# Topics in Graphical Models

$$P(X_1, \ldots X_n) = \prod_{i=1}^{n} P(X_i \mid pa(X_i))$$

- ## Representation
  - Which joint probability distributions does a graphical model represent?

    conditional independence
    - local Markov assump$^n$
    - D-separation
    - Markov Blanket

- ## Inference
  - How to answer questions about the joint probability distribution?
    - Marginal distribution of a node variable
    - Most likely assignment of node variables

    $P(X_1, X_3)$ or $P(X_2 \mid X_4)$

    Variable elimination

    $\max_{X_1, X_3} P(X_1, X_3)$

- ## Learning
  - How to learn the parameters and $\boxed{\text{structure}}$ of a graphical model?

    $D = \{ X_1^{(i)} \ldots X_n^{(i)} \}_{i=1}^{m}$ — m data points

# Max Likelihood score for graph structure

$D = \{ X_1^{(j)} \dots X_n^{(j)} \}_{j=1}^m$  m - data points  n - # variables

ML score for graph structure $\mathcal{G}$

$\hat{\theta}_{\mathcal{G}}$ — MLE estimates of parameters given G   or MAP

$$\arg\max_{\mathcal{G}} \log \hat{P}(\mathcal{D} \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = \arg\max_{\mathcal{G}} \sum_{i=1}^n \hat{I}(X_i, \mathbf{Pa}_{X_i})$$

$\hat{P}(D \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = \prod_{j=1}^m \hat{P}(X_1 \dots X_n) = \prod_{i=1}^n \hat{P}(X_i^{(j)} \mid pa(X_i))$

$$\log \hat{P}(\mathcal{D} \mid \hat{\theta}_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log \hat{P}\left( x_i^{(j)} \mid \mathbf{x}_{\mathbf{Pa}_{X_i}}^{(j)} \right)$$

$$= m \sum_{i=1}^n \sum_{x_i} \sum_{\mathbf{x}_{\mathbf{Pa}_{X_i}}} \hat{P}(x_i, \mathbf{x}_{\mathbf{Pa}_{X_i}}) \log \hat{P}\left( x_i \mid \mathbf{x}_{\mathbf{Pa}_{X_i}} \right)$$

$$= -m \sum_{i=1}^n \hat{H}(X_i \mid \mathbf{Pa}_{X_i})$$

$H(z) = -\sum_z p(z) \log p(z)$

$$= m \sum_{i=1}^n [\hat{I}(X_i, \mathbf{Pa}_{X_i}) - \hat{H}(X_i)]$$

Doesn't depend on graph structure $\mathcal{G}$

# ML score is Decomposable

- Log data likelihood

$$\log \widehat{P}(\mathcal{D} \mid \widehat{\theta}_{\mathcal{G}}, \mathcal{G}) = m \sum_{i=1}^{n} \left[ \widehat{I}(X_i, \mathbf{Pa}_{X_i}) - \widehat{H}(X_i) \right]$$
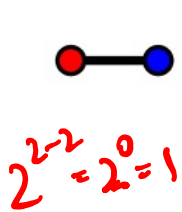
- Decomposable score:

  – Decomposes over families in BN (node and its parents)

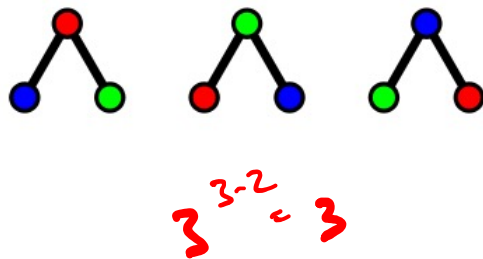  – Will lead to significant computational efficiency!!!

