# CS 540 Introduction to Artificial Intelligence
## **Unsupervised Learning II**

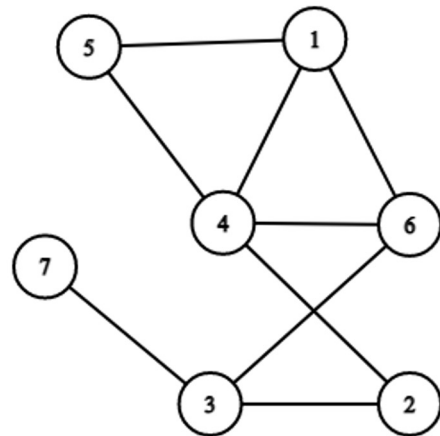University of Wisconsin-Madison
Fall 2023

# Unsupervised Learning II Outline

- Finish up Other Clustering Types
  - Graph-based clustering, graph cuts, spectral clustering
- Unsupervised Learning: Visualization
  - t-SNE: algorithm, examples, vs. PCA
- Unsupervised Learning: Density Estimation
  - Kernel density estimation: high-level intro

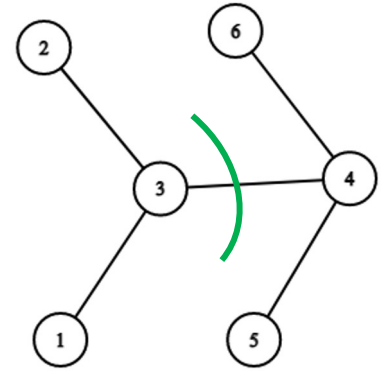# Other Types of Clustering

**Graph**-based/proximity-based

- Recall: Graph G = (V,E) has vertex set V, edge set E.
  - Edges can be weighted or unweighted
  - Edges encode **similarity** between vertices**:**
$$w_{ij} = \text{sim}(v_i, v_j)$$

- Don't need to KEEP vectors for each v.
  - Only keep the edges (possibly weighted)

# Graph-Based Clustering

**Want:** partition V into $V_1$ and $V_2$

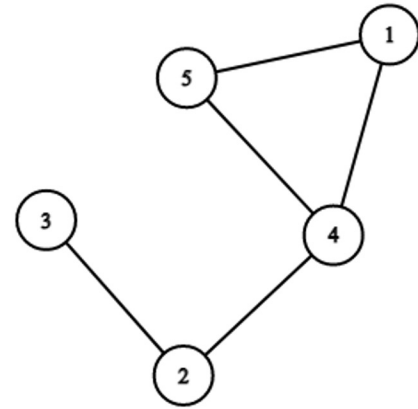- Implies a graph "cut"
- One idea: minimize the **weight** of the cut



$$W(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

$$\text{cut}(A_1, \ldots, A_k) := \frac{1}{2} \sum_{i=1}^{k} W(A_i, \overline{A_i}).$$
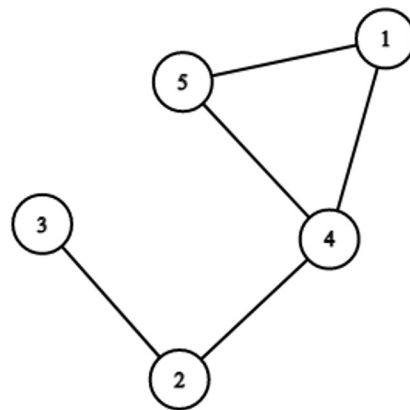
# Graph-Based Clustering

**How do we compute these?**

- Hard problem → heuristics
  - Greedy algorithm
  - "Spectral" approaches

- Spectral clustering approach:
  - **Adjacency** matrix $A_{ij} = w_{ij}$

$$A = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

# Spectral Clustering

- Spectral clustering approach:
  - **Adjacency** matrix
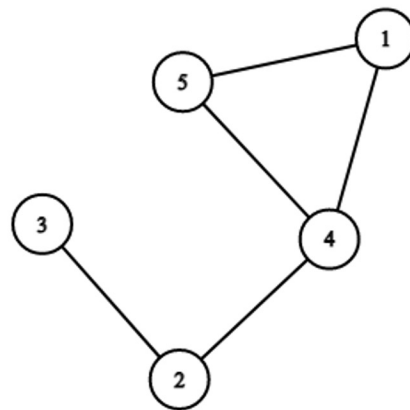  - **Degree** matrix $D_{ii} = \sum_{j=1}^{n} A_{ij}$

