



香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 13

# Dimensionality Reduction

Junxian He  
Oct 22, 2024

# Midterm Exam

Tomorrow (Oct 24), 1:20pm-2:40pm, one A4-size double-sided cheetsheet is allowed  
(either printing or handwriting is fine)

We have two rooms for the exam for sparse seat plans:

1. For SIS ID ending with an even digit: Room 2303
2. For SIS ID ending with an odd digit: Room 2504

# High-Dimensional Data

- High-Dimensions = Lot of Features

Document classification

Features per document =

thousands of words/unigrams

millions of bigrams, contextual  
information

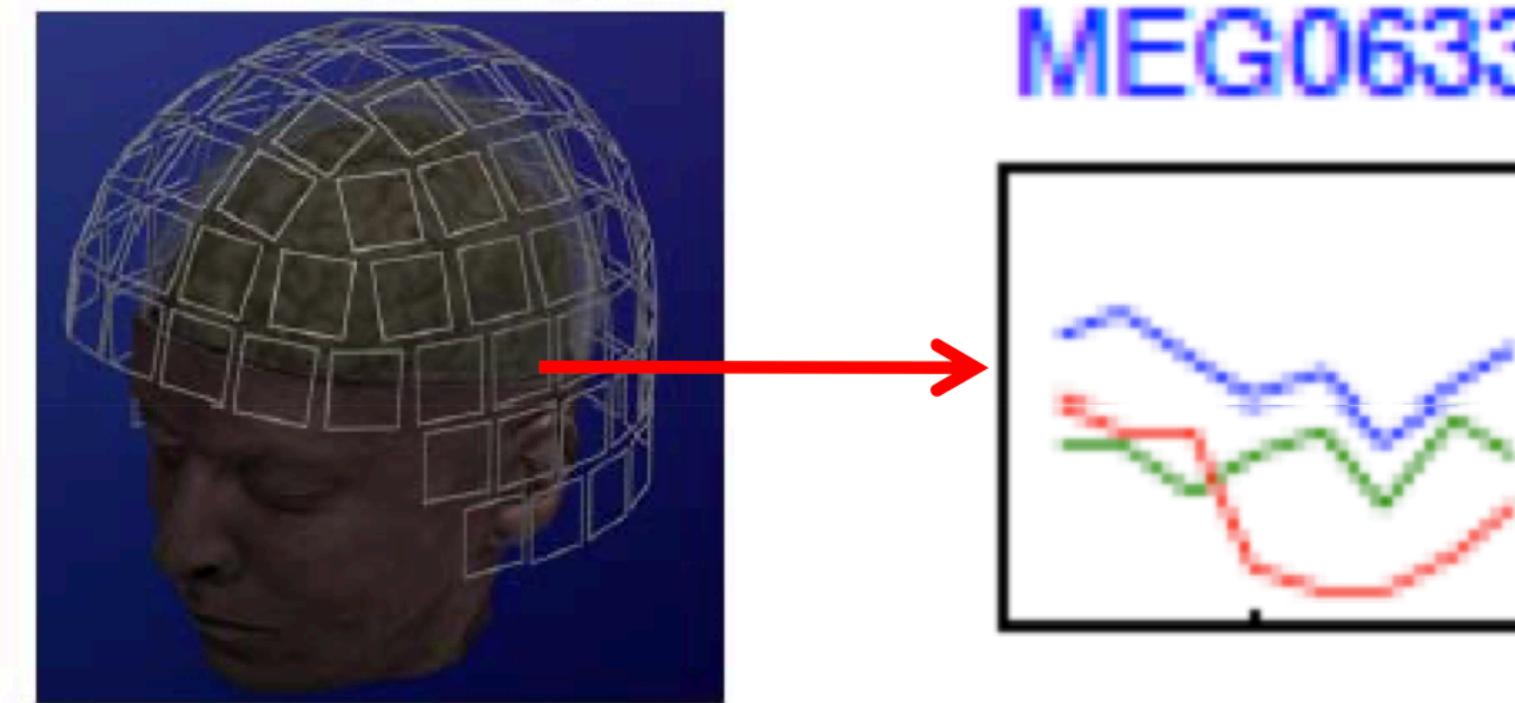


# High-Dimensional Data

- High-Dimensions = Lot of Features

## MEG Brain Imaging

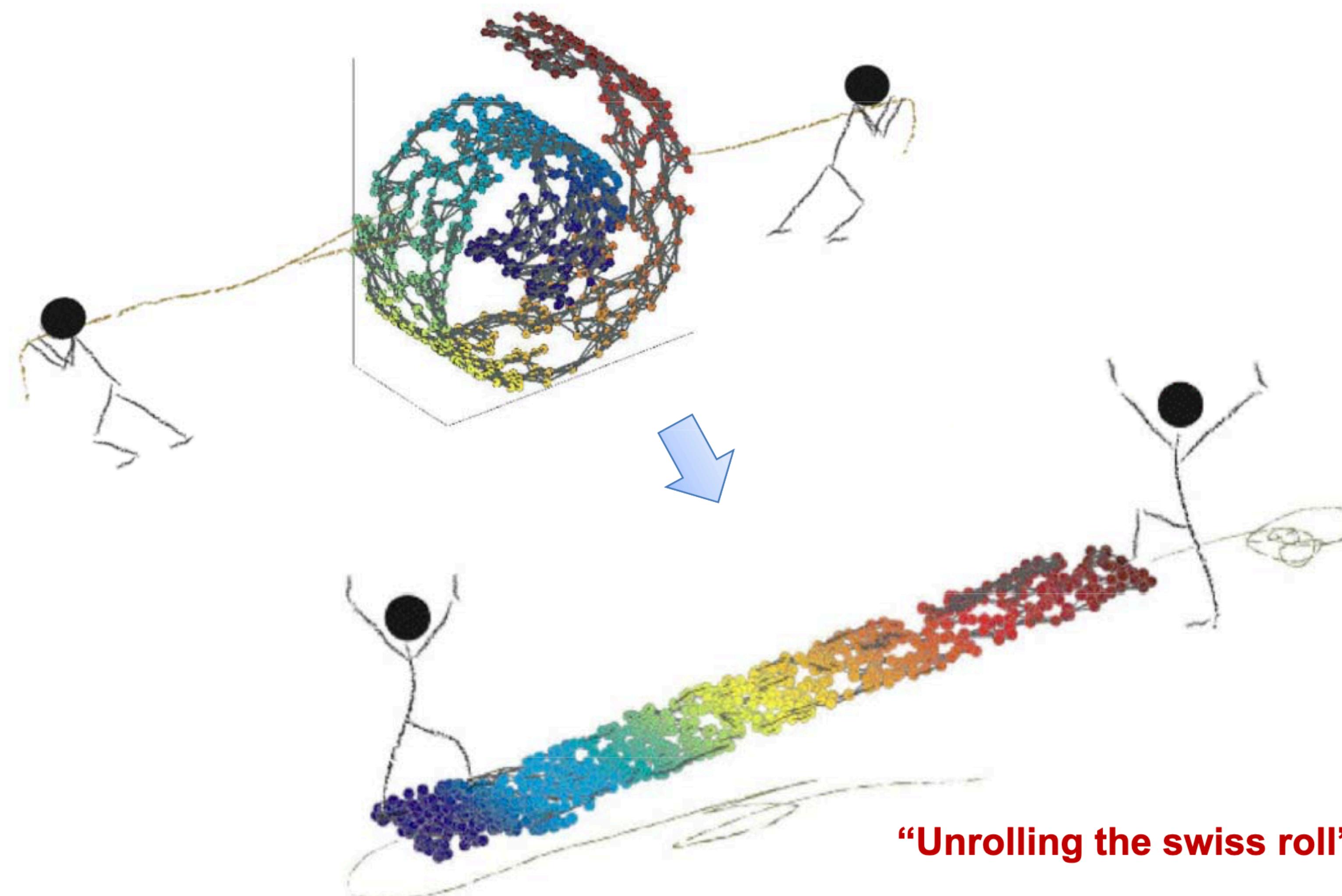
120 locations x 500 time points  
x 20 objects



# Curse of Dimensionality

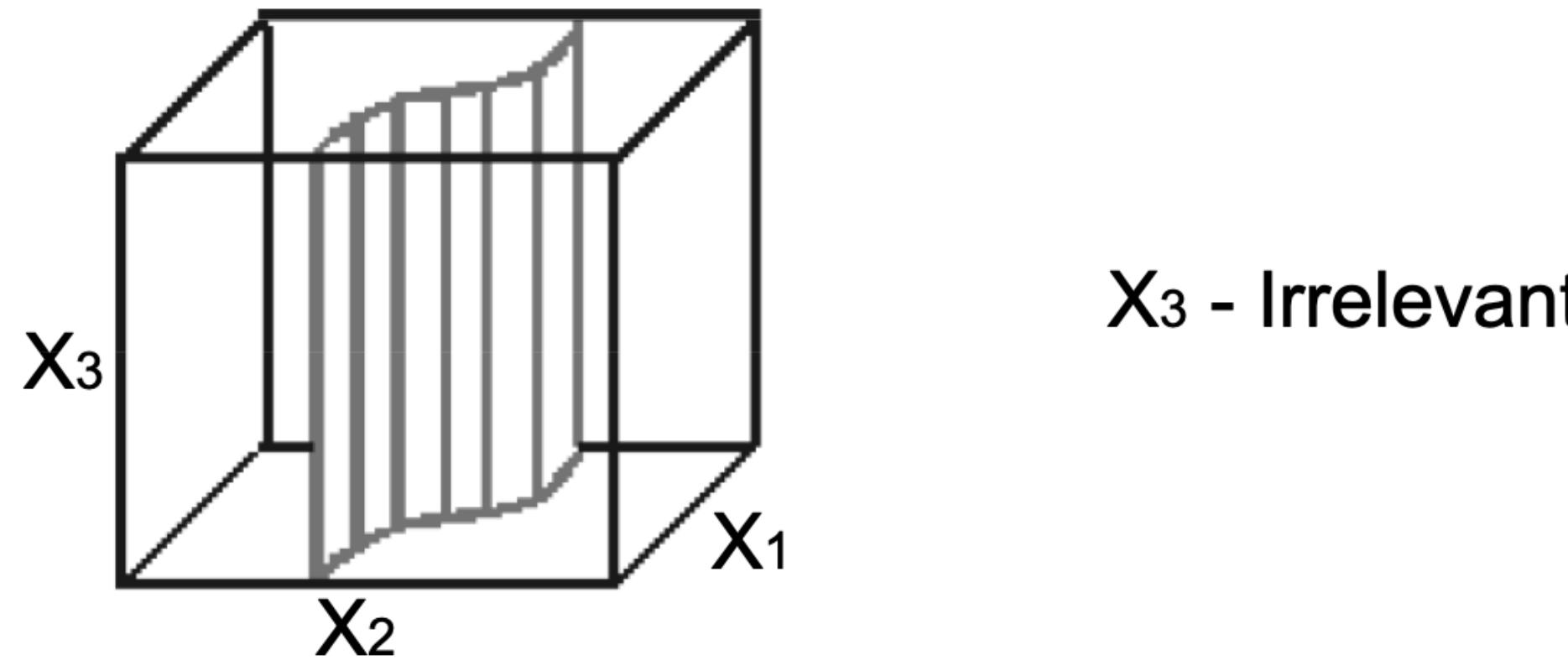
- Why are more features bad?
  - Redundant features (not all words are useful to classify a document)  
more noise added than signal
  - Hard to store and process data (computationally challenging)
  - Hard to interpret and visualize
  - Complexity of decision rule tends to grow with # features

