30.08.2024

# Statistical Methods in AI (CS7.403)

Lecture-8: Clustering (k-means, Gaussian Mixture Models)

Ravi Kiran (ravi.kiran@iiit.ac.in)

https://ravika.github.io

@vikataravi

Center for Visual Information Technology (CVIT)

IIIT Hyderabad

Repeat until convergence: {

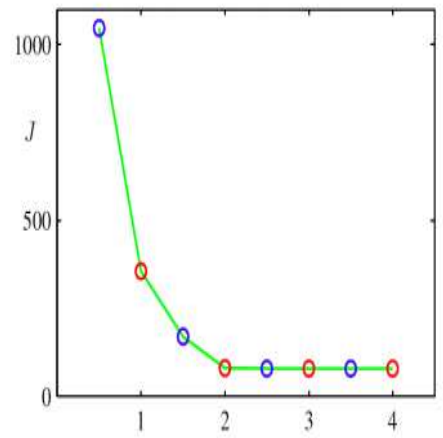For every $i$, set

$$c^{(i)} := \arg\min_j \|x^{(i)} - \mu_j\|^2.$$

For each $j$, set

$$\mu_j := \frac{\sum_{i=1}^{m} 1\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{m} 1\{c^{(i)} = j\}}.$$

}

$$\sum_{k=1}^{K} \sum_{i=1}^{n_k} \|x_{ki} - \mu_k\|^2$$

- Whenever an assignment is changed, the sum squared distances $J$ of data points from their assigned cluster centers is reduced.

- Whenever a cluster center is moved, $J$ is reduced.

- Test for convergence: If the assignments do not change in the assignment step, we have converged (to at least a local minimum).



- K-means cost function after each E step (blue) and M step (red). The algorithm has converged after the third M step

- The objective $J$ is non-convex (so coordinate descent on $J$ is not guaranteed to converge to the global minimum)



$$se = \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left\| x_{ki} - \mu_k \right\|^2$$

1. Initialize **cluster centroids** $\mu_1, \mu_2, \ldots, \mu_k \in \mathbb{R}^n$ randomly.

2. Repeat until convergence: {

   For every $i$, set
   $$c^{(i)} := \arg\min_j \|x^{(i)} - \mu_j\|^2.$$

   For each $j$, set
   $$\mu_j := \frac{\sum_{i=1}^{m} 1\{c^{(i)} = j\} x^{(i)}}{\sum_{i=1}^{m} 1\{c^{(i)} = j\}}.$$
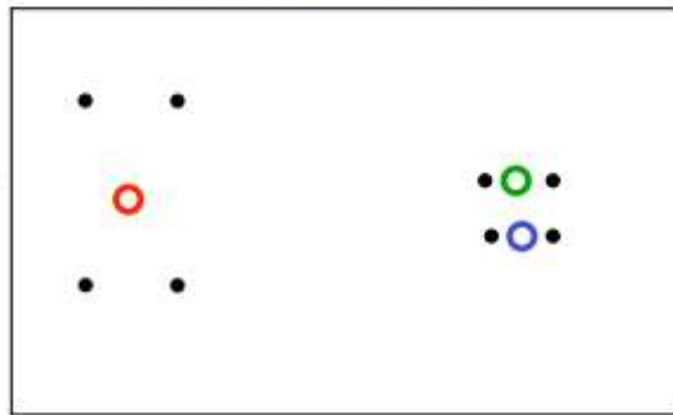
   }

- The objective $J$ is non-convex (so coordinate descent on $J$ is not guaranteed to converge to the global minimum)

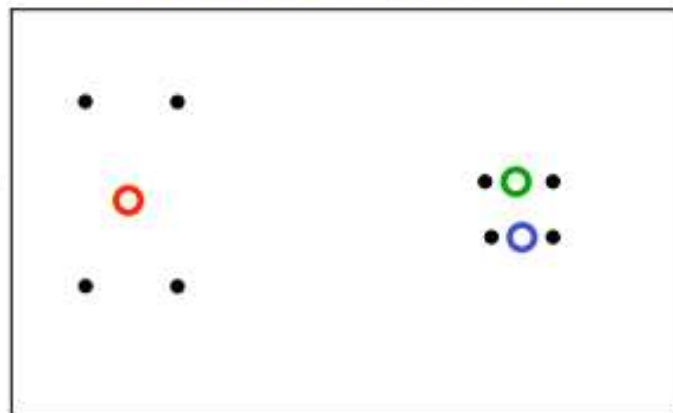- There is nothing to prevent k-means getting stuck at local minima.