$$D = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix} \quad A = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$
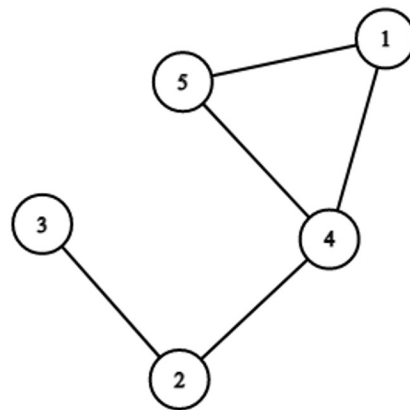
# Spectral Clustering

- Spectral clustering approach:
  - 1. Compute **Laplacian L** = **D** − **A**
  
  (Important tool in graph theory)



$$L = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix} - \begin{bmatrix} 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 & -1 & -1 \\ 0 & 2 & -1 & -1 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ -1 & -1 & 0 & 3 & -1 \\ -1 & 0 & 0 & -1 & 2 \end{bmatrix}$$

**Degree Matrix**  **Adjacency Matrix**  **Laplacian**

# Spectral Clustering
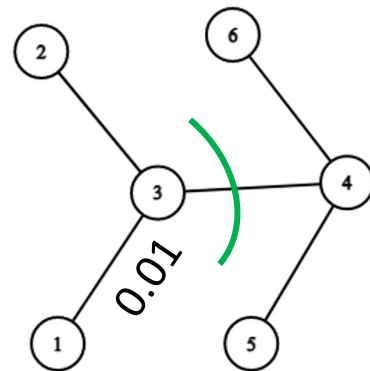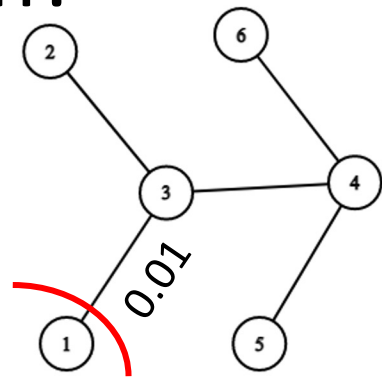
- Spectral clustering approach:
  - 1. Compute **Laplacian L = D − A**
  - 1a (optional): compute normalized Laplacian:
    $$L = I − D^{-1/2}AD^{-1/2}, \quad or \quad L = I − D^{-1}A$$
  - 2. Compute $j$ **smallest** eigenvectors of **L**
  - 3. Set $U$ to be the $n$ x $j$ matrix with $u_1, \ldots, u_j$ as columns. Take the $n$ rows formed as points.
  - 4. Run k-means on the representations.

# Why normalized Laplacian?

**Want:** partition V into $V_1$ and $V_2$

- Implies a graph "cut"
- One idea: minimize the **weight** of the cut
  - Downside: might only get cut of one node
  - Need: "**balanced**" cut

# Why Normalized Laplacian?

**Want:** partition V into $V_1$ and $V_2$

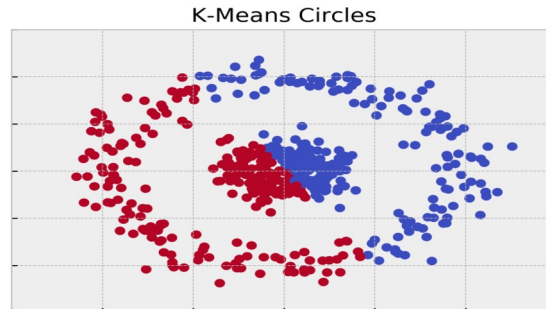- Just minimizing weight is not always a good idea.
- We want **balance!**

$$\text{Ncut}(A_1, \ldots, A_k) := \frac{1}{2} \sum_{i=1}^{k} \frac{W(A_i, \overline{A_i})}{\text{vol}(A_i)}$$

