

lecture 7/01 : Features (recap) NO CLASS Friday July 4th

- PCA new feat = linear combination (original feat)
- Feature selection using L1-Ridge
- Bootstrap margins, margin per feature

~~BINARY Margin~~ Margin (x) = $y \cdot \text{classifier score} = y \cdot \frac{F(x)}{\text{Score}}$

$y \in \{-1\}$

if $y=+1$: want large score $> 0 \Leftrightarrow$ large margin

if $y=-1$: want large neg score $\leq 0 \Leftrightarrow$ large margin

if $F(x) \approx 0 \Leftrightarrow$ low margin \Leftrightarrow small confidence

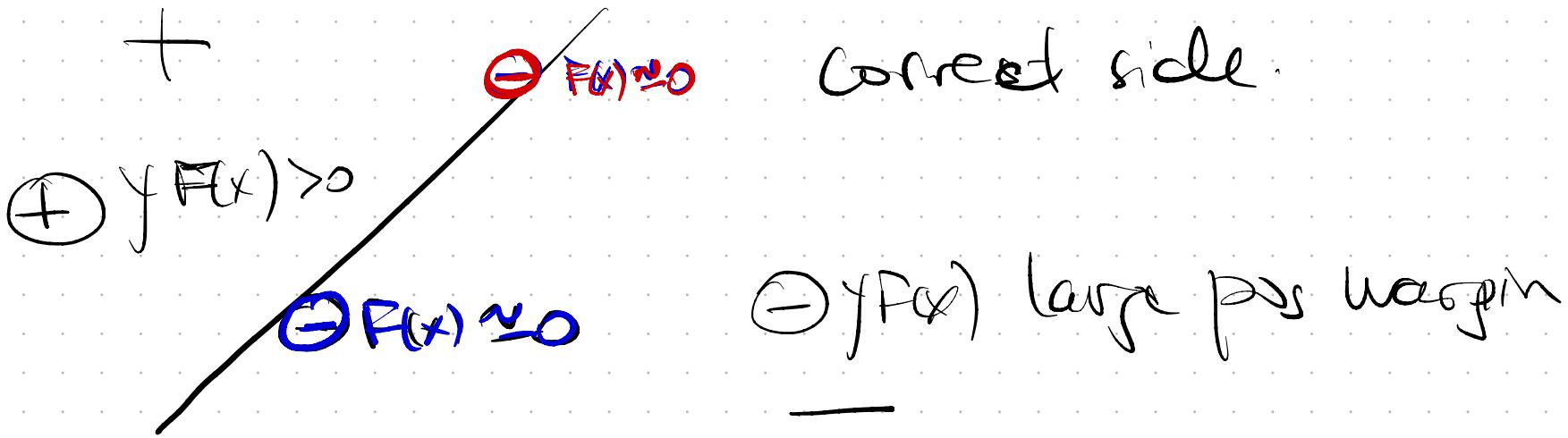
if $yF(x) < 0 \Leftrightarrow$ mistake classification

if margin $\ll 0 \Leftrightarrow$ mistake with high confidence

Quality of classifier \leftrightarrow large (pos) margins
margin better measured than accuracy

" how far
we predict

from 50-50
confidence



Increasing margin \leftrightarrow increase confidence.

GB, AdaBoost : increase margins $YF(x)$ during training
even for points classified correctly
already $YF(x) \neq 0$.