A bad local optimum

- The objective $J$ is non-convex (so coordinate descent on $J$ is not guaranteed to converge to the global minimum)

- There is nothing to prevent k-means getting stuck at local minima.

- We could try many random starting points

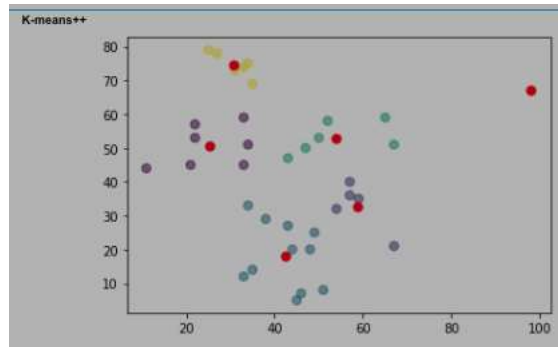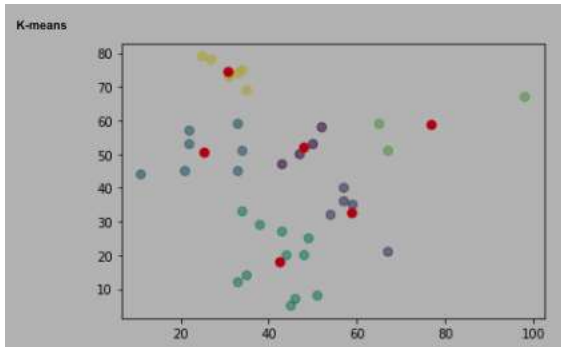$$\sum_{k=1}^{K}\sum_{i=1}^{n_k}\left\|x_{ki}-\mu_k\right\|^2$$

A bad local optimum

# K-means++: Improving K-means initialization

- Common way to improve k-means - smart initialization!

- General idea - try to get good coverage of the data.

- k-means++ algorithm:

    1. Pick the first center randomly
    2. For all points $\mathbf{x}^{(n)}$ set $d^{(n)}$ to be the distance to closest center.
    3. Pick the new center to be at $\mathbf{x}^{(n)}$ with probability proportional to $d^{(n)2}$
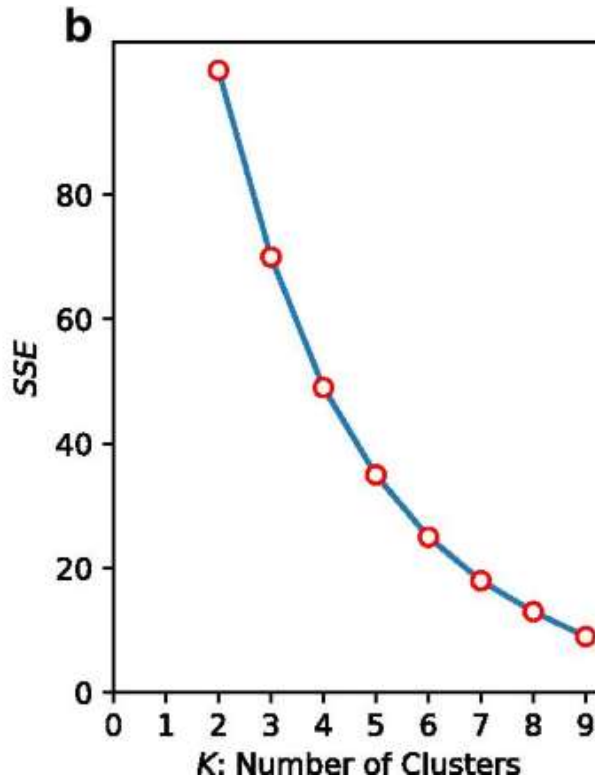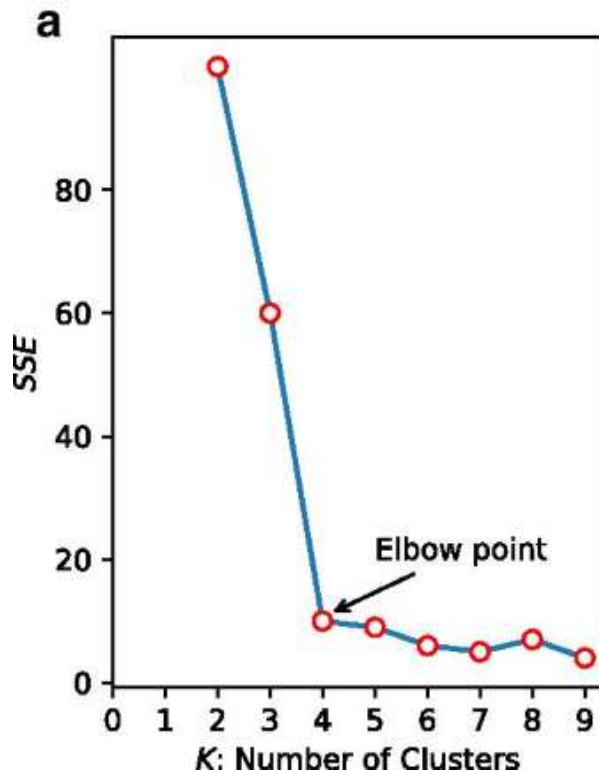    4. Repeat steps 2+3 until you have k centers

# K-means++: Improving K-means initialization

- Common way to improve k-means - smart initialization!

- General idea - try to get good coverage of the data.

- k-means++ algorithm:

  1. Pick the first center randomly
  2. For all points $\mathbf{x}^{(n)}$ set $d^{(n)}$ to be the distance to closest center.
  3. Pick the new center to be at $\mathbf{x}^{(n)}$ with probability proportional to $d^{(n)2}$
  4. Repeat steps 2+3 until you have k centers

# How to choose k ?

$$se = \sum_{k=1}^{K} \sum_{i=1}^{n_k} \left\| x_{ki} - \mu_k \right\|^2$$



**a**

SSE

Elbow point

*K*: Number of Clusters

**b**

SSE

*K*: Number of Clusters

# Regularization

- Penalize "overly" large or "overly" small clusters
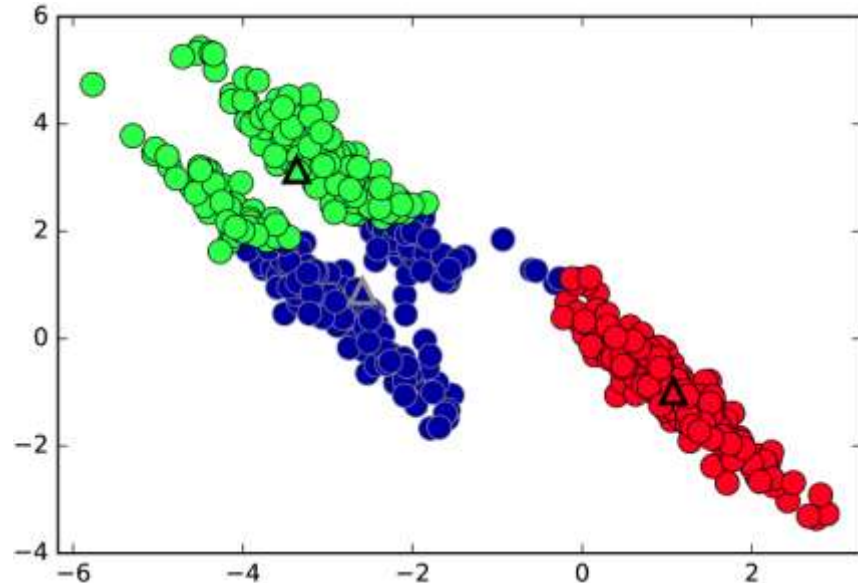
# K-mediods

- Squared Euclidean distance loss function of K-means not robust.

