Machine Learning Land

# Machine Learning: Unsupervised Learning

**Konstantin Tretyakov**
http://kt.era.ee

**IFI Summer School 2014**

STACC Software Technology and Applications Competence Center

BIIT

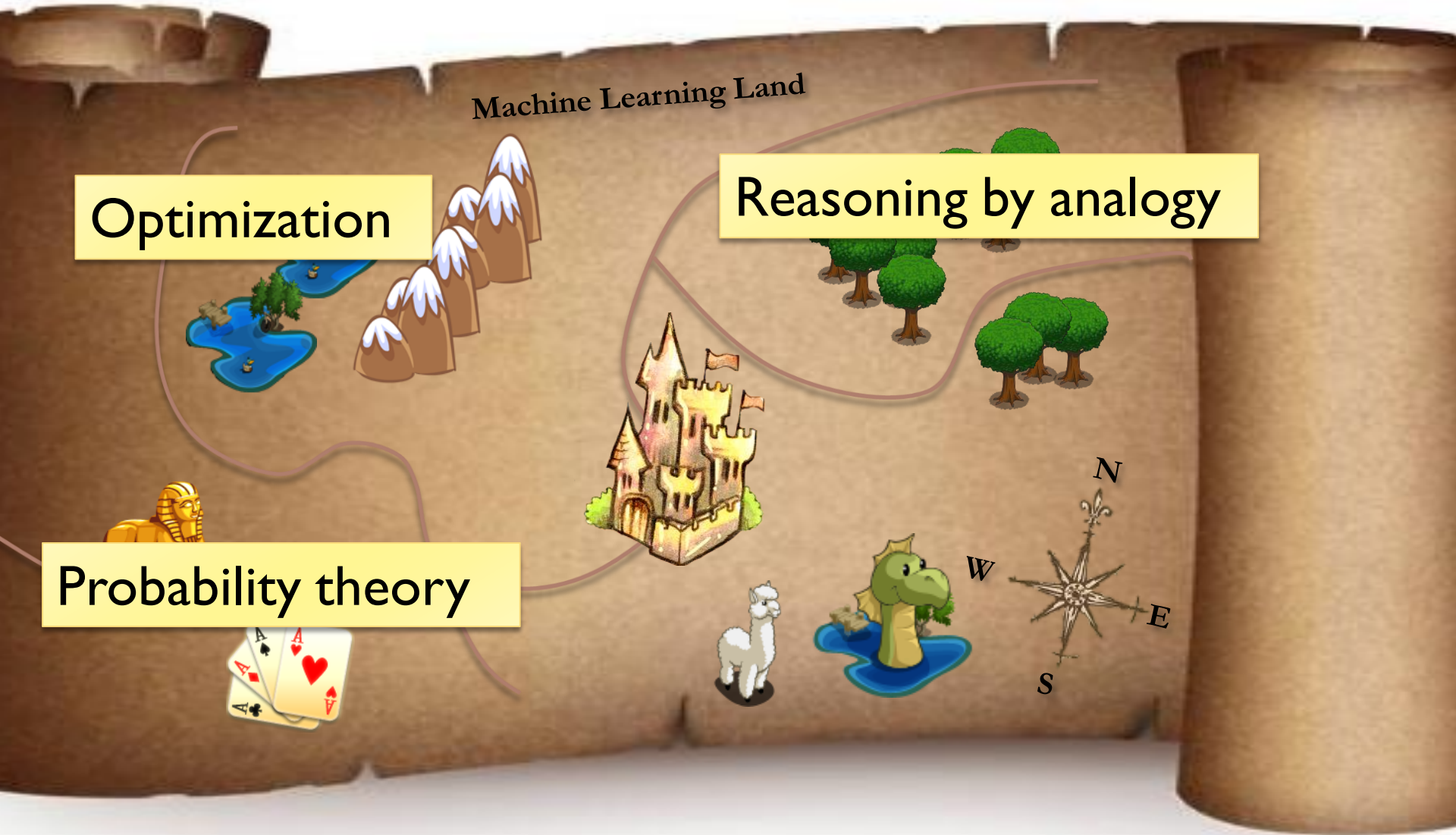TARTU ÜLIKOOL · UNIVERSITAS TARTUENSIS · 1632

# So far…

# So far…



Machine Learning Land

Optimization

Reasoning by analogy

Probability theory

# So far…

Optimization

Fermat's theorem
Gradient methods
Batch & On-line

MLE,
MAP,
Bayesian Estimation
Risk Optimization

Probability theory

Machine Learning

Reasoning by analogy

k-NN,
Kernel methods

N

W

E

S

# So far...

Optimization

Fermat's theorem
Gradient [method]
Batch & [O...]

MLE,
MAP,
Bayesian Estimation
Risk Optimization

Probability theory

Machine Learning

Reasoning by analogy
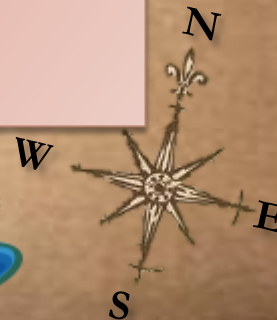
k-NN,

[...]ods

k-NN, Decision tree learning
Linear regression (OLS, RR),
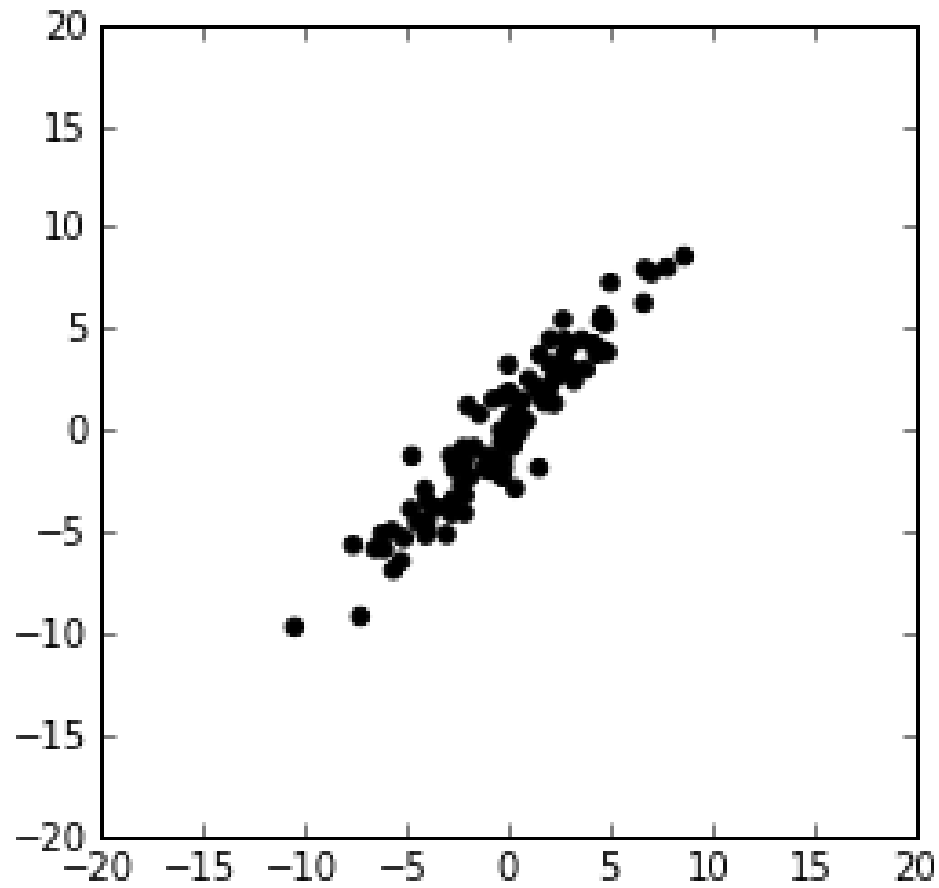Linear classification (SVM, Perceptron),
Bayes & Naïve Bayes
Kernel-<xxx>