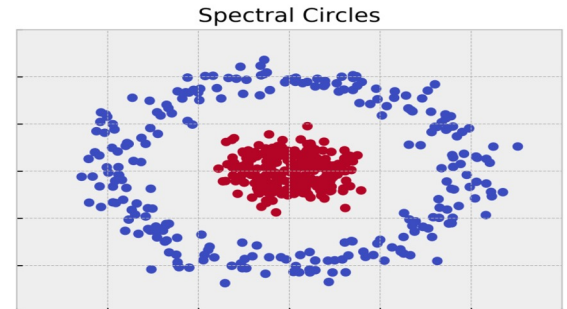$$\text{vol}(A) = \sum_{i \in A} \text{degree}(i)$$

# Spectral Clustering

**Q**: Why do this?

- 1. graph induces an "effective resistance distance" , similar to shortest path distance but also considers how many paths there are
- 2. Can handle intuitive separation (Euclidean dist can't!)



Credit: William

# Break & Quiz

**Q 1.1**: We have two datasets: a social network dataset $S_1$ which shows which individuals are friends with each other along with image dataset $S_2$.

What kind of clustering can we do? Assume we do not make additional data transformations.

- A. k-means on both $S_1$ and $S_2$
- B. graph-based on $S_1$ and k-means on $S_2$
- C. k-means on $S_1$ and graph-based on $S_2$
- D. hierarchical on $S_1$ and graph-based on $S_2$

# Break & Quiz

**Q 1.1**: We have two datasets: a social network dataset $S_1$ which shows which individuals are friends with each other along with image dataset $S_2$.

What kind of clustering can we do? Assume we do not make additional data transformations.

- A. k-means on both $S_1$ and $S_2$
- **B. graph-based on $S_1$ and k-means on $S_2$**
- C. k-means on $S_1$ and graph-based on $S_2$
- D. hierarchical on $S_1$ and graph-based on $S_2$

# Break & Quiz

**Q 1.1**: We have two datasets: a social network dataset $S_1$ which shows which individuals are friends with each other along with image dataset $S_2$.
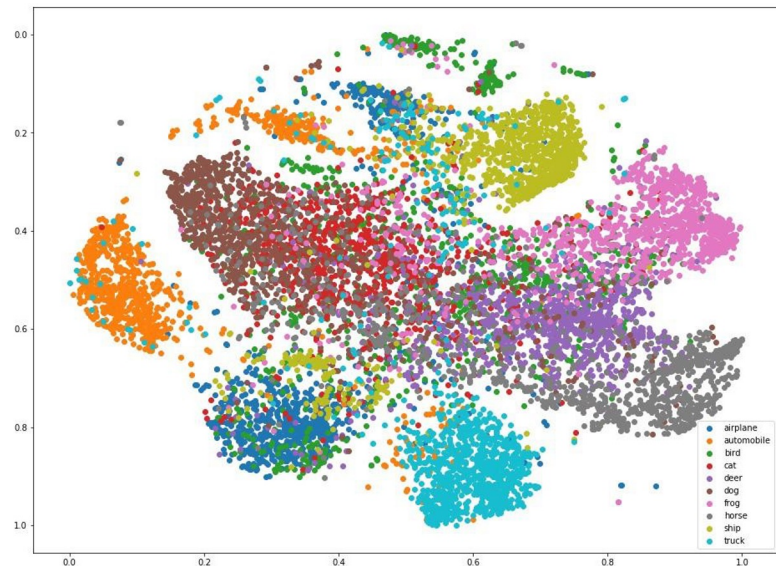
What kind of clustering can we do? Assume we do not make additional data transformations.

- A. k-means on both $S_1$ and $S_2$ **(No: can't do k-means on graph)**
- **B. graph-based on $S_1$ and k-means on $S_2$**
- C. k-means on $S_1$ and graph-based on S **(Same as A)**
- D. hierarchical on $S_1$ and graph-based on $S_2$ **(No: $S_2$ is not a graph)**

# Unsupervised Learning Beyond Clustering

Data analysis, dimensionality reduction, etc

- Already talked about PCA.
- Note: PCA can be used for visualization, but not specifically designed for it.
- Some algorithms are **specifically** for visualization.



Philip Slingerland

# Dimensionality Reduction & Visualization

Typical dataset: MNIST
- Handwritten digits 0-9
  - 60,000 images (small by ML standards)
  - 28×28 pixel (784 dimensions)
  - Standard for image experiments

- Dimensionality reduction?
  - Reducing dimensionality to 2-3 dimensions allows people to visualize data points and their relationships.