# Dimensionality Reduction

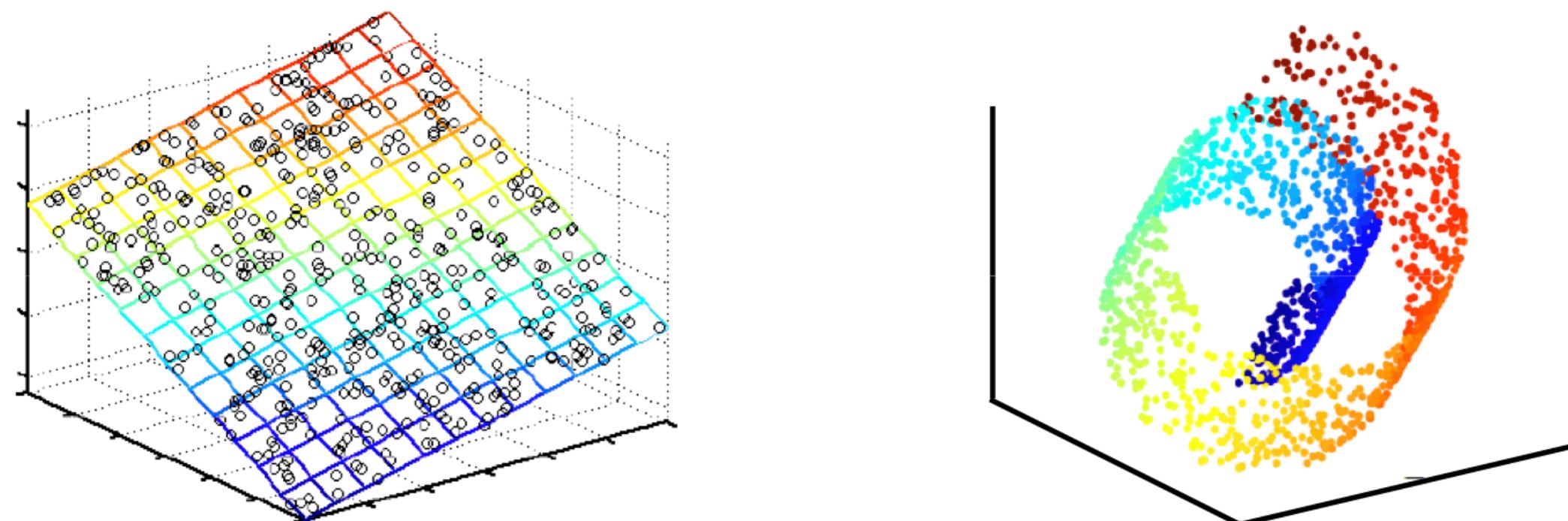


# Dimensionality Reduction

- Feature Selection – Only a few features are relevant to the learning task



- Latent features – Some linear/nonlinear combination of features provides a more efficient representation than observed features



# Latent Feature Extraction

Combinations of observed features provide more efficient representation, and capture underlying relations that govern the data

E.g. Ego, personality and intelligence are hidden attributes that characterize human behavior instead of survey questions

Topics (sports, science, news, etc.) instead of documents

- Linear

Principal Component Analysis (PCA)

Factor Analysis

Independent Component Analysis (ICA)

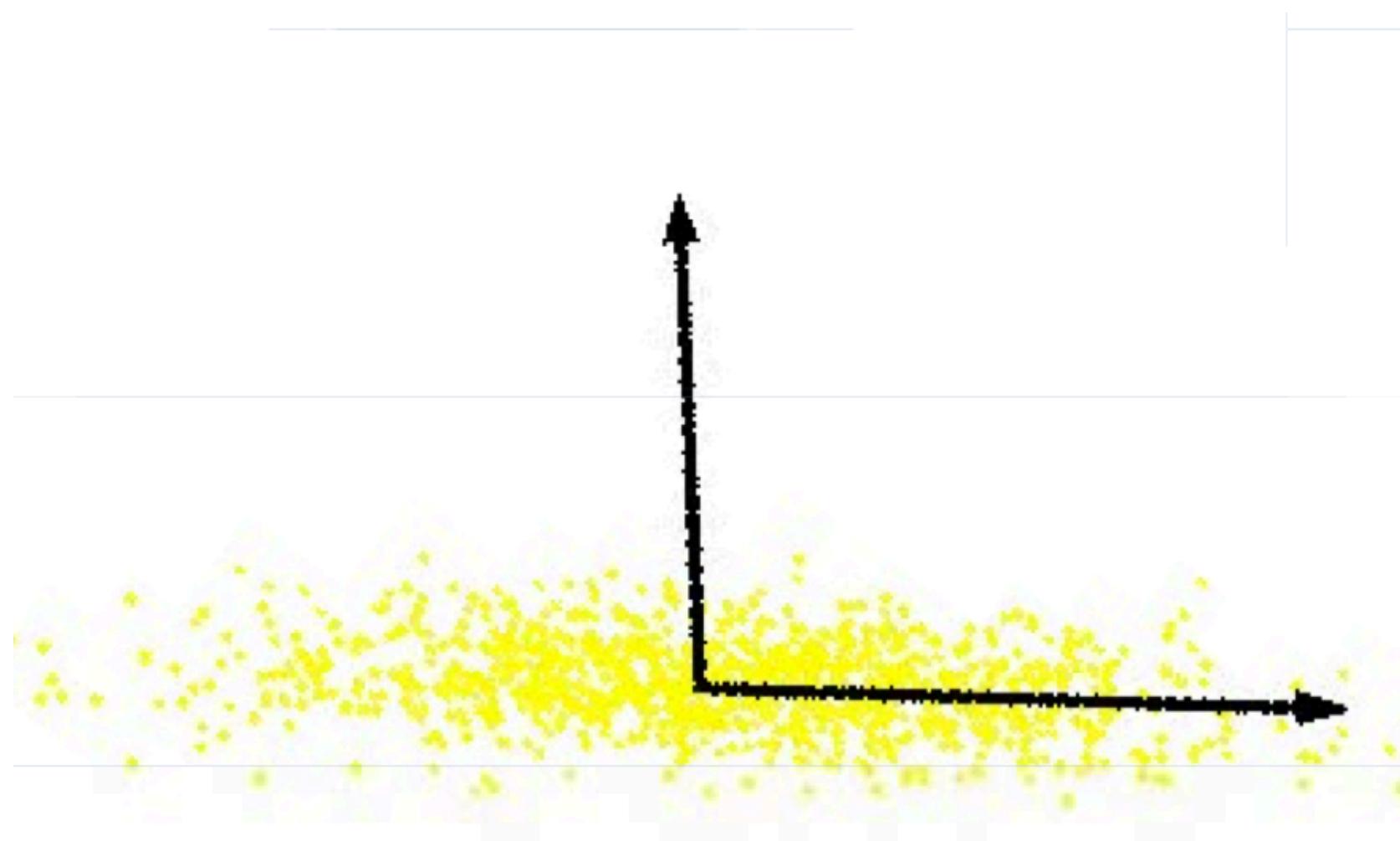
- Nonlinear

ISOMAP

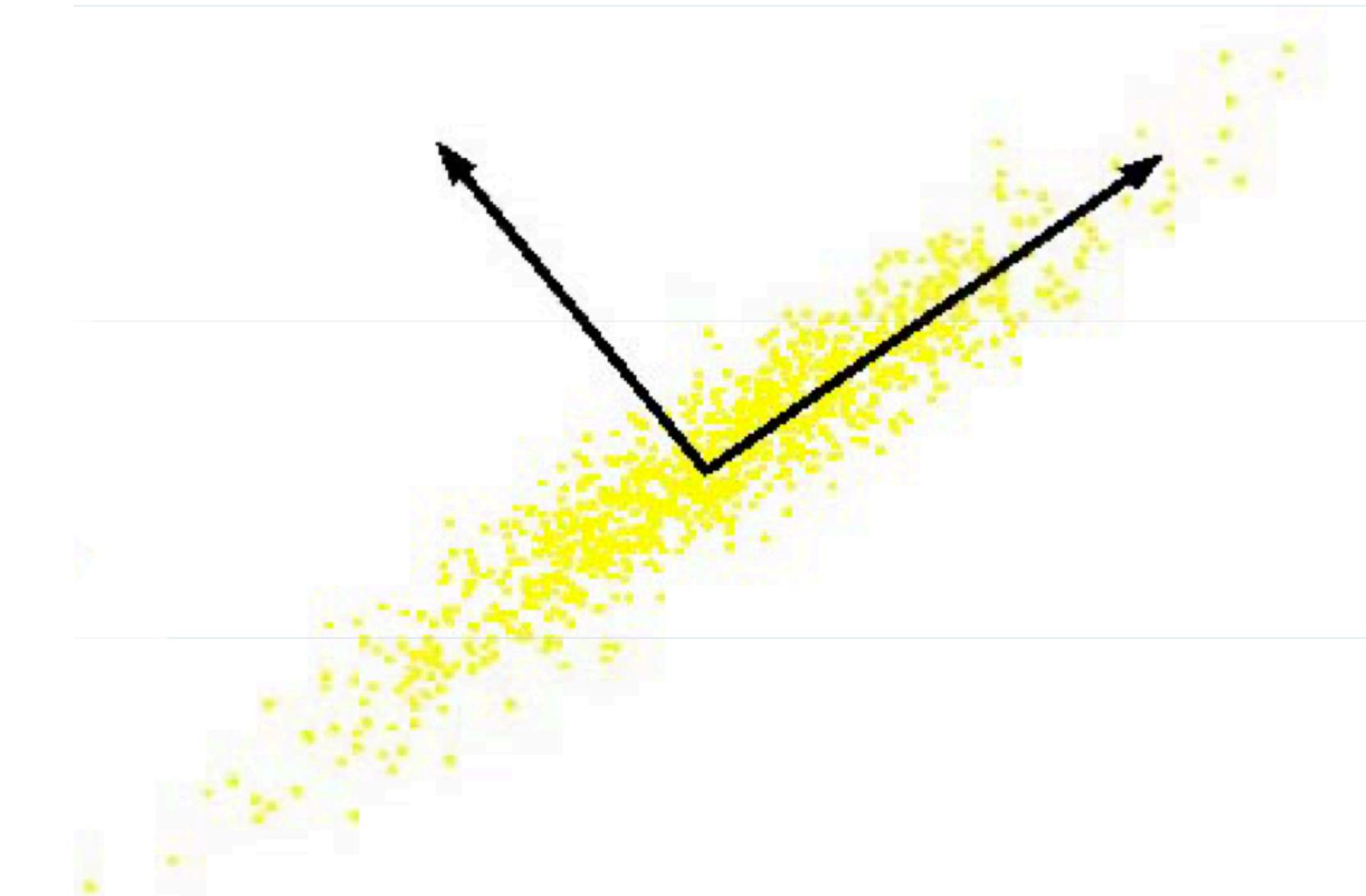
Local Linear Embedding (LLE)

