# Machine Learning from Data

Lecture 19: Spring 2021

# Today's Lecture
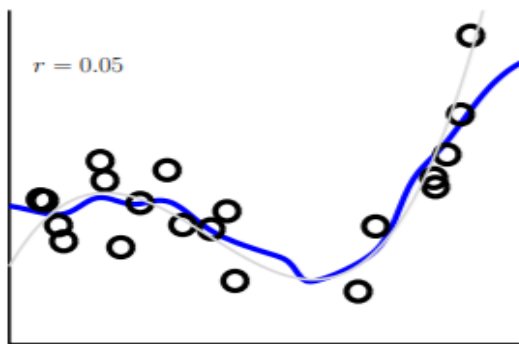
- Unsupervised Learning
- K-means clustering
- Probability Density Estimation
- Gaussian Mixture Models

# Radial Basis Functions

## Nonparametric RBF

$$g(\mathbf{x}) = \sum_{n=1}^{N} \left( \frac{\alpha_n(\mathbf{x})}{\sum_{m=1}^{N} \alpha_m(\mathbf{x})} \right) \cdot y_n$$

$$\alpha_n(\mathbf{x}) = \phi \left( \frac{\|\mathbf{x} - \mathbf{x}_n\|}{r} \right) \qquad \text{(bump on } \mathbf{x})$$
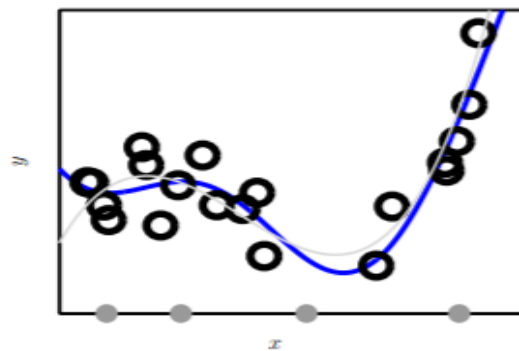


$r = 0.05$

No Training

## Parametric $k$-RBF-Network

$$h(\mathbf{x}) = w_0 + \sum_{j=1}^{k} w_j \cdot \phi \left( \frac{\|\mathbf{x} - \boldsymbol{\mu}_j\|}{r} \right)$$

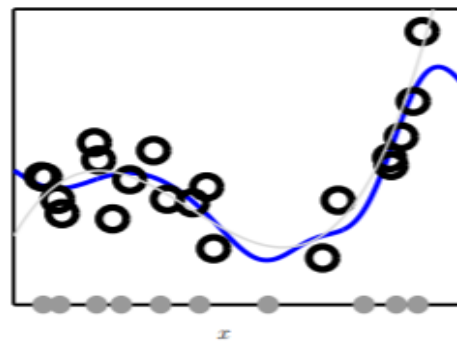$$= \mathbf{w}^{\mathbf{T}} \boldsymbol{\Phi}(\mathbf{x})$$

(bump on $\boldsymbol{\mu}_j$)

linear model given $\boldsymbol{\mu}_j$

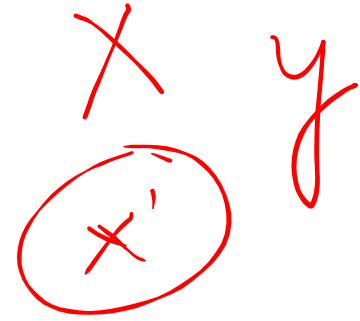choose $\boldsymbol{\mu}_j$ as centers of $k$-clusters of data



$k = 4, r = \frac{1}{k}$ $\qquad\qquad$ $k = 10$, regularized

# Unsupervised Learning

- Preprocessor to organize the data for supervised learning:

    Organize data for faster nearest neighbor search

    Determine centers for RBF bumps.

- Important to be able to organize the data to identify patterns.
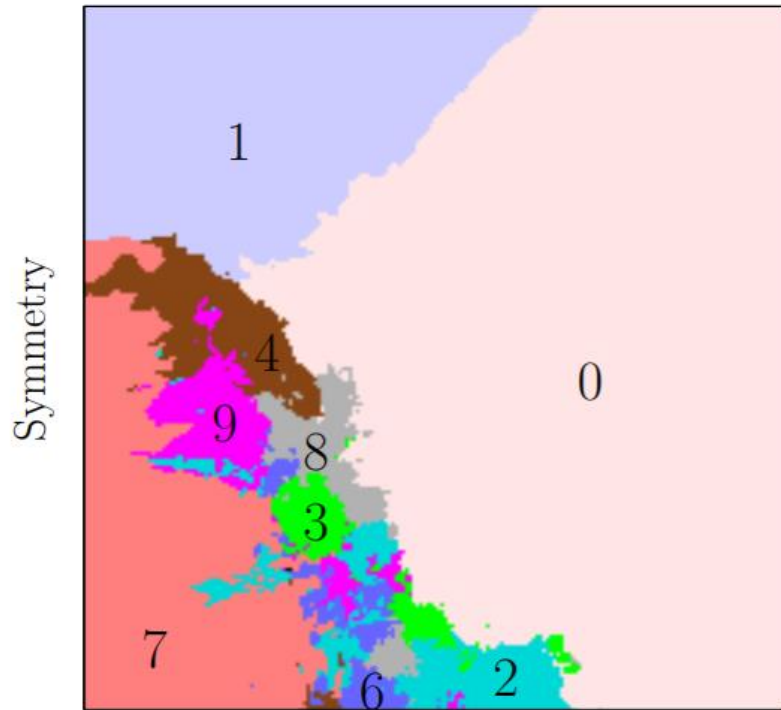
    Learn the patterns in data, e.g. the patterns in a language before getting into a supervised setting.