*N*

*W*

*E*

*S*

# Next

Machine Learning

Optimization

Fermat's theorem
Gradient
Batch & O

Reasoning by analogy

k-NN,

k-NN, Decision tree learning
Linear regression (OLS, RR),
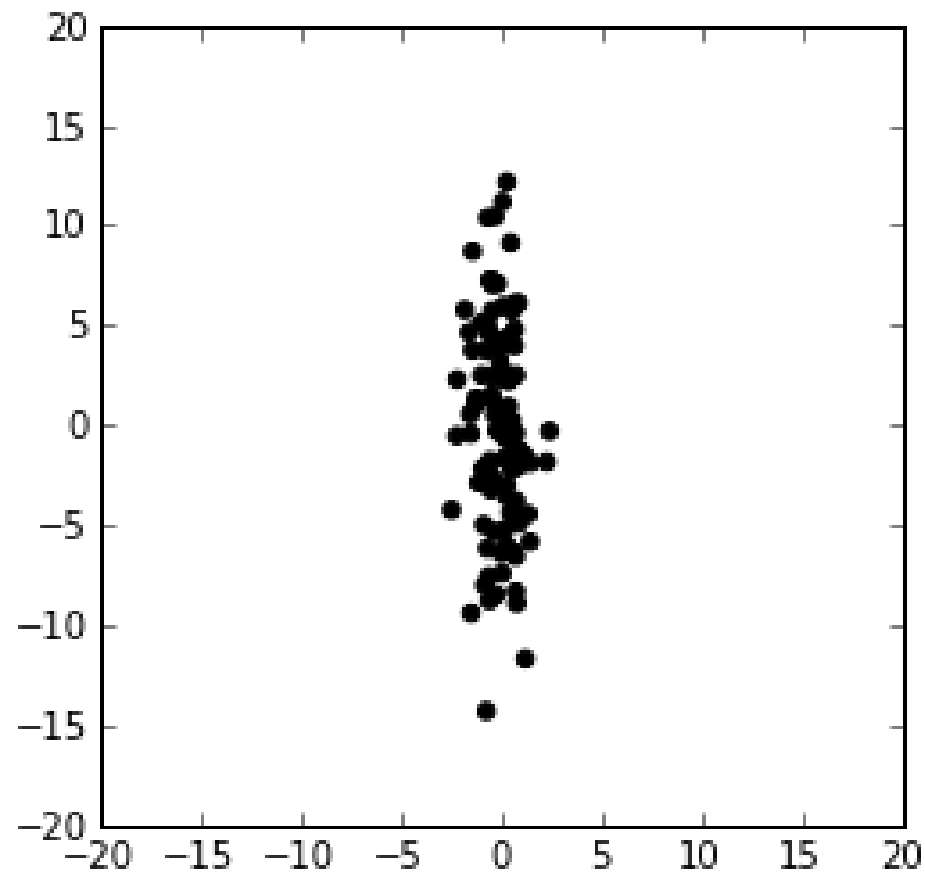Linear classification (SVM, Perceptron),
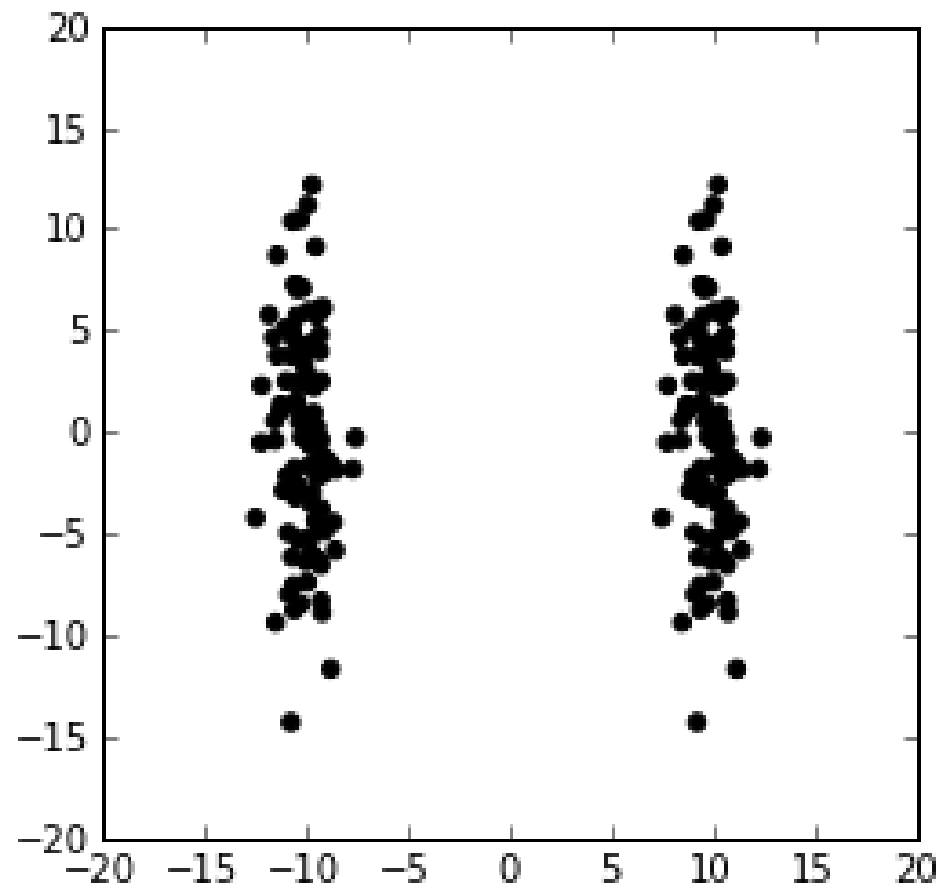Bayes & Naïve Bayes
Kernel-<xxx>

MLE,
MAP,
Bayesian Estimation
Risk Optimization

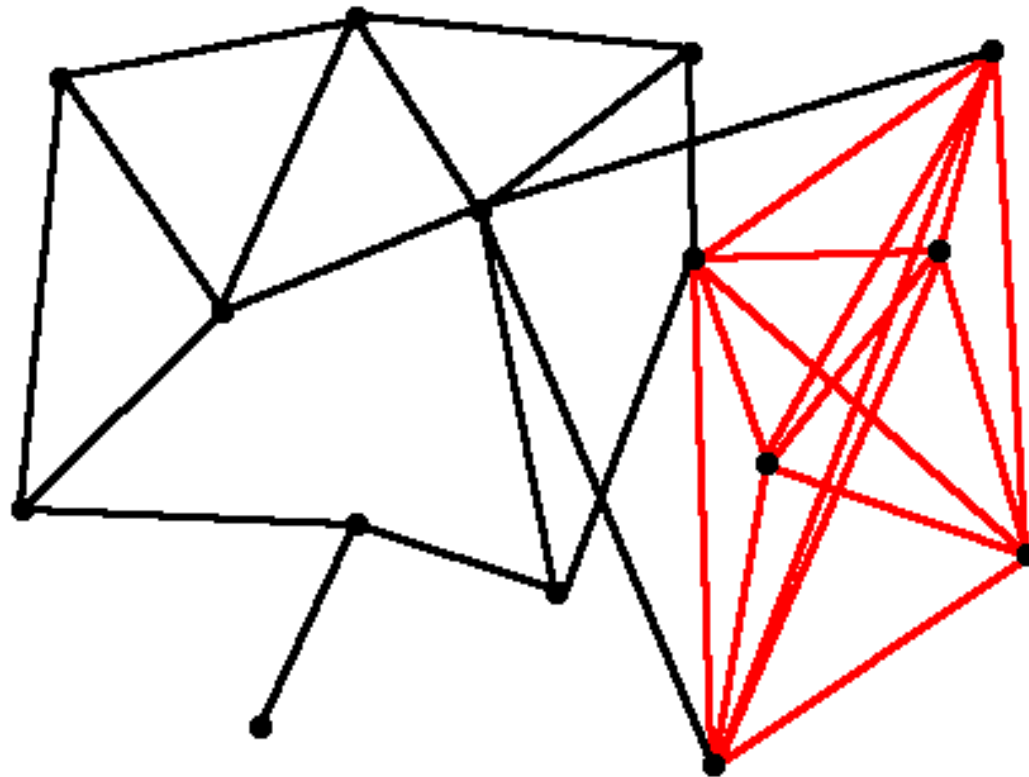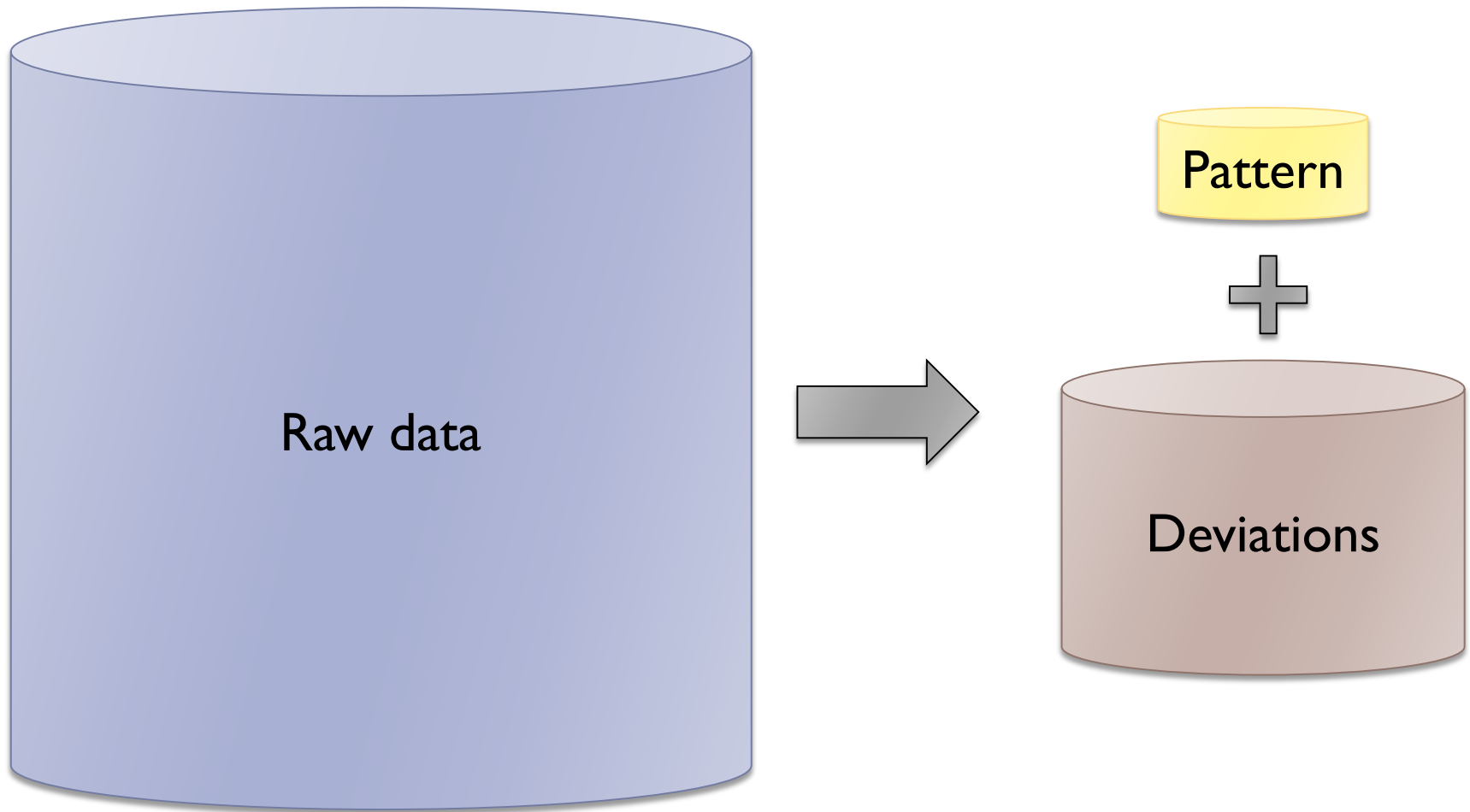Unsupervised
Learning