# How many trees are there?

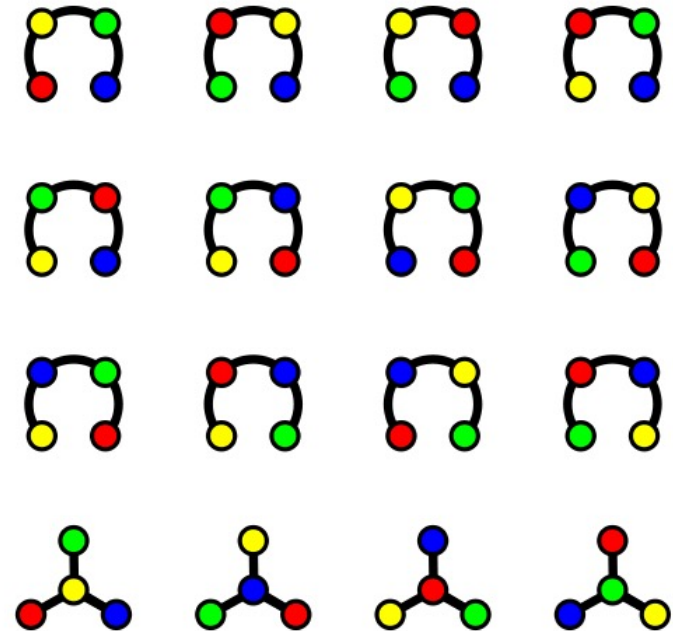- Trees – every node has at most one parent
- $n^{n-2}$ possible trees (Cayley's Theorem)
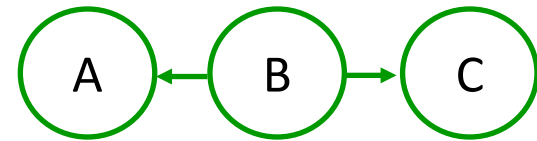
$n=4$

$n=2$

$n=3$

$2^{2-2} = 2^0 = 1$

$3^{3-2} = 3$

$4^{4-2} = 16$
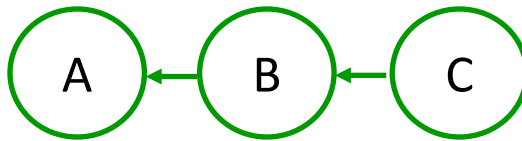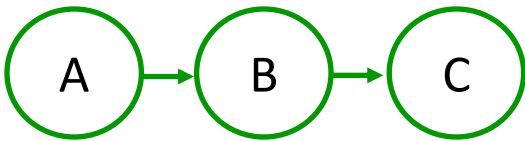
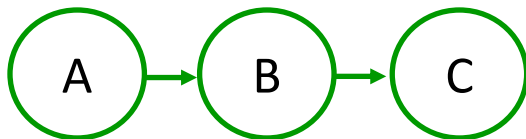Nonetheless – Efficient optimal algorithm finds best tree!

# Scoring a tree

$$\arg\max_{\mathcal{G}} \log \widehat{P}(\mathcal{D} \mid \widehat{\theta}_{\mathcal{G}}, \mathcal{G}) = \arg\max_{\mathcal{G}} \sum_{i=1}^{n} \widehat{I}(X_i, \mathbf{Pa}_{X_i}) \quad \checkmark$$
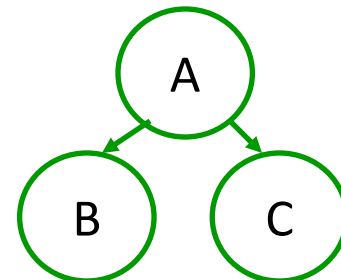
Equivalent Trees (same score):   I(A,B) + I(B,C)



Score provides indication of structure:



I(A,B) + I(B,C)        I(A,B) + I(A,C) ✔

# Chow-Liu algorithm

$I(X_i, pa(X_i))$

- For each pair of variables $X_i, X_j$
  - Compute empirical distribution: $\hat{P}(x_i, x_j) = \dfrac{\text{Count}(x_i, x_j)}{m}$ ✓
  - Compute mutual information:

