

Machine Learning for Data Science (CS4786)

Lecture 8

Mixture Models, Dimensionality Reduction

Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2017fa/>

TOWARDS HARD GAUSSIAN MIXTURE MODEL

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$, ellipsoids $\hat{\Sigma}_j^0$ and initial proportions π^0 randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1})^\top (\hat{\Sigma}_j^{m-1})^{-1} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1}) - \log(\pi_j^{m-1})$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t \quad \hat{\Sigma}_j^m = \frac{1}{|C_j^m|} \sum_{t \in C_j^m} (\mathbf{x}_t - \hat{\mathbf{r}}_j^m)(\mathbf{x}_t - \hat{\mathbf{r}}_j^m)^\top \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

TOWARDS HARD GAUSSIAN MIXTURE MODEL

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$, ellipsoids $\hat{\Sigma}_j^0$ and initial proportions π^0 randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1})^\top (\hat{\Sigma}_j^{m-1})^{-1} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1}) - \log(\pi_j^{m-1})$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t \quad \hat{\Sigma}_j^m = \frac{1}{|C_j^m|} \sum_{t \in C_j^m} (\mathbf{x}_t - \hat{\mathbf{r}}_j^m)(\mathbf{x}_t - \hat{\mathbf{r}}_j^m)^\top \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

TOWARDS HARD GAUSSIAN MIXTURE MODEL

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$, ellipsoids $\hat{\Sigma}_j^0$ and initial proportions π^0 randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1})^\top (\hat{\Sigma}_j^{m-1})^{-1} (\mathbf{x}_t - \hat{\mathbf{r}}_j^{m-1}) - \log(\pi_j^{m-1})$$
$$d(\mathbf{x}_t, C_j)$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t \quad \hat{\Sigma}_j^m = \frac{1}{|C_j|} \sum_{t \in C_j} (\mathbf{x}_t - \hat{\mathbf{r}}_j^m)(\mathbf{x}_t - \hat{\mathbf{r}}_j^m)^\top \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

GENERAL HARD MIXTURE MODEL

- For all $j \in [K]$, initialize π^0 and parameters $\theta_1, \dots, \theta_K$ randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \operatorname{argmin}_{j \in [K]} d(\mathbf{x}_t, \theta_j) - \log(\pi_j^{m-1})$$

- 2 For each $j \in [K]$, set new representative as

$$\text{compute } \theta_j \text{ for cluster } C_j \quad \& \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

Multivariate Gaussian

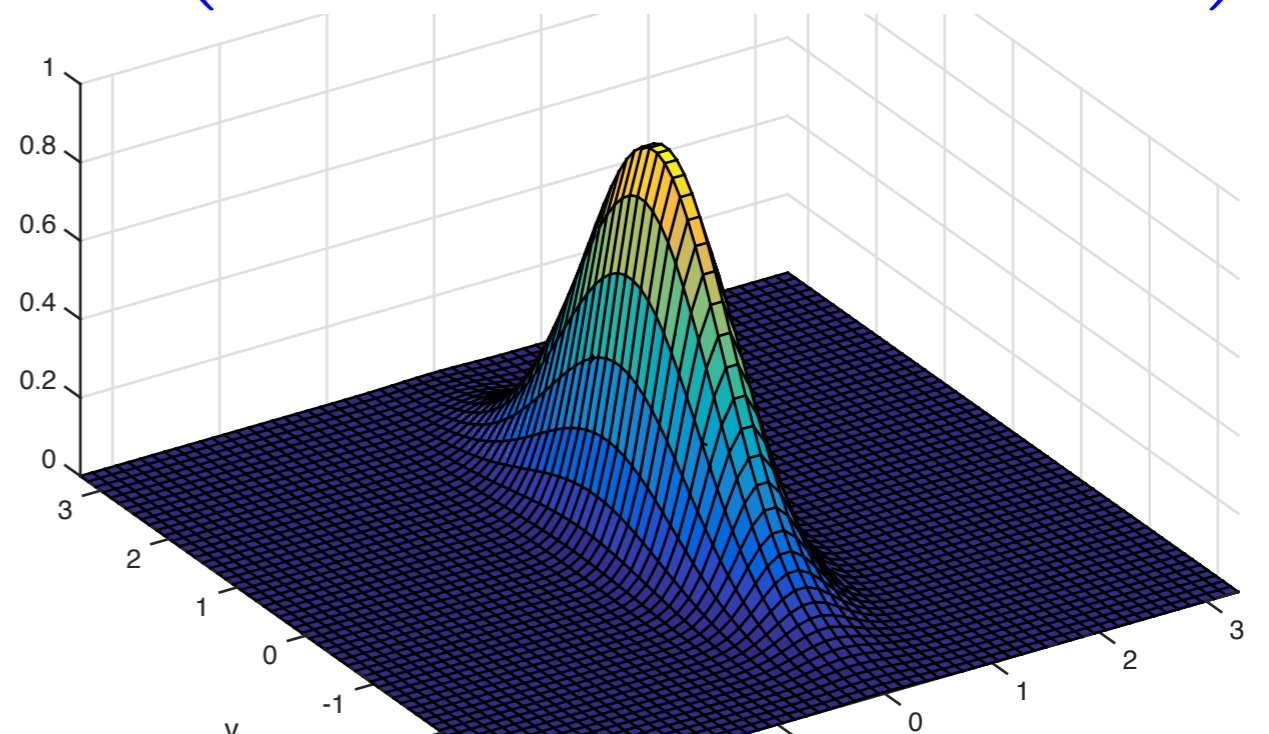
- Two parameters:
 - Mean $\mu \in \mathbb{R}^d$
 - Covariance matrix Σ of size $d \times d$

$$p(x; \mu, \Sigma) = (2\pi)^{d/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$

Multivariate Gaussian

- Two parameters:
 - Mean $\mu \in \mathbb{R}^d$
 - Covariance matrix Σ of size $d \times d$

$$p(x; \mu, \Sigma) = (2\pi)^{d/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$



HARD GAUSSIAN MIXTURE MODEL

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$, ellipsoids $\hat{\Sigma}_j^0$ and initial proportions π^0 randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \arg \max_{j \in [K]} p(\mathbf{x}_t, \hat{\mathbf{r}}_j^{m-1}, \hat{\Sigma}_j^{m-1}) \times \pi^m(j)$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{1}{|\hat{C}_j^m|} \sum_{\mathbf{x}_t \in \hat{C}_j^m} \mathbf{x}_t \quad \hat{\Sigma}_j^m = \frac{1}{|C_j^m|} \sum_{t \in C_j^m} (\mathbf{x}_t - \hat{\mathbf{r}}_j^m)(\mathbf{x}_t - \hat{\mathbf{r}}_j^m)^\top \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

GENERAL HARD MIXTURE MODEL

- For all $j \in [K]$, initialize π^0 and parameters $\theta_1, \dots, \theta_K$ randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$\hat{c}^m(\mathbf{x}_t) = \arg \max_{j \in [K]} p(\mathbf{x}_t, \theta_j) \times \pi_j^{m-1}$$

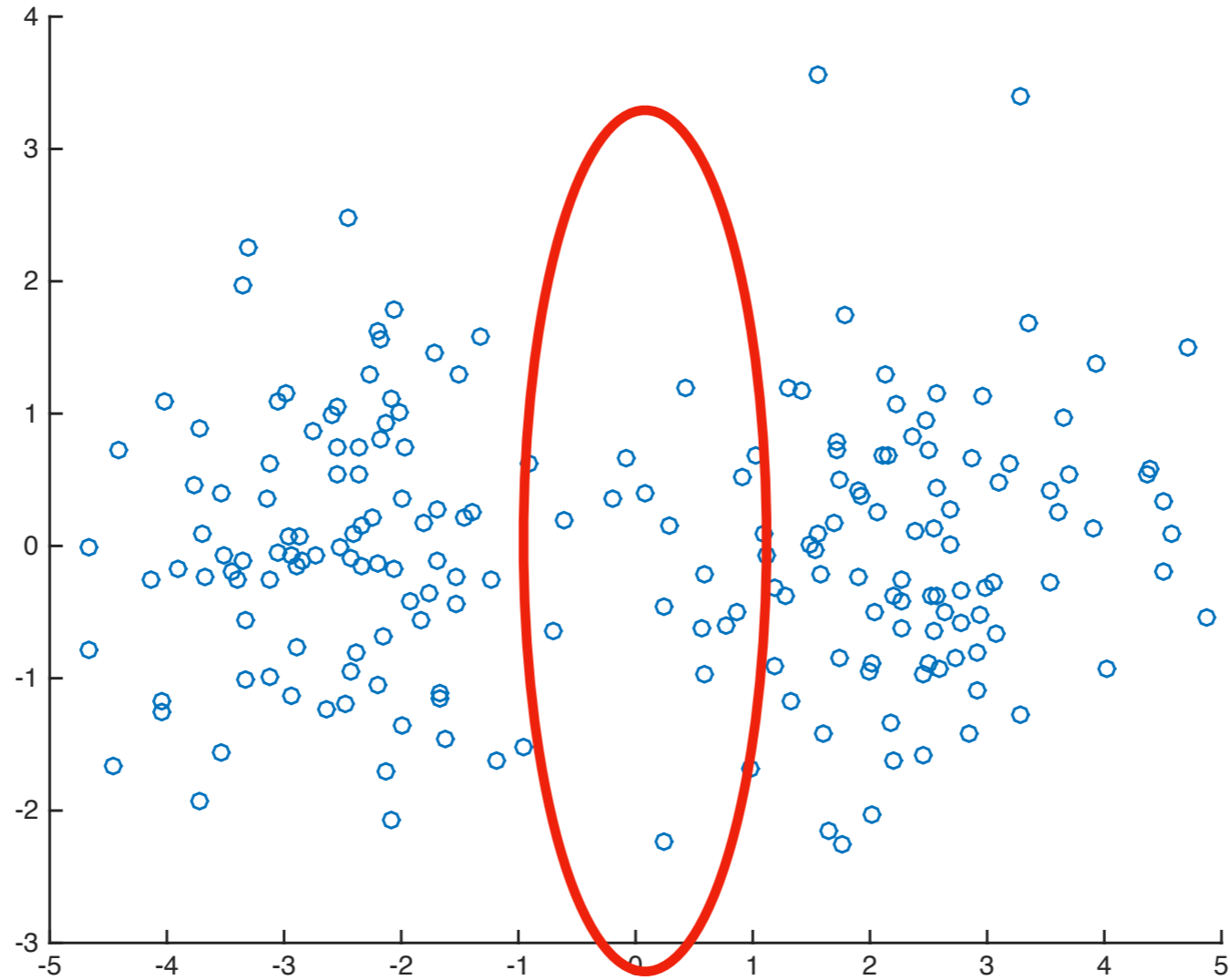
- 2 For each $j \in [K]$, set new representative as

$$\text{compute } \theta_j \text{ for cluster } C_j \quad \& \quad \pi_j^m = \frac{|C_j^m|}{n}$$