Probability theory

*N*

*E*

*S*

AATAACGGCCCGATGAGGAAACGAACGGTCGCACT
AAGATGAGACATGTCCCGAAAGGTGCATAAGTTAT
GGACGAAAAACTTTCTTCGCCCTTTGATGTGCCCC
AGCGCGGGATGAGGATCAGCCCCCGCATTAGTTCA
ATATGCGAGCTTTCGCGCTCGGAAAGGGCAATAAA
GCGACGGCCCCGATGAGGGGTGTTACTAGATTGGA
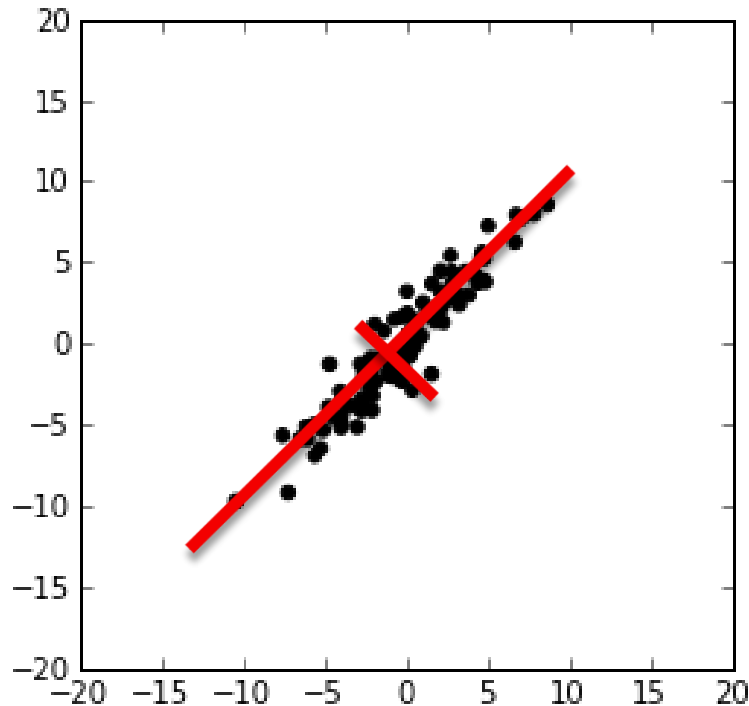TGGGTGGTTCAGATCTCGGCTTACCCCCTTTATCA
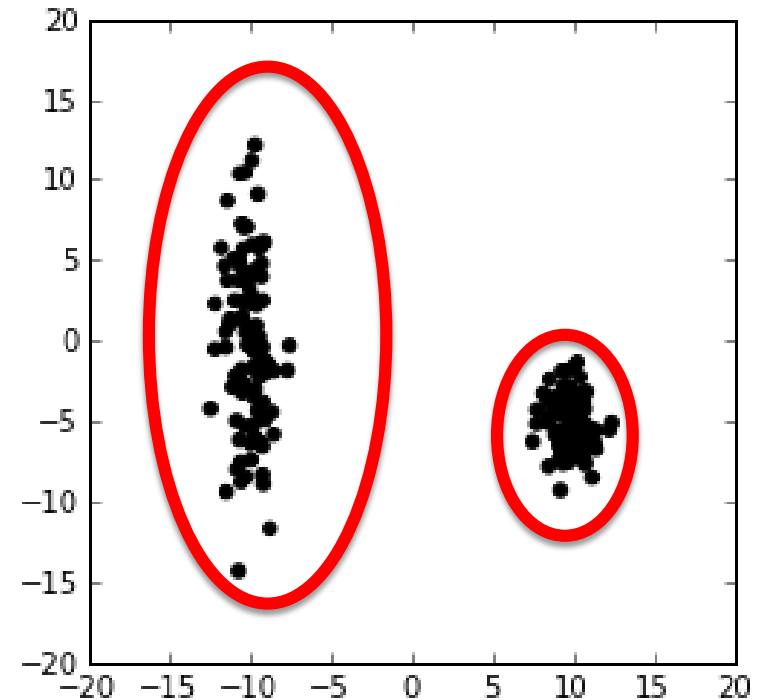ACCCTGCTACAGACTCGTTGAGAATGCTACGGATC

# Data Mining

Raw data

Pattern

+

Deviations

# Unsupervised learning patterns

**Decomposition**
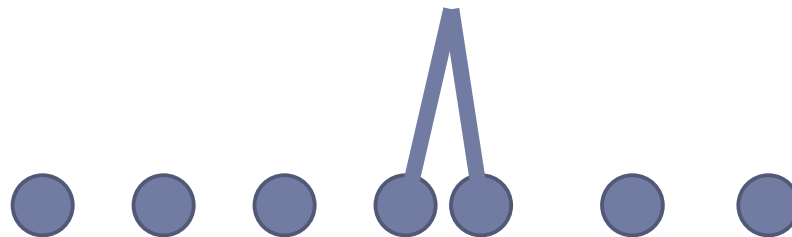
**Clustering**

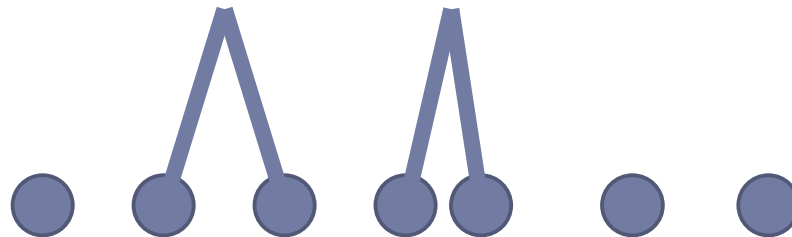# Quiz

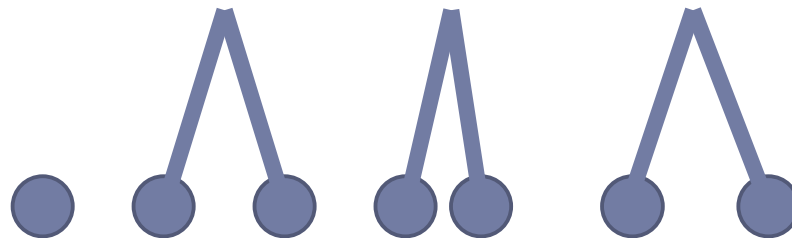▸ Why would one need clustering?

# Hierarchical clustering

# Hierarchical clustering
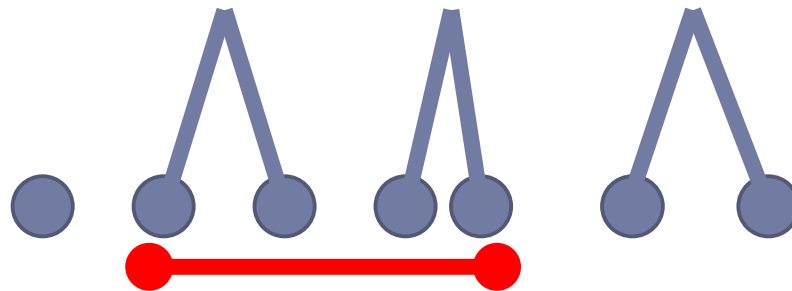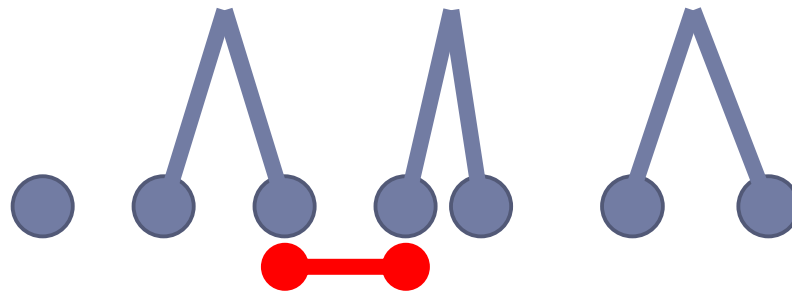
# Hierarchical clustering

# Hierarchical clustering

Complete linkage

Single linkage

# Hierarchical clustering

Average linkage

Ward linkage

$$\sigma^2$$

# Hierarchical clustering

# Hierarchical clustering

# Hierarchical clustering

# Partitional vs Hierarchical



Partitional clustering finds a fixed number of clusters

Hierarchical clustering creates a series of clusterings contained in each other

# K-means

# K-means

# K-means

# K-means

# K-means

# K-means

# K-means



$$\text{argmin}_{c_1,\dots,c_K} \sum_i \left\| x_i - c_{\textbf{closest\_to}(i)} \right\|^2$$

# K-means

$$\text{argmin}_{c_1,\ldots,c_K} \sum_i \left\| x_i - c_{\text{closest\_to}(i)} \right\|^2$$

- **Need to find cluster centers $c_k$.**
$$c_1 = ?, c_2 = ?, \ldots, c_K = ?$$

# K-means

$$\mathrm{argmin}_{c_1,\ldots,c_K} \sum_i \left\| x_i - c_{\mathbf{closest\_to}(i)} \right\|^2$$

- **Need to find cluster centers $c_k$.**

$$c_1 = ?, c_2 = ?, \ldots, c_K = ?$$

- **Introduce *latent variables* (one for each $x_i$)**

$$a_i = \mathbf{closest\_cluster\_center}(i)$$

$$a_1 = ?, a_2 = ?, a_3 = ?, \ldots, a_n = ?$$

# K-means

$$\mathrm{argmin}_{c_1,\dots,c_K} \sum_i \left\| x_i - c_{\mathbf{closest\_to}(i)} \right\|^2$$

- **For fixed $c_k$ we can find optimal $a_i$**

- **For fixed $a_i$ we can find optimal $c_k$.**

- **Iterate to convergence.**

# Fuzzy vs Hard



Each object belongs to each cluster with some weight (the weight can be zero)

Each object belongs to exactly one cluster

IFI Summer School. June 2014

# Gaussian Mixture Modeling

$$X \sim [N(\boldsymbol{\mu}_1, \sigma_1^2) \quad \text{or} \quad N(\boldsymbol{\mu}_2, \sigma_2^2)]$$

**Given $X$, estimate $\mu_i, \sigma_i^2$**

# Gaussian Mixture Modeling

$$X \sim [N(\boldsymbol{\mu}_1, \sigma_1^2) \quad \text{or} \quad N(\boldsymbol{\mu}_2, \sigma_2^2)]$$

**Given $X$, estimate $\boldsymbol{\mu}_i, \sigma_i^2$**

$\Longrightarrow$ **MLE**

# Gaussian Mixture Modeling

$$X \sim [N(\boldsymbol{\mu}_1, \sigma_1^2) \ \text{ or } \ N(\boldsymbol{\mu}_2, \sigma_2^2)]$$

**Given $X$, estimate $\mu_i, \sigma_i^2$**

$\Rightarrow$ **MLE**

$\Rightarrow$ **Expectation-Maximization (EM)**

# SKLearn's Clustering

```
from sklearn.cluster
import
    Ward,
    KMeans,
    DBScan,
    MeanShift,
    SpectralClustering,
    AffinityPropagation
```

# SKLearn's Clustering

**from** `sklearn.cluster`

**import**

`Ward,`
`KMeans,`
`DBScan,`
`MeanShift,`

**Use feature vectors**

**Use distance matrix**

`SpectralClustering,`
`AffinityPropagation`

# Quiz

▸ Fuzzy clustering means that _____

▸ K-means finds a set of cluster centers, which have the smallest _____

▸ K-means can get stuck in a local minimum (Y/N)?

