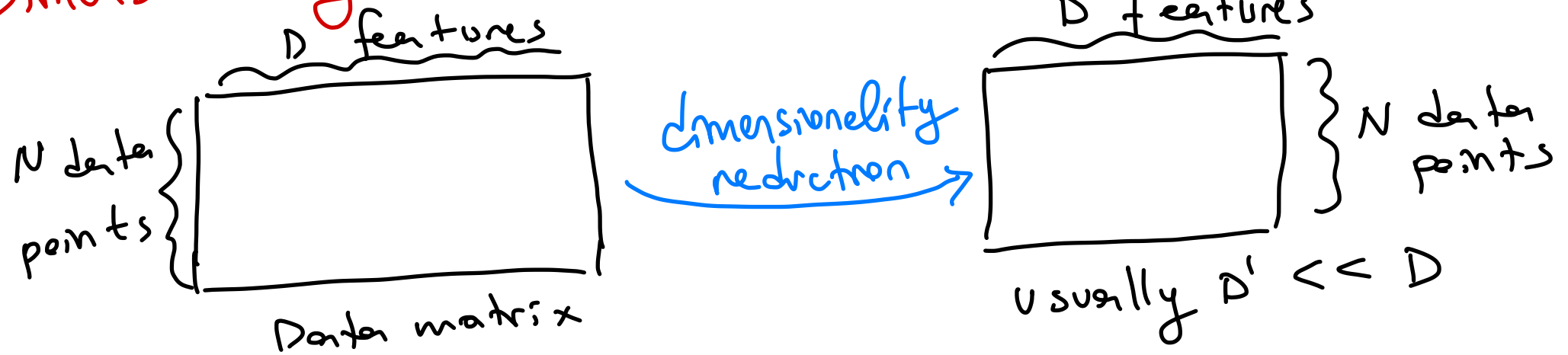


Dimensionality Reduction



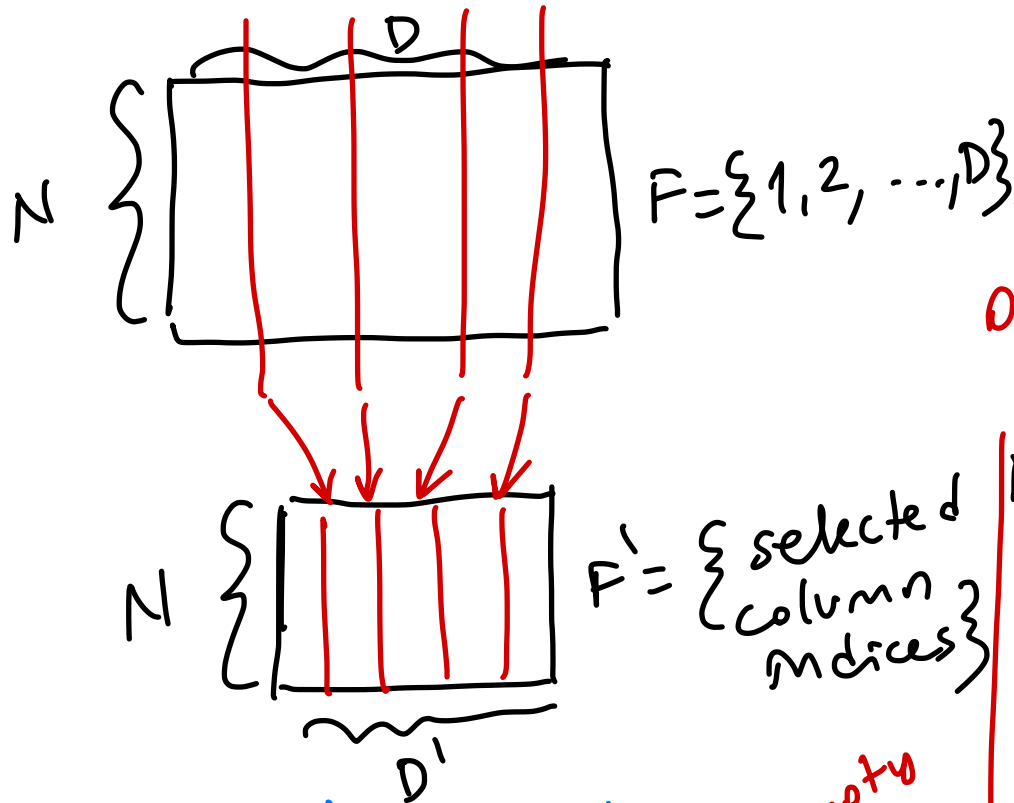
Reasons

- ① To reduce computational complexity
- ② To reduce storage complexity
- ③ To reduce data acquisition cost
- ④ To increase robustness
- ⑤ To increase interpretability
- ⑥ To enable visualization (when $D'=2$ or $D'=3$)

⋮

Feature Selection

$\mathcal{X} = \{x_i\}_{i=1}^N$ where $x_i \in \mathbb{R}^D$
 We will select a subset of
