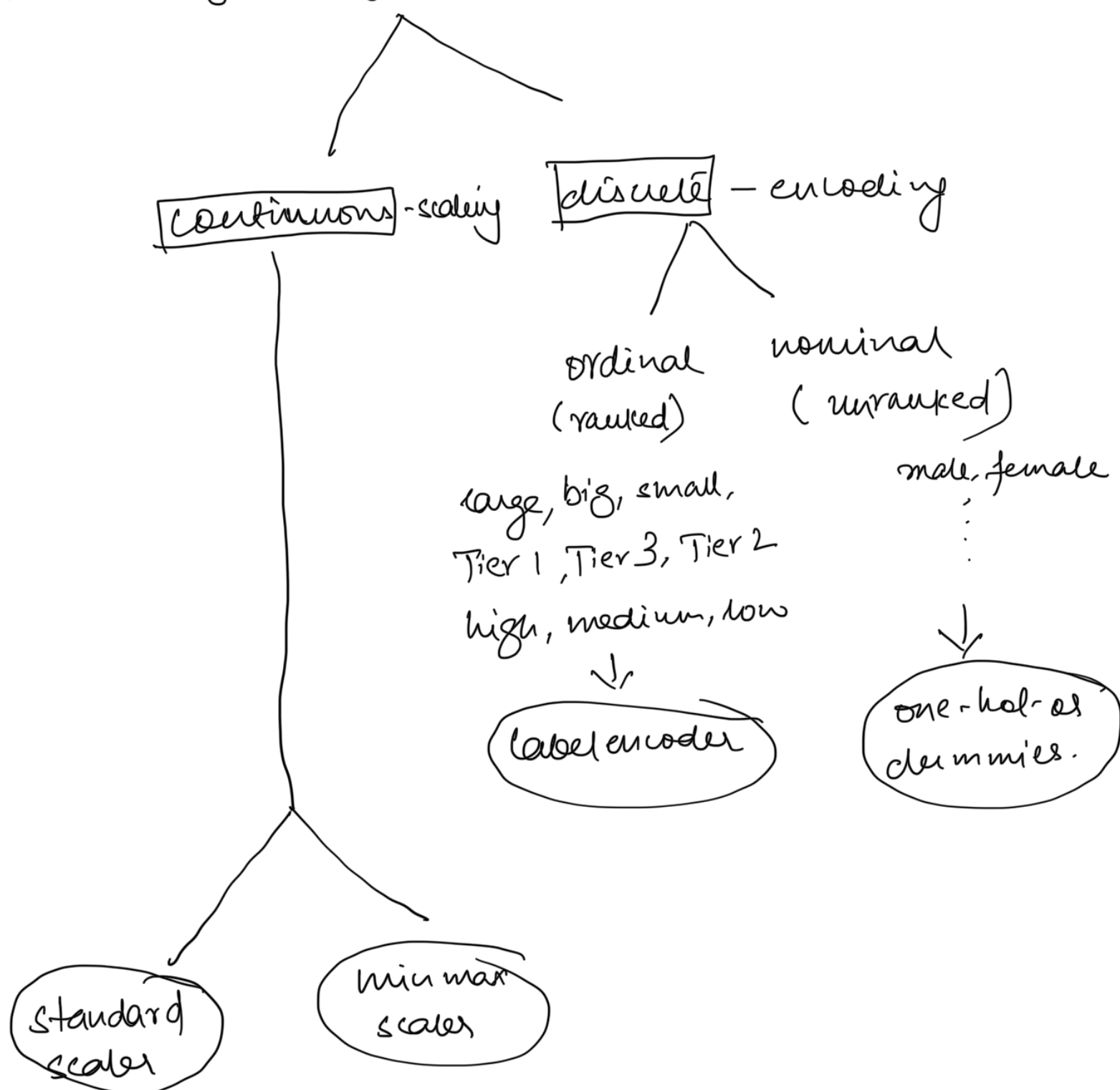


Feature engineering :-



$$\bar{x} = \frac{x - \mu}{\sigma}$$

$$\begin{cases} \text{mean} = 0 \\ \text{variance} = 1 \end{cases}$$

$$\bar{x} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

$$\begin{cases} \text{max} = 1, \text{min} = 0 \end{cases}$$

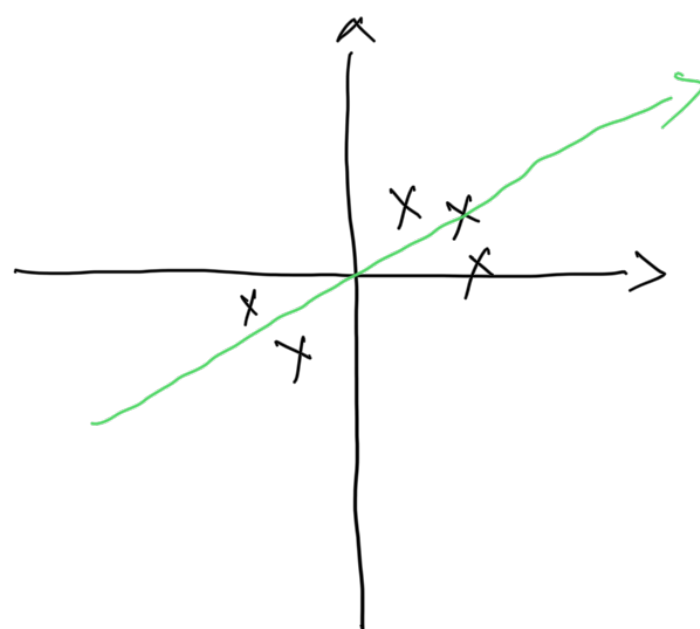