# Unsupervised learning patterns

**Decomposition**

**Clustering**

# Canonical basis



$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \alpha \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \beta \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

# Alternative basis

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$
$$= \alpha \begin{pmatrix} 0.9 \\ -0.1 \end{pmatrix} + \beta \begin{pmatrix} 0.1 \\ 0.9 \end{pmatrix}$$

# Alternative basis



$$\binom{x_1}{x_2}$$
$$= \alpha \binom{0.4}{-0.6} + \beta \binom{0.6}{0.4}$$

# Linear Decomposition

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \ldots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \cdots + \alpha_m \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$

# Linear Decomposition

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \cdots + \alpha_m \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$

# Linear Decomposition

$$
\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 0.0 \\ 0.1 \\ 0.1 \\ 0.2 \\ \vdots \\ 0.0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0.3 \\ 0.2 \\ 0.2 \\ 0.1 \\ \vdots \\ 0.3 \end{pmatrix} + \alpha_m \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.0 \\ \vdots \\ 0.0 \end{pmatrix}
$$

# Linear Decomposition

$$
\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 0.0 \\ 0.1 \\ 0.1 \\ 0.2 \\ \vdots \\ 0.0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0.3 \\ 0.2 \\ 0.2 \\ 0.1 \\ \vdots \\ 0.3 \end{pmatrix} + \alpha_m \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.0 \\ \vdots \\ 0.0 \end{pmatrix}
$$

# Linear Decomposition

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \ldots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + \cdots + \alpha_m \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix}$$
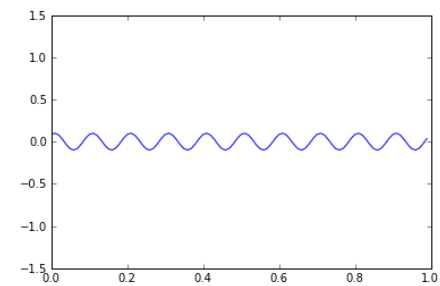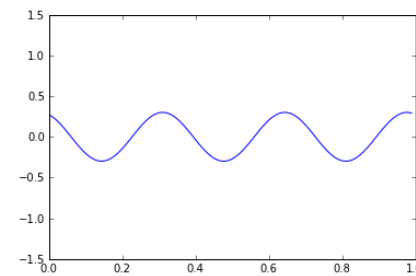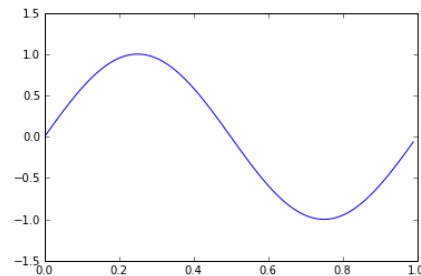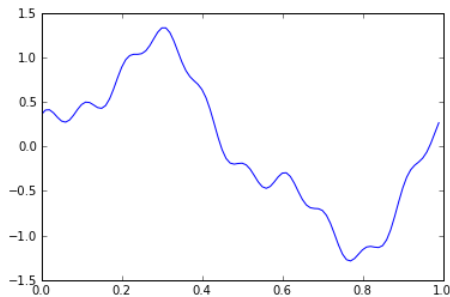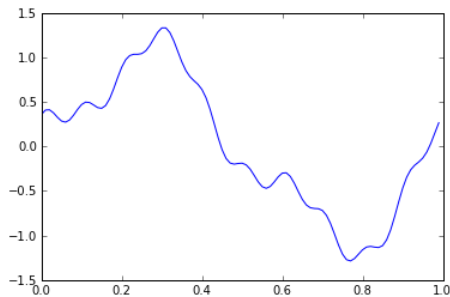
# Linear Decomposition

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 0.0 \\ 0.1 \\ 0.1 \\ 0.2 \\ \vdots \\ 0.0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0.3 \\ 0.2 \\ 0.2 \\ 0.1 \\ \vdots \\ 0.3 \end{pmatrix} + \alpha_m \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.0 \\ \vdots \\ 0.0 \end{pmatrix}$$

# Linear Decomposition

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \dots \\ x_{100000} \end{pmatrix} = \alpha_1 \begin{pmatrix} 0.0 \\ 0.1 \\ 0.1 \\ 0.2 \\ \vdots \\ 0.0 \end{pmatrix} + \alpha_2 \begin{pmatrix} 0.3 \\ 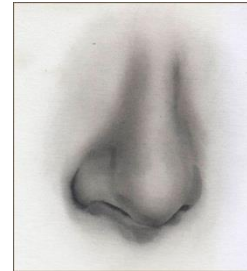0.2 \\ 0.2 \\ 0.1 \\ \vdots \\ 0.3 \end{pmatrix} + \alpha_m \begin{pmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.0 \\ \vdots \\ 0.0 \end{pmatrix}$$

# Linear Decomposition

$$x = \alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_m v_m$$

# Linear Decomposition

$$x = \alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_m v_m$$

$$x = \begin{pmatrix} \vdots & \vdots & \cdots & \vdots \\ v_1 & v_2 & & v_m \\ \vdots & \vdots & \cdots & \vdots \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_m \end{pmatrix}$$

# Linear Decomposition

$$x = \alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_m v_m$$

$$x = V\alpha$$

# Linear Decomposition

$$x = \alpha_1 v_1 + \alpha_2 v_2 + \cdots + \alpha_m v_m$$

$$x = V\alpha$$

$$\alpha = ?$$

# Linear Decomposition

$$x = \alpha_1 \boldsymbol{v_1} + \alpha_2 \boldsymbol{v_2} + \cdots + \alpha_m \boldsymbol{v_m}$$

$$\boldsymbol{x} = \boldsymbol{V\alpha}$$
$$\boldsymbol{\alpha} = \boldsymbol{V^+ x}$$

# How do we find a good basis?

# Linear projection

# Linear projection

# Linear projection

# Linear projection

# Idea: Maximize projection variance

▸ For a point $x_i$ and a unit basis vector $v$ the length of projection of $x_i$ onto $v$ is given by

$$p = \langle v, x_i \rangle = v^T x_i$$

# Projection variance

$$p_i = v^T x_i$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_v^2 = \frac{1}{n} \sum_i (p_i - \overline{p})^2$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_v^2 = \frac{1}{n} \sum_i (p_i - \overline{p})^2$$

$$\boldsymbol{v} = \text{argmax}_{\boldsymbol{v}} \, \sigma_v^2$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_v^2 = \frac{1}{n} \sum_i (p_i - \overline{p})^2$$

Pre-center data, so that $\overline{p} = \boldsymbol{v}^T \overline{\boldsymbol{x}} = 0$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_{\boldsymbol{v}}^2 = \frac{1}{n} \sum_i (p_i)^2$$

Pre-center data, so that $\overline{p} = \boldsymbol{v}^T \overline{\boldsymbol{x}} = 0$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_v^2 = \frac{1}{n} \sum_i (p_i)^2 = \frac{1}{n} \|\boldsymbol{p}\|^2$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_{\boldsymbol{v}}^2 = \frac{1}{n} \sum_i (p_i)^2 = \frac{1}{n} \|\boldsymbol{p}\|^2$$

$$= \frac{1}{n} \|\boldsymbol{X}\boldsymbol{v}\|^2$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_{\boldsymbol{v}}^2 = \cdots$$

$$= \frac{1}{n}\|\boldsymbol{X}\boldsymbol{v}\|^2 = \frac{1}{n}(\boldsymbol{X}\boldsymbol{v})^T(\boldsymbol{X}\boldsymbol{v})$$

$$= \frac{1}{n}\boldsymbol{v}^T \boldsymbol{X}^T \boldsymbol{X}\boldsymbol{v} = \boldsymbol{v}^T \boldsymbol{\Sigma}\boldsymbol{v}$$

# Projection variance

$$p_i = v^T x_i$$

$$\sigma_v^2 = v^T \Sigma v$$

# Projection variance

$$p_i = \boldsymbol{v}^T \boldsymbol{x}_i$$

$$\sigma_{\boldsymbol{v}}^2 = \boldsymbol{v}^T \boldsymbol{\Sigma} \boldsymbol{v}$$