- Use L1 loss function $J = \sum_{i=1}^{n} \sum_{k=1}^{K} r_{ik} \|x_i - \mu_k\|_1$ instead of squared Euclidean distance.

- Use an iterative procedure as before.
  - Prototype is the median of the points assigned to a cluster.

# K-means: Additional issues

- 'Hard' assignments
- Euclidean ➜ Favours 'spherical' clusters of equal 'contribution'
- Sensitive to initialization
- Sensitive to outliers
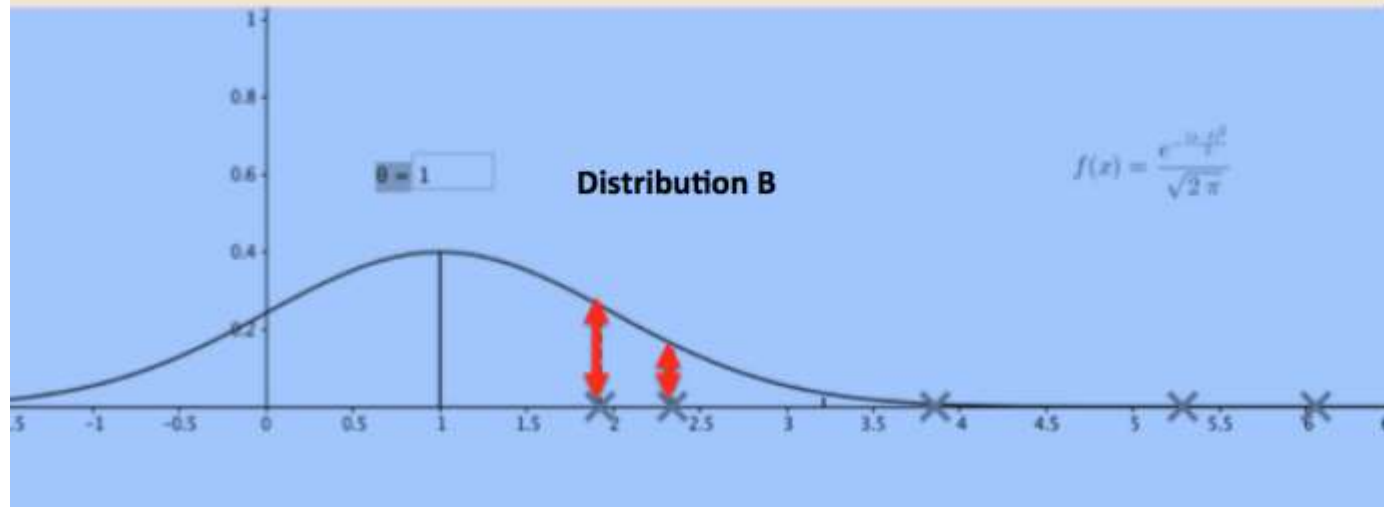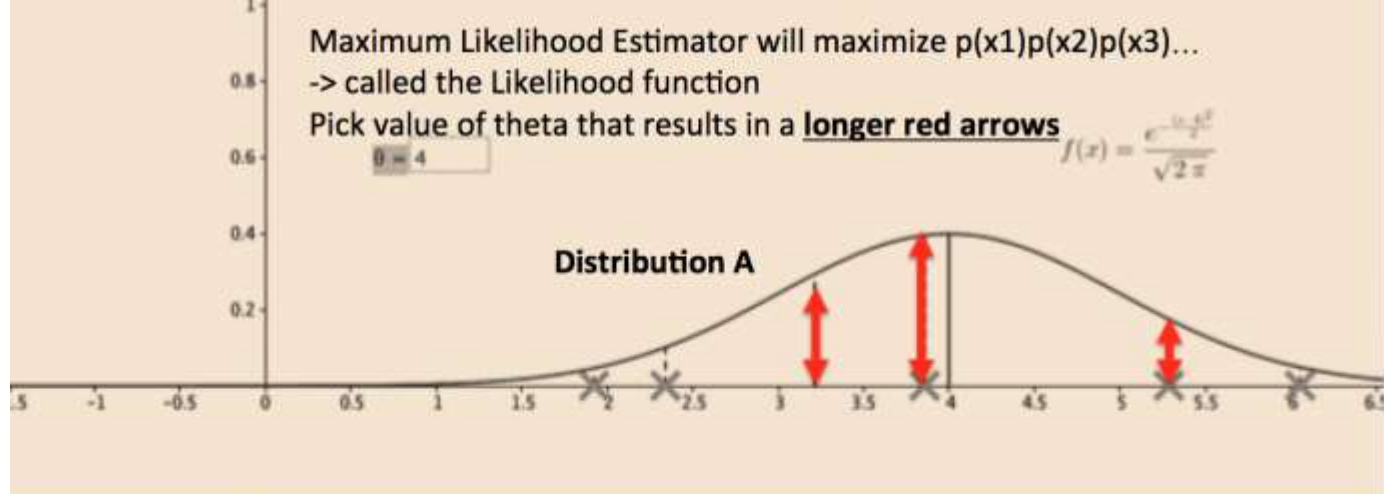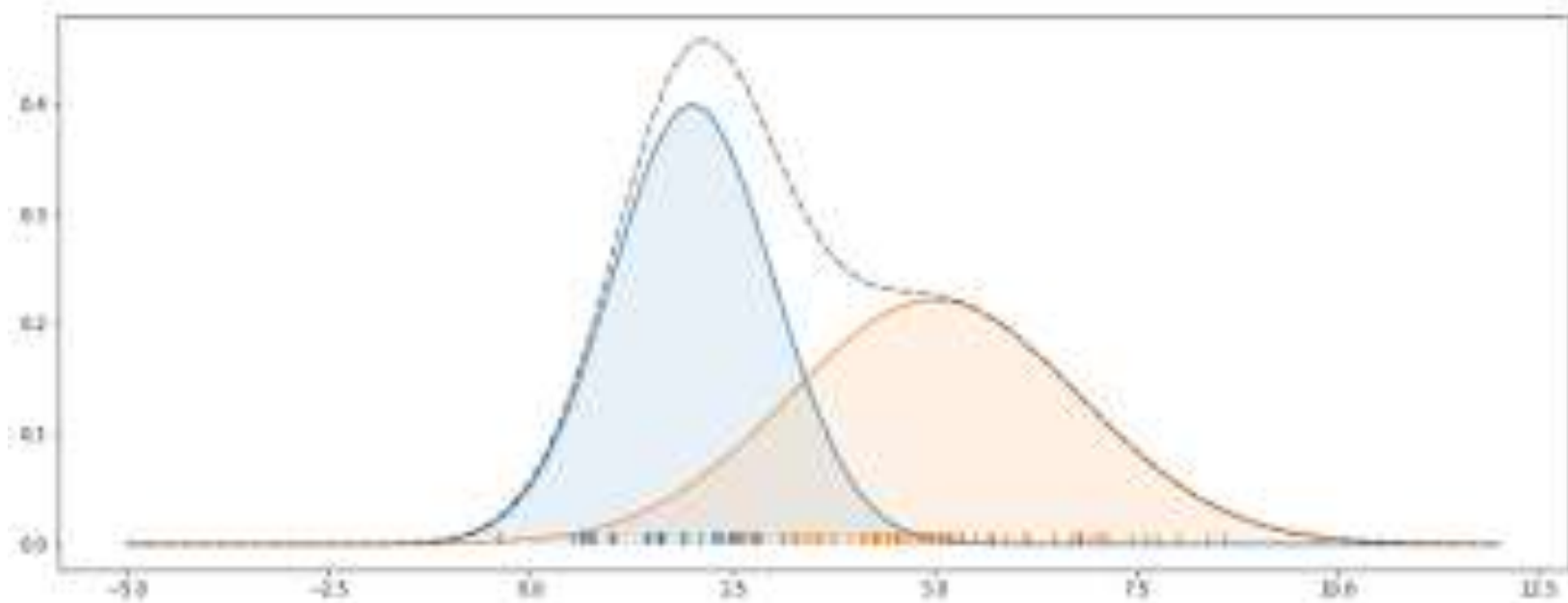
# Maximum Likelihood Estimation