# Dimensionality Reduction & Visualization

## Run PCA on MNIST

- PCA is a linear mapping,
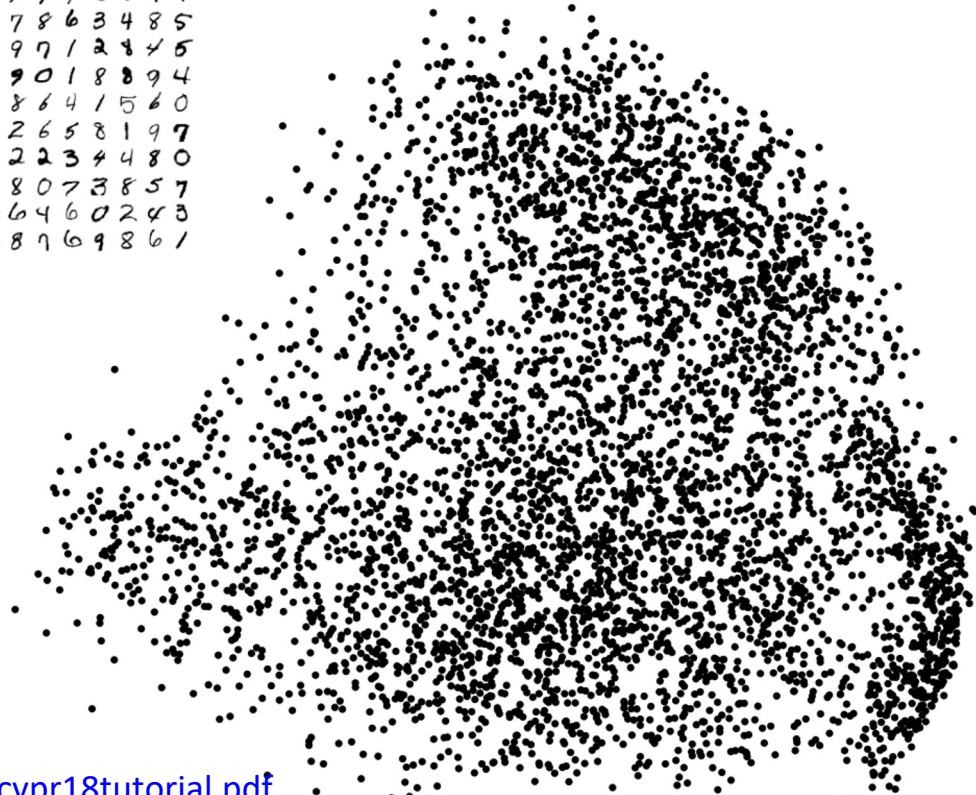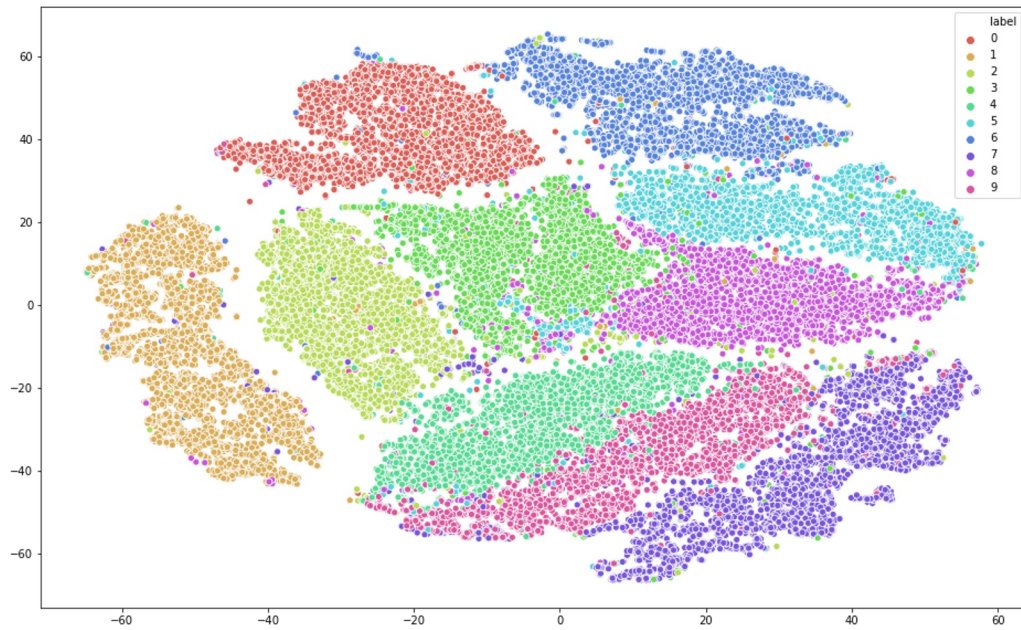  (can be restrictive)



