constrained optimization  
(support vector machines)

Define Lagrangian

$$\mathcal{L} = S_0^T \Sigma_X S_0 + \lambda_0 (1 - S_0^T S_0)$$

compute derivatives wrt  $s_0$   
and  $\lambda_0$

$$\lambda_0 : s_0^T s_0 = 1$$

$$\Sigma_x s_0 = \lambda_0 s_0$$

optimal solution for  $s_0$  is  
given by the eigenvector  
of  $\Sigma_x$ .

First principal component,  
To find the second principal  
component;

$$s_1^T s_0 = 1$$

$$\begin{aligned} E [s_0 x x^T s_1^T] &= s_0 \Sigma_x s_1^T \\ &= \lambda_0 s_0 s_1^T = 0 \end{aligned}$$

$$\mathcal{L} = s_1^T \Sigma_x s_1 + \lambda_1 (1 - s_1^T s_1) + \gamma s_1^T s_0$$

$$s.t. \quad s_1^T s_1 = 1 \quad s_1^T s_0 = 0$$

Taking derivatives wrt  $s_1, \lambda_1$   
&

$$\gamma = 0$$

$$\Sigma_x s_1 = \lambda_1 s_1 \quad . \quad \text{The best choice}$$

if  $\lambda_1 = \text{var}(y_1)$

We can construct remaining by induction.

SVD

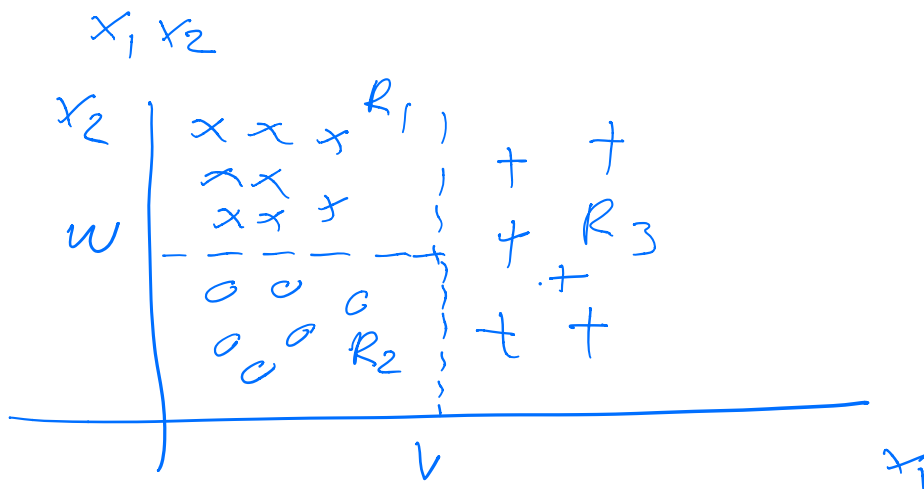
$$X \in \mathbb{R}^{n \times p}$$

$$S = \begin{bmatrix} \sigma_0 & & \\ & \ddots & \\ & & \sigma_{(p)} \\ & & & 0 \dots 0 \end{bmatrix}$$