- 3 $m \leftarrow m + 1$

Demo

Pitfall of Hard Assignment



(SOFT) GAUSSIAN MIXTURE MODEL

- For all $j \in [K]$, initialize cluster centroids $\hat{\mathbf{r}}_j^0$ and ellipsoids $\hat{\Sigma}_j^0$ randomly and set $m = 1$
- Repeat until convergence (or until patience runs out)
 - 1 For each $t \in \{1, \dots, n\}$, set cluster identity of the point

$$Q_t^m(j) \propto p(\mathbf{x}_t, \hat{\mathbf{r}}_j^{m-1}, \hat{\Sigma}_j^{m-1}) \times \pi_j^{m-1}(j)$$

- 2 For each $j \in [K]$, set new representative as

$$\hat{\mathbf{r}}_j^m = \frac{\sum_{t=1}^n Q_t(j) \mathbf{x}_t}{\sum_{t=1}^n Q_t(j)} \quad \hat{\Sigma}_j^m = \frac{\sum_{t=1}^n Q_t(j) (\mathbf{x}_t - \hat{\mathbf{r}}_j^m) (\mathbf{x}_t - \hat{\mathbf{r}}_j^m)^\top}{\sum_{t=1}^n Q_t(j)}$$

$$\pi_j^m = \frac{\sum_{t=1}^n Q_t(j)}{n}$$

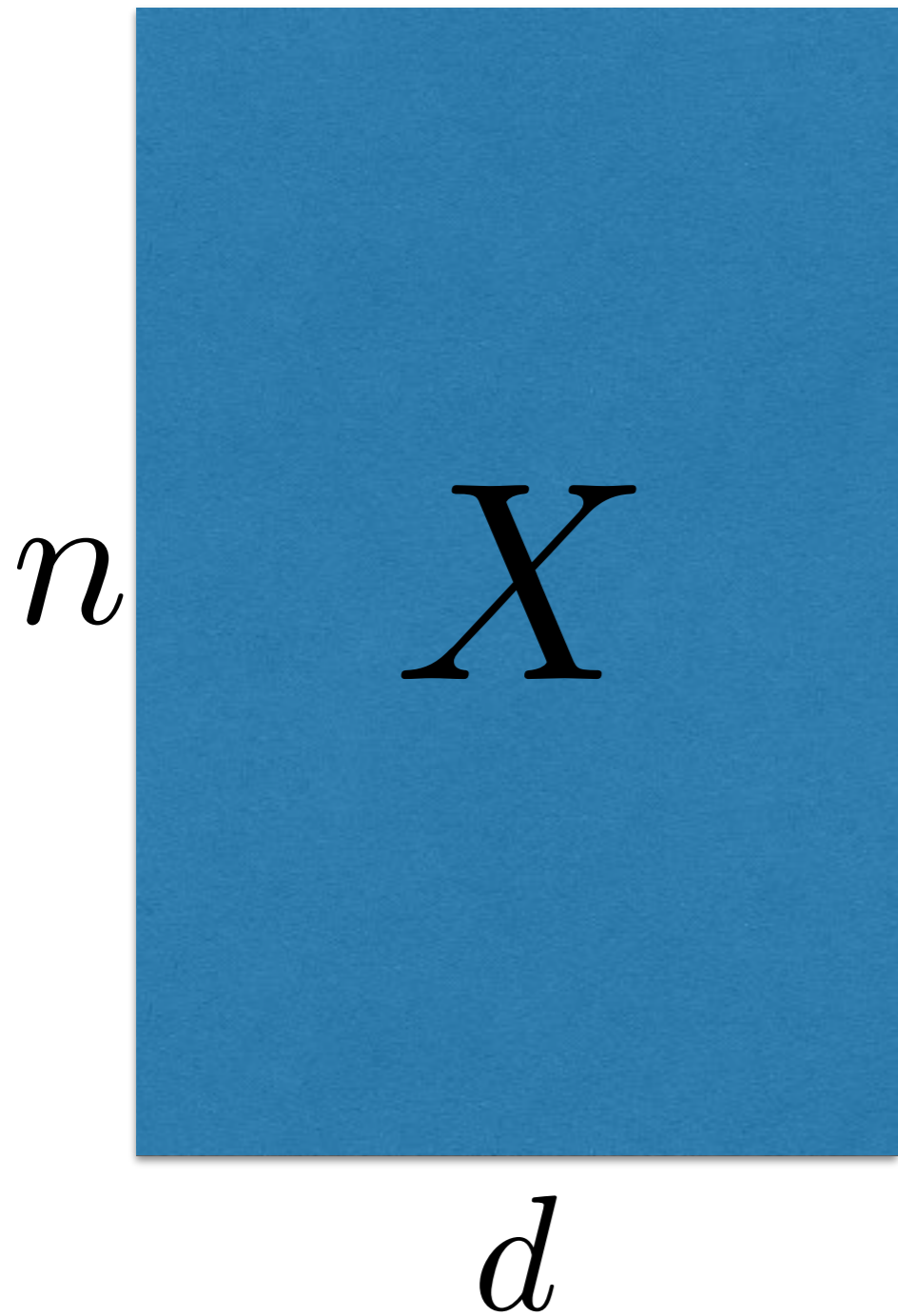
- 3 $m \leftarrow m + 1$

How to choose K

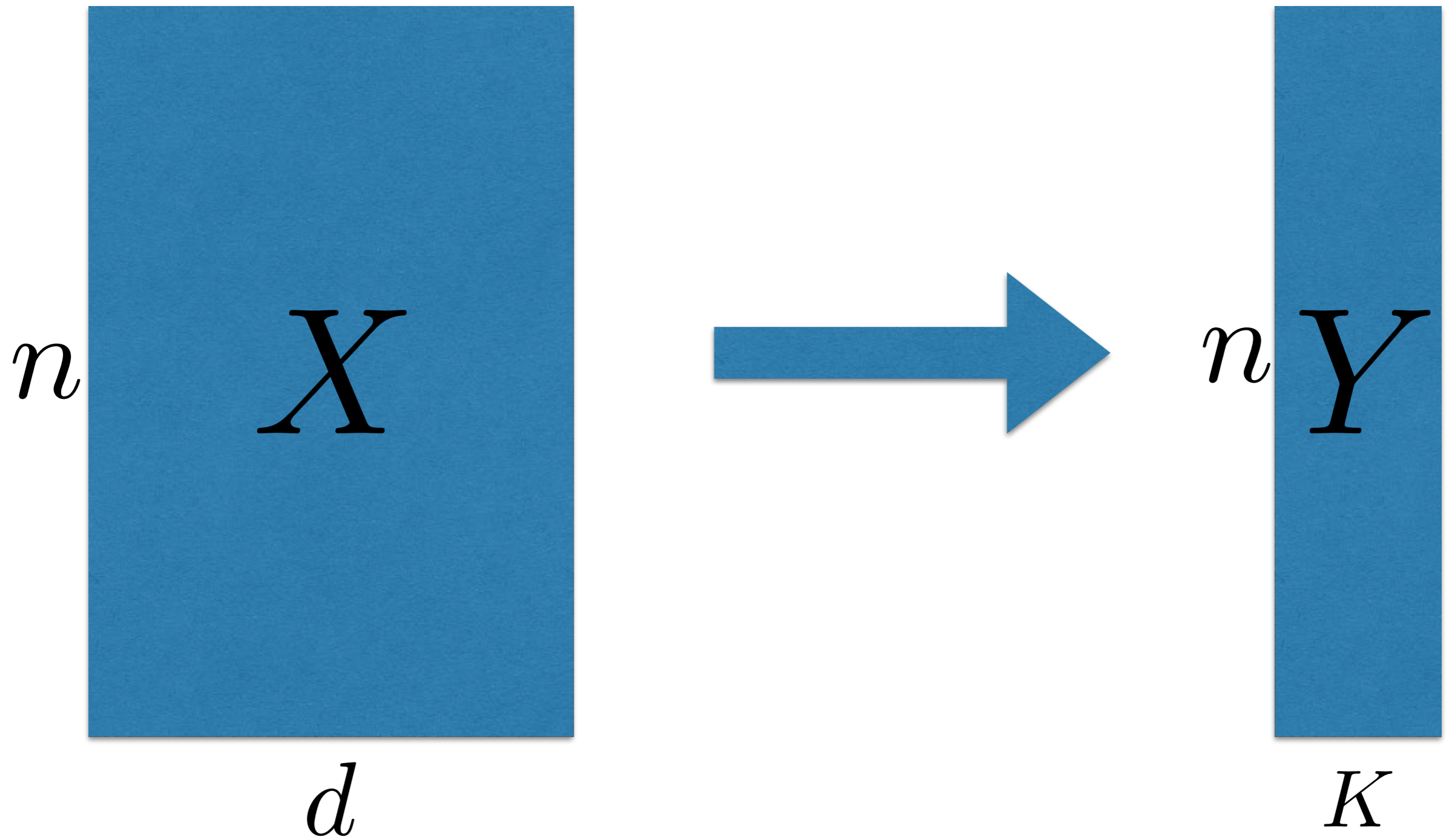
- Elbow method:
 - plot Objective versus K , typically it monotonically decreases.
 - Pick point where there is a kink (explanation in variance is not as much)
 - Intuition: look at rate of change
- Add to objective penalty $+ p(K)$ and minimize, where p increases with K
 - intuition we prefer smaller clusters
 - Use prior knowledge to pick p
 - (AIC, BIC etc can be seen to be specific cases)