Image source:
http://deeplearning.csail.mit.edu/slide_cvpr2018/laurens_cvpr18tutorial.pdf

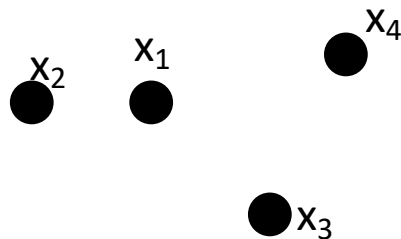# Visualization: **T-SNE**

Typical dataset: MNIST

- **T-SNE**: project data into just 2 dimensions

- Try to maintain structure

- MNIST Example

- **Input**: $x_1, x_2, ..., x_n$

- **Output**: 2D/3D $y_1, y_2, ..., y_n$

# **T-SNE** Algorithm: Step 1

How does it work? Two steps

- **1.** Turn vectors into probability pairs
- **2**. Turn pairs back into **(lower-dim)** vectors

**Intuition**: probability that $x_i$ would pick $x_j$ as its neighbor under a Gaussian probability

# T-SNE Examples

- Examples: (from Laurens van der Maaten)
- **Movies**:
  https://lvdmaaten.github.io/tsne/examples/netflix_tsne.jpg

# T-SNE Examples

- Examples: (from Laurens van der Maaten)
- **NORB**:
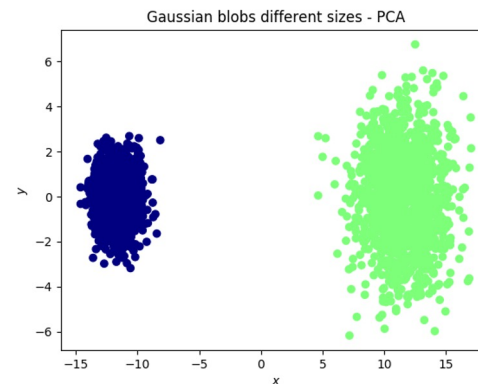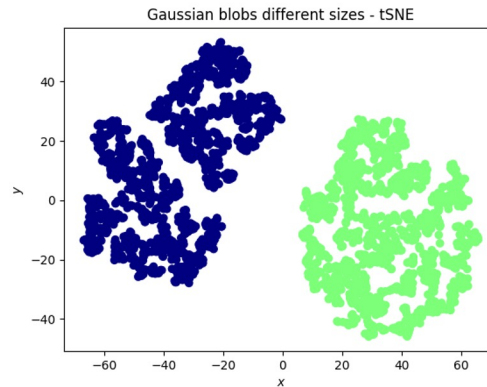  https://lvdmaaten.github.io/tsne/examples/norb_tsne.jpg

# Visualization: **T-SNE**

t-SNE vs PCA?

- "Local" vs "Global"

- Lose information in t-SNE
  - not a bad thing necessarily

- Downstream use

Good resource/credit:

https://www.thekerneltrip.com/statistics/tsne-vs-pca/



Gaussian blobs different sizes - tSNE



Gaussian blobs different sizes - PCA

# Break & Quiz

**Q 2.1**: Can we do t-SNE on NLP (words) or graph datasets?

- A. Never
- B. Yes, after running PCA on them
- C. Yes, after mapping them into $R^d$ (ie, embedding)
- D. Yes, after running hierarchical clustering on them

# Break & Quiz

**Q 2.1**: Can we do t-SNE on NLP (words) or graph datasets?

- A. Never
- B. Yes, after running PCA on them
- **C. Yes, after mapping them into R$^d$ (ie, embedding)**
- D. Yes, after running hierarchical clustering on them