$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$
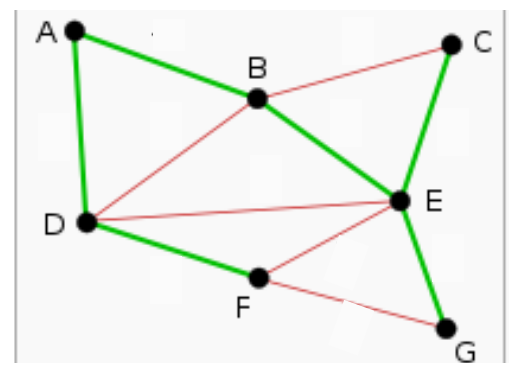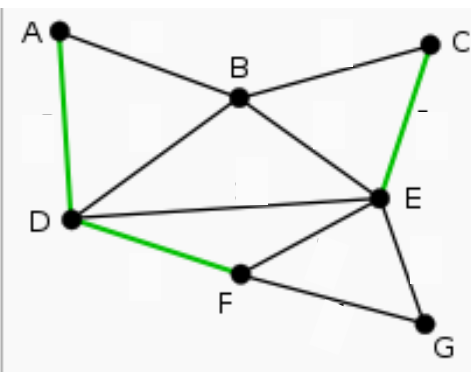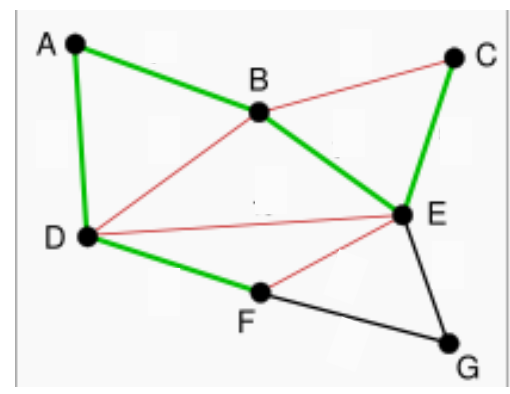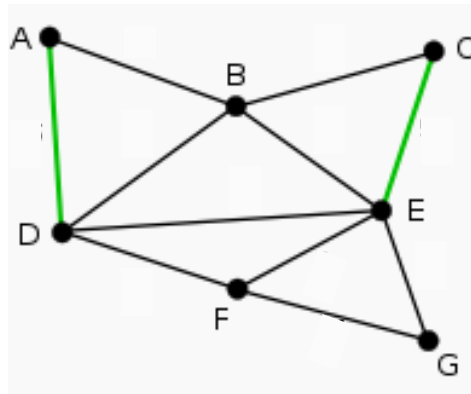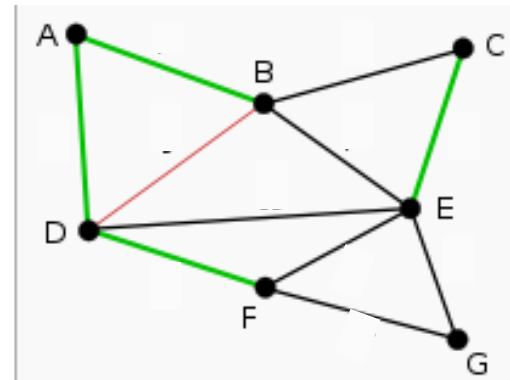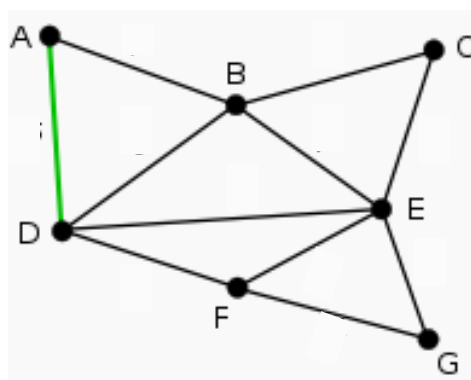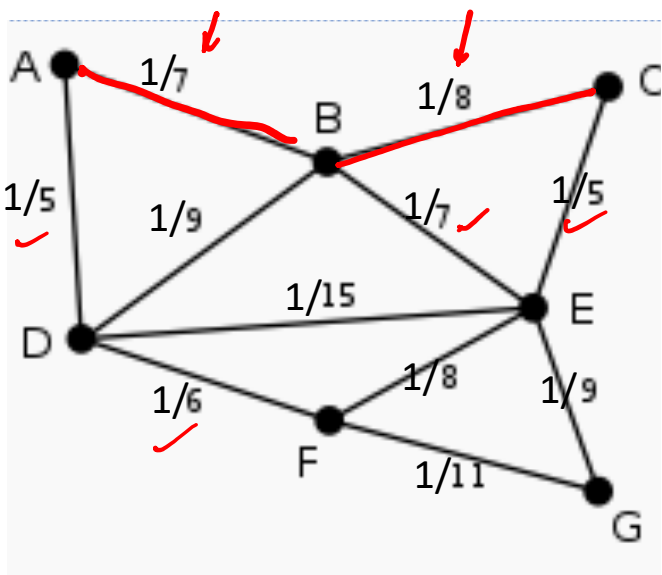
- Define a graph
  - Nodes $X_1, \ldots, X_n$
  - Edge (i,j) gets weight $\hat{I}(X_i, X_j)$