Dimensionality Reduction :-

Dimension \rightarrow No. of features
input shape

- 1) Feature selection.
- 2) PCA X (Feature summarization).

Principal Component Analysis :-

$$2d \rightarrow 1d$$



- ① Scale to origin
- ② Best fit line.

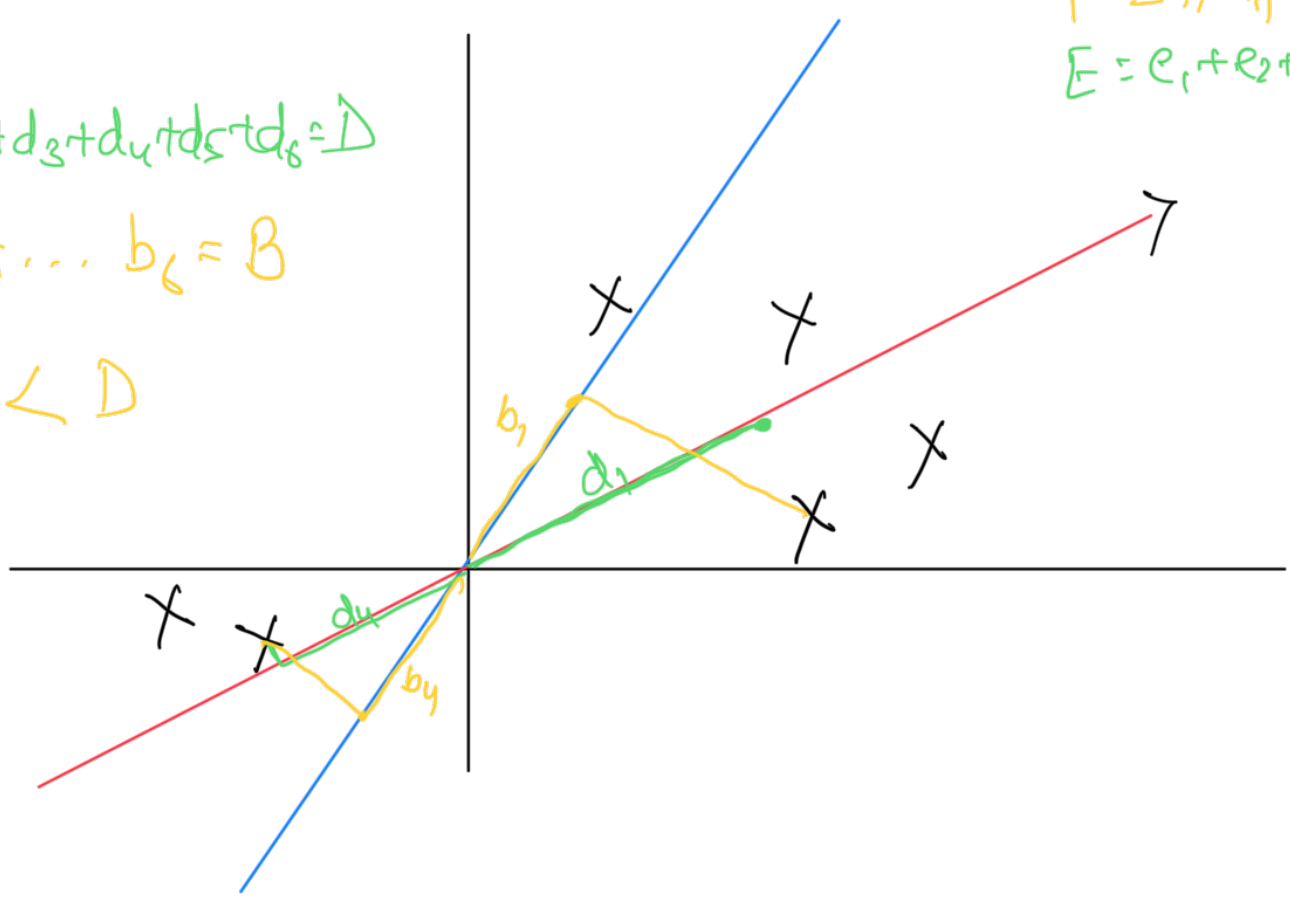
$$d_1 + d_2 + d_3 + d_4 + d_5 + d_6 = D$$