    amazon.com organizes books into categories

# Clustering Digits
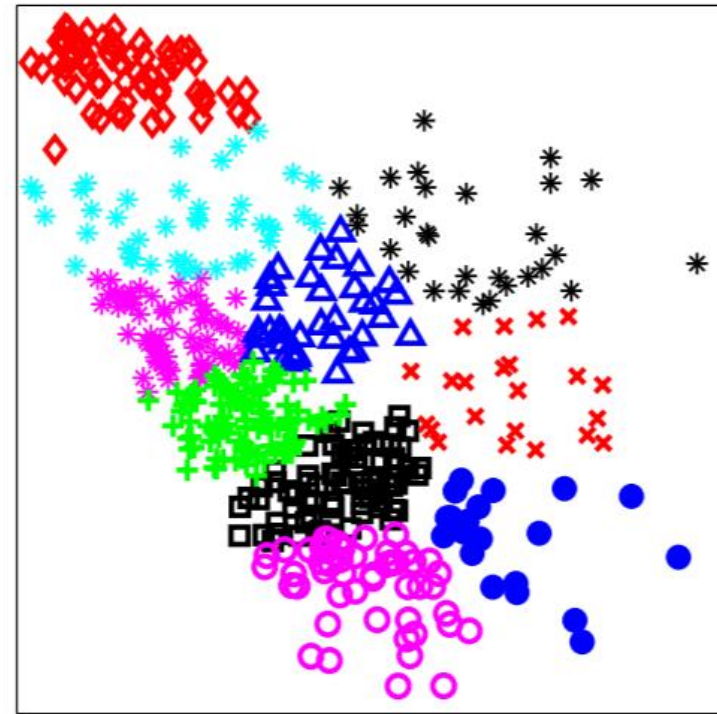


21-NN rule, 10 Classes

*labels*

10 Clustering of Data

# Clustering [Unsupervised]

$$D = [x_1, x_2 \ldots x_N]$$

→ Assume we are given
k-subsets.

→ $S_1 \ S_2 \ldots \ S_k$

Cluster data into
some subsets

↓

We need to know
how many?

i) Disjoint (Non overlapping)

ii) Complete (cover entire dataset)

$$S_i \cap S_j = \phi \quad i \neq j$$

$$\bigcup_{i=1}^{k} S_i = D$$

Trivial

OUTPUT

What makes clustering good?

1) Tightly bound — within cluster ✓

2) Well separated — between cluster.

$$S_1 \, S_2 \cdots S_K \implies \mu_j$$

$$E_j = \sum_{x_i \in S_j} \| x_i - \mu_j \|^2 \longrightarrow \text{Spread of cluster.}$$

$$E = \sum_{j=1}^{K} E_j \longrightarrow E \simeq \text{small}$$

Minimize $\longrightarrow$

$$E = \sum_{n=1}^{N} \| x_n - \mu(x_n) \|^2$$

K-means clustering error.

Output $\longrightarrow$ $S_1$ $S_2$ $\cdots$ $S_K$ $\Big\}$ NP-hard
$\quad\quad\quad\quad$ $\mu_1$ $\mu_2$ $\quad\quad$ $\mu_K$ $\quad\quad\uparrow$

Alternating Algorithm / Lloyd's

i) $\mu_1, \mu_2 \cdots \mu_K$ (given) $\longrightarrow$ greedily.

$\quad\quad\quad\quad\quad$ ii) $S_1$ $S_2$ $\cdots$ $S_K$ (given)



$\mu_{in}$

$x_n$

$E_j = \sum_{x \in S_j} \|x - \mu_j\|^2$

$\mu_j^* \quad \boxed{\mu_j^* = \dfrac{\sum_{x \in S_j} x}{|S_j|}}$

# Lloyd's Algorithm

1) Find the centers $\mu_1 \; \mu_2 \; \text{---} \; \mu_k$ . How?
   $\rightarrow$ Randomly selected. (first step)
   OR $\rightarrow$ Mean vector in $S_j$

2) Assign data to centers greedily.
   $$\arg\min \|x_n - \mu_j\|^2$$

Repeat Until convergence.

# Convergence

Change $\longrightarrow$ Error $\longrightarrow$ Improving.