- Optimal tree BN
  - Compute maximum weight spanning tree (e.g. Prim's, Kruskal's algorithm O(nlog n)) ✓
  - Directions in BN: pick any node as root, breadth-first-search defines directions

# Chow-Liu algorithm example

$\{A^{(j)} \ldots G^{(j)}\}_{j=1}^{m}$

# Scoring general graphical models

- Graph that maximizes ML score -> complete graph!
- Information never hurts

  $H(A|B) \geq H(A|B,C)$ ✓

- Adding a parent always increases ML score

  $I(A,B,C) \geq I(A,B)$ ✓

- The more edges, the fewer independence assumptions, the higher the likelihood of the data, but will overfit…

- Why does ML for trees work?

  Restricted model space – tree graph

# Regularizing

- Model selection
  - Use MDL (Minimum description length) score
  - BIC score (Bayesian Information criterion)
- Still NP –hard

  **Theorem**: The problem of learning a BN structure with at most $d$ parents is NP-hard for any (fixed) $d>1$ (Note: tree d=1)

- Mostly heuristic (exploit score decomposition)
- Chow-Liu: provides best tree approximation to any distribution.
- Start with Chow-Liu tree. Add, delete, invert edges. Evaluate BIC score

# What you should know

- Learning BNs ← *directed graphical model*
  - Maximum likelihood or MAP learns parameters
  - ML score
    - Decomposable score
    - Information theoretic interpretation (Mutual information)
  - Best tree (Chow-Liu)
  - Other BNs, usually local search with BIC score
    *regularized ML score*

# Unsupervised Learning

Aka Learning without labels

$$y$$

$$P(X_1 \cdots X_n)$$

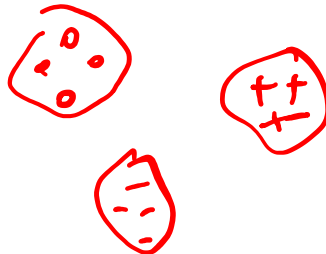➢ Learning and inference using probability distributions & densities

    MLE/MAP ✓

    Graphical models ✓

➢ Dimensionality Reduction

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \longrightarrow \tilde{X} = \begin{bmatrix} \tilde{x}_1 \\ \vdots \\ \tilde{x}_d \end{bmatrix}$$

$$d \ll D$$

➢ Clustering

# Dimensionality Reduction
# PCA

Aarti Singh

Machine Learning 10-701
April 19, 2023

Slides Courtesy: Tom Mitchell, Eric Xing, Lawrence Saul

# High-Dimensional data

- High-Dimensions = Lot of Features

Document classification

Features per document =

   thousands of words/unigrams

   millions of bigrams, contextual

   information

Surveys - Netflix

  480189 users x 17770 movies

| | movie 1 | movie 2 | movie 3 | movie 4 | movie 5 | movie 6 |
|---|---|---|---|---|---|---|
| Tom | 5 | ? | ? | 1 | 3 | ? |
| George | ? | ? | 3 | 1 | 2 | 5 |
| Susan | 4 | 3 | 1 | ? | 5 | 1 |
| Beth | 4 | 3 | ? | 2 | 4 | 2 |