$$b_1 + b_2 + \dots + b_6 = B$$

$$B < D$$

$$F = f_1 + f_2 + f_3 + f_4 + f_5 + f_6$$

$$E = e_1 + e_2 + e_3 + e_4 + e_5 + e_6$$



$$\sqrt{E} < \sqrt{F}$$

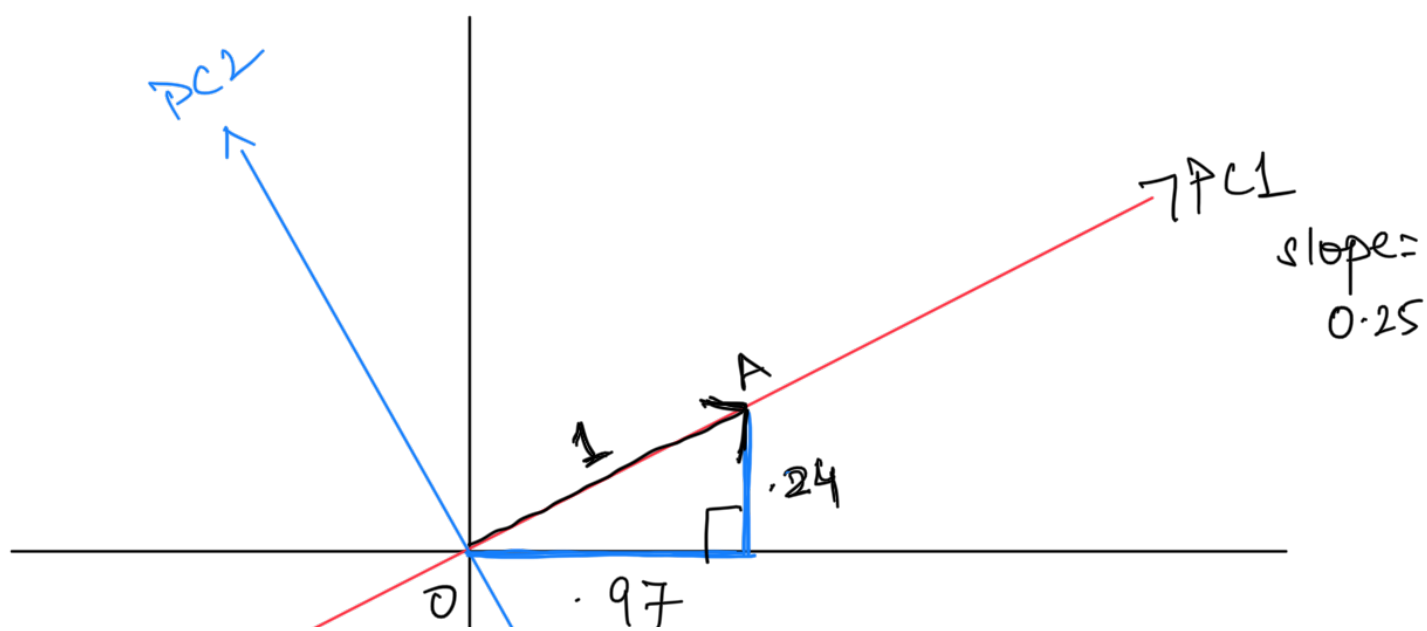
$$SS(E) < SS(F)$$

A best fit line is one

- a. whose SS(E) is less.
 $E < F \therefore$ Red is best fit line.
- b. where (SSD) is more (variance)
 $B < D \therefore$ Red is best fit line.