 $\{1, 2, \dots, D\}$



of possible subsets
 $= 2^D - 1 - 1$
 (empty set)
 (full set)

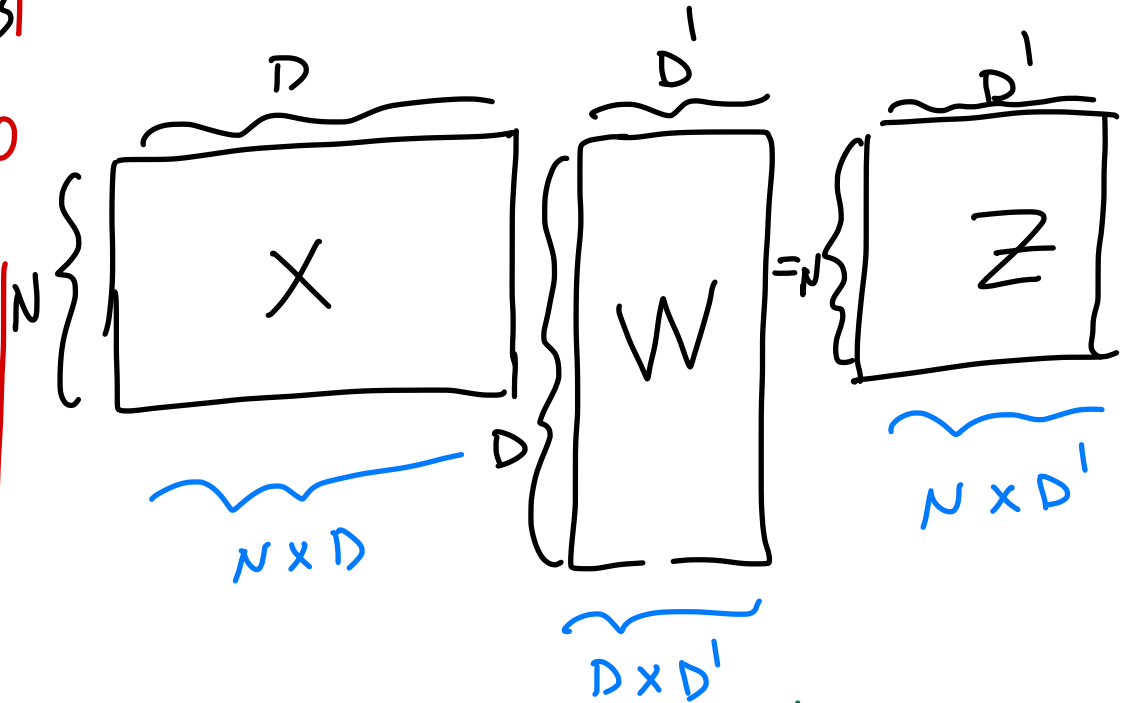
Feature Extraction

$\mathcal{X} = \{x_i\}_{i=1}^N$ where $x_i \in \mathbb{R}^D$

$x_i \in \mathbb{R}^D \xrightarrow{\text{projection}} z_i \in \mathbb{R}^{D'}$

$$z_i = W^T \cdot x_i$$

$[D \times 1] \quad [D \times D] \quad [D \times 1]$



$W \Rightarrow$ model parameters

① Forward Selection

- $F' = \emptyset$ (empty set)