DIMENSIONALITY REDUCTION

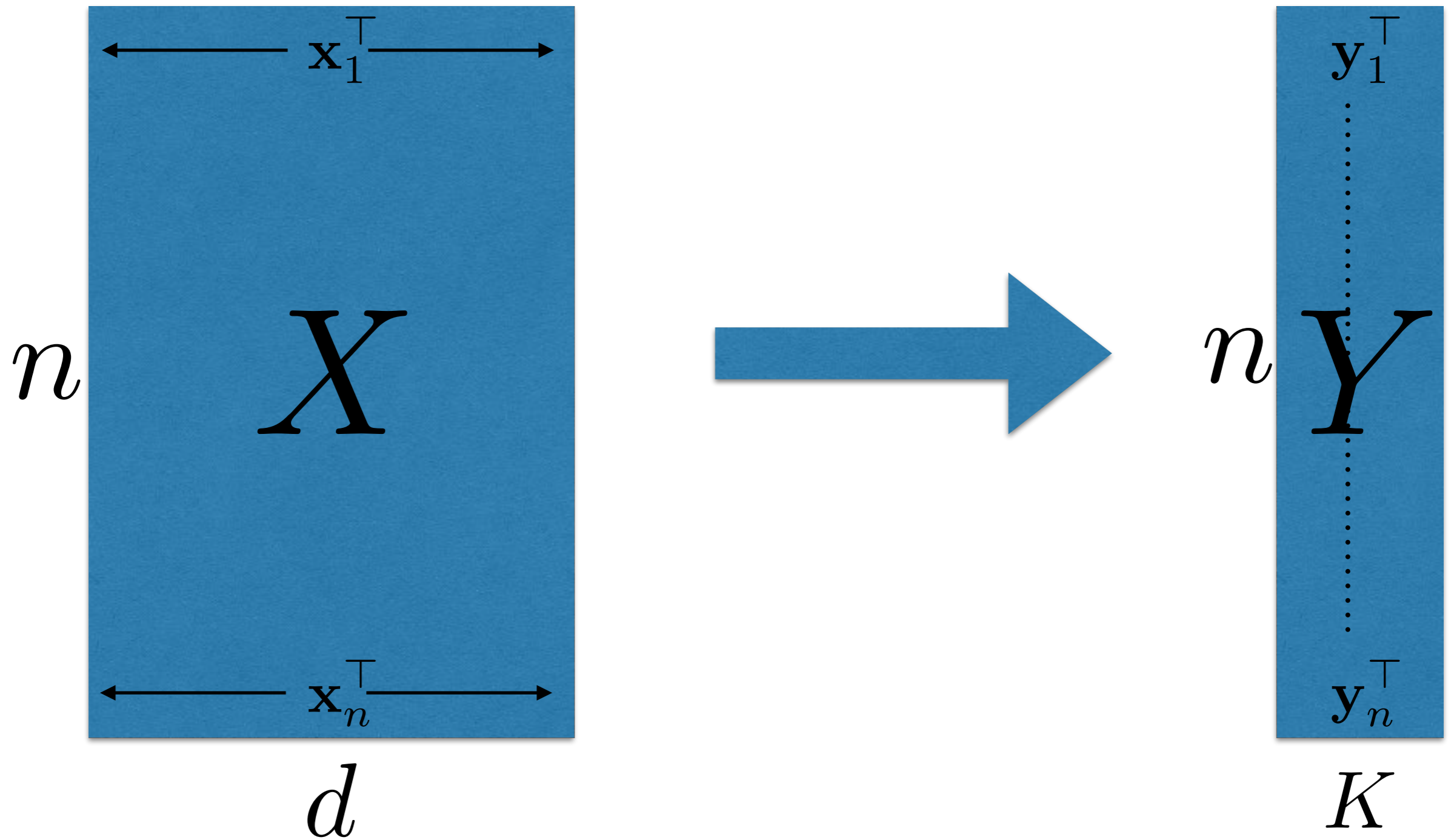
DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION

- You are provided with n data points each in \mathbb{R}^d
- Goal: Compress data into n points in \mathbb{R}^K where $K \ll d$
 - Retain as much information about the original data set
 - Retain desired properties of the original data set

WHY DIMENSIONALITY REDUCTION?

- For computational ease
 - As input to supervised learning algorithm
 - Before clustering to remove redundant information and noise
- Data compression & Noise reduction
- Data visualization

DIMENSIONALITY REDUCTION

Given feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, compress the data points into low dimensional representation $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^K$ where $K \ll d$

DIMENSIONALITY REDUCTION

Desired properties:

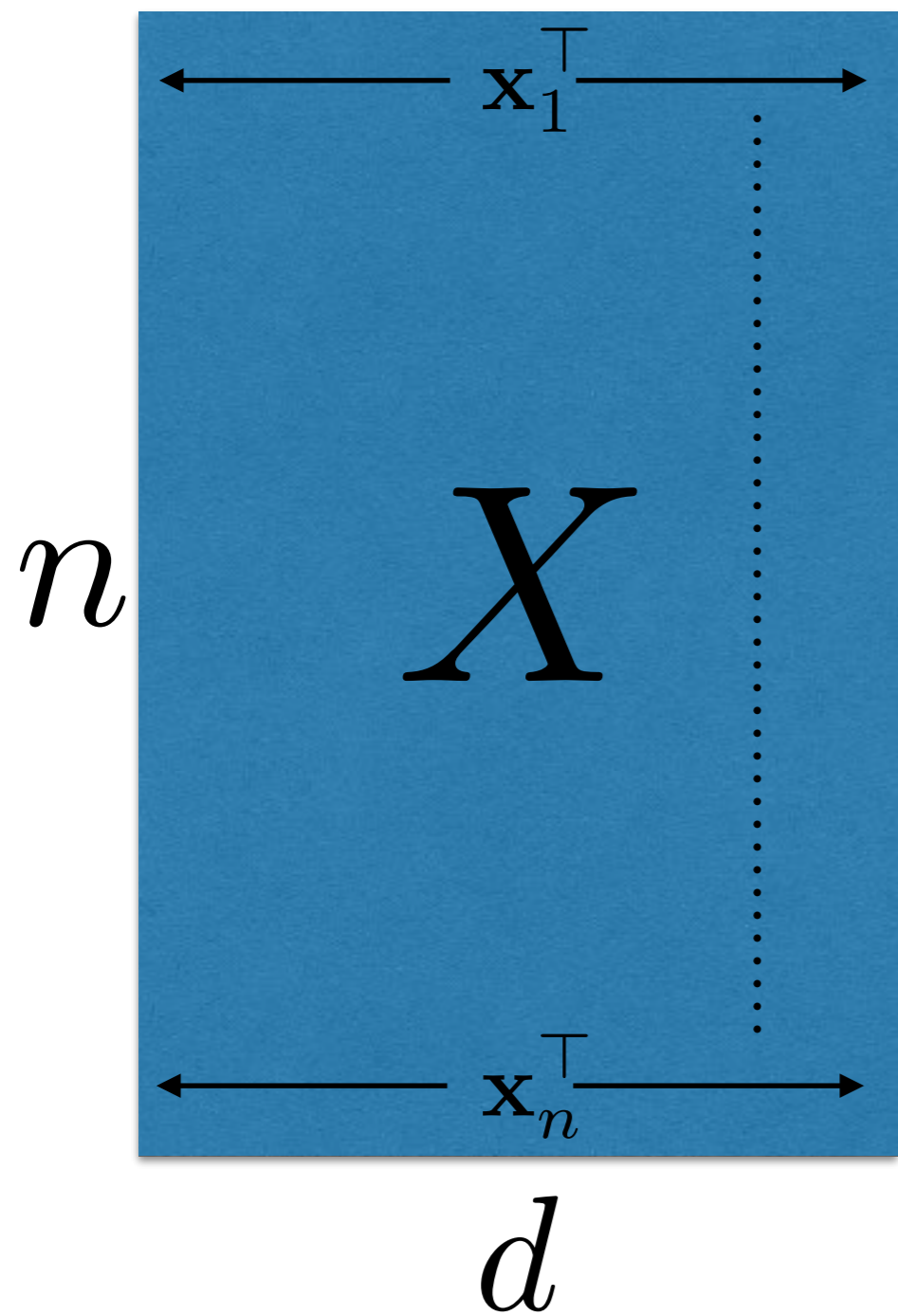
- ① Original data can be (approximately) reconstructed
- ② Preserve distances between data points
- ③ “Relevant” information is preserved
- ④ Noise is reduced

DIM REDUCTION: LINEAR TRANSFORMATION

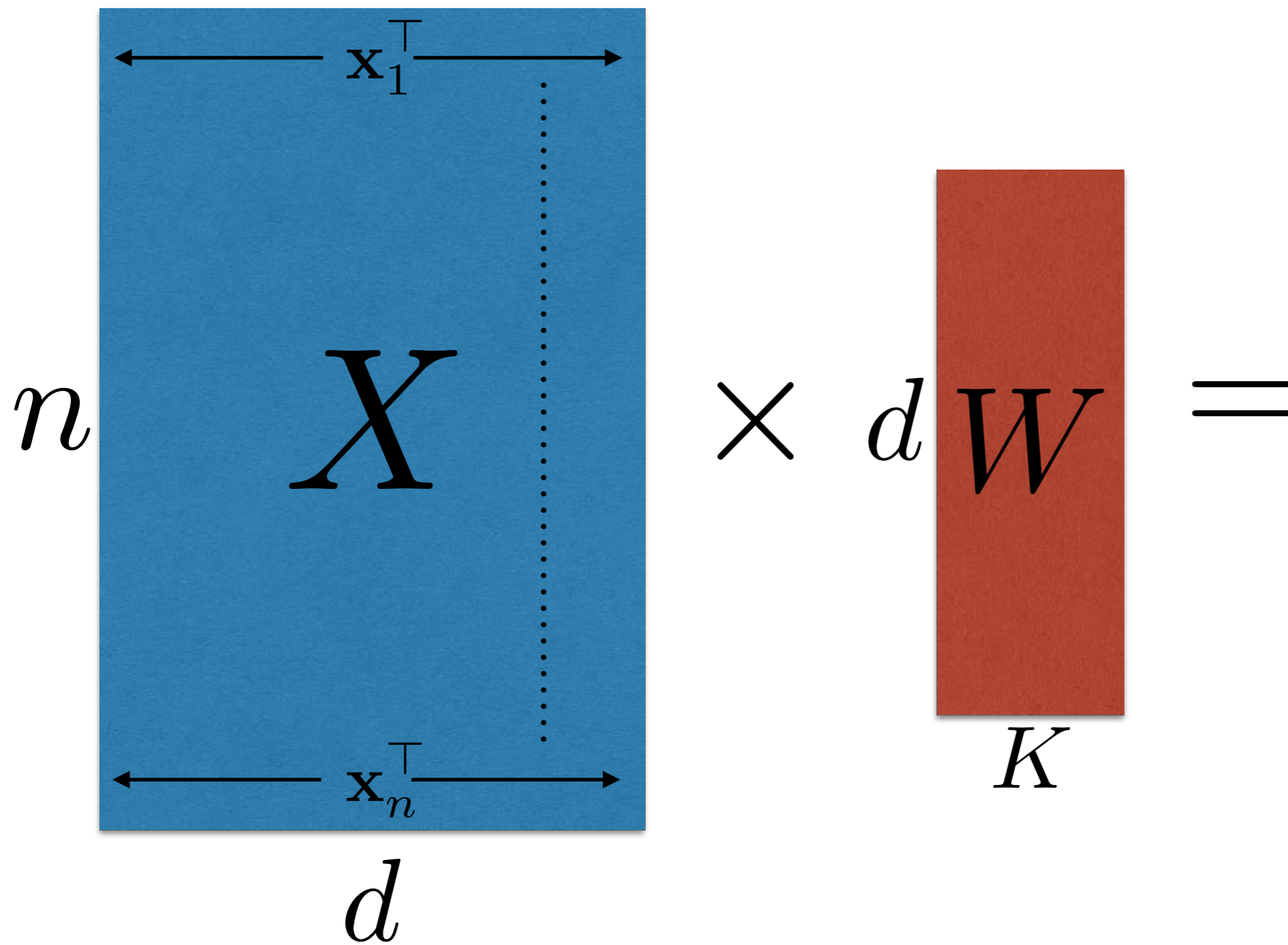
- Pick a low dimensional subspace
- Project linearly to this subspace
- Subspace retains as much information