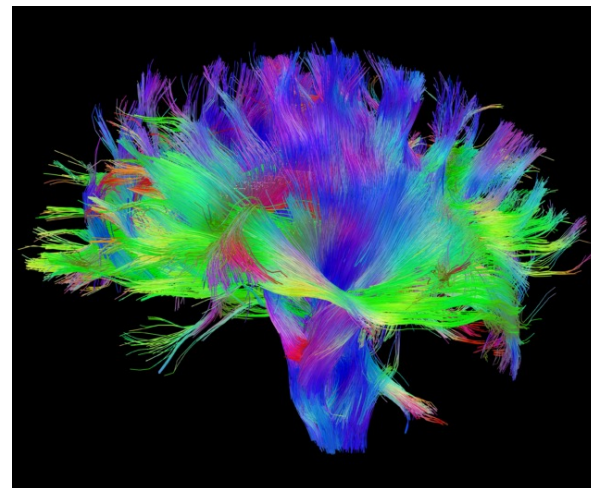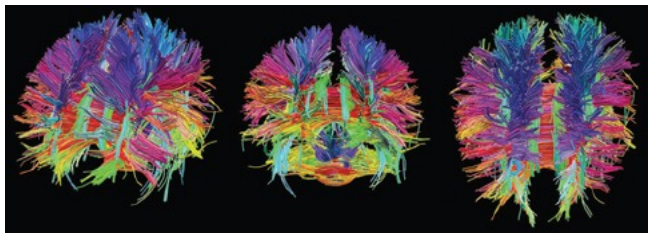3

# High-Dimensional data

- High-Dimensions = Lot of Features

High resolution images
  millions of pixels

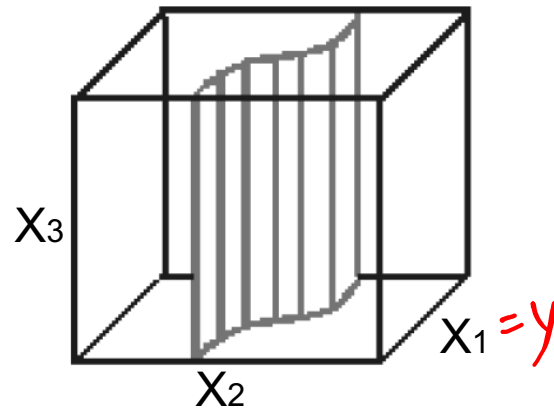Diffusion scans of Brain
  300,000 brain fibers

# 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 interpret and visualize

  – Hard to store and process data (computationally challenging)

  – Complexity of decision rule tends to grow with # features. Hard to learn complex rules as it needs more data (statistically challenging)

# Dimensionality Reduction

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



$X_3$ - Irrelevant

$l_1$ penalty
'''
lasso

$y \leftarrow w_2 x_2 + w_3 x_3$

$X_3$

$X_2$

$X_1 = y$

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



$X_3$

$X_2$

$X_1$

$X_3$

$X_2$

$X_1$

6

# Latent Features

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

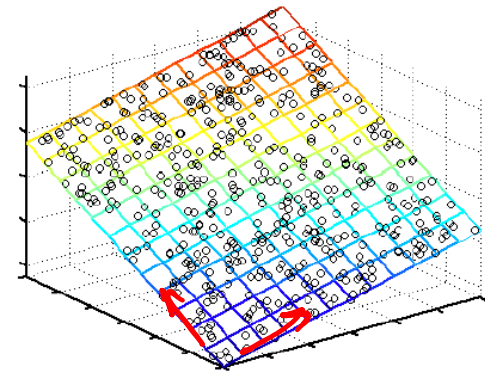           Topics (sports, science, news, etc.) instead of documents

Often may not have physical meaning

- Linear

    **Principal Component Analysis (PCA)** ✓

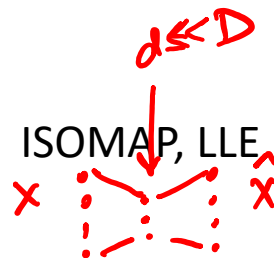    Factor Analysis

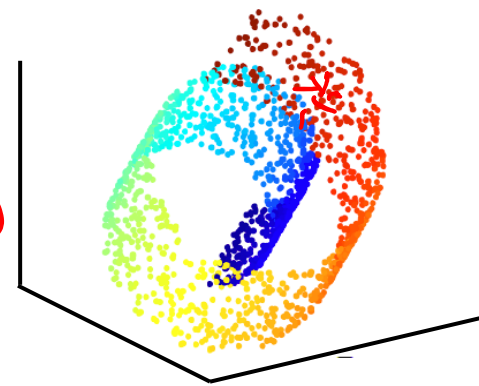    Independent Component Analysis (ICA)