Laplacian Eigenmaps

# Principal Component Analysis (PCA)



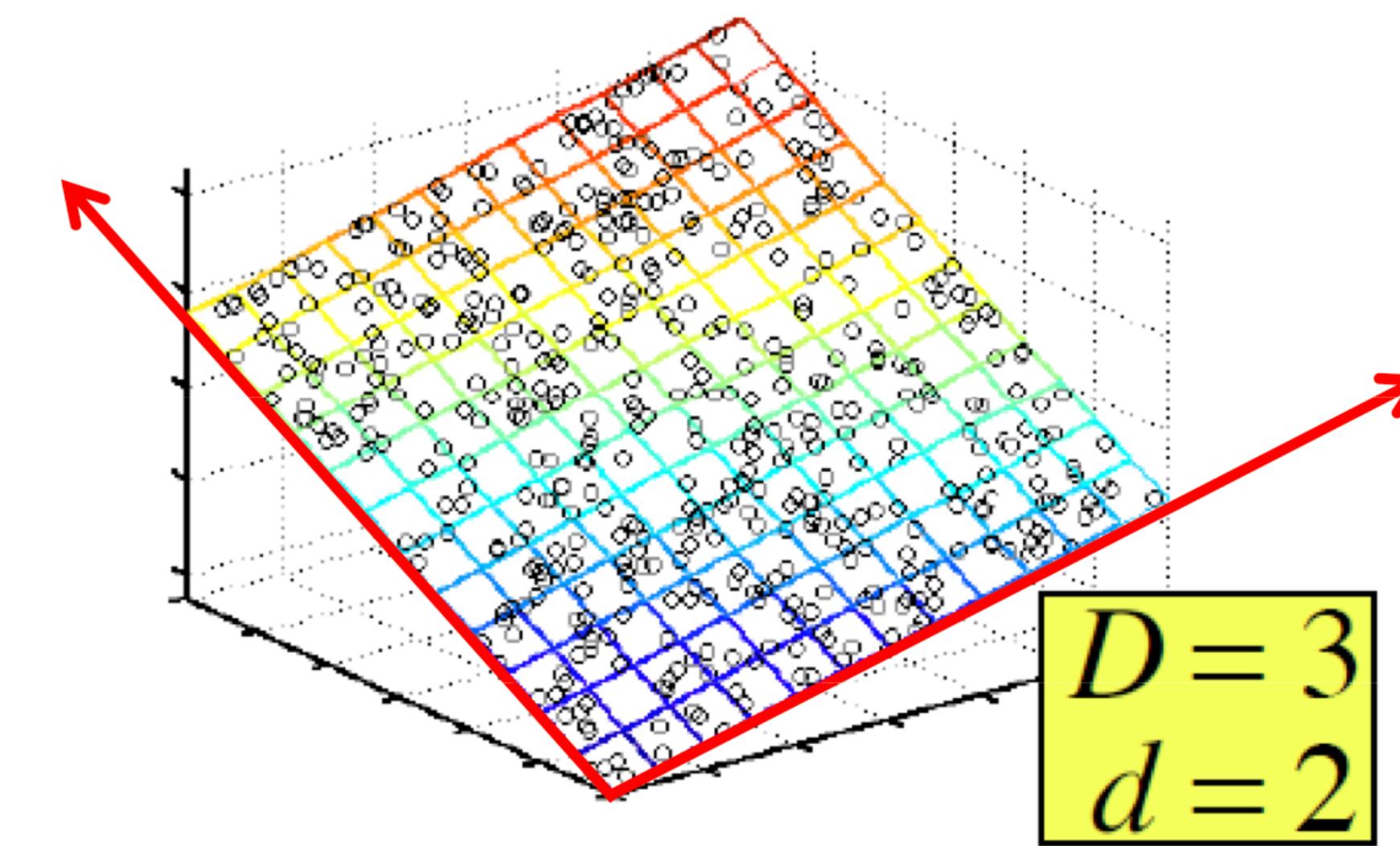
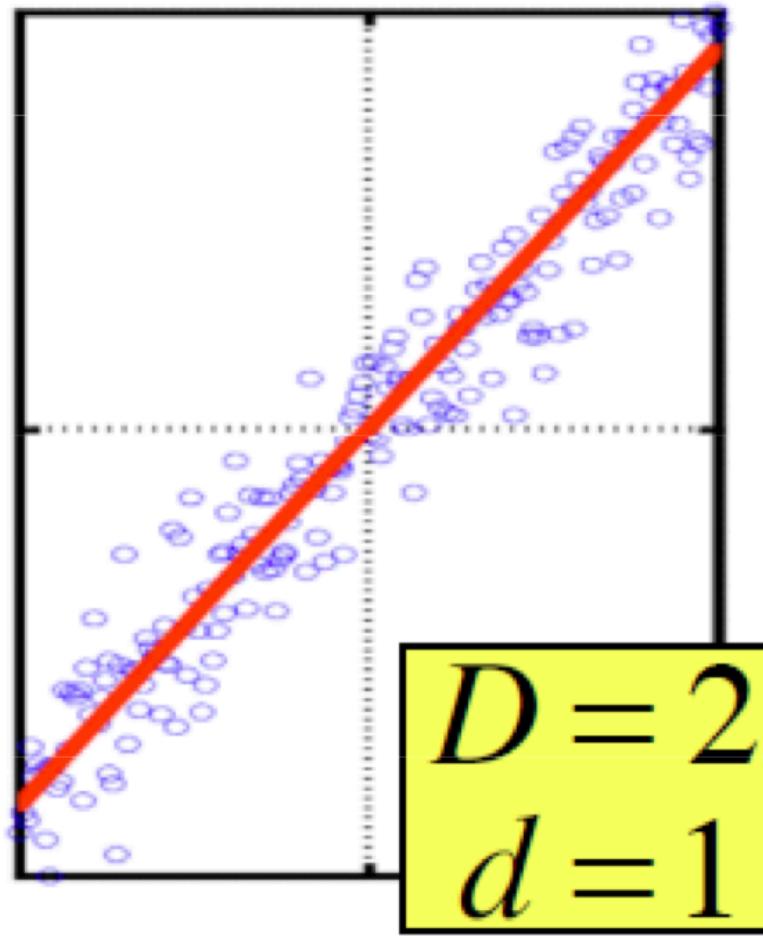
Only one relevant feature



Both features become relevant

Can we transform the features so that we only need to preserve one latent feature? Find linear projection so that projected data is uncorrelated.