clement
boosting
Alg

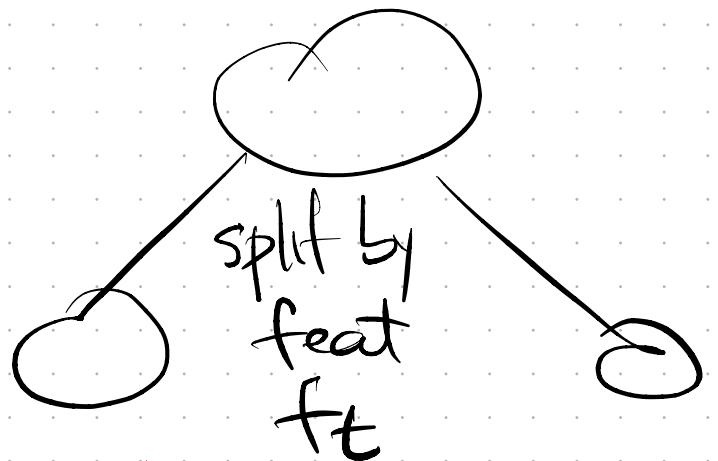
$$\text{margin} = e^{YF(x)}$$

$$\frac{1}{1 + e^{-YF(x)}}$$

Margin Analysis GB: $F(x) = \sum_{t=1}^T h_t(x)$

Adaboost: $F(x) = \sum_{t=1}^{T+1} \alpha_t h_t(x)$

h_t = dec. stump



Reorder the summation per feature

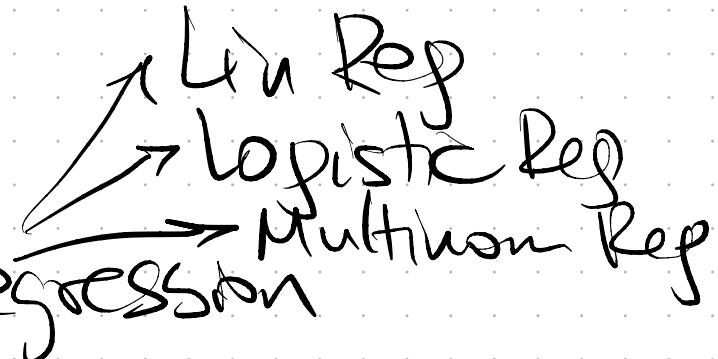
$$F(x) = \left[\sum_{f_t=f_1} \alpha_t h_t(x) \cdot y \right] + \left[\sum_{f_t=f_2} \alpha_t h_t(x) \cdot y \right] + \left[\sum_{f_t=f_3} \alpha_t h_t(x) \cdot y \right] + \dots + \left[\sum_{f_t=f_D} \alpha_t h_t(x) \cdot y \right]$$

margin_{f₁}(x) margin_{f₂}(x) margin_{f₃}(x) ... margin_{f_D}(x)

\Rightarrow contribution of each feature to the classifier

Note online, optional PB HW4.

HW4 Feature Selection using L1-reg Regression



- Use very large L_1 penalty train Regression
 - ⇒ many Feat. end up with 0 or close-to-0 wtf.
 - L_1 penalty \Rightarrow sparsity. (acc not that good, irrelevant)

- Select $G = 100$ features by $| \text{coef} |$ largest 200

- Train (real) classifier without L_1 penalty.

ex L₂-reg Repession

ex NNet

ex. Naive Bayes, GDA,

not good: Dec Tree, because D.T. have "natural" sparsity by # splits

PCA Recap ① rotates data to align dims with axis

Σ = sigma = covar - dim are sorted by $\text{var}[\text{dim}]$

Σ = summation - each dim = projector = linear comb (original feat)

② Selects first R dim \Rightarrow losing $\text{var}[\cdot]$ information

$\mu = \sum_{i=1}^N x_i$ existing on dim [R+1 : D]

data centered $X = X - \mu$. Each column = feat has 0-mean.

$$\sum_{DXD} = \text{covar}(X) = \frac{1}{N} \begin{matrix} X^T \\ DXN \end{matrix} \begin{matrix} X \\ NxD \end{matrix} =$$