$$d_1^2 + d_2^2 + d_3^2 + \dots + d_n^2 = SSD$$

This Best fit line is PC1

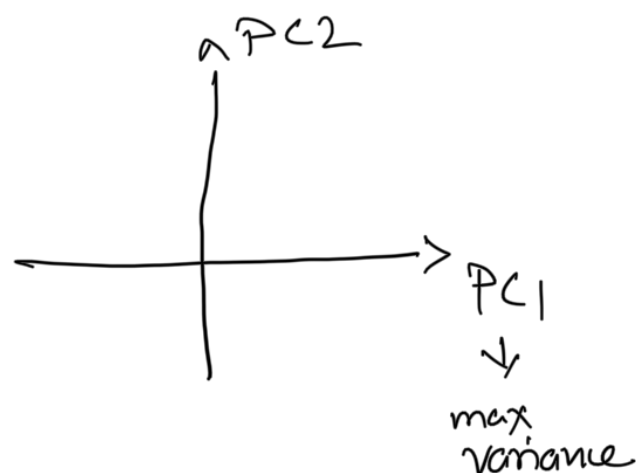
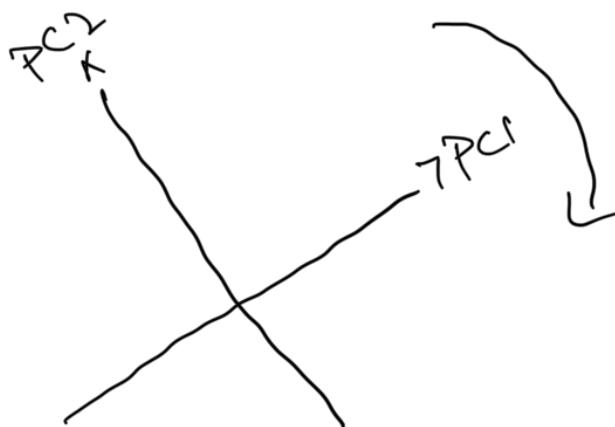


$4x + 1y \Rightarrow$ linear combination.

$\vec{OA} \rightarrow$ unit vector
or
Eigenvector of PC1

and

$$d_1^2 + d_2^2 + \dots + d_n^2 = \text{eigenvalue of PC1}$$



$$\text{variation of } \underline{PC1} = \frac{d_1^2 + \dots + d_n^2}{n-1} = \underline{15}$$

$$\boxed{15/18\%}$$

$$\text{" } \underline{PC2} = \frac{b_1^2 + \dots + b_n^2}{n-1} = \underline{3}$$

$$\boxed{3/18\%}$$

10 dim \rightarrow 10 PCs. Plot variance ratio

$$\begin{bmatrix} 0.28 & 0.15 & 0.10 & 0.08 & \dots \end{bmatrix}_{1 \times 10}$$