# Principal Component Analysis (PCA)

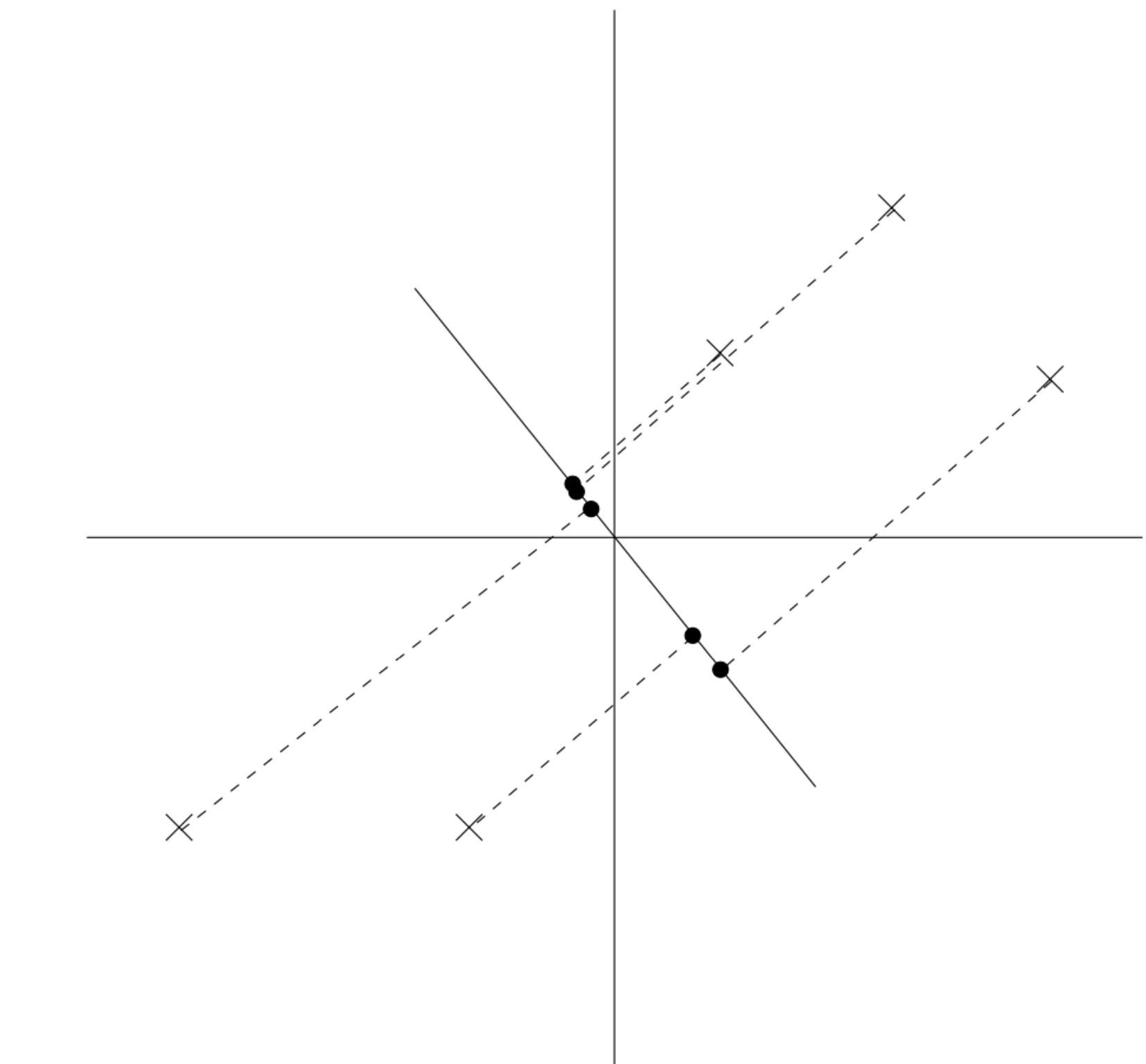
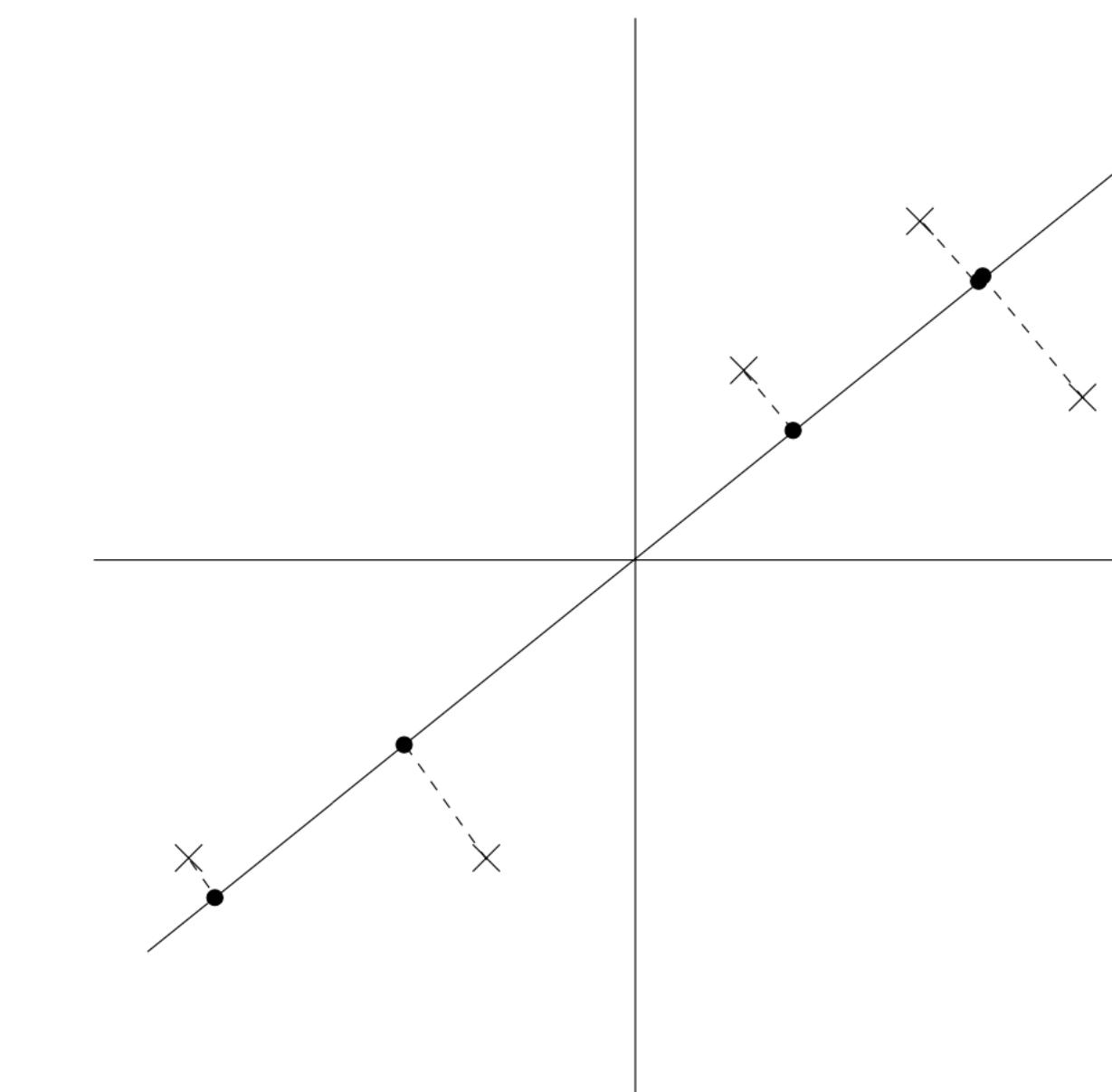
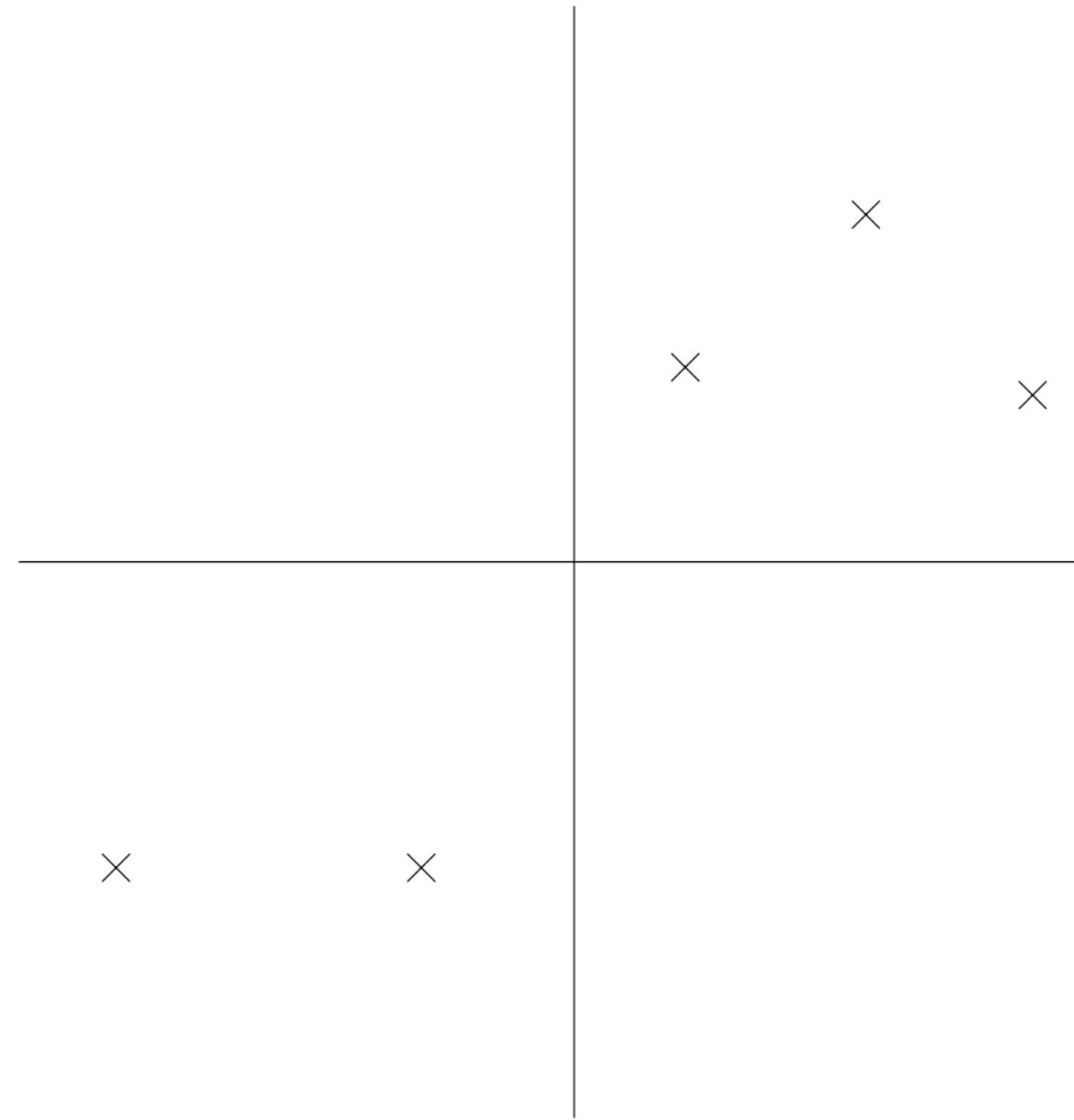


**Assumption:** Data lies on or near a low  $d$ -dimensional linear subspace.

Axes of this subspace are an effective representation of the data

Identifying the axes is known as Principal Components Analysis, and can be obtained by Eigen or Singular value decomposition

# Principal Component Analysis (PCA)



Project the data onto different directions

Which projection is better?

We want the low-dim features that can discriminate the data the most

# Normalizing Data

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}$$
$$\mu_j = \frac{1}{n} \sum_{i=1}^n x_j^{(i)}$$
$$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n (x_j^{(i)} - \mu_j)^2$$