DIM REDUCTION: LINEAR TRANSFORMATION

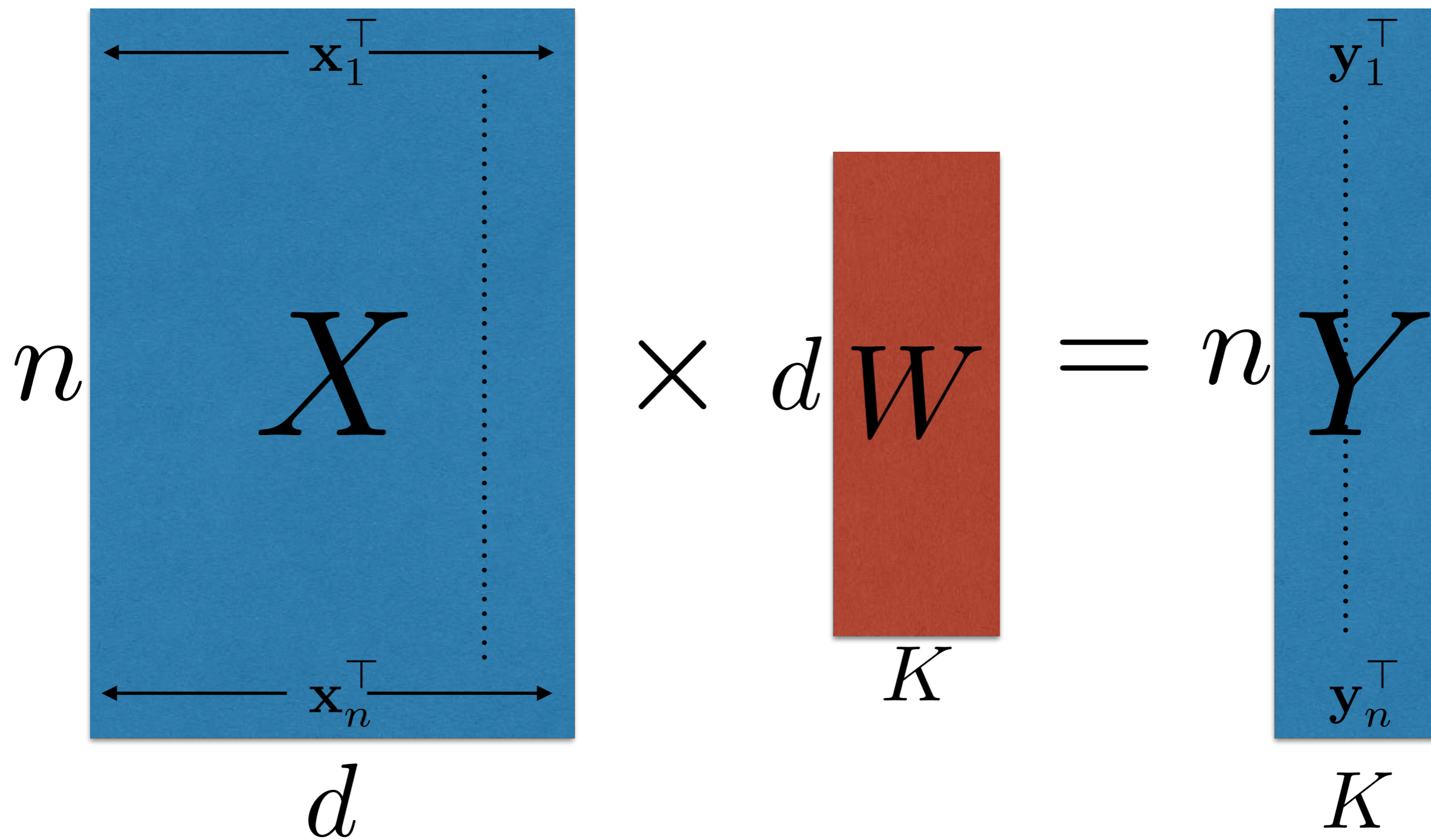
DIM REDUCTION: LINEAR TRANSFORMATION



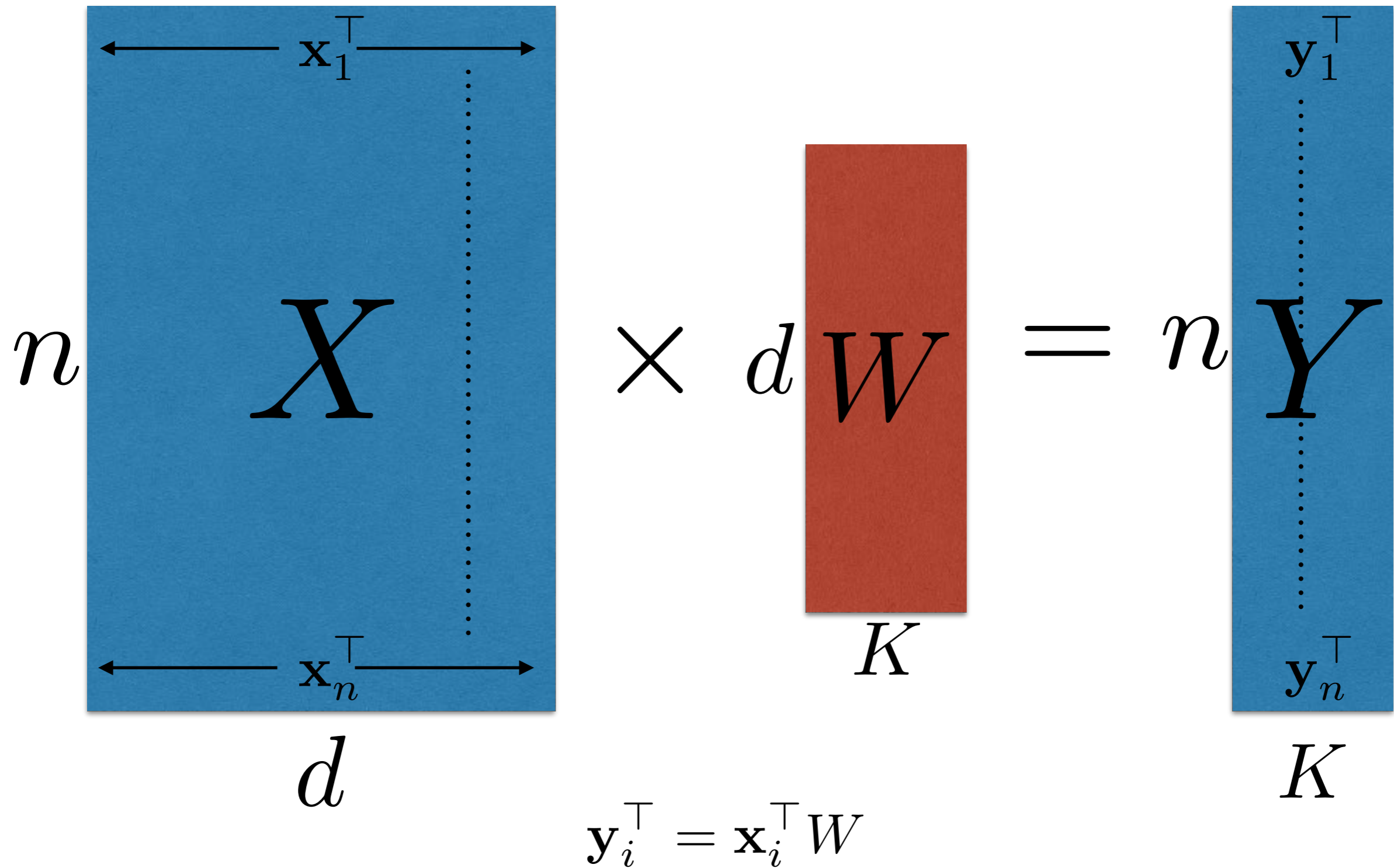
DIM REDUCTION: LINEAR TRANSFORMATION



DIM REDUCTION: LINEAR TRANSFORMATION



DIM REDUCTION: LINEAR TRANSFORMATION



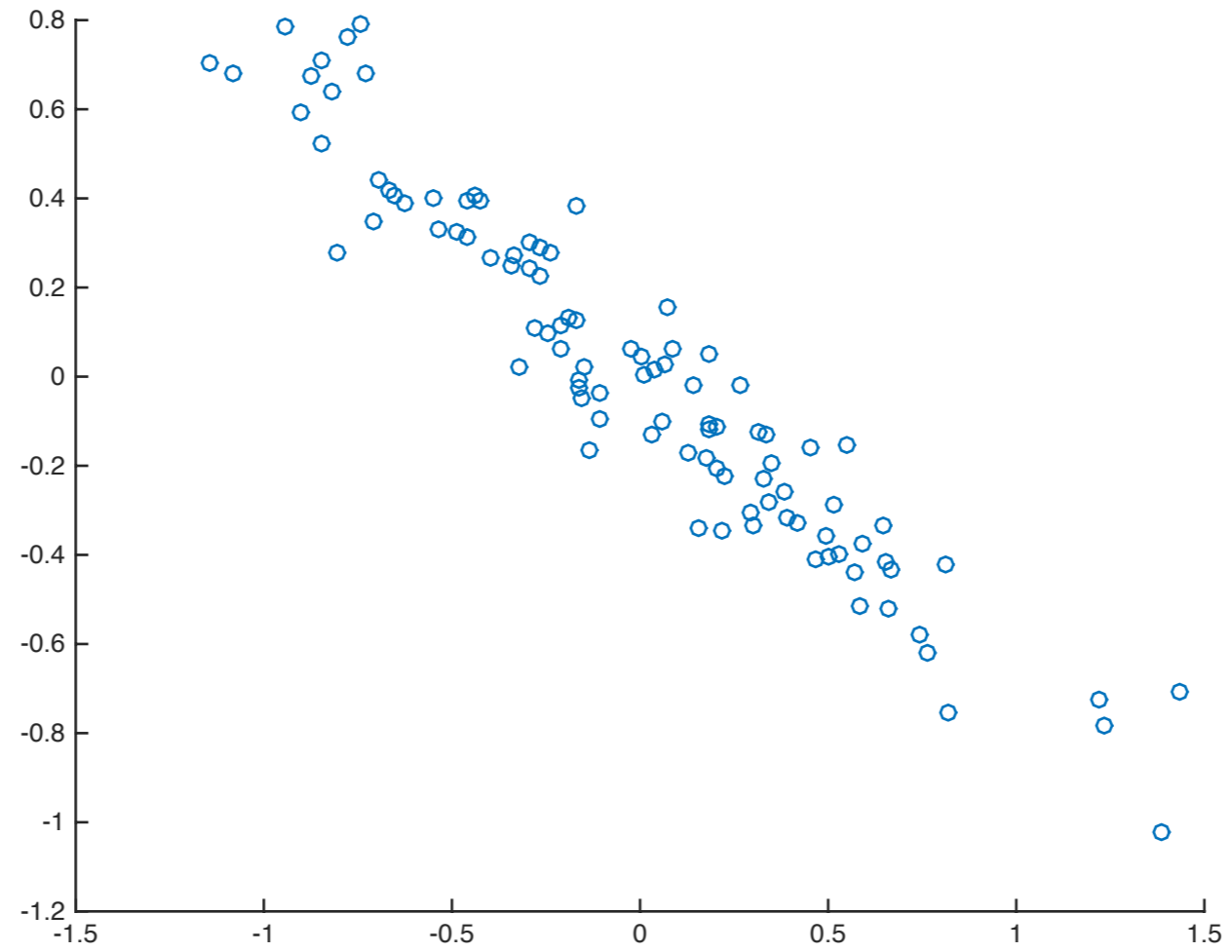
Example: Students in classroom



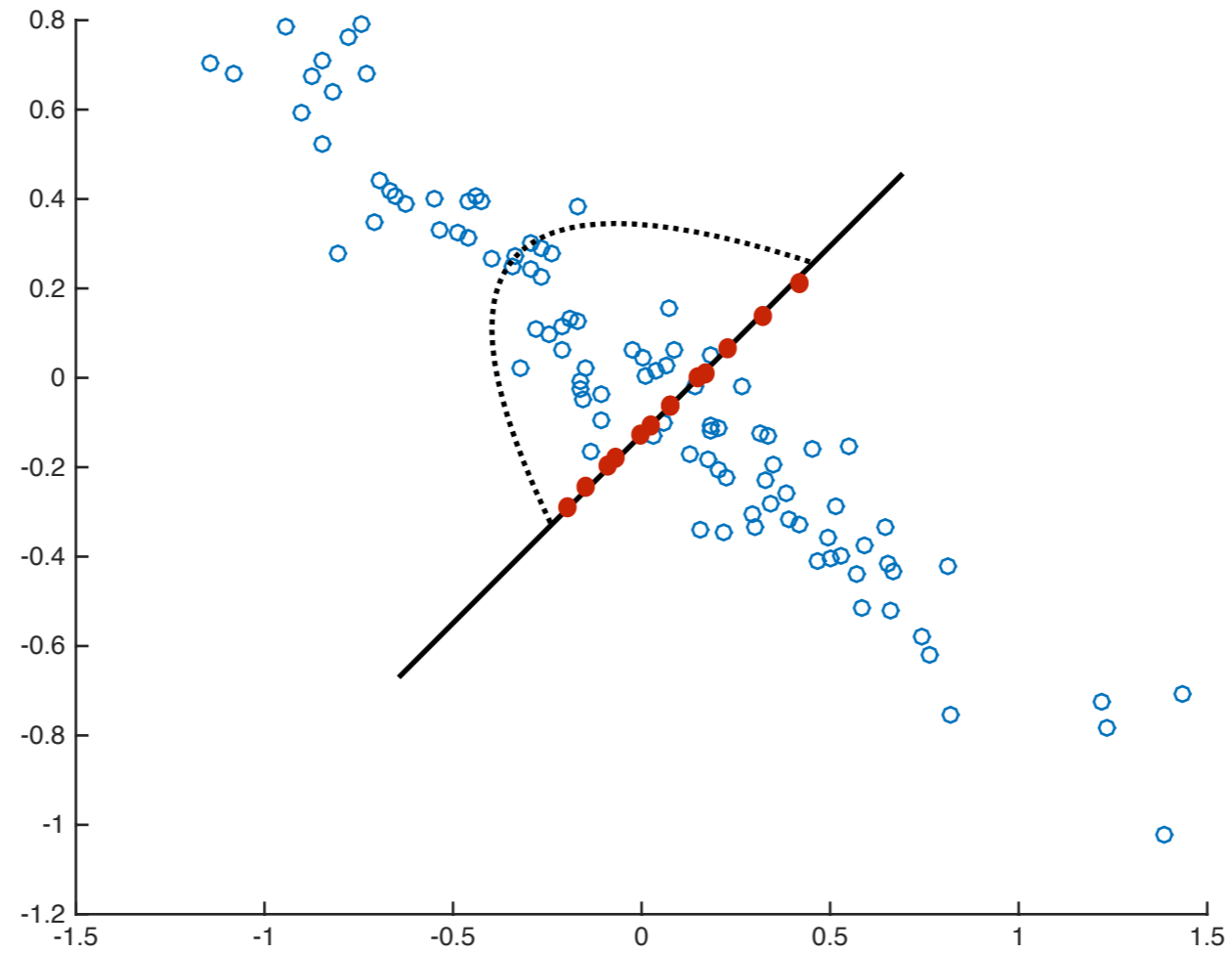