**Data covariance matrix**
$$X^T X$$

# Objective function

$$\text{argmax}_{v} \ v^T \Sigma v$$

$$s.t. \|v\|^2 = 1$$

# Optimization

$$\text{argmax}_v \; \boldsymbol{v}^T \boldsymbol{\Sigma} \boldsymbol{v}$$

$$s.t. \|\boldsymbol{v}\|^2 = 1$$

**Method of Lagrange multipliers…**

$$\boldsymbol{\Sigma} \boldsymbol{v} = \lambda \boldsymbol{v}$$

# Optimization

$$\text{argmax}_v \; v^T \Sigma v$$

$$s.t. \|v\|^2 = 1$$

**Method of Lagrange multipliers…**

$$\Sigma v = \lambda v$$
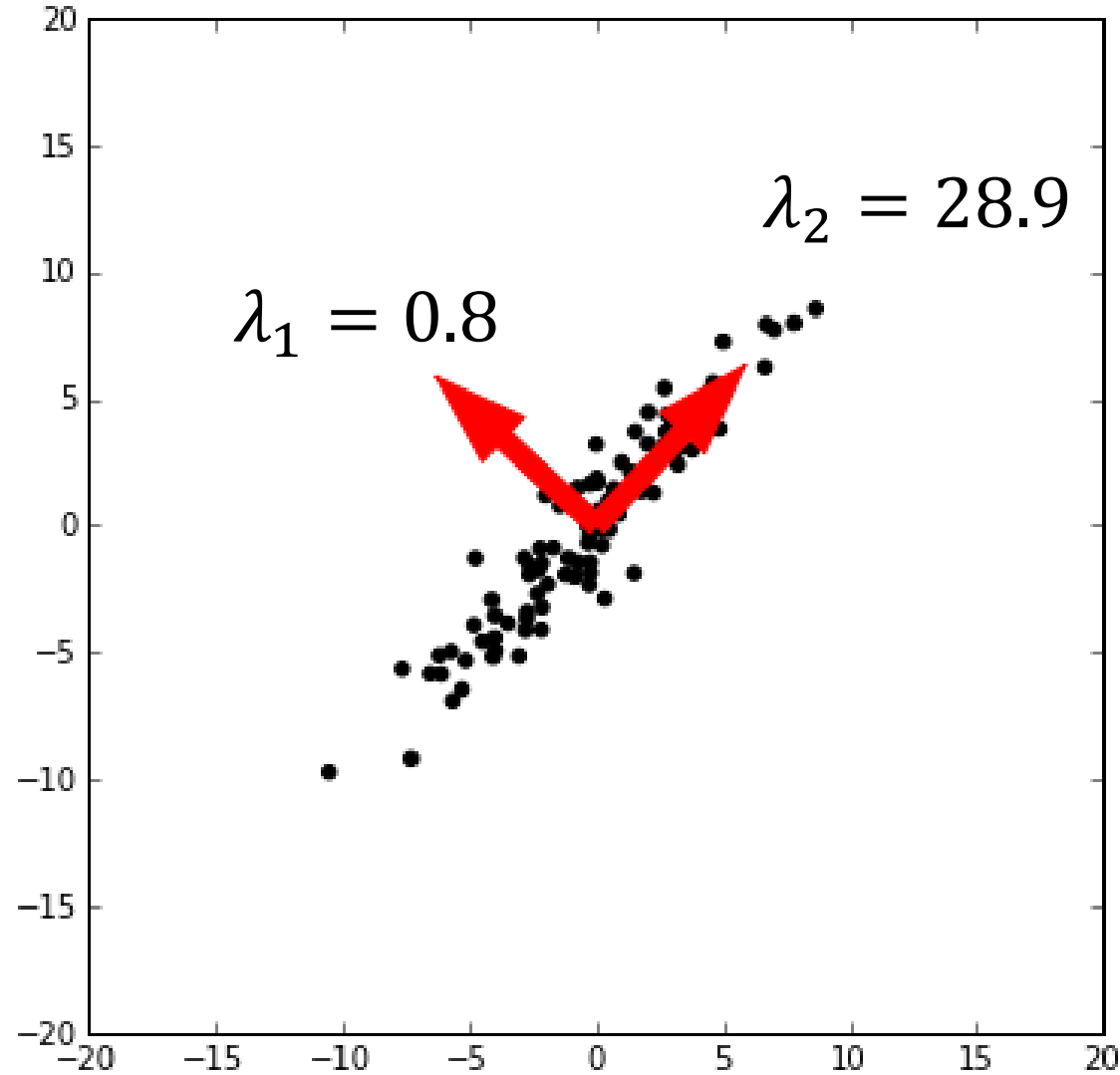
**Eigenvector of $\Sigma$**     **Eigenvalue**

# Example

```
Xc = X - mean(X, axis=0)

Sigma = Xc.T * Xc / n

                (= cov(Xc, rowvar=0))

lambdas, vs = eigh(Sigma)
```

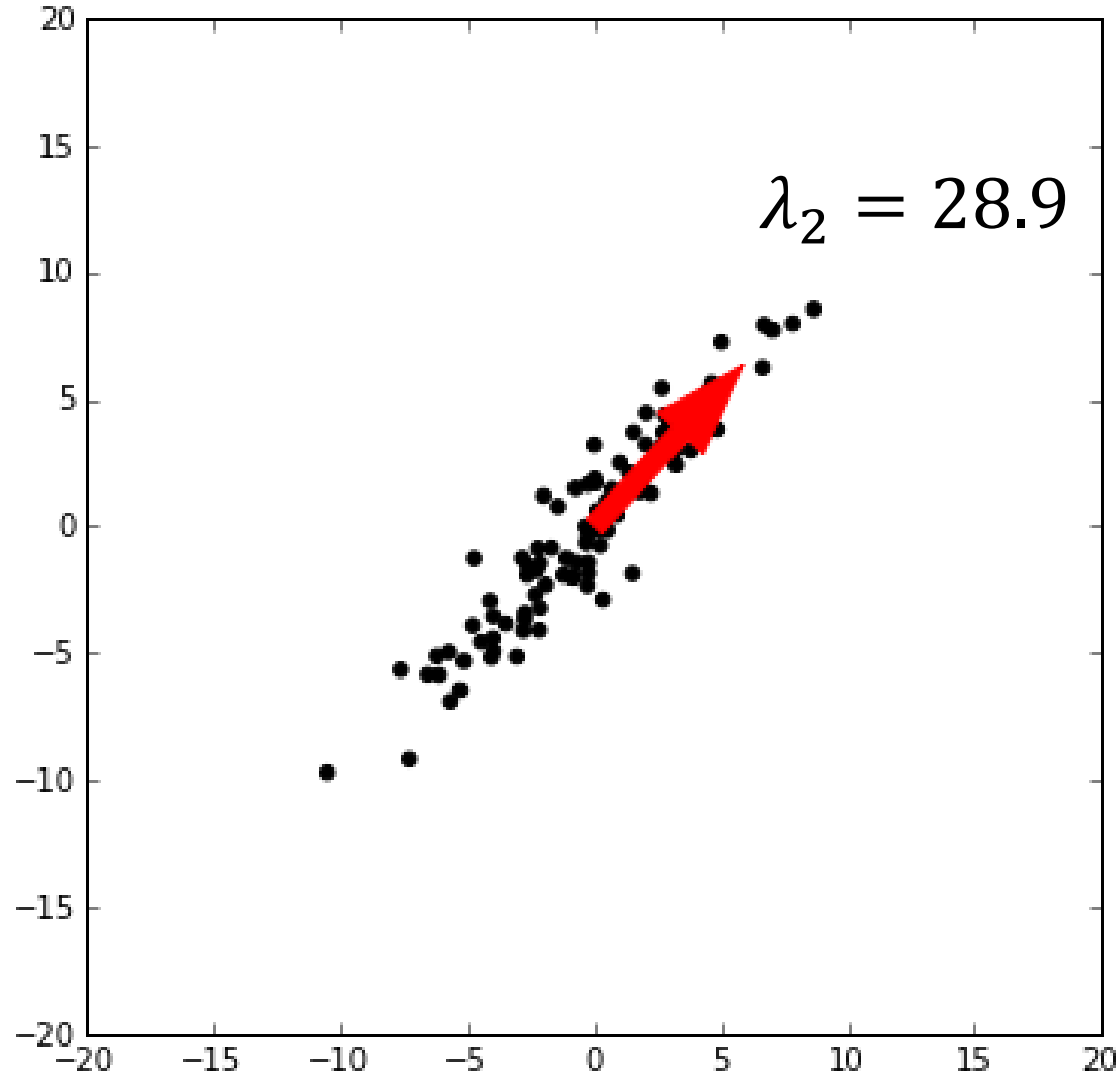# Example

# Example



$$\lambda_2 = 28.9$$

$$\lambda_1 = 0.8$$

$$\sigma_i^2 = \boldsymbol{v}_i^T \boldsymbol{\Sigma} \boldsymbol{v}_i = \boldsymbol{v}_i^T \lambda_i \boldsymbol{v}_i = \lambda_i \|\boldsymbol{v}_i\|^2 = \lambda_i$$