Different features may have different scales

After normalization, each feature has 0 mean and variance 1

# Principal Component Analysis (PCA)

Let  $v$  be the principal component

Find vector that maximizes sample variance of projection

$$\frac{1}{n} \sum_{i=1}^n (v^T x_i)^2 = \frac{1}{n} v^T X X^T v$$

$$\max_v v^T X X^T v \quad \text{s.t.} \quad v^T v = 1$$

Lagrangian:  $\max_v v^T X X^T v - \lambda(v^T v - 1)$

$$\partial/\partial v = 0$$

$$(X X^T - \lambda I)v = 0$$

$$\Rightarrow (X X^T)v = \lambda v$$

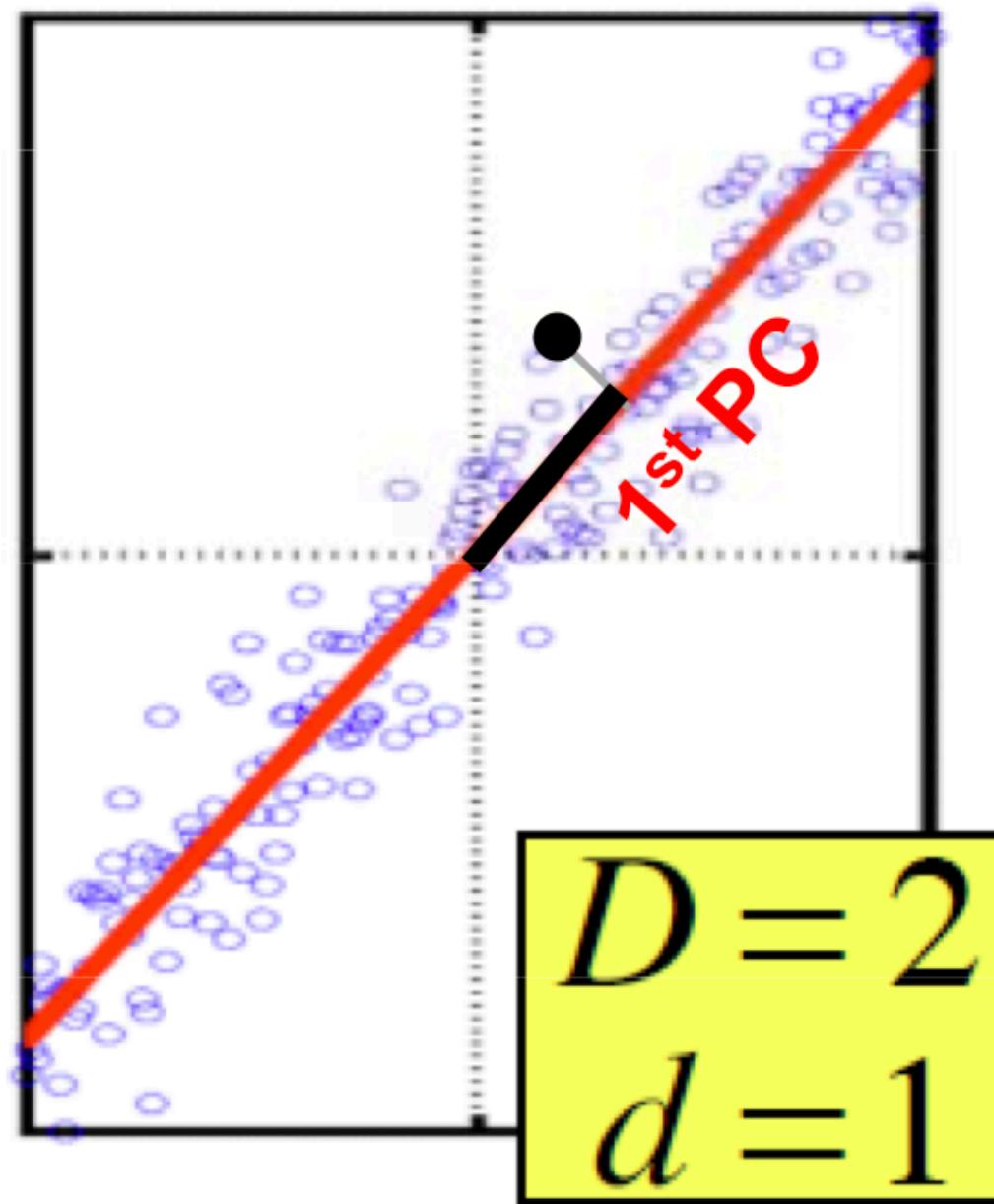
Definition of eigenvectors

## K-dimensional Cases

If we project our data into a k-dimensional subspace ( $k < d$ ), we should choose  $v_1, v_2, \dots, v_k$  to be the top  $k$  eigenvectors of  $XX^T$

For symmetric matrices, eigenvectors for distinct eigenvalues are orthogonal

# Principal Component Analysis (PCA)



Principal Components (PC) are orthogonal directions that capture most of the variance in the data

1<sup>st</sup> PC – direction of greatest variability in data

Projection of data points along 1<sup>st</sup> PC discriminate the data most along any one direction

# Principal Component Analysis (PCA)

Sample variance of projection =  $\mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda$

**Thus, the eigenvalue  $\lambda$  denotes the amount of variability captured along that dimension.**

The 1<sup>st</sup> Principal component  $v_1$  is the eigenvector of the sample covariance matrix  $\mathbf{X} \mathbf{X}^T$  associated with the largest eigenvalue  $\lambda_1$

The 2<sup>nd</sup> Principal component  $v_2$  is the eigenvector of the sample covariance matrix  $\mathbf{X} \mathbf{X}^T$  associated with the second largest eigenvalue  $\lambda_2$

And so on ...

# Computing the Principal Components (PCs)

Eigenvectors are solutions of the following equation:

$$(\mathbf{X}\mathbf{X}^T)\mathbf{v} = \lambda\mathbf{v} \quad (\mathbf{X}\mathbf{X}^T - \lambda\mathbf{I})\mathbf{v} = 0$$

Non-zero solution  $\mathbf{v} \neq 0$  possible only if

$$\det(\mathbf{X}\mathbf{X}^T - \lambda\mathbf{I}) = 0$$

We can compute the eigenvalues from this equation

This is a  $D^{\text{th}}$  order equation in  $\lambda$ , can have at most  $D$  distinct solutions (roots of the characteristic equation)

Once eigenvalues are computed, solve for eigenvectors (Principal Components) using

$$(\mathbf{X}\mathbf{X}^T - \lambda\mathbf{I})\mathbf{v} = 0$$

# Another Interpretation

**Minimum Reconstruction Error:** PCA finds vectors  $v$  such that projection on to the vectors yields minimum MSE reconstruction

$$\frac{1}{n} \sum_{i=1}^n \|x_i - (v^T x_i)v\|^2$$

# Dimensionality Reduction using PCA

The eigenvalue  $\lambda$  denotes the amount of variability captured along that dimension.

Zero eigenvalues indicate no variability along those directions => data lies exactly on a linear subspace

Only keep data projections onto principal components with non-zero eigenvalues, say  $v_1, \dots, v_d$  where  $d = \text{rank}(XX^T)$

## Original Representation

data point

$$x_i = [x_i^1, x_i^2, \dots, x_i^D]$$

(D-dimensional vector)

## Transformed representation projections

$$[v_1^T x_i, v_2^T x_i, \dots, v_d^T x_i]$$

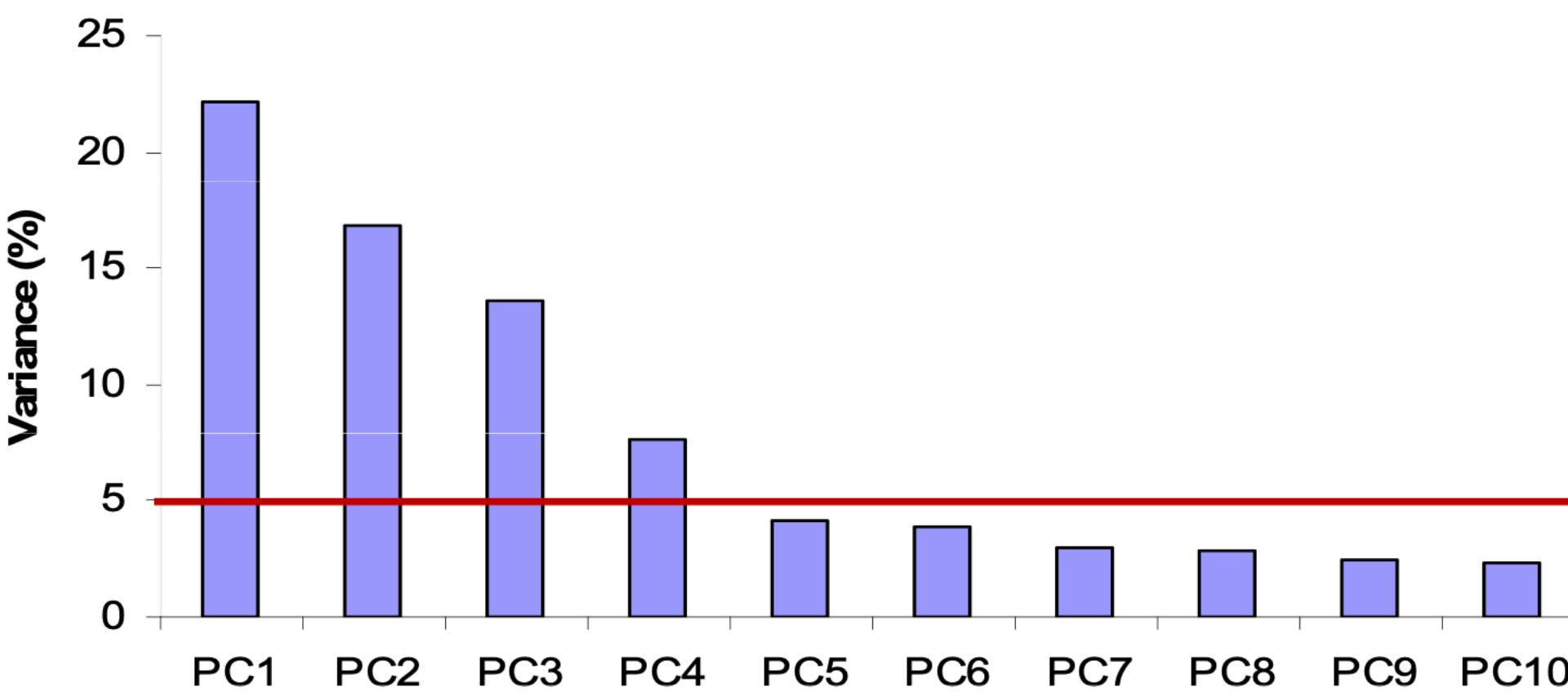
(d-dimensional vector)

# Dimensionality Reduction using PCA

Usually data lies near a linear subspace, as noise introduces small variability

Only keep data projections onto principal components with **large** eigenvalues

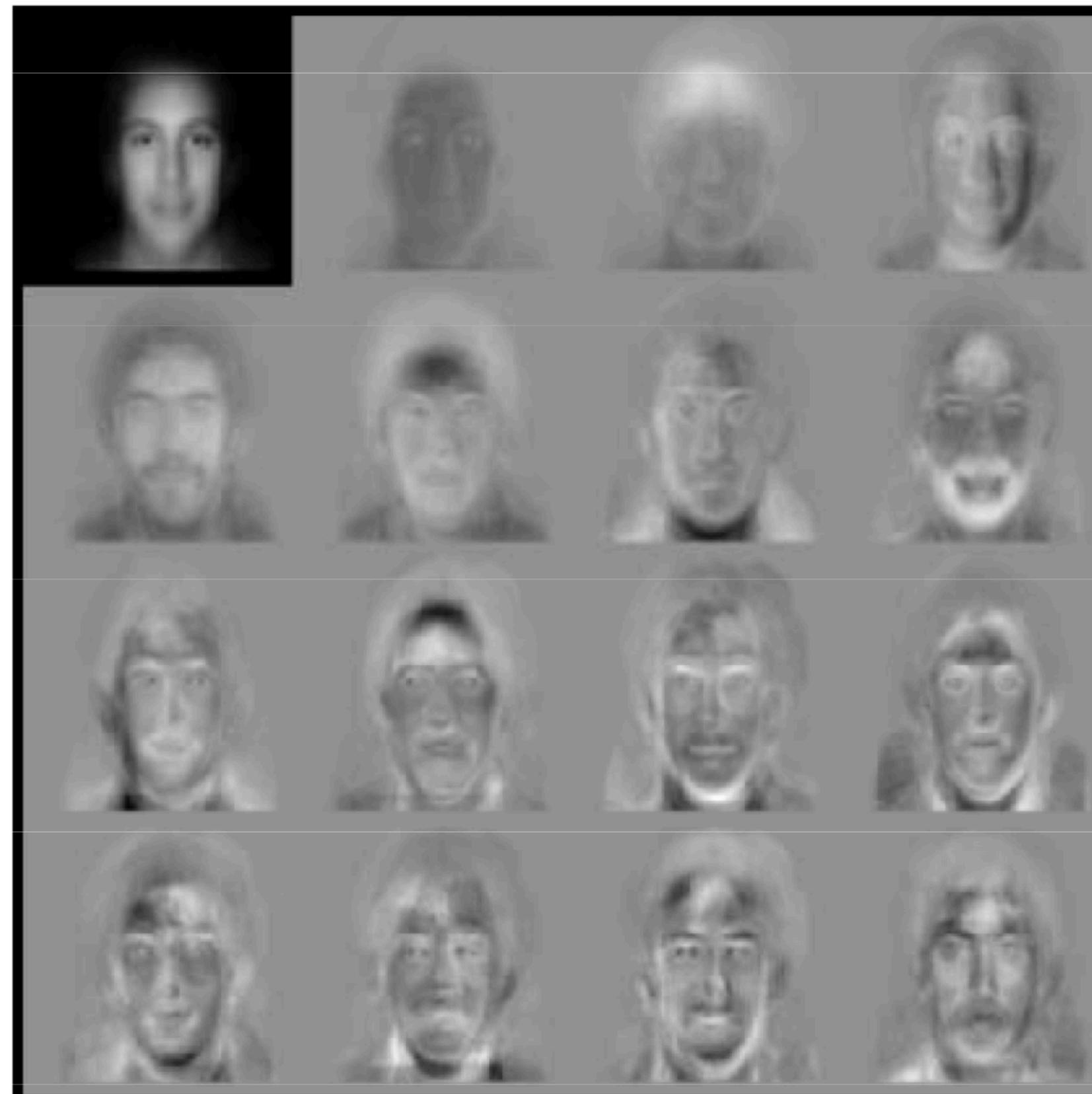
Can *ignore* the components of lesser significance.