Maximum Likelihood Estimator will maximize p(x1)p(x2)p(x3)...
-> called the Likelihood function
Pick value of theta that results in a **longer red arrows**

θ = 4

$$f(x) = \frac{e^{-\frac{1}{2}x^2}}{\sqrt{2\pi}}$$

**Distribution A**

θ = 1

**Distribution B**

$$f(x) = \frac{e^{-\frac{1}{2}x^2}}{\sqrt{2\pi}}$$

https://pmirla.github.io/2016/07/20/maximum-likelihood-explanation.html

# Mixture Models

- Data distribution $p(x)$ assumed to be a **weighted sum** of $K$ distributions

$$p(x) = \sum_{k=1}^{K} \pi_k p(x|\boldsymbol{\theta}_k)$$

  where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1$, $\pi_k \geq 0$ (intuitively, $\pi_k$ is the proportion of data generated by the $k$-th distribution)

- Each component distribution $p(x|\boldsymbol{\theta}_k)$ represents a "cluster" in the data
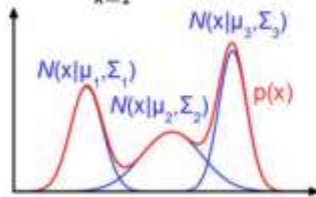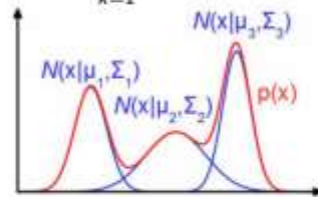
# Mixture Models

- Data distribution $p(x)$ assumed to be a **weighted sum** of $K$ distributions

$$p(x) = \sum_{k=1}^{K} \pi_k p(x|\theta_k)$$

where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1, \quad \pi_k \geq 0$ (intuitively, $\pi_k$ is the proportion of data generated by the $k$-th distribution)

- Each component distribution $p(x|\theta_k)$ represents a "cluster" in the data
- **Gaussian Mixture Model (GMM):** component distributions are Gaussians

$$p(x) = \sum_{k=1}^{K} \pi_k \mathcal{N}(x|\mu_k, \Sigma_k)$$

# Mixture Models

- Data distribution $p(x)$ assumed to be a **weighted sum** of $K$ distributions

$$p(x) = \sum_{k=1}^{K} \pi_k p(x|\boldsymbol{\theta}_k)$$

where $\pi_k$'s are the **mixing weights**: $\sum_{k=1}^{K} \pi_k = 1$, $\pi_k \geq 0$ (intuitively, $\pi_k$ is the proportion of data generated by the $k$-th distribution)
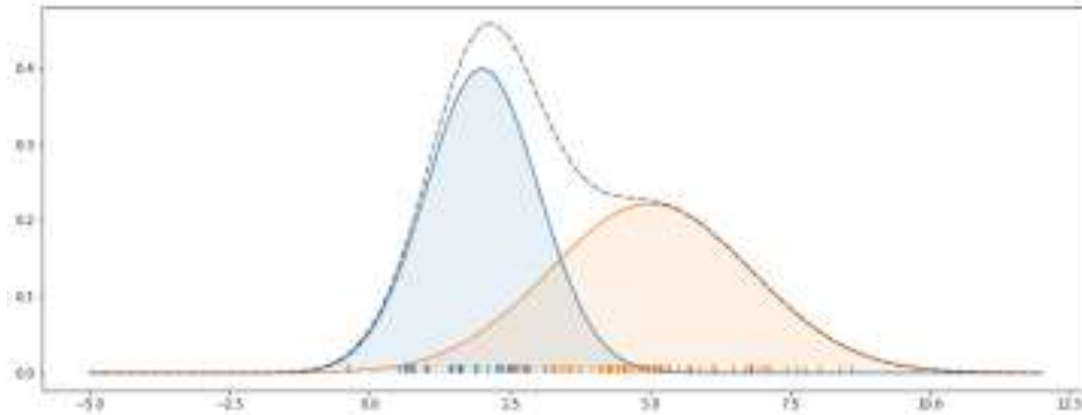
- Each component distribution $p(x|\boldsymbol{\theta}_k)$ represents a "cluster" in the data