Example: Students in classroom



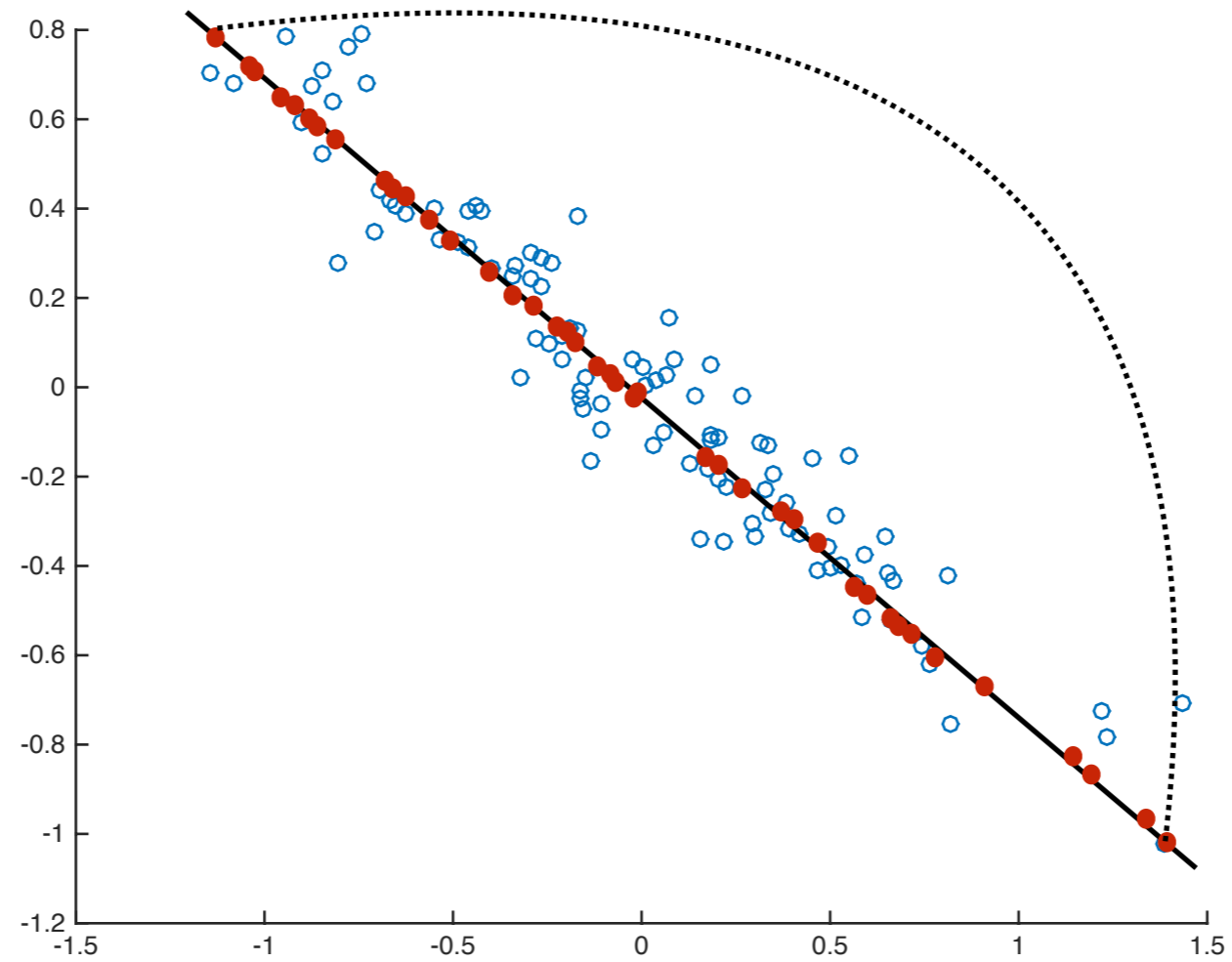
PCA: VARIANCE MAXIMIZATION



PCA: VARIANCE MAXIMIZATION



PCA: VARIANCE MAXIMIZATION



DIM REDUCTION: LINEAR TRANSFORMATION

Prelude: reducing to 1 dimension

DIM REDUCTION: LINEAR TRANSFORMATION

Prelude: reducing to 1 dimension

x_1 ●

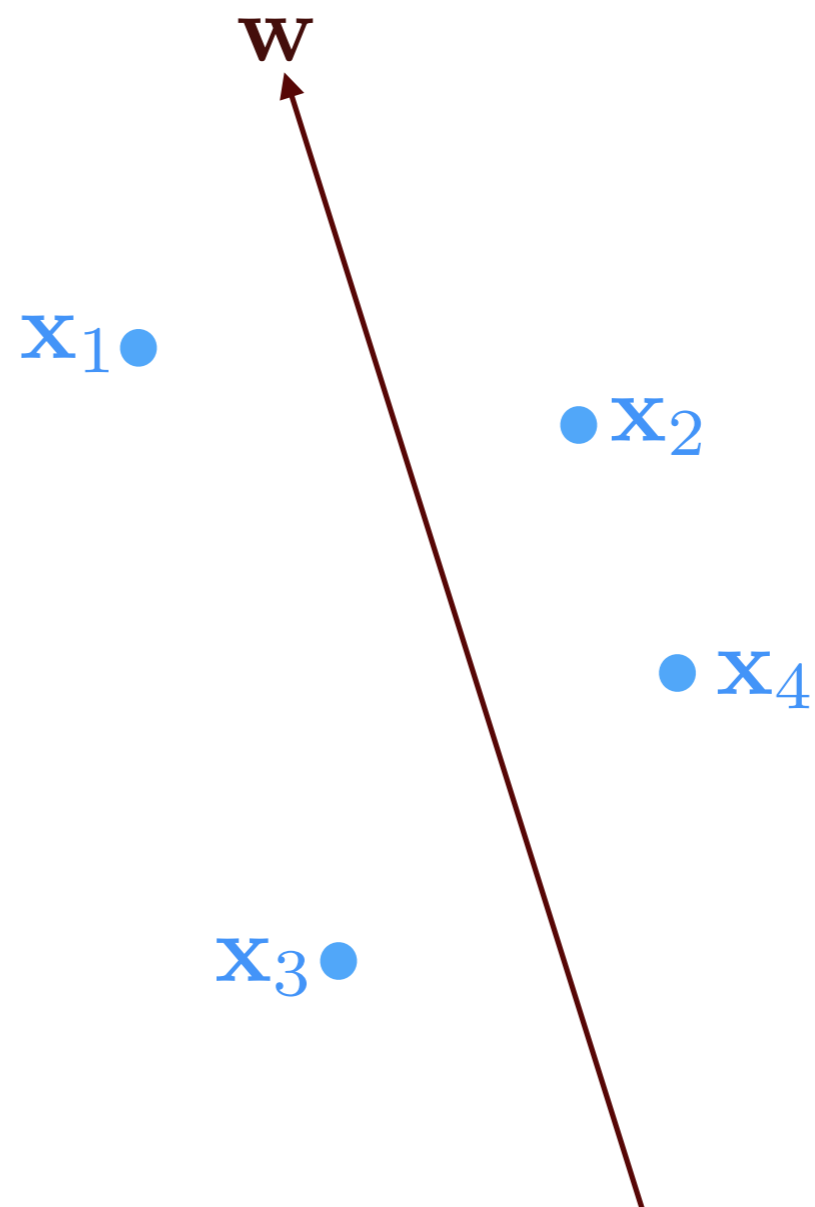
● x_2

● x_4

x_3 ●