You might lose some information, but if the eigenvalues are small, you don't lose much

It is not lossless compression

# Example: faces



**Eigenfaces  
from 7562  
images:**

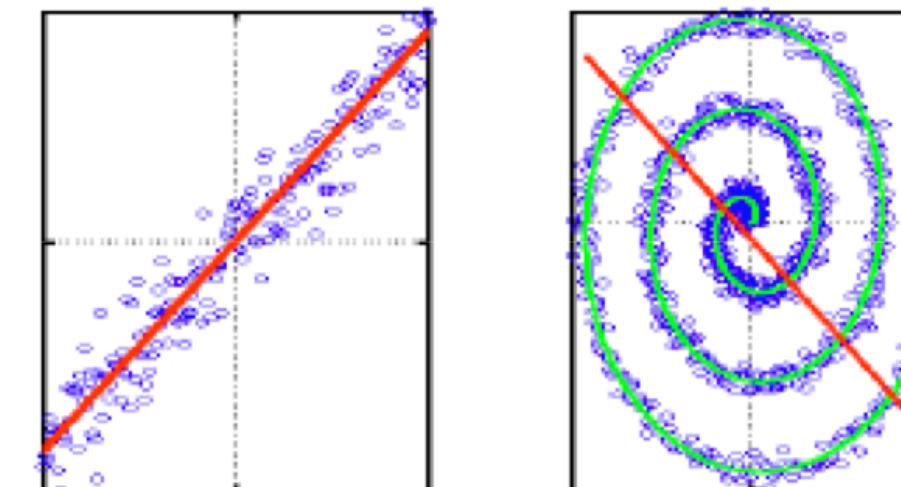
**top left image  
is linear  
combination  
of rest.**

Sirovich & Kirby (1987)  
Turk & Pentland (1991)

# Properties of PCA

- **Strengths**

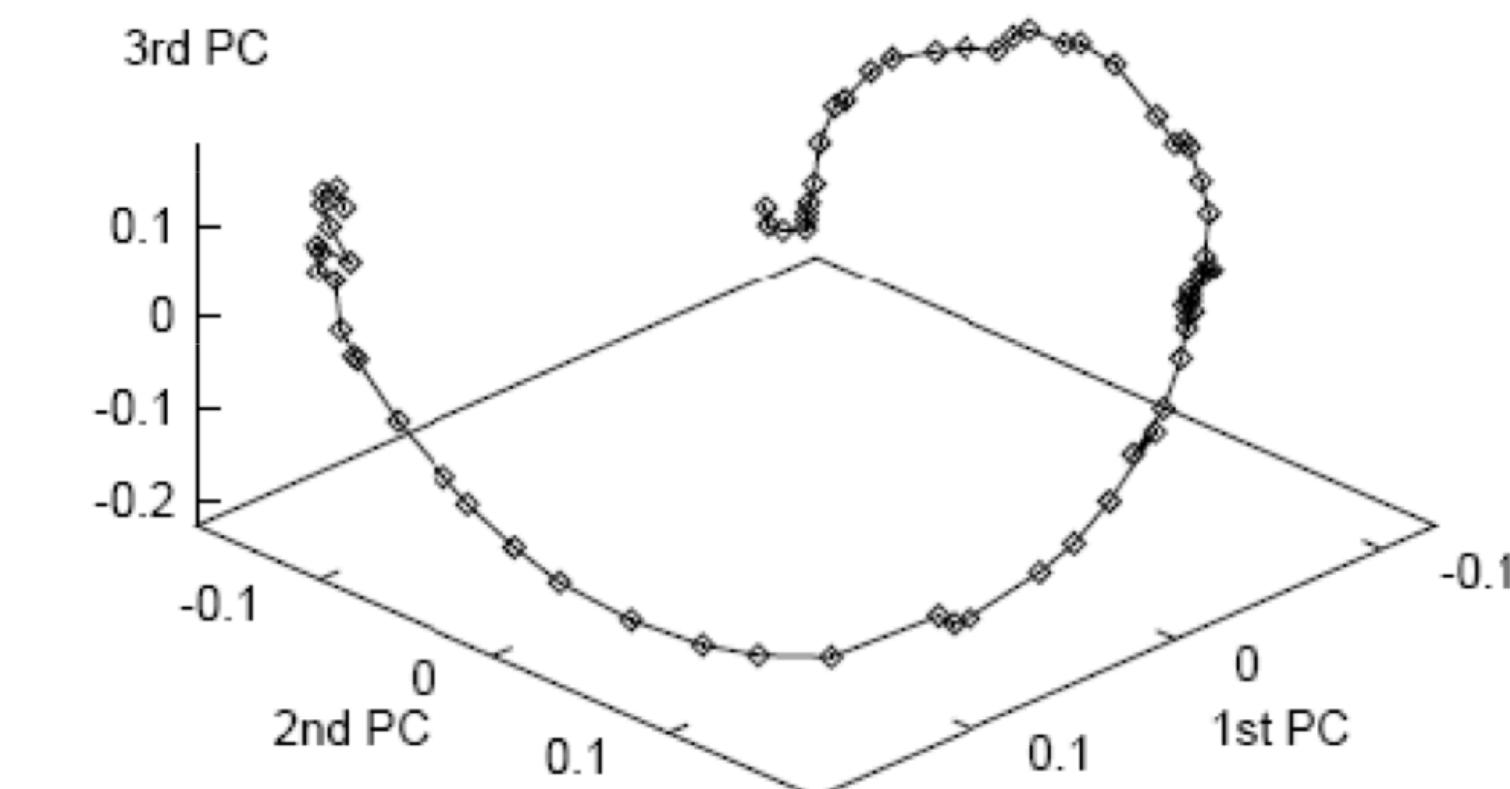
- Eigenvector method
- No tuning parameters
- Non-iterative
- No local optima



- **Weaknesses**

- Limited to second order statistics
- Limited to linear projections

Nonlinear example





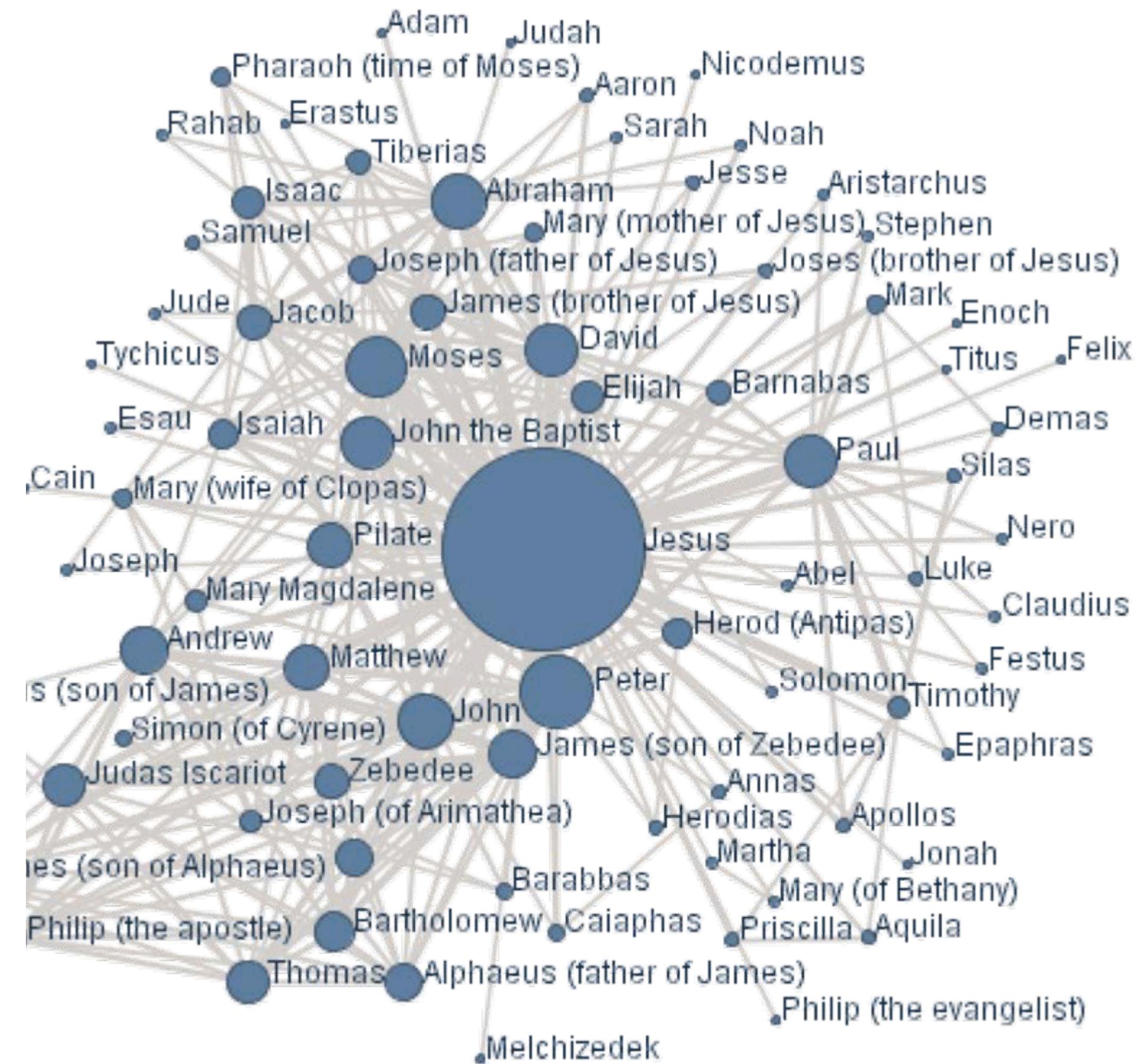
香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 14