> > > > >

$$(0.28 + 0.15 + 0.10 + 0.08)$$

$$n\text{-components} = 0.80$$

\downarrow
10 PCs it will give $n \text{ PCs} < 10$
or > 0.80

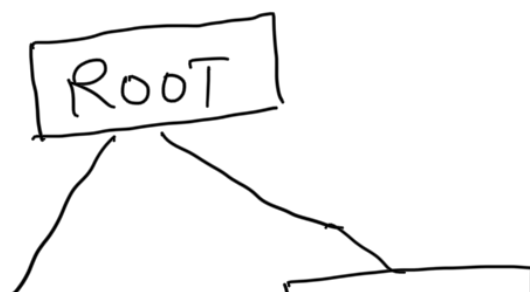
$$\sum_{i=1}^n x_i^2 \approx 0.00$$

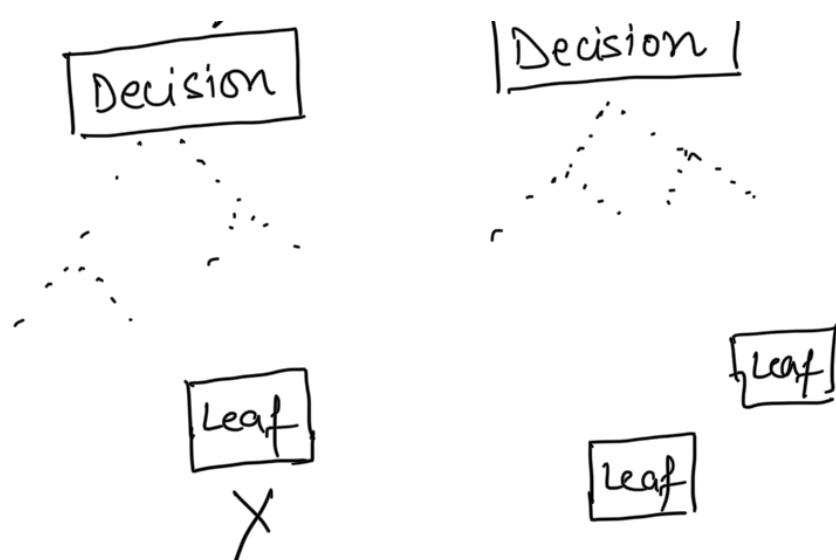
- ① Scale origin
- ② find best fit line
min error or
max var ✓
- ③ PC1 is best fit line
- ④ Normalise PC1
↳ ev of PC1
↳ value of PC1
- ⑤ PC2 orthogonal to PC1
- ⑥ Rotate the PC1 and PC2
- ⑦ Calculate variance ratio
 $PC1 > PC2$
- ⑧ Using a threshold select no. of PC's.

Total dim > no. of PC's

2 → 1
20 → 5/4/7
250 → 2/20/100...

Tree based methods.



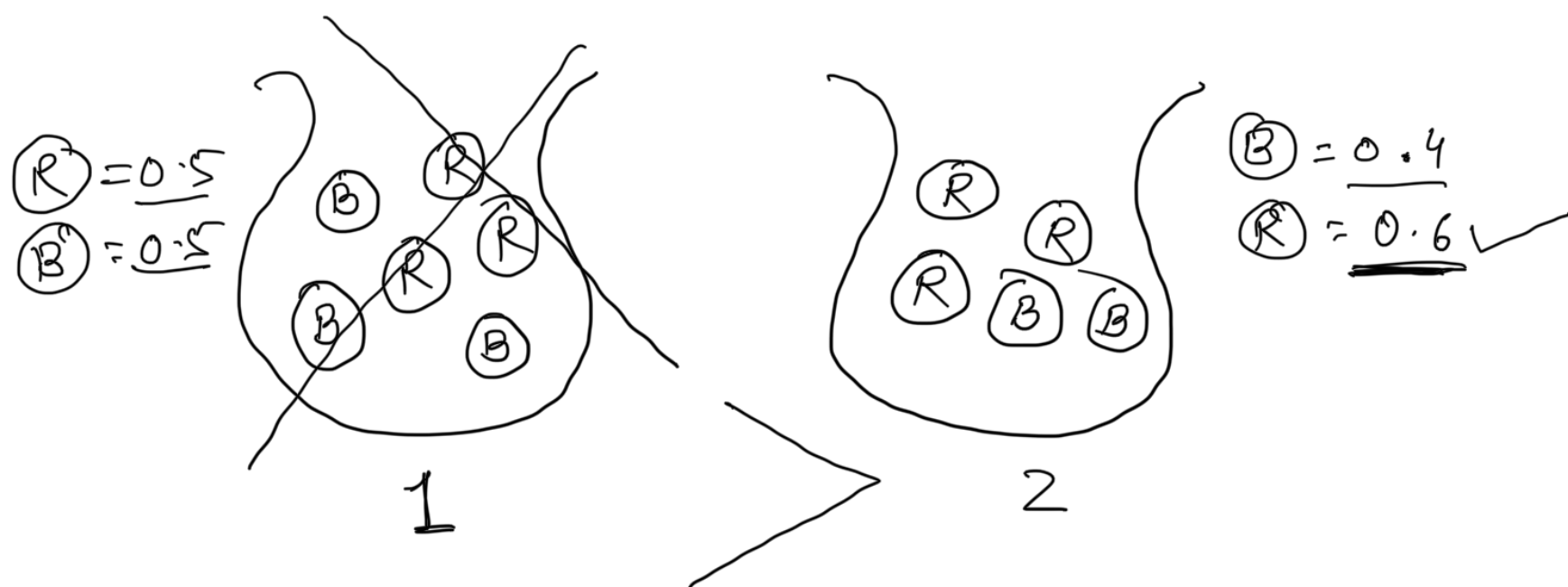


① What are these nodes and how do we split??

② The order of split...