- At each iteration, find the best new feature to be added to F'

$$d^* = \arg \min_d \text{Error}(F' \cup d)$$

↓
union

- Add d^* to F' if $\text{Error}(F' \cup d) < \text{Error}(F')$

$F' = \emptyset$

$t=1 \Rightarrow$ 1 2 3 4 5 6 $\Rightarrow F' = \{1\}$

$t=2 \Rightarrow$ X $\{1,2\}$ $\{1,3\}$ $\{1,4\}$ $\{1,5\}$ $\{1,6\}$ $\Rightarrow F' = \{1,4\}$

if $E(\{1,4\}) < E(\{1\}) \Rightarrow \text{YES}$

$t=3 \Rightarrow$ X $\{1,4,2\}$ $\{1,4,3\}$ X $\{1,4,5\}$ $\{1,4,6\}$ $\Rightarrow F' = \{1,4,5\}$

if $E(\{1,4,5\}) < E(\{1,4\}) \Rightarrow \text{YES}$

$t=4 \Rightarrow$ X $\{1,4,5,2\}$ $\{1,4,5,3\}$ X X $\{1,4,5,6\}$ $\Rightarrow F' = \{1,4,5\}$

if $E(\{1,4,5,2\}) < E(\{1,4,5\}) \Rightarrow \text{NO}$

of subsets = $2^6 - 1 - 1 = 62$ models

we trained only 18 models.

separate set of data points other than the training data

↓
STOP

② Backward Elimination

- $F' = F$ (full set)
- At each iteration, find the best feature to be removed from F' . $d^* = \arg \min_d \text{Error}(F'/d)$ ↪ set difference
- Remove d^* from F' if $\text{Error}(F'/d) < \text{Error}(F')$

$$F' = \{1, 2, 3, 4, 5, 6\}$$

$$t=1 \Rightarrow \underbrace{\{2, 3, 4, 5, 6\}}_{\text{remove 1}} \quad \underbrace{\{1, 3, 4, 5, 6\}}_{\text{remove 2}} \quad \underbrace{\{1, 2, 4, 5, 6\}}_{\text{remove 3}} \quad \underbrace{\{1, 2, 3, 5, 6\}}_{\text{remove 4}} \quad \underbrace{\{1, 2, 3, 4, 6\}}_{\text{remove 5}} \quad \underbrace{\{1, 2, 3, 4, 5\}}_{\text{remove 6}}$$

if $\text{Error}(\{1, 3, 4, 5, 6\}) < \text{Error}(\{1, 2, 3, 4, 5, 6\}) \Rightarrow \text{YES}$

$$F' = \{1, 3, 4, 5, 6\}$$

$$t=2 \Rightarrow \underbrace{\{3, 4, 5, 6\}} \quad \times \quad \underbrace{\{1, 4, 5, 6\}} \quad \underbrace{\{1, 3, 5, 6\}} \quad \underbrace{\{1, 3, 4, 6\}} \quad \underbrace{\{1, 3, 4, 5\}}$$

if $\text{Error}(\{1, 4, 5, 6\}) < \text{Error}(\{1, 3, 4, 5, 6\}) \Rightarrow \text{NO}$

$$F' = \{1, 3, 4, 5, 6\}$$

↓
STOP

We trained only 12 models.