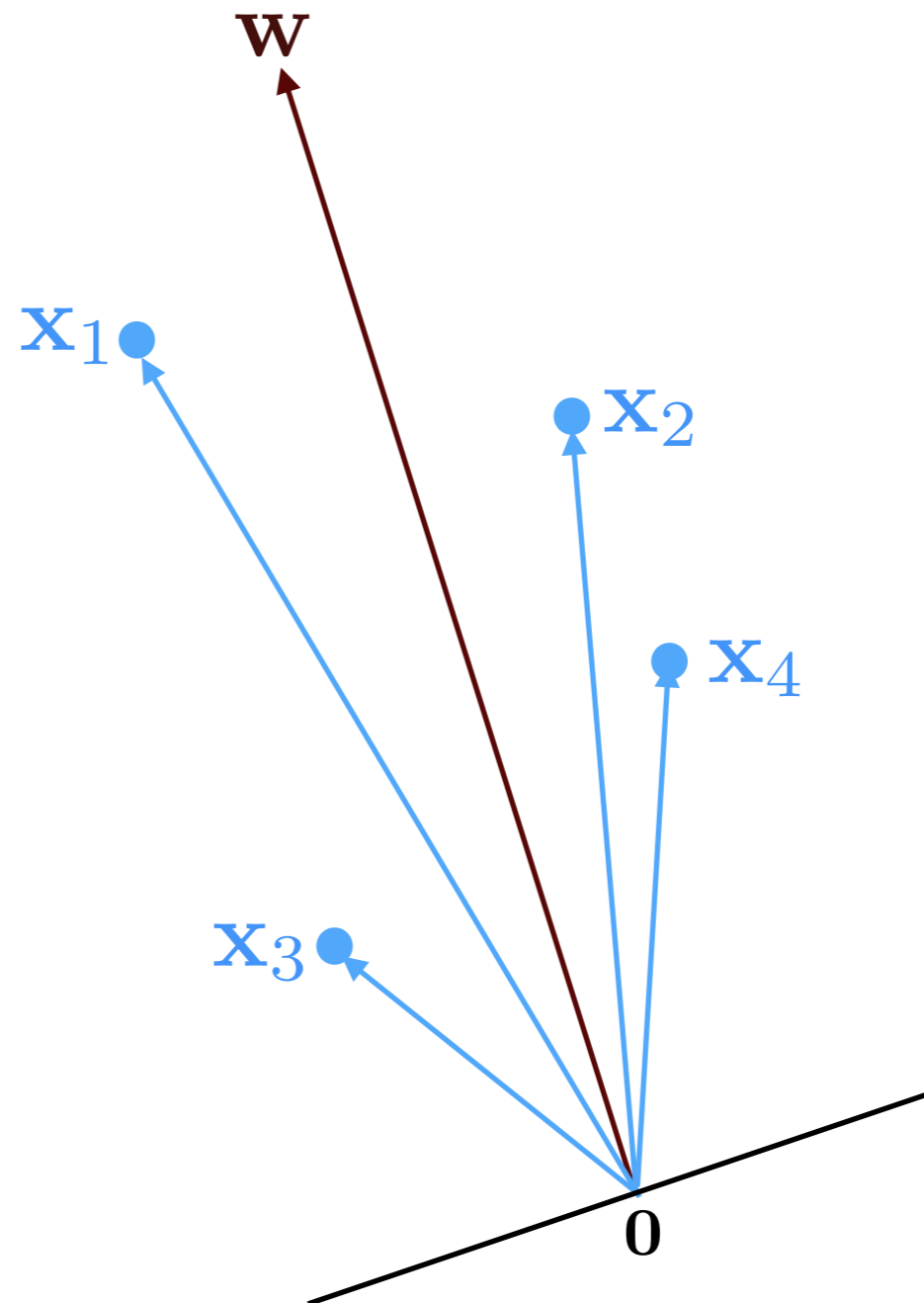
DIM REDUCTION: LINEAR TRANSFORMATION

Prelude: reducing to 1 dimension



DIM REDUCTION: LINEAR TRANSFORMATION

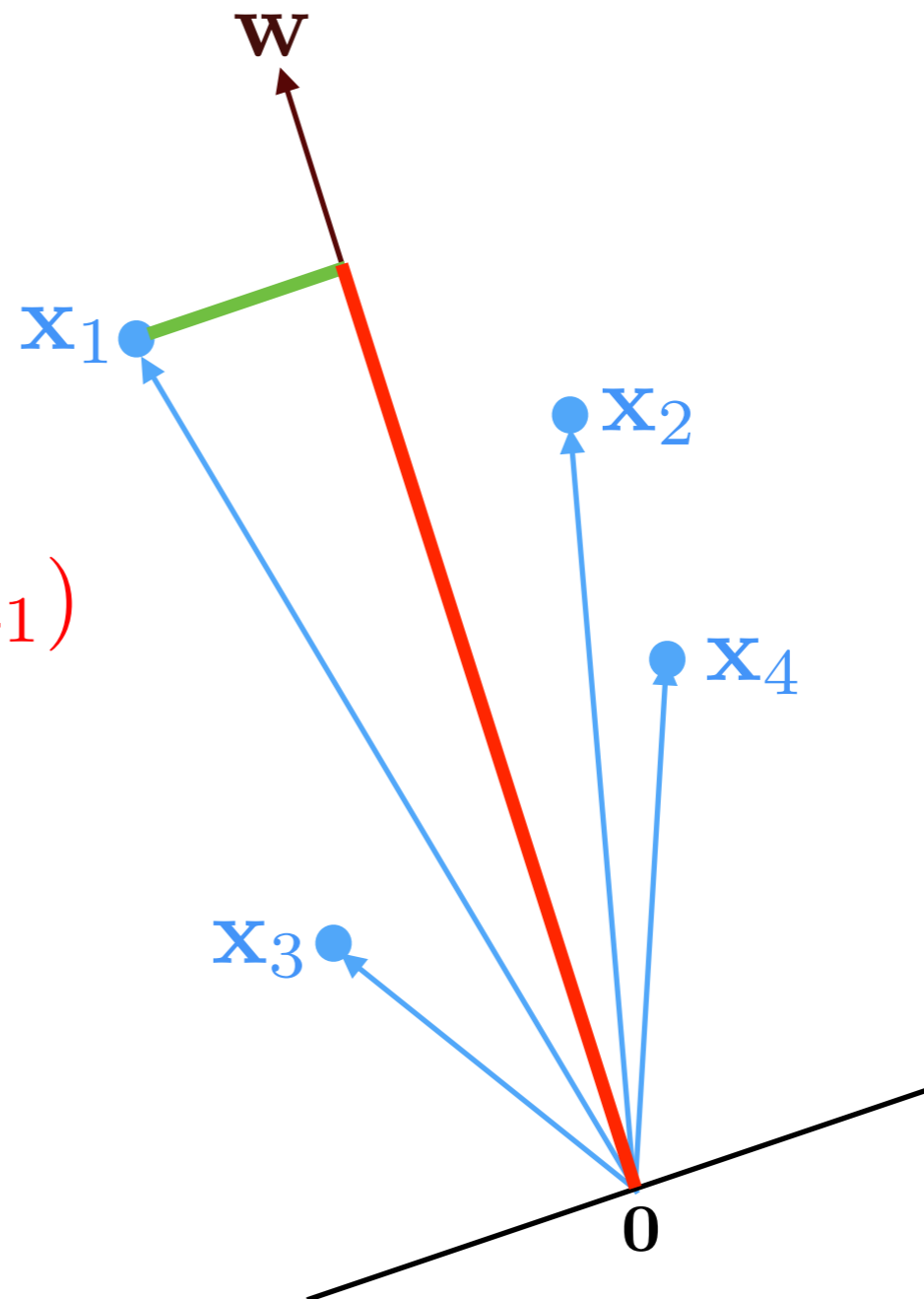
Prelude: reducing to 1 dimension



DIM REDUCTION: LINEAR TRANSFORMATION

Prelude: reducing to 1 dimension

$$y_1 = \mathbf{w}^T \mathbf{x}_1 = \|\mathbf{x}_1\| \cos(\angle \mathbf{w} \mathbf{x}_1)$$



PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\mathbf{w}_1 = \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2\end{aligned}$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w}\end{aligned}$$

PCA: VARIANCE MAXIMIZATION

- Pick directions along which data varies the most
- First principal component:

$$\begin{aligned}\mathbf{w}_1 &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top \mathbf{x}_t - \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top \mathbf{x}_t \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \left(\mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) \right)^2 \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \frac{1}{n} \sum_{t=1}^n \mathbf{w}^\top (\mathbf{x}_t - \boldsymbol{\mu}) (\mathbf{x}_t - \boldsymbol{\mu})^\top \mathbf{w} \\ &= \arg \max_{\mathbf{w}: \|\mathbf{w}\|_2=1} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}\end{aligned}$$

$\boldsymbol{\Sigma}$ is the covariance matrix

Review

- Review covariance
- Review Eigen vectors

Covariance Matrix

- Its a $d \times d$ matrix, $\Sigma[i, j]$ measures “covariance” of features i and j

$$\Sigma[i, j] = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t[i] - \mu[i])(\mathbf{x}_t[j] - \mu[j])$$

PCA: VARIANCE MAXIMIZATION

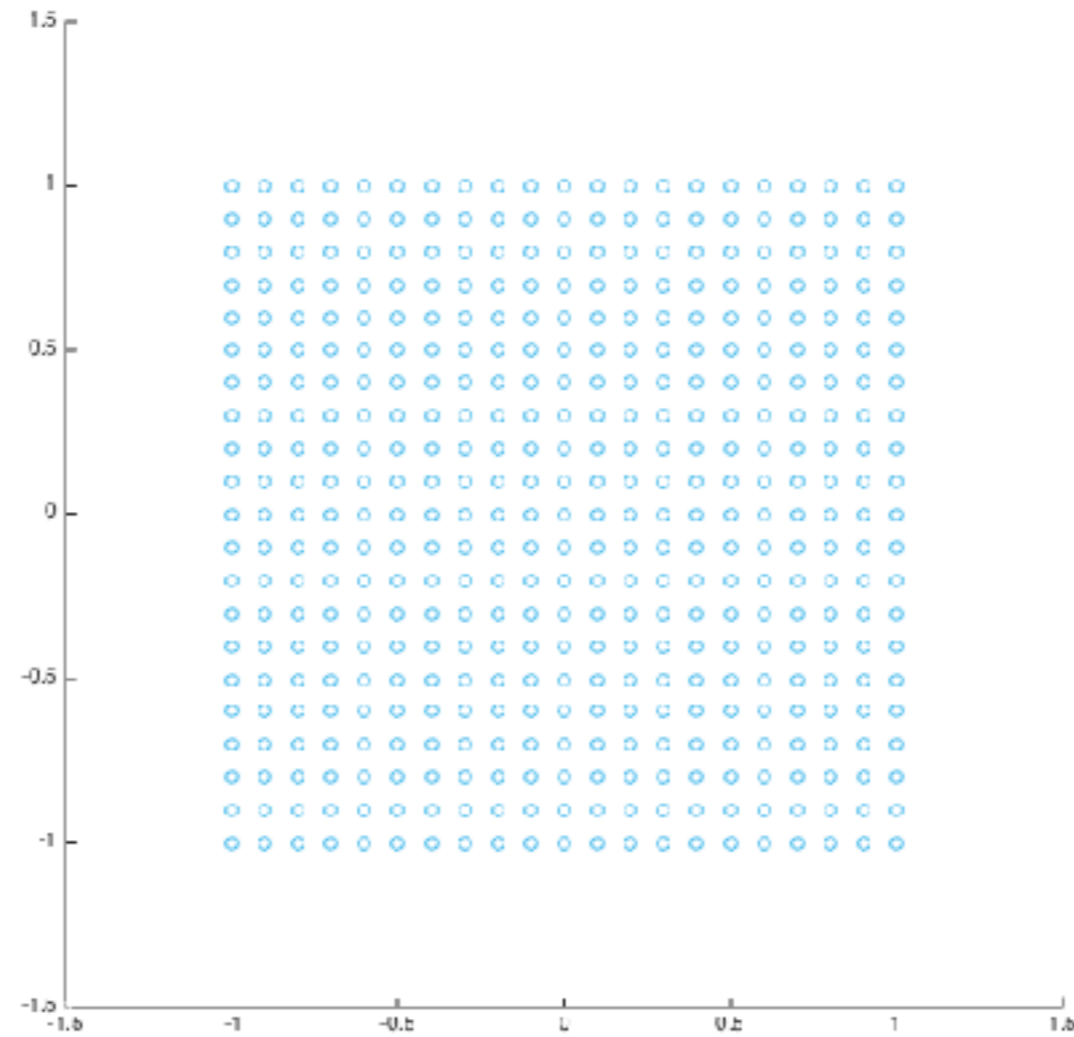
Covariance matrix:

$$\Sigma = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^\top$$

- Its a $d \times d$ matrix, $\Sigma[i, j]$ measures “covariance” of features i and j

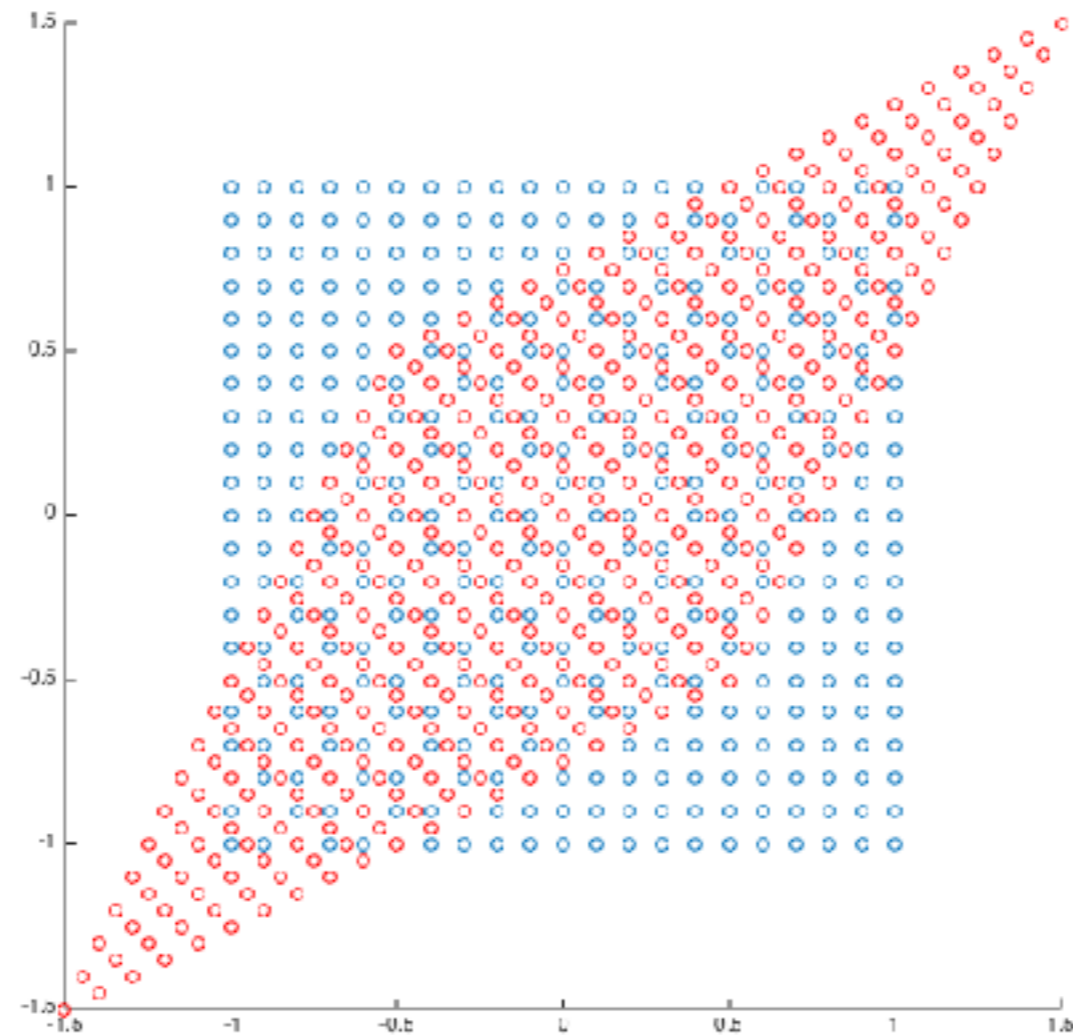
$$\Sigma[i, j] = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t[i] - \mu[i])(\mathbf{x}_t[j] - \mu[j])$$

What are Eigen Vectors?