Entropy

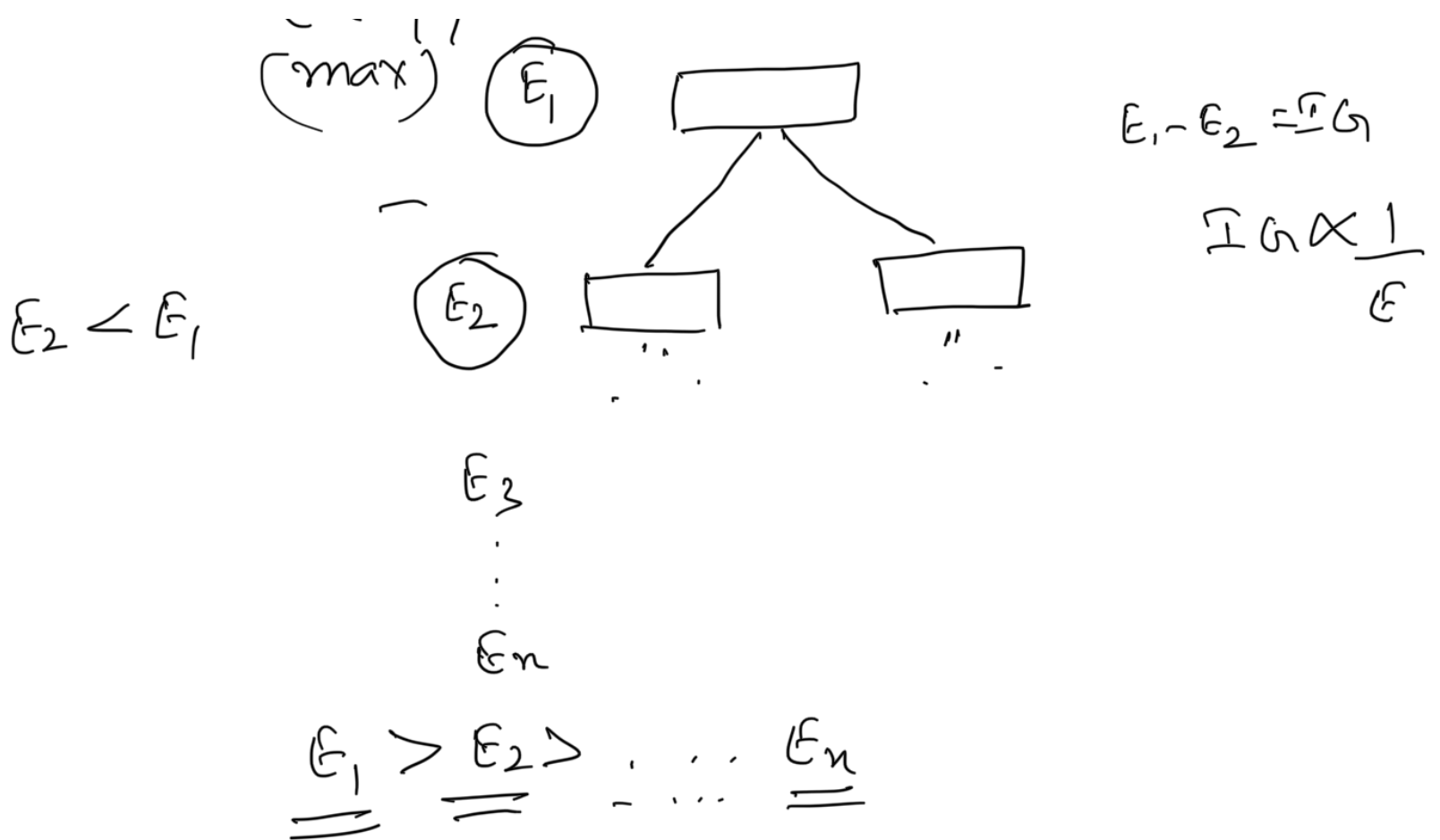
measure of impurity
uncertainty



① Choose any bag, any color (Red/Blue)

② If you draw a ball of the chosen color from the chosen bag you win 10,000\$. Which color, bag??

Entropy



$$H = - \sum_{i=1}^n P_i \log_2 P_i$$

Initial Entropy (9+, 5-)

① ②

$$H = - \sum_{i=1}^2 P_i \log_2 P_i$$

$$= - [P_1 \log_2 P_1 + P_2 \log_2 P_2]$$

$$= - \left[\frac{9}{9+5} \log_2 \frac{9}{9+5} + \frac{5}{9+5} \log_2 \frac{5}{9+5} \right]$$

$$= - \left[\frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} \right]$$

$$= 0.940$$

$$H = 0.940$$

Hum

windy.