- Nonlinear

    Kernel PCA

    Laplacian Eigenmaps, ISOMAP, LLE

    Autoencoders

$d \ll D$

$MSE(x \hat{x})$

$x \quad \hat{x}$

# Principal Component Analysis (PCA)

$D = 2$
$d = 1$

$D = 3$
$d = 2$

When 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

# Data for PCA

Data $X = [x_1, x_2, \ldots, x_n]$ where each data point $x_i$ is D-dimensional vector

X is D x n matrix

Assume data are centered i.e. sample mean $\dfrac{1}{n}\displaystyle\sum_{i=1}^{n} x_i = 0$

What if data is not centered?

①     Subtract off sample mean from each data point ✓

Since data matrix is centered, sample covariance matrix can be written as

②     $$S = \frac{1}{n}XX^\top$$

$\mathbb{E}\left[(z_i - \mathbb{E}z)(z_i - \mathbb{E}z)^\top\right]$

$\dfrac{1}{n}\sum_{i=1}^{n} \quad 0 \quad 0$

$\dfrac{1}{n}\sum_{i=1}^{n} z_i$

# Principal Component Analysis (PCA)



$$D = 2$$
$$d = 1$$

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

$1^{st}$ PC – direction of greatest variability in data

Projection of data points along $1^{st}$ PC discriminate the data most along any one direction

Take a data point $x_i$ (D-dimensional vector)

Projection of $x_i$ onto the $1^{st}$ PC $v$ is $v^T x_i$

# Principal Component Analysis (PCA)



$$D = 2$$
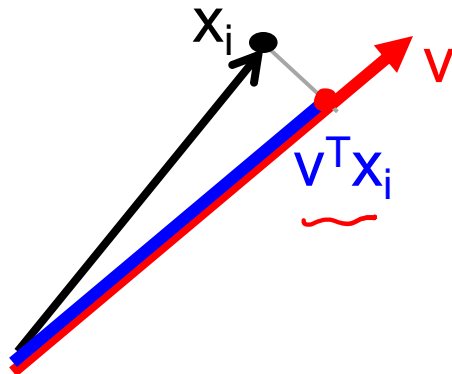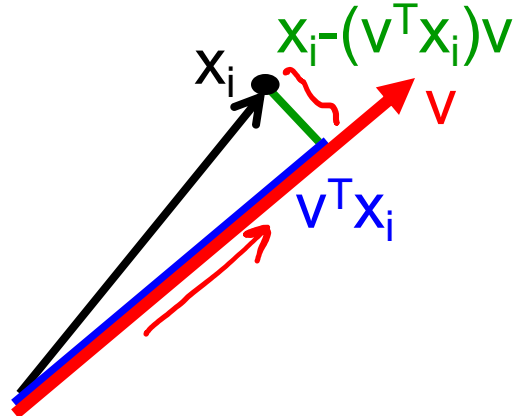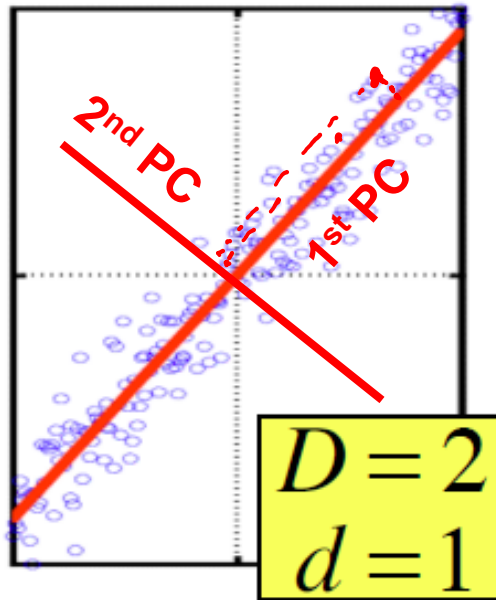$$d = 1$$

$$x_i - (v^T x_i)v$$
$$v$$
$$x_i$$
$$v^T x_i$$

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

1st PC – direction of greatest variability in data

2nd PC – Next orthogonal (uncorrelated) direction of greatest variability

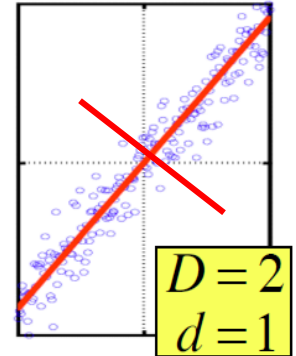(remove all variability in first direction, then find next direction of greatest variability)