Principal Component Analysis (PCA)

- PCA is a feature extraction algorithm

$$x \in \mathbb{R}^D \quad z \in \mathbb{R}^{D'} \quad W \in \mathbb{R}^{D \times D'}$$

2 → 1
D → D'

We would like to find the direction that maximizes the variance.

$$\begin{aligned} \text{VAR}(z) &= \text{VAR}(W^T \cdot x) \\ &= W^T \cdot \text{VAR}(x) \cdot W \\ &= W^T \cdot \Sigma_x \cdot W \end{aligned}$$

optimize these

$\begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \rightarrow$ Covariance matrix of original data points

$$\begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}$$

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{N1} & x_{N2} \end{bmatrix}$$

$$W = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

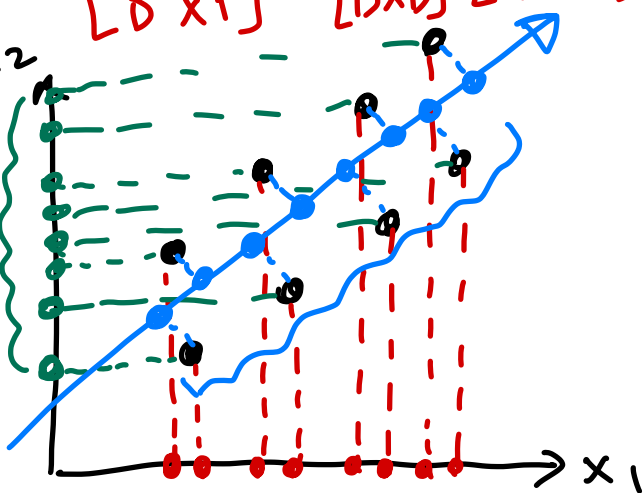
$$= \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{N1} \end{bmatrix}$$

$$X = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{N1} & x_{N2} \end{bmatrix}$$

$$W = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{N2} \end{bmatrix}$$

range of the projected points



range of the projected data points

maximize $\text{VAR}(z) = w^T \cdot \Sigma_x \cdot w$

assume w^* is the optimum solution.

$\tilde{w} = 2 \cdot w^*$

$(\tilde{w})^T \Sigma_x (\tilde{w}) = (2w^*)^T \cdot \Sigma_x (2w^*) > (w^*)^T \cdot \Sigma_x (w^*)$

$\left. \begin{aligned} 2x_1 + 3x_2 + 7 &= 0 \\ 4x_1 + 6x_2 + 14 &= 0 \end{aligned} \right\}$ are the same lines.

maximize $w^T \cdot \Sigma_x \cdot w$

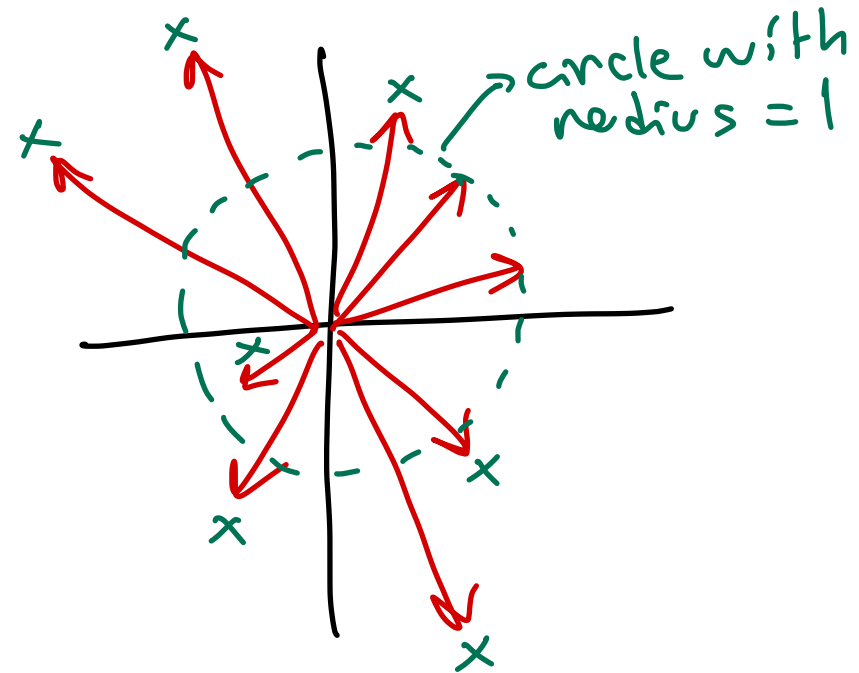
subject to: $\|w\|_2^2 = 1 \quad] \alpha$