# Break & Quiz

**Q 2.1**: Can we do t-SNE on NLP (words) or graph datasets?

- A. Never **(No: too strong)**
- B. Yes, after running PCA on them **(No: can't run PCA on words or graphs directly. Need vectors)**
- **C. Yes, after mapping them into $R^d$ (ie, embedding)**
- D. Yes, after running hierarchical clustering on them **(No: hierarchical clustering gives us a graph)**

# Short Intro to Density Estimation

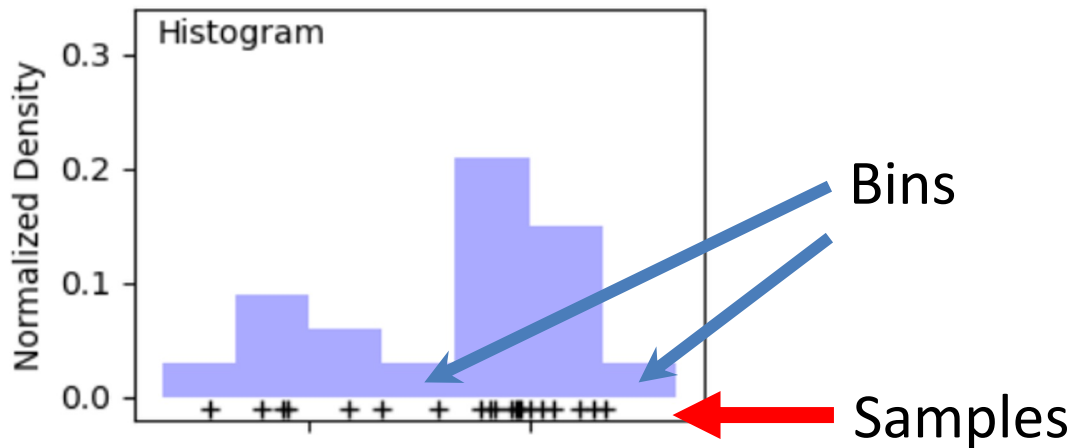Goal: given samples $x_1$, ..., $x_n$ from some distribution $P$, estimate P.

- Compute statistics (mean, variance)
- Generate samples from P
- Run inference



Zach Monge

# Simplest Idea: Histograms

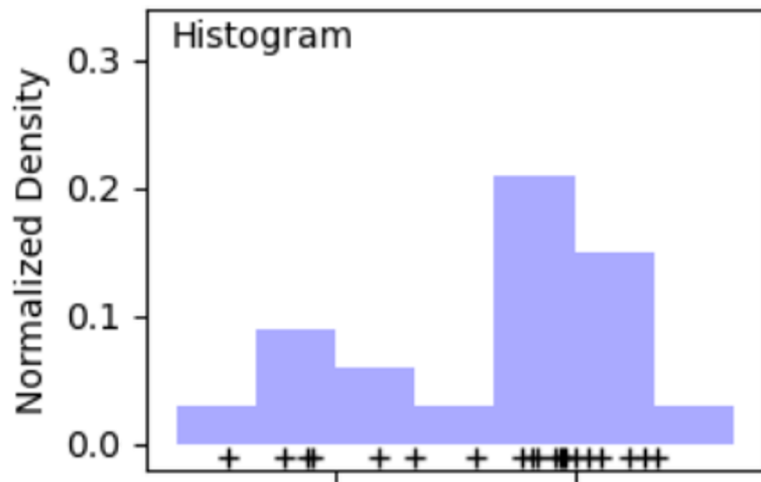Goal: given samples $x_1$, ..., $x_n$ from some distribution $P$, estimate P.



Define bins; count # of samples in each bin, normalize

# Simplest Idea: Histograms

Goal: given samples $x_1, …, x_n$ from some distribution $P$, estimate P.

**Downsides:**

i) High-dimensions: most bins are empty.

ii) Not continuous.

iii) How to choose bins?

# Kernel Density Estimation

Goal: given samples $x_1, \ldots, x_n$ from some distribution $P$, estimate P.

**Idea**: represent density as combination of "kernels"

$$f(x) = \frac{1}{nh} \sum_{i=1}^{n} K\left(\frac{x - x_i}{h}\right)$$

Center at each point

Kernel function: often Gaussian

Width parameter

# Kernel Density Estimation

**Idea**: represent density as combination of kernels

- "Smooth" out the histogram