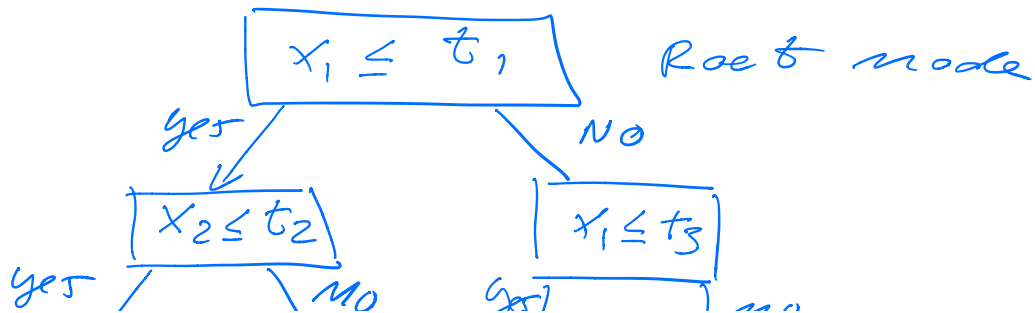
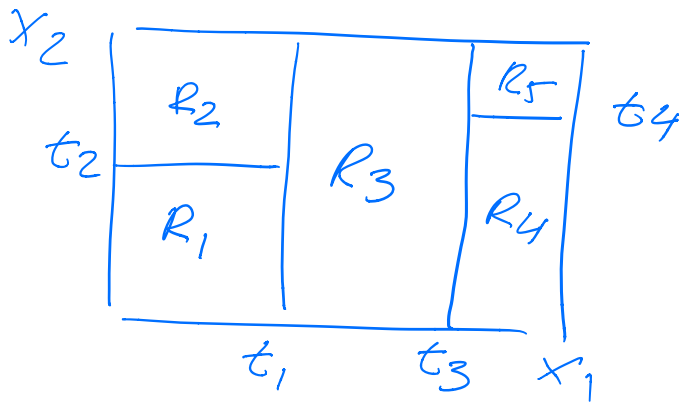
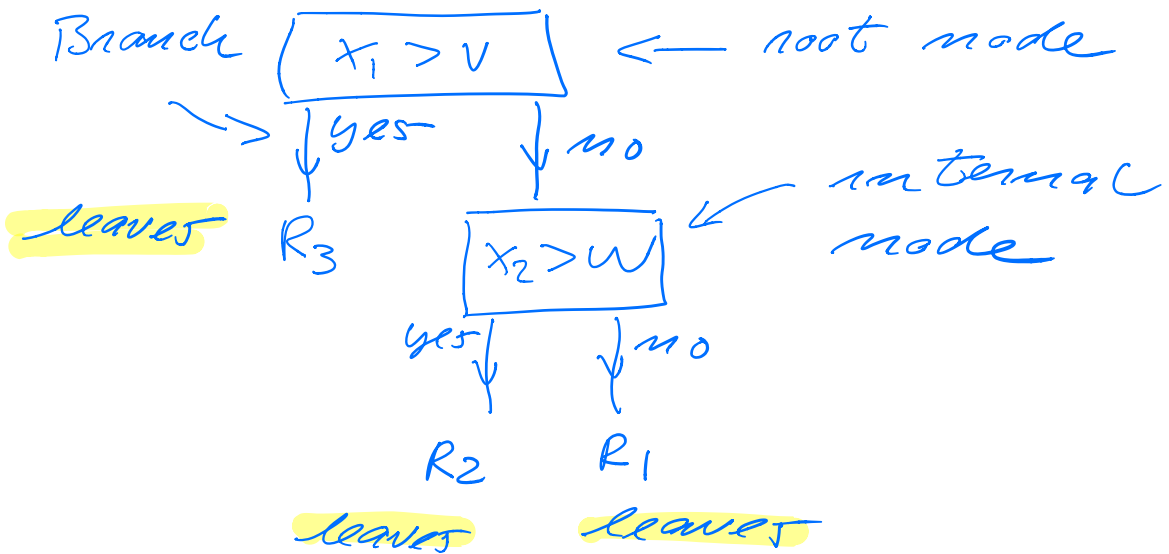
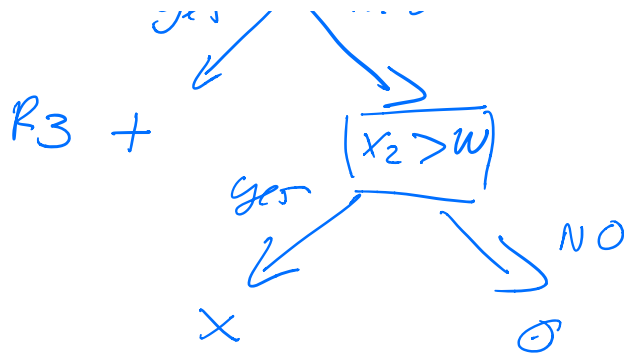
$$X = \underbrace{U}_{n \times n} \underbrace{S}_{n \times p} \underbrace{V^T}_{p \times p}$$

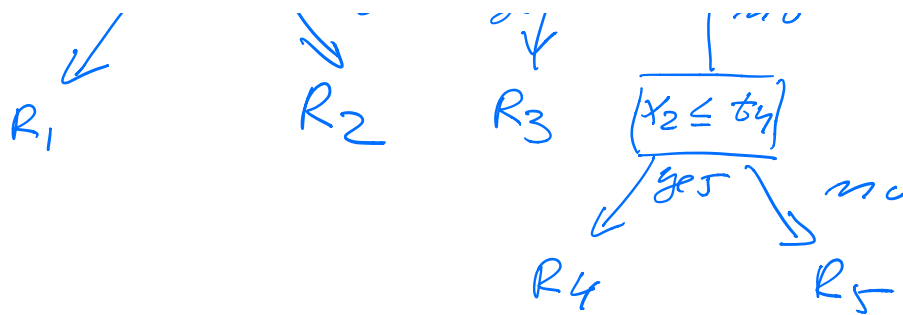
Trees, Forests, jungles and more

Decision Trees, Regression:



if  $x_1 > v$   
 then ...





## Algorithm for regression case

- i) Divide the space into sets of possible values for  $x_1, x_2, \dots, x_p$  into  $K$  distinct and non-overlapping regions  $R_1, R_2, \dots, R_K$
- (ii) For every observation that falls into region  $R_i$ , the prediction is given by the mean value of the observations in  $R_i$

Example:  $R_1, R_2$

if mean of  $R_1 = \mu_1 = 5$   
and  $R_2 = \mu_2 = 10$

an observation that falls in  $R_1$  gets a prediction of 5  
For  $R_2$ , the prediction is 10

..

+ +

How do we construct  
 $R_1, R_2 \dots R_K$ ?

Basic approach is to find  
boxes where we minimize

$$\sum_{j=1}^K \sum_{i \in R_j} (y_i - \bar{y}_{R_j})^2$$

$\bar{y}_{R_j}$  = mean value within  
the  $j$ th-box,

Time consuming,

Rather: start with all observations  
and define single Region. Then  
successively split into smaller  
regions.

- Define a cut point  $s$   
Define Regions

$$R_1(j, s) \{x \mid x_j < s\} \text{ and } \underbrace{\{x \mid x_j \geq s\}}_{R_2(j, s)}$$

- MINIMIZE

$$\sum_{j \in R_1} (y_j - \bar{y}_{R_1})^2 + \sum_{j \in R_2} (y_j - \bar{y}_{R_2})^2$$