# Example



$$\lambda_2 = 28.9$$

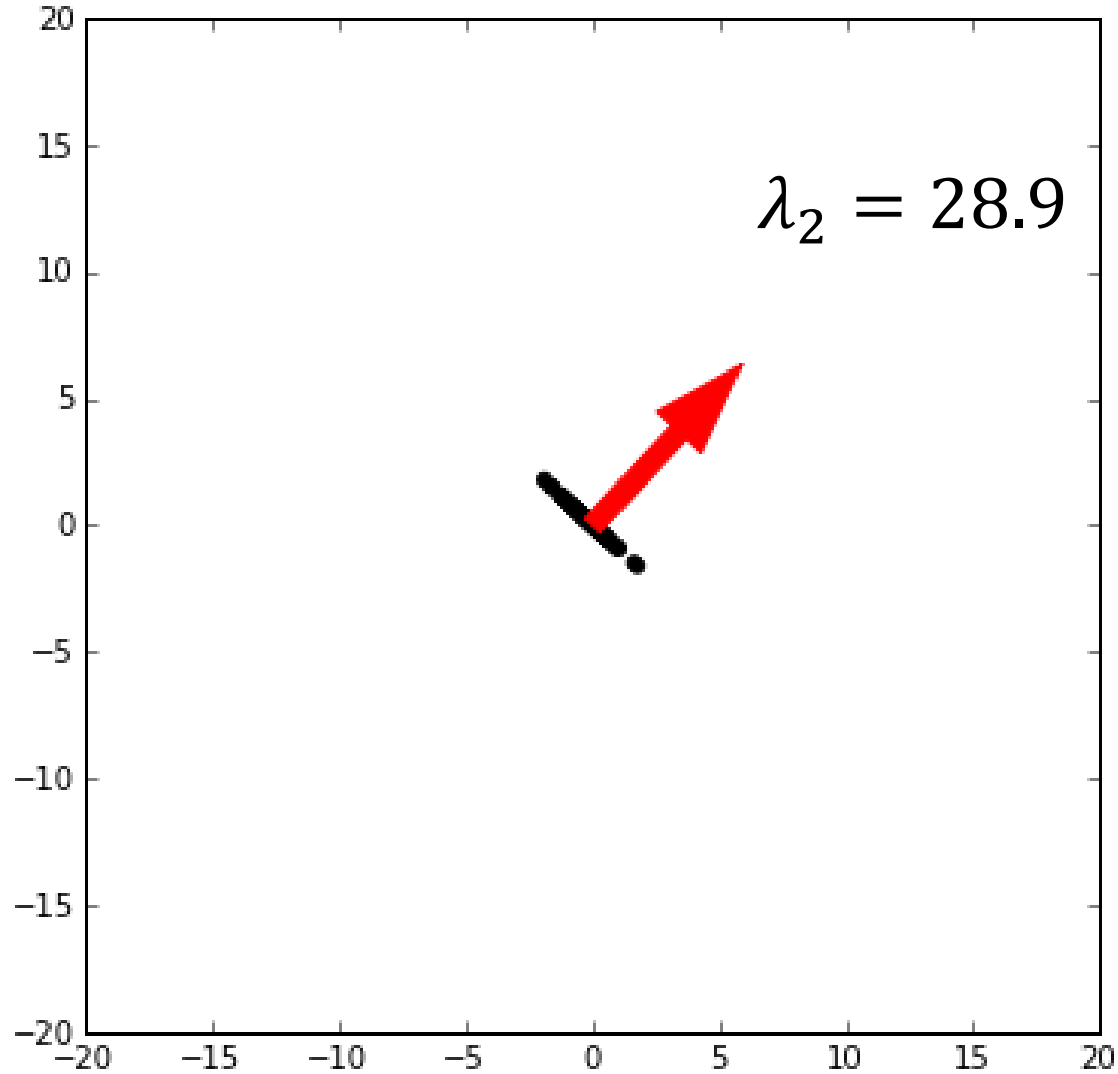# Example



$$\lambda_2 = 28.9$$

# Example



$\lambda_2 = 28.9$

$\lambda_1 = 0.8$
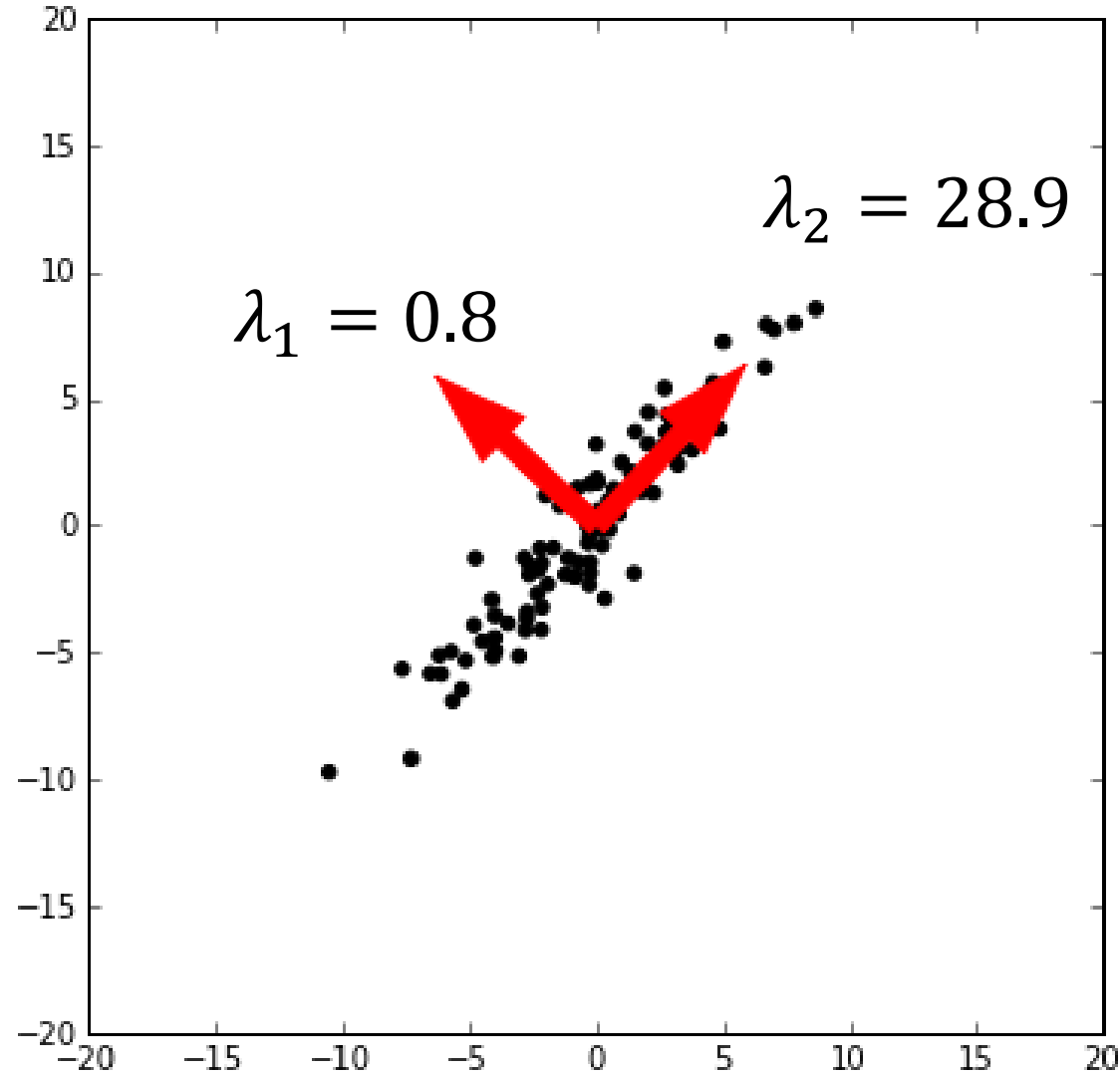
# Principal Components Analysis

**Principal components** are the **eigenvectors** of the **covariance matrix**.

$$V, \lambda = \mathbf{eig}(\Sigma)$$
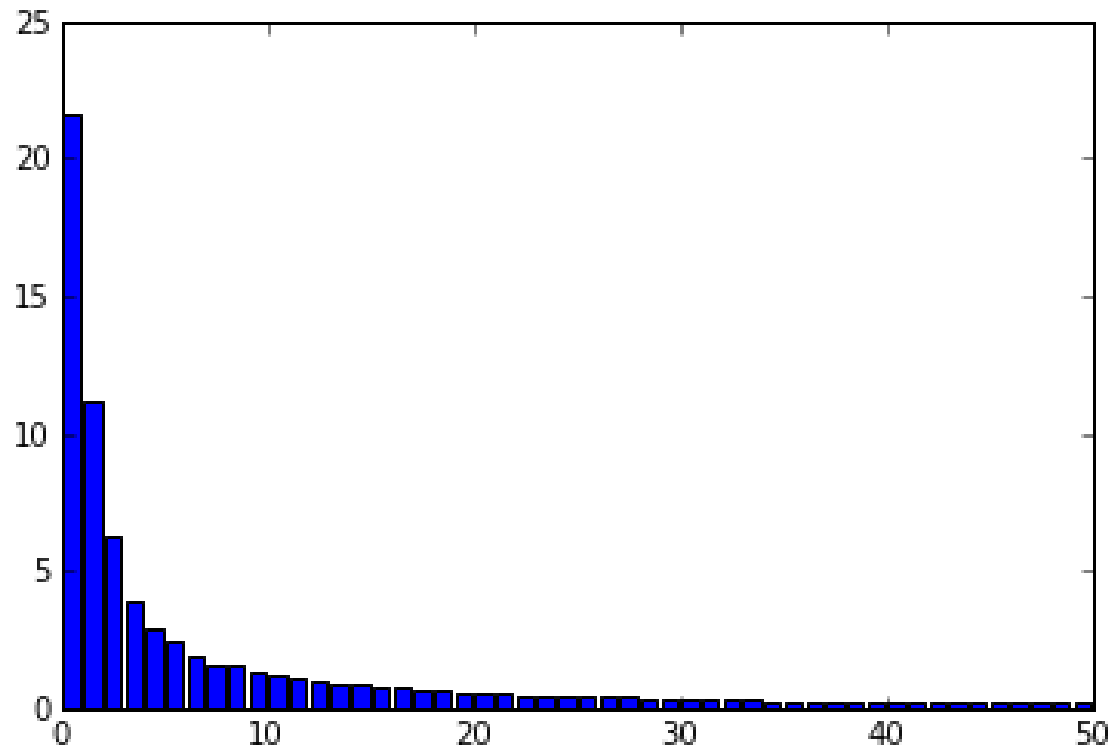
# Principal Components Analysis

**Principal components** are the **eigenvectors** of the **covariance matrix**.
$$V, \lambda = \mathbf{eig}(\Sigma)$$

For each PC, the corresponding eigenvalue $\lambda_i$ shows the **amount of variance explained** by the component.