# Probabilistic Graphical Models

# What Are Graphical Models?

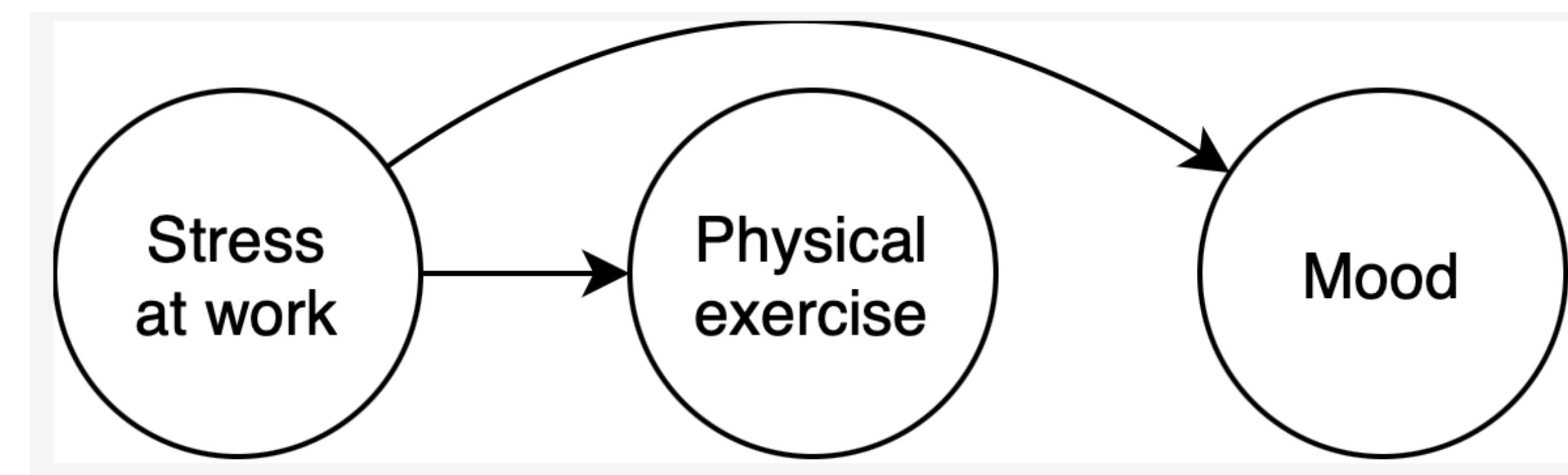
- Informally, a GM is just a graph representing **relationship** among random variables
    - Nodes: random variables (features, not examples)
    - Edges (or absence of edges): relationship
  - Looks simple!
    - But detail matters, as always.
    - What exactly do we mean by **relationship**?



# Relationship between two random variables

- Many types of relationships exist:
  - X and Y are correlated
  - X and Y are dependent
  - X and Y are independent
  - X and Y are partially correlated given Z
  - X and Y are conditionally dependent given Z
  - X and Y are conditionally independent given Z
  - X causes Y
  - Y causes X
- ...

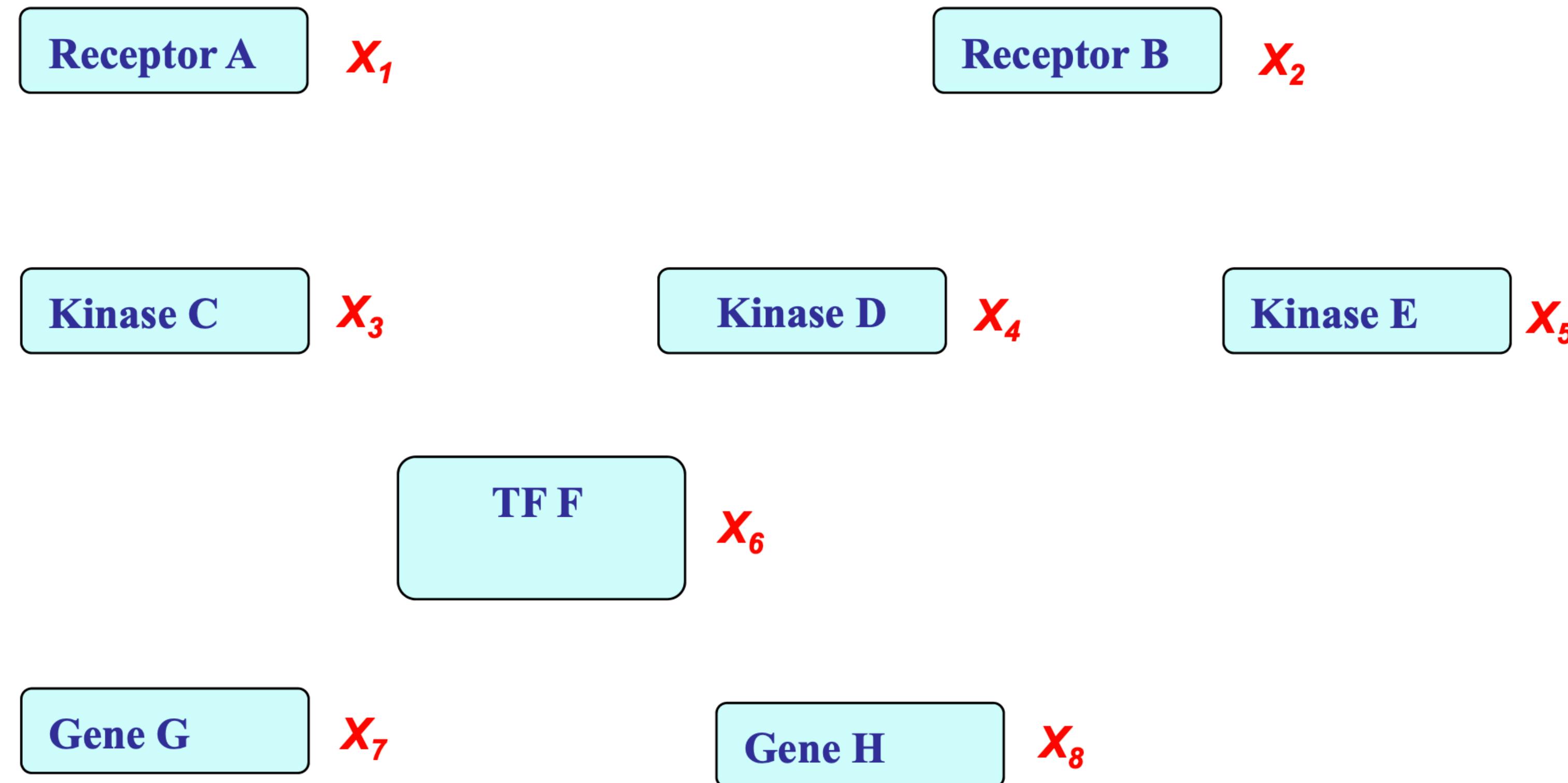
Correlation does not imply causation



# What is a Graphical Model?

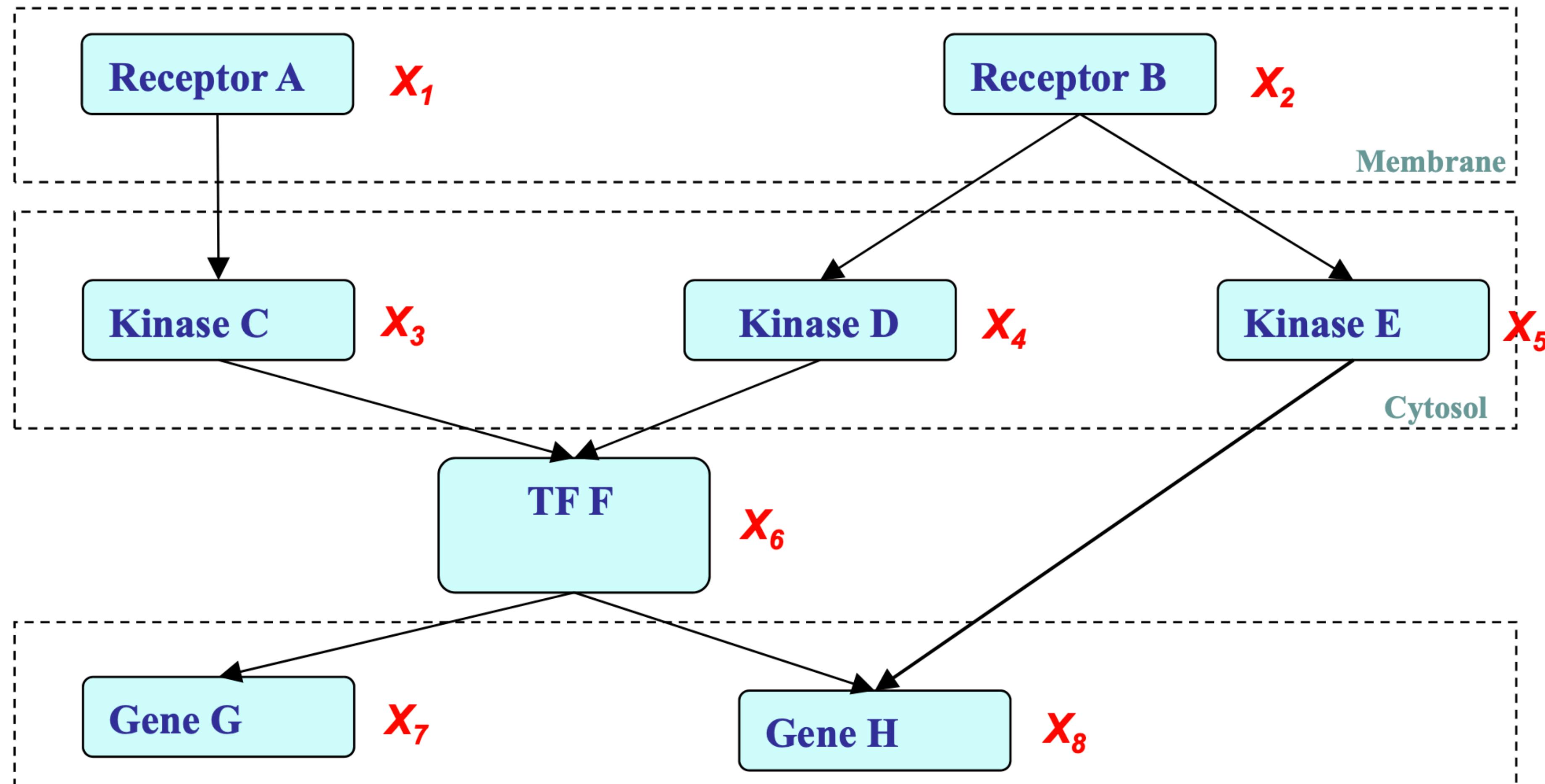
Graphical model represents a multivariate distribution in High-D space