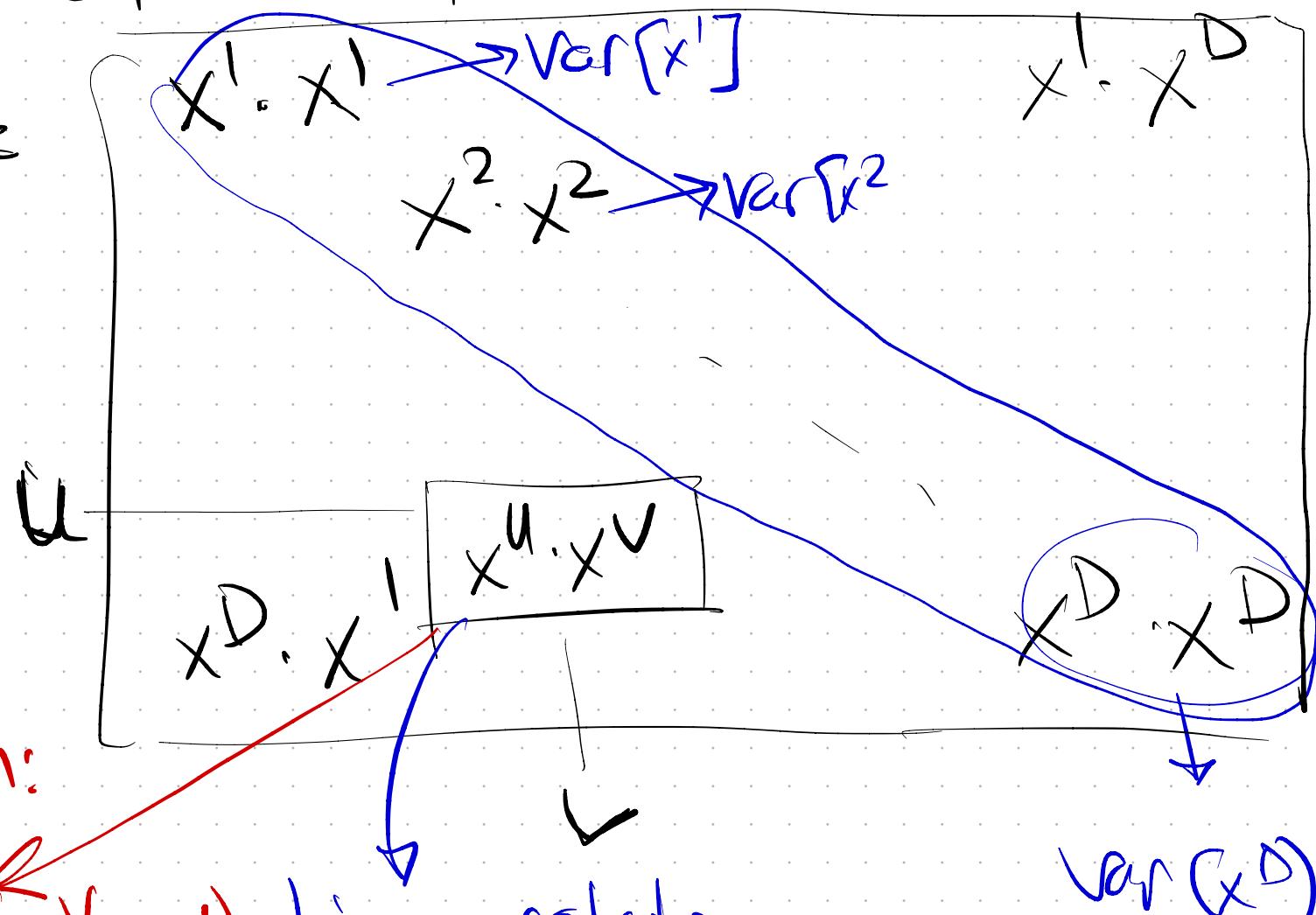
orthonormal vectors

$$e_1, e_2, \dots, e_D$$

$$\|e_i\|=1$$

$$\|e_i \cdot e_j\|=0 \text{ (perpendicular)}$$

$$\sum_{i=1}^N (x_i^u - \mu^u)(x_i^v - \mu^v) \quad \text{lin correlation (feat } u \text{, feat } v)$$