# Principal Components Analysis

## Eigenvalue spectrum of $\Sigma$

# Principal Components Analysis

Data projection onto PC $i$:
$$p = Xv_i$$

Data projection onto multiple PCs:
$$X_{\mathrm{proj}} = XV_*$$

Data reconstruction from PC coordinates:
$$X_{\mathrm{proj}}V_*^T = X$$

# SKLearn's PCA

```
from sklearn.decomposition
                    import PCA


model = PCA(n_components=2)
model.fit(X)
X_t = model.transform(X)


model.components_[1,:]
```
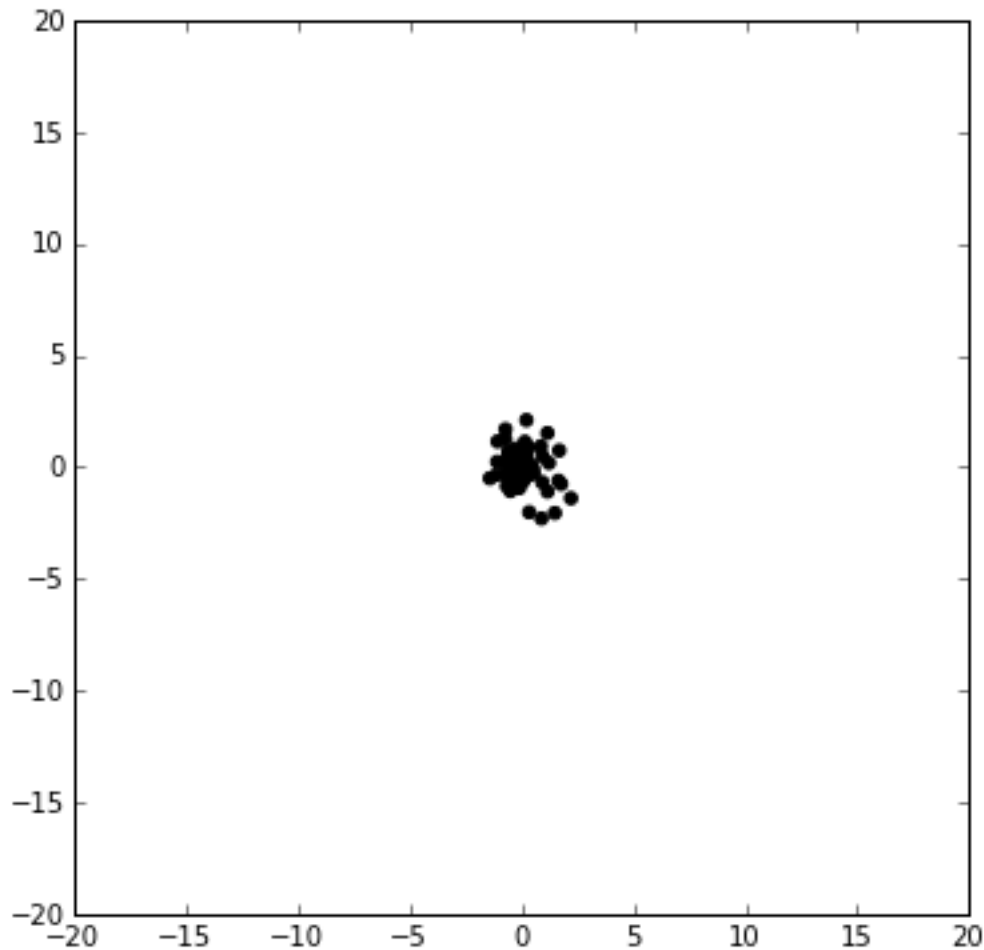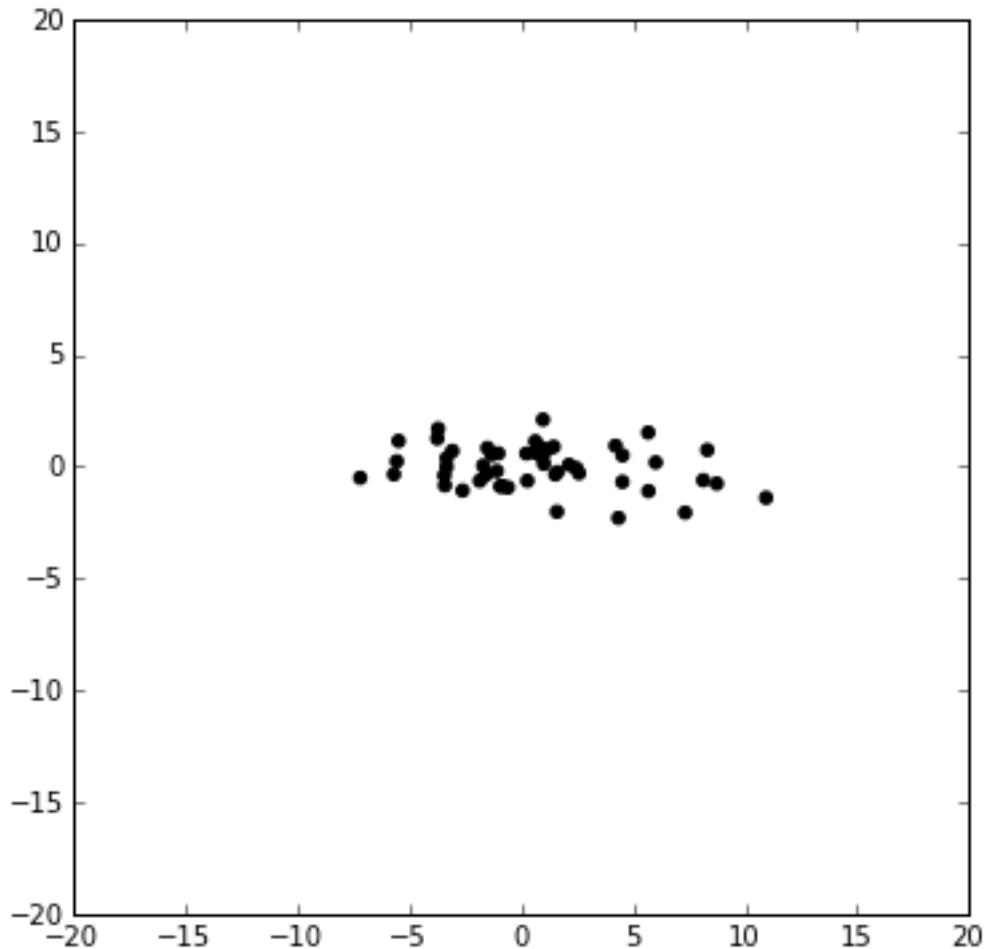
# SKLearn's PCA

```
from sklearn.decomposition
import
    PCA,
    SparsePCA,
    ProbabilisticPCA,
    KernelPCA,
    FastICA,
    NMF,
    DictionaryLearning,
    ...
```
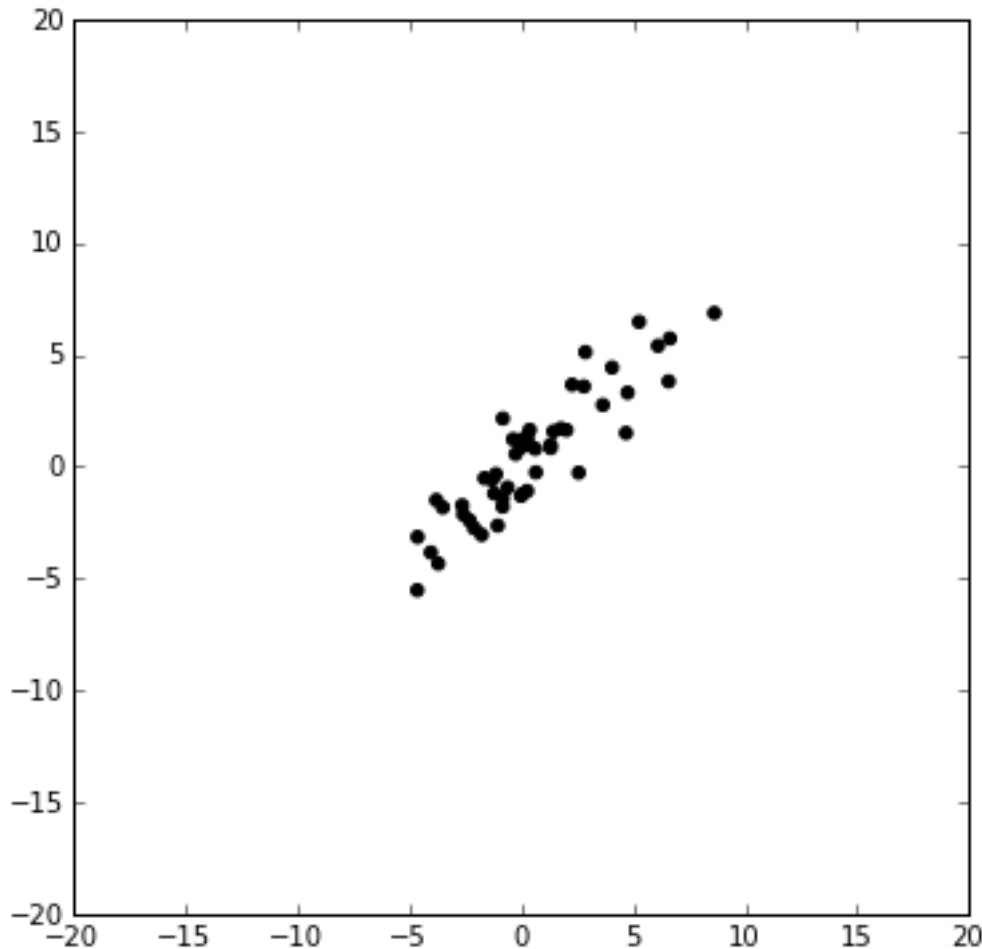
# PCA: Geometric intuition



$$\boldsymbol{X} \sim N(0,1)$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