- **Gaussian Mixture Model (GMM):** component distributions are Gaussians

$$p(x) = \sum_{k=1}^{K} \pi_k \mathcal{N}(x|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



- Mixture models used in many data modeling problems, e.g.,
  - Unsupervised Learning: Clustering (+density estimation)
  - Supervised Learning: Mixture of Experts models

# Mixture Models



A GMM represents a **distribution** as

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

with $\pi_k$ the mixing coefficients, where:

$$\sum_{k=1}^{K} \pi_k = 1 \quad \text{and} \quad \pi_k \geq 0 \quad \forall k$$
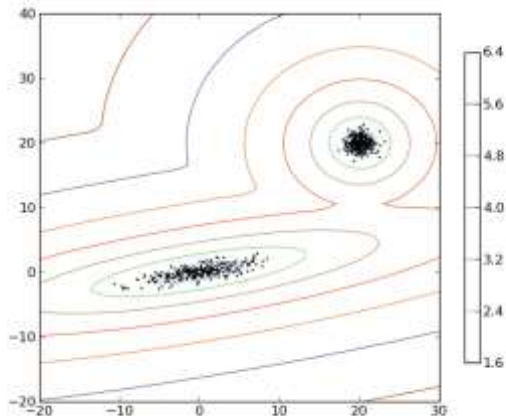
# Mixture Models

Most common mixture model: Gaussian mixture model (GMM)

- A GMM represents a **distribution** as

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

with $\pi_k$ the mixing coefficients, where:

$$\sum_{k=1}^{K} \pi_k = 1 \quad \text{and} \quad \pi_k \geq 0 \quad \forall k$$



http://scikit-learn.sourceforge.net/0.5/auto_examples/gmm/plot_gmm_pdf.html

- Can think of the data $\{x_1, x_n, \dots, x_N\}$ using a "generative story"
  - For each example $x_n$, first choose its cluster assignment $z_n \in \{1, 2, \dots, K\}$ as

  $$z_n \sim \text{Multinoulli}(\pi_1, \pi_2, \dots, \pi_K) \quad \text{aka "categorical"}$$

  - Now generate $x$ from the Gaussian with id $z_n$

  $$x_n | z_n \sim \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$$

# Resources

- Textbook
  - PRML (Bishop) – Chapter 9: 9.1,9.2,9.3.2
  - Pattern Classification (Duda, Hart, Stork)
    - 10.4.3,10.6.1,10.7.1,10.7.2,10.8,10.10
- Videos
  - https://www.youtube.com/watch?v=REypj2sy_5U&list=PLBv09BD7ez_4e9LtmK626Evn1ion6ynrt
  - https://www.youtube.com/watch?v=rVfZHWTwXSA

- Blog posts/Lecture Notes
  - https://www.cse.iitk.ac.in/users/piyush/courses/pml_winter16/slides_lec7.pdf
  - https://see.stanford.edu/materials/aimlcs229/cs229-notes8.pdf
  - https://www.cs.toronto.edu/~jlucas/teaching/csc411/lectures/lec15_16_handout.pdf
  - https://mbernste.github.io/posts/gmm_em/
  - https://www.ritchievink.com/blog/2019/05/24/algorithm-breakdown-expectation-maximization/