$$L_P = w^T \cdot \Sigma_x w - \alpha (\|w\|_2^2 - 1)$$

$$= w^T \cdot \Sigma_x w - \alpha (w^T w - 1)$$

$$\frac{\partial L_P}{\partial w} = 2 \cdot \Sigma_x \cdot w - 2 \cdot \alpha \cdot w = 0$$

$$\cancel{2} \cdot \Sigma_x \cdot w = \cancel{2} \cdot \alpha \cdot w$$



$Ax = b$

- A : matrix
- x : vector
- b : vector
- \rightarrow : scalar
- x : eigenvector
- b : eigenvalue

D eigenvalues

$$\alpha_1, \alpha_2, \dots, \alpha_D \Rightarrow \alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_D$$

w^* \Rightarrow the eigenvector that corresponds to the largest eigenvalue (α_1) [the first eigenvector]

Exercise: If $D' = 2 \Rightarrow$ we need to pick the first two eigenvectors. Prove that.

$$W^* = \begin{bmatrix} | & | \\ w_1 & w_2 \\ | & | \end{bmatrix}$$

first eigenvector

second eigenvector

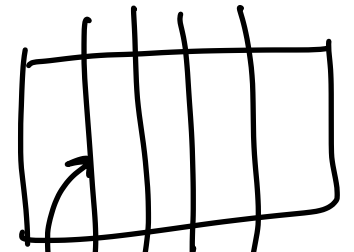
PCA Algorithm

Step 1: Calculate Σ_x .

Step 2: Find first D' eigenvectors of Σ_x .

eigenvectors that correspond to D' largest eigenvalues

Projection Step: $\underbrace{z_i}_{D \times 1} = \underbrace{W^T}_{D \times D} \cdot \underbrace{(x_i - \hat{\mu})}_{D \times 1} \quad \forall i$