And so on …

# Principal Component Analysis (PCA)

Let $v_1, v_2, \ldots, v_d$ denote the principal components

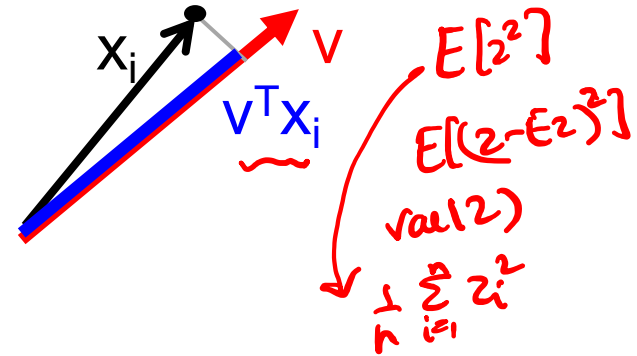Orthogonal and unit norm $\quad v_i^T v_j = 0 \quad i \neq j$ ✓

$\quad v_i^T v_i = 1$ ✓

Find vector that maximizes sample variance of projection

$D = 2$
$d = 1$

$$\frac{1}{n}\sum_{i=1}^{n}(\mathbf{v}^T\mathbf{x}_i)^2 = \frac{\mathbf{v}^T \overset{1 \times d}{\mathbf{X}}\overset{d \times d}{\mathbf{X}^T}\overset{d \times 1}{\mathbf{v}}}{n}$$

$x_i$ $\quad$ v $\quad$ $v^T x_i$

$E[z^2]$

$E[(z - Ez)^2]$

$var(z)$

$\frac{1}{n}\sum_{i=1}^{n} z_i^2$

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

$\sum v_i^2 = 1$

Poll:

Convex $\qquad$ Non-convex set

$z^T M z \geq 0$

$\|v\| \leq 1$

$v^T v \leq 1$

$v_2$ $\quad$ $v_1$

➤ Is this a convex optimization problem?

# Principal Component Analysis (PCA)
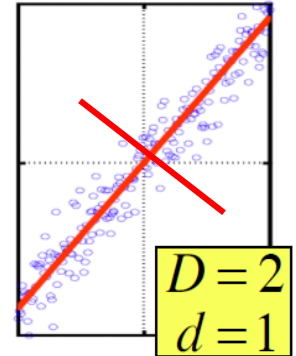
Let $v_1, v_2, \ldots, v_d$ denote the principal components

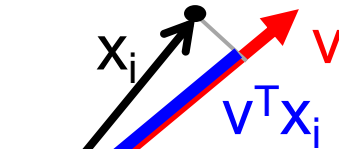Orthogonal and unit norm $\quad v_i^T v_j = 0 \quad i \neq j$

$$v_i^T v_i = 1$$



$D = 2$
$d = 1$

Find vector that maximizes sample variance of projection

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



$x_i$ $\quad$ $v$

$v^T x_i$

$$\max_{\mathbf{v}} \ \mathbf{v}^T\mathbf{X}\mathbf{X}^T\mathbf{v} \quad \text{s.t.} \quad \mathbf{v}^T\mathbf{v} = 1 \qquad \lambda$$

Lagrangian: $\max_{\mathbf{v}} \mathbf{v}^T\mathbf{X}\mathbf{X}^T\mathbf{v} - \lambda\mathbf{v}^T\mathbf{v}$

Wrap constraints into the objective function

$evec(XX^T)$

$$2XX^Tv - 2\lambda v = 0$$

$$\partial/\partial\mathbf{v} = 0 \qquad (\mathbf{X}\mathbf{X}^T - \lambda\mathbf{I})\mathbf{v} = 0 \qquad \Rightarrow \boxed{(\mathbf{X}\mathbf{X}^T)\mathbf{v} = \lambda\mathbf{v}}$$

Sample var when projecting on v

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

$eval(XX^T)$

13