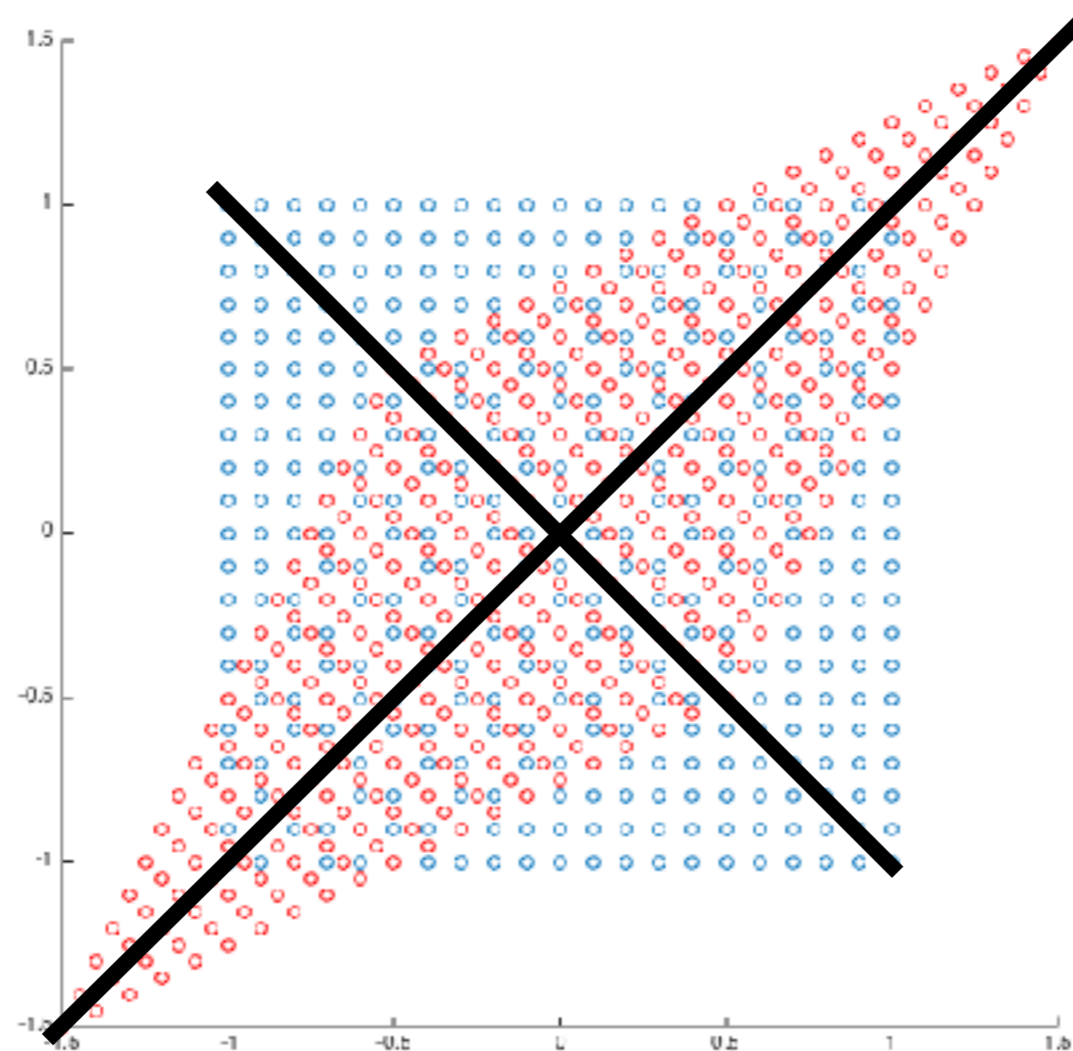
What are Eigen Vectors?

$$x \mapsto Ax$$



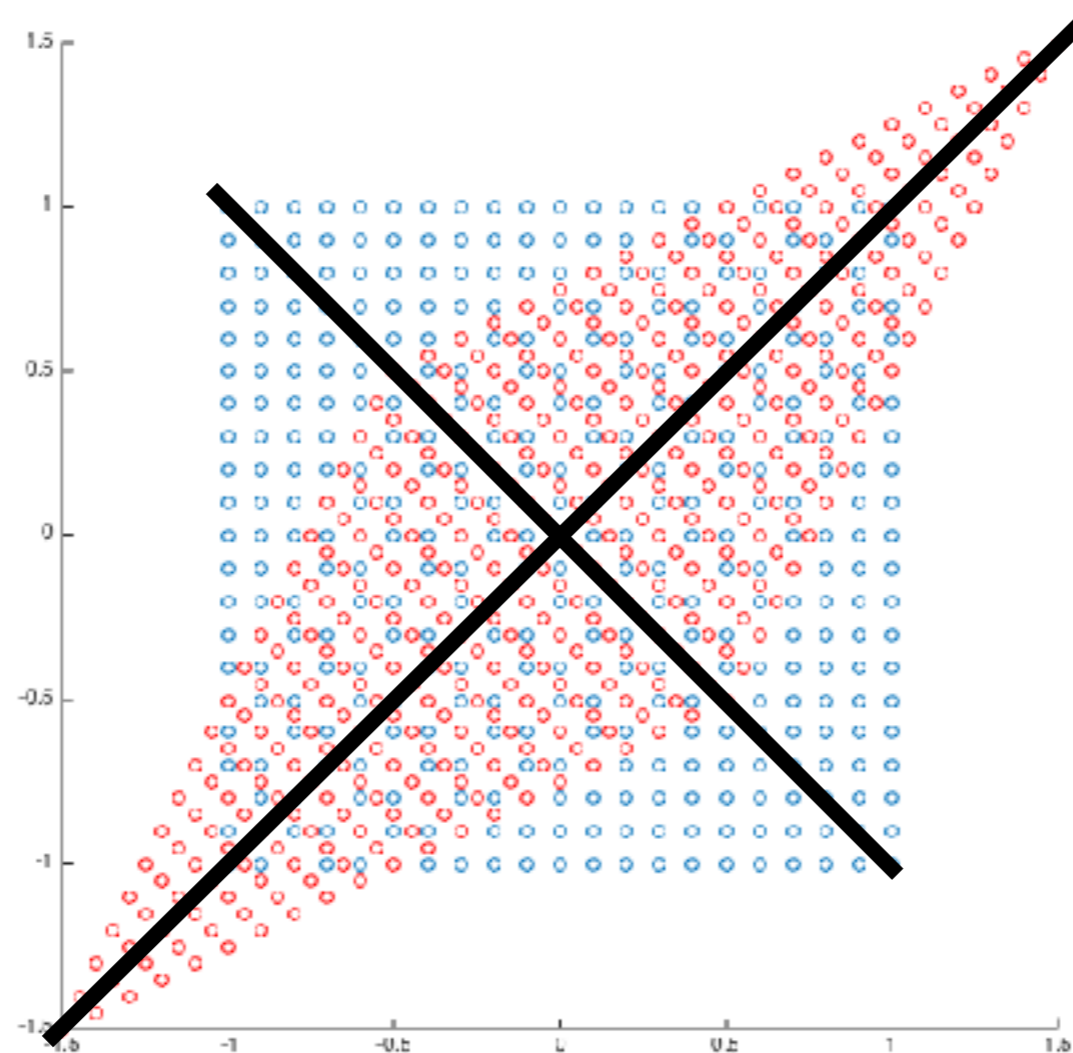
What are Eigen Vectors?

$$x \mapsto Ax$$



What are Eigen Vectors?

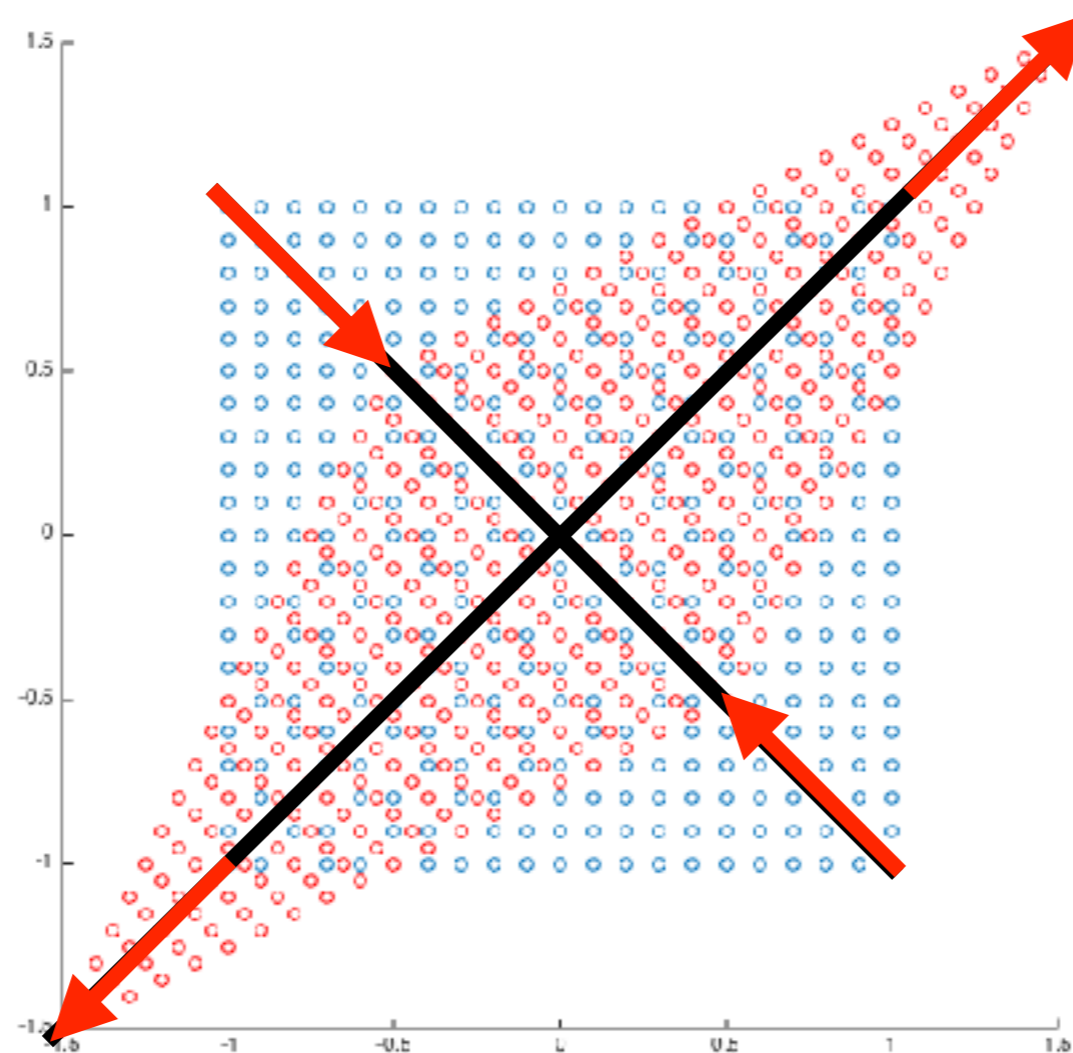
$$x \mapsto Ax$$



$$Ax = \lambda x$$

What are Eigen Vectors?

$$x \mapsto Ax$$



$$Ax = \lambda x$$