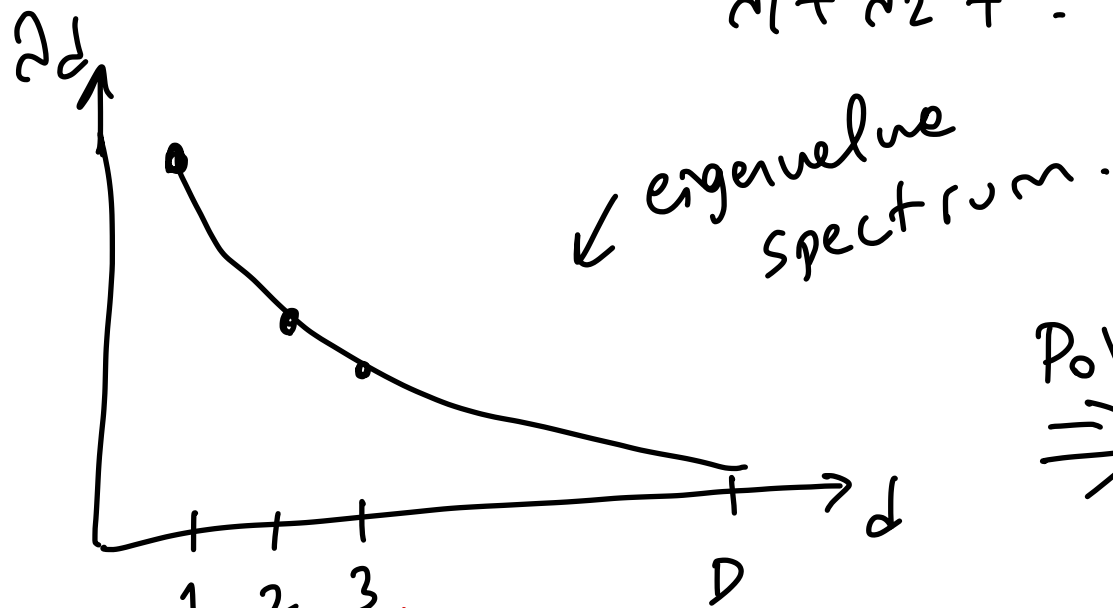


$$W = \begin{bmatrix} | & \dots & | \\ w_1 & \dots & w_{D'} \\ | & \dots & | \end{bmatrix}$$

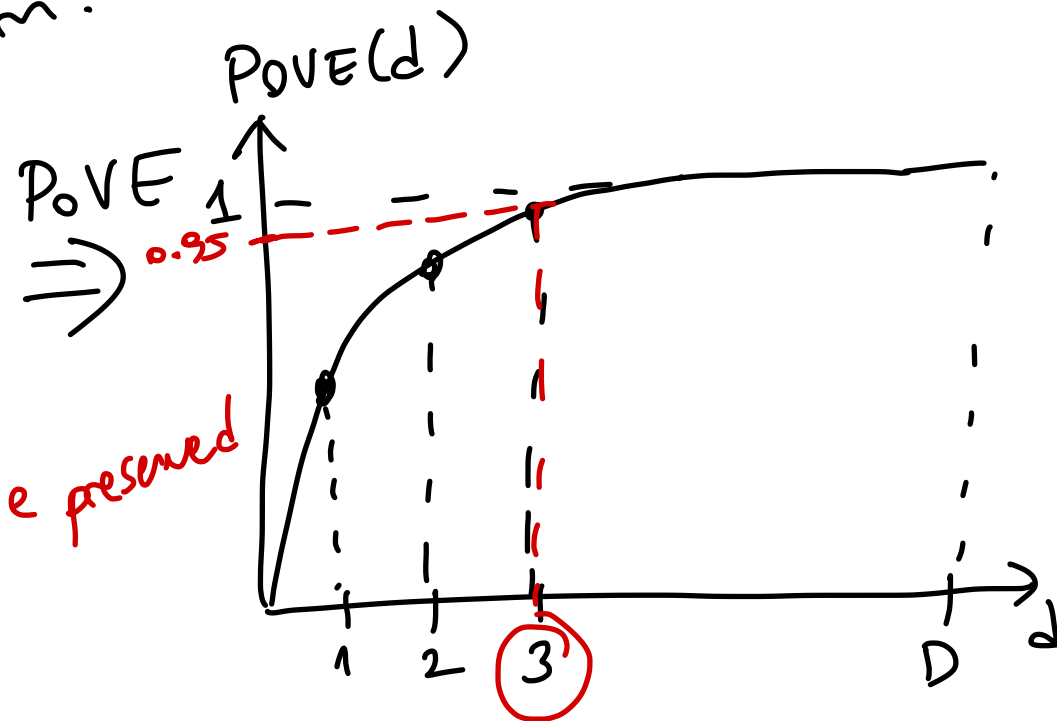
How to pick D' ? Proportion of Variance Explained (PoVE)

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_D \geq 0$$

$$\text{PoVE}(D') = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_{D'}}{\lambda_1 + \lambda_2 + \dots + \lambda_D}$$



Rule of thumb:
at least 95% of variance should be preserved



3 dimensions should be enough.