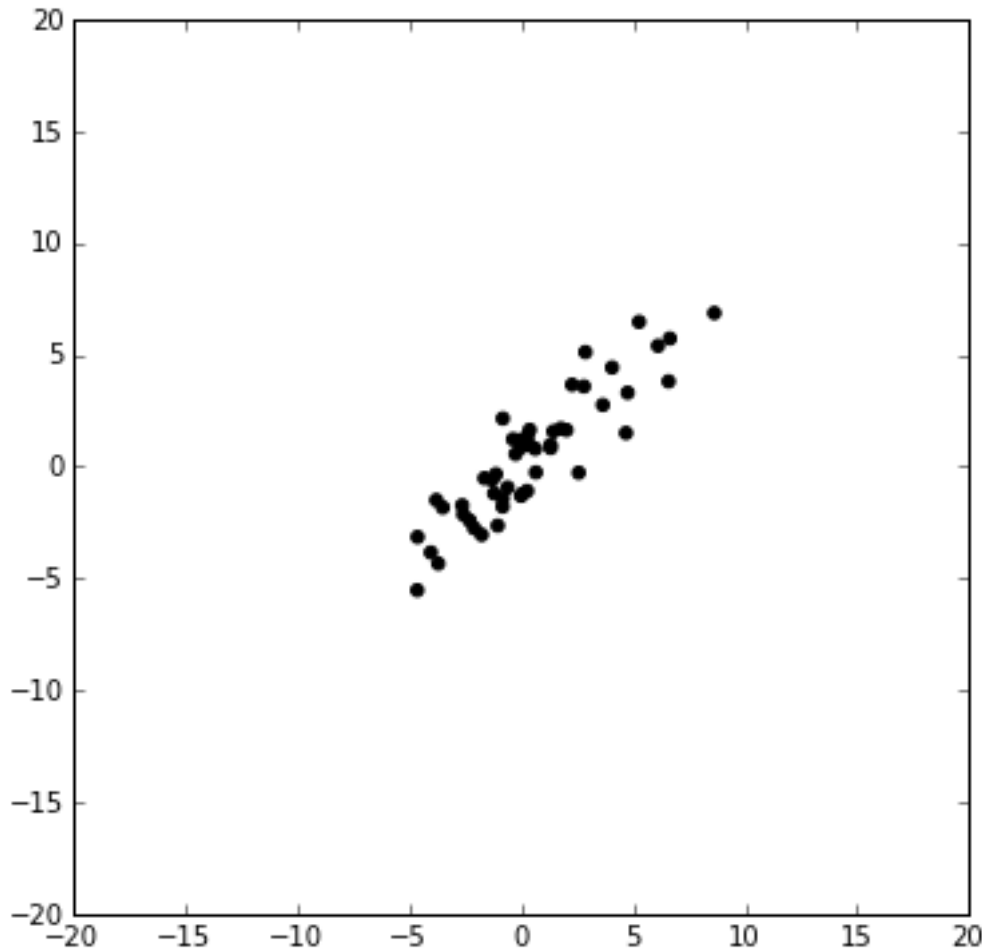$$X' = X \begin{pmatrix} 5 & 0 \\ 0 & 0.9 \end{pmatrix}$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X' = X \begin{pmatrix} 5 & 0 \\ 0 & 0.9 \end{pmatrix}$$

$$X'' = X' \begin{pmatrix} \cos 0.8 & -\sin 0.8 \\ \sin 0.8 & \cos 0.8 \end{pmatrix}$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

$$(X'')^T(X'') = (XDR)^T(XDR)$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

$$(X'')^T (X'')$$
$$= (XDR)^T (XDR)$$
$$= R^T D^T X^T X D R$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

$$(X'')^T (X'')$$
$$= (XDR)^T (XDR)$$
$$= R^T D^T X^T X D R$$
$$= R^T D^2 R$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

$$\Sigma = V \Lambda V^T$$
$$(V = R^T, \Lambda = D^2)$$

$$(X'')^T(X'')$$
$$= (XDR)^T(XDR)$$
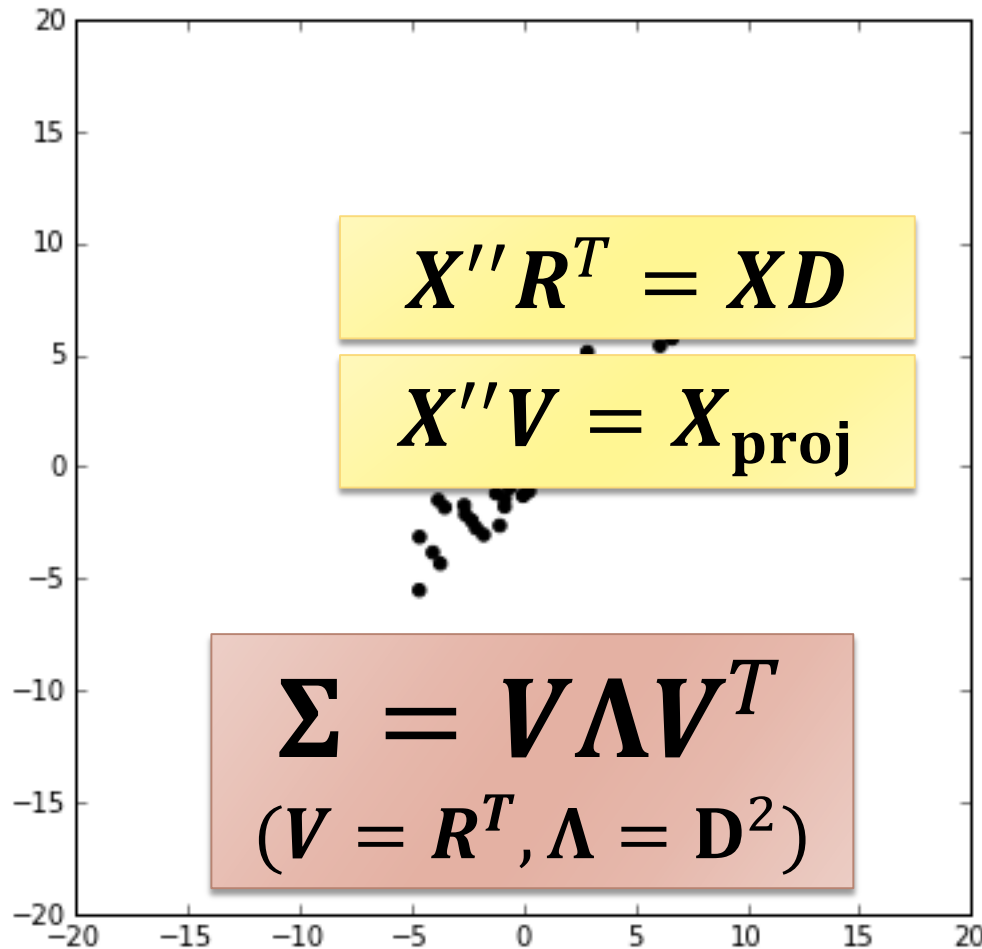$$= R^T D^T X^T X D R$$
$$= R^T D^2 R$$

# PCA: Geometric intuition



$$X \sim N(0,1)$$

$$X'' = X \cdot D \cdot R$$

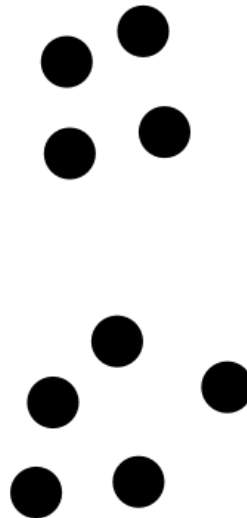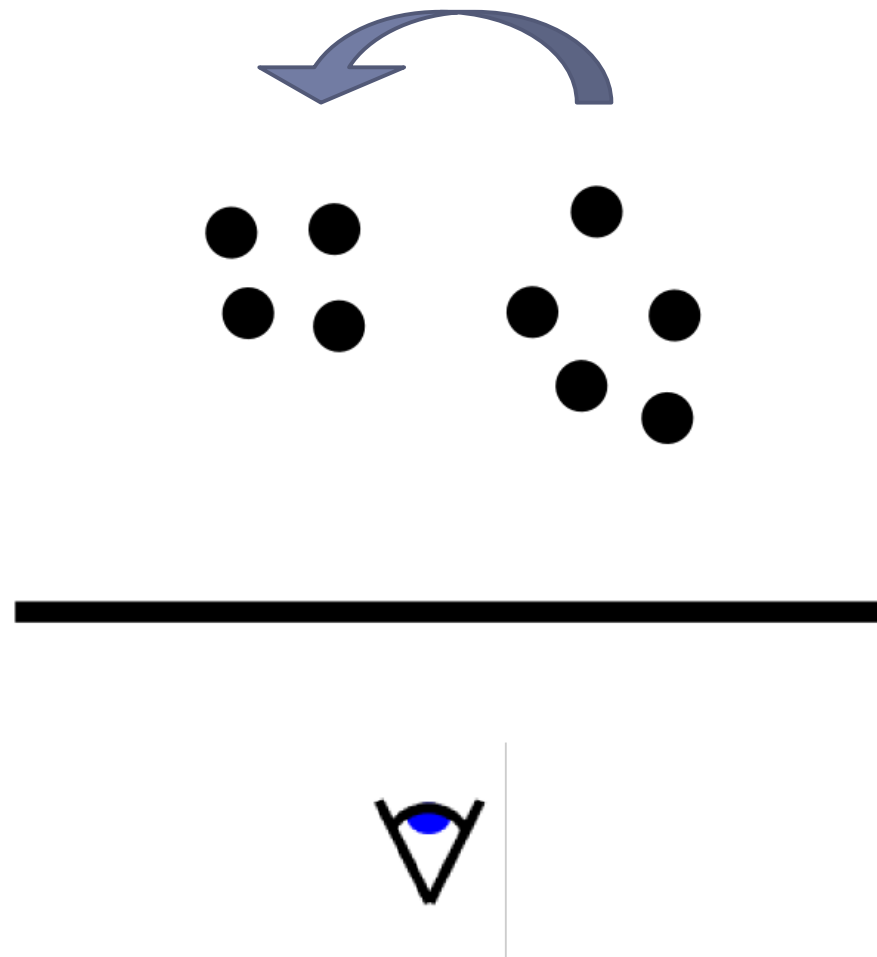$$\begin{aligned}(X'')^T(X'') \\ = (XDR)^T(XDR) \\ = R^T D^T X^T XDR \\ = R^T D^2 R\end{aligned}$$

Boxes on plot:

$$X'' R^T = XD$$

$$X'' V = X_{\text{proj}}$$

$$\Sigma = V\Lambda V^T$$
$$(V = R^T, \Lambda = D^2)$$

# Quiz

▸ Principal components are _____ of the _____ matrix.

▸ Eigenvalue spectrum shows how much _____ is explained by each _____.

▸ If $\Sigma = V\Lambda V^{T}$, then

$$X_{\text{proj}} = \underline{\quad\quad}$$