A possible world for cellular signal transduction:



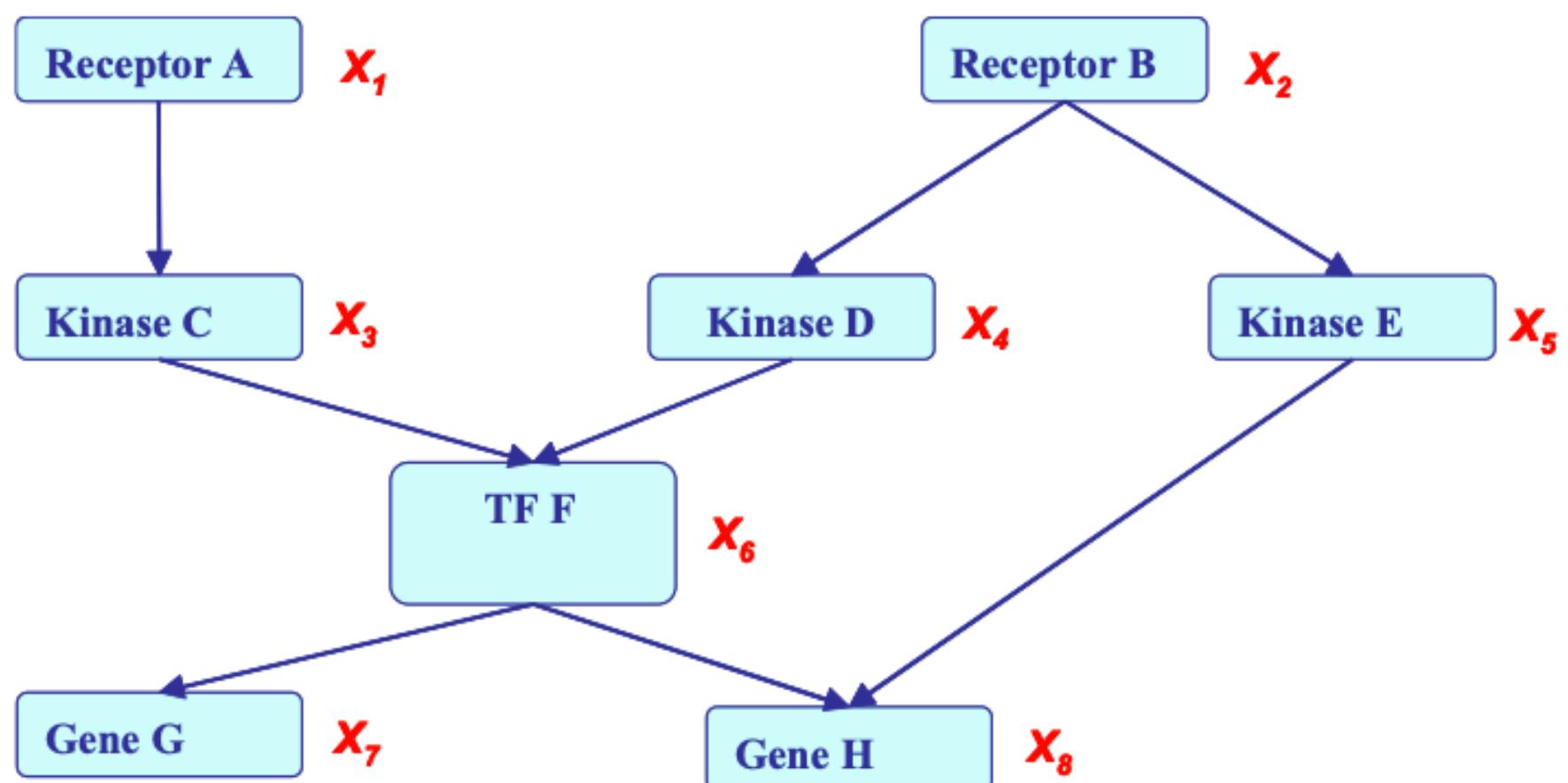
# Structure Simplifies Representation

Dependencies among variables



# Probabilistic Graphical Models

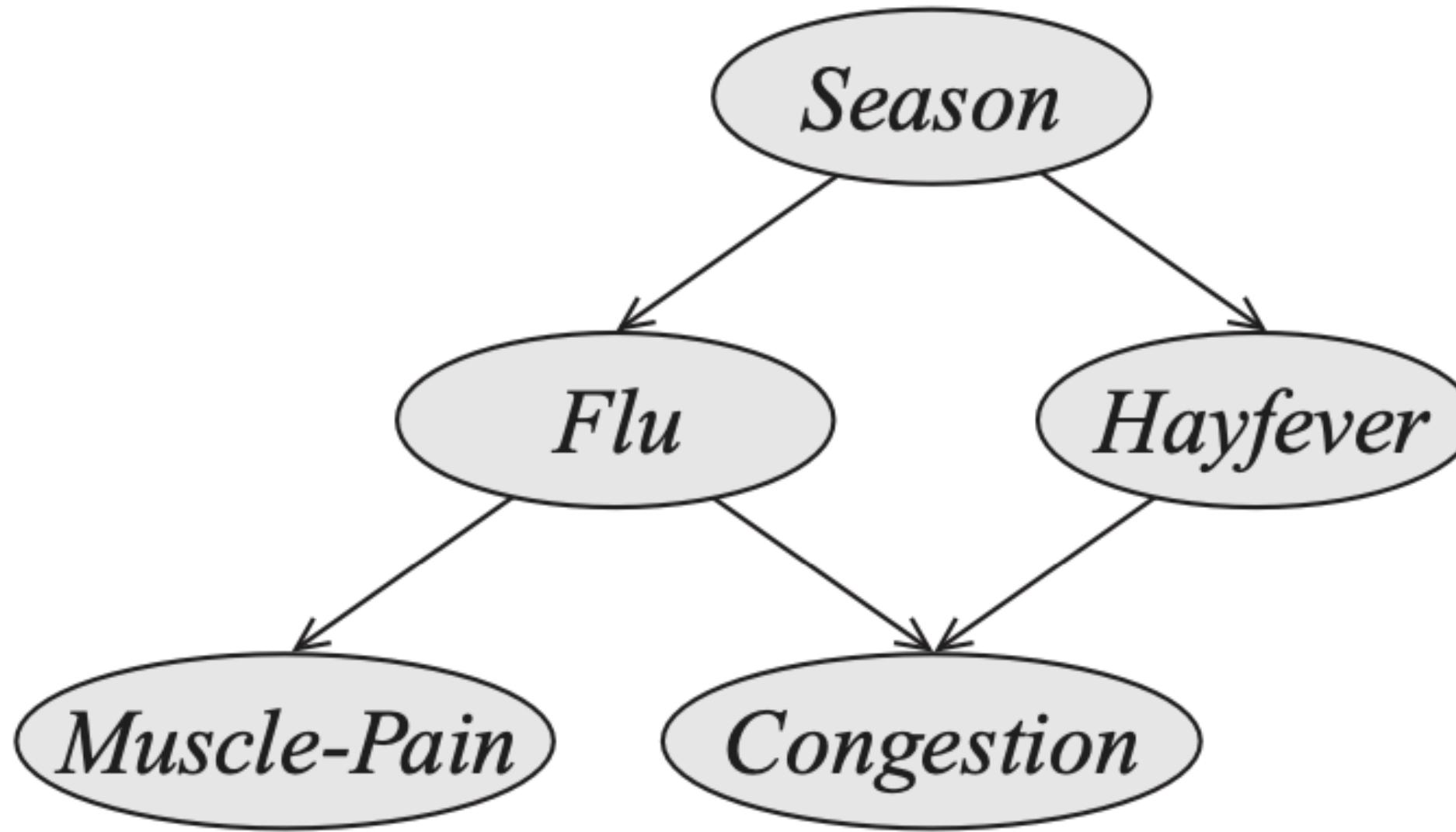
- If  $X_i$ 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{aligned} & P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ = & P(X_1) P(X_2) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_2) \\ & P(X_6 | X_3, X_4) P(X_7 | X_6) P(X_8 | X_5, X_6) \end{aligned}$$

Stay tune for what are these independencies!

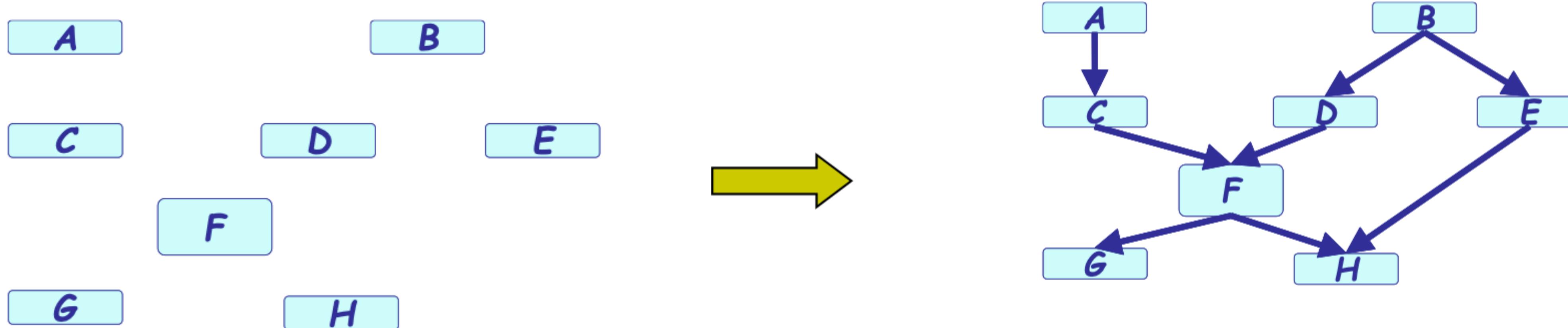
# Another Example



$$P(\text{Congestion} \mid \text{Flu}, \text{Hayfever}, \text{Season}) = P(\text{Congestion} \mid \text{Flu}, \text{Hayfever});$$

# What is a PGM After All

It is a smart way to write/specify/compose/design exponentially-large probability distributions without paying an exponential cost, and at the same time endow the distributions with **structured semantics**



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$P(X_{1:8}) = P(X_1)P(X_2)P(X_3 | X_1X_2)P(X_4 | X_2)P(X_5 | X_2)$$

$$P(X_6 | X_3, X_4)P(X_7 | X_6)P(X_8 | X_5, X_6)$$

More formal definition:

It refers to a family of distributions on a set of random variables that are compatible with all the probabilistic independence propositions encoded by a graph that connects these variables

**Probabilistic Graphical Model is a  
graphical language to express  
conditional independence**

**Thank You!**  
**Q & A**