- Σ = nice math matrix \Rightarrow orthonormal decomposition
 - sym
 - pos semidef
 - form $A^T A$
- $(e_1, \lambda_1) (e_2, \lambda_2)$... - (e_D, λ_D) eigen vector-values
 • $e\Sigma = \lambda e$ "e does not change direction mult by Σ "

$$\Sigma_{\text{covar}} = E \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{bmatrix} E^T$$

Eigen vectors = cols Eigen values Eigenvec = rows

Notes: Why $\lambda \geq 0$
 $\lambda_1 > \lambda_2 > \dots > \lambda_D > 0$

① Rotation PCA dim are obtained by projecting x on e_1, e_2, e_D

first dim

$$X \cdot e_1 = X_{\text{new}}^1$$

~~NXD DX1 NX1~~

2nd dim

$$X \cdot e_2 = X_{\text{new}}^2$$

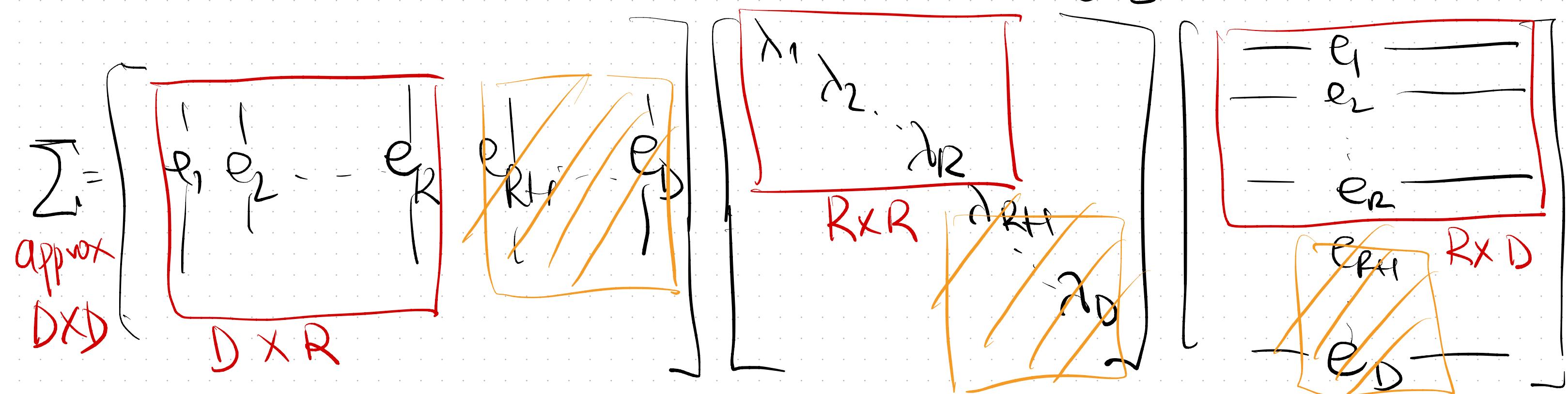
$\Rightarrow X_{\text{new}}^1 = \text{lin combination}$
 of orig feat with coef
 given by $e_1 = [e_1^1 \ e_1^2 \ \dots \ e_1^D]$

② R dim R << D (for ex D=1000 R=20)

$$X \cdot \begin{bmatrix} e_1^T & e_2^T & \dots & e_{20}^T \end{bmatrix} = X_{\text{new}}$$

$N \times D$ $D \times 20$ $N \times 20$

Theory: $R=20$ first PCA-dim \Leftrightarrow approx Σ with first 20 (eigenvectors) eigen vals



Rank deficient at most R

Information lost $\propto \sum_{\text{approx}} \geq \sum$ depends on magnitude
of ignored eigenvalues.

- $\lambda_{R+1} \approx \lambda_{R+2} \dots \approx \lambda_D$ close to 0 $\Rightarrow \sum_{\text{approx}}$ very good

$R=? \Leftarrow$ don't want miss on large eigenvalues.

- if $\lambda_{R+1} = \lambda_{R+2} = \dots = \lambda_D = 0$ (actually 0) $\Rightarrow \sum_{\text{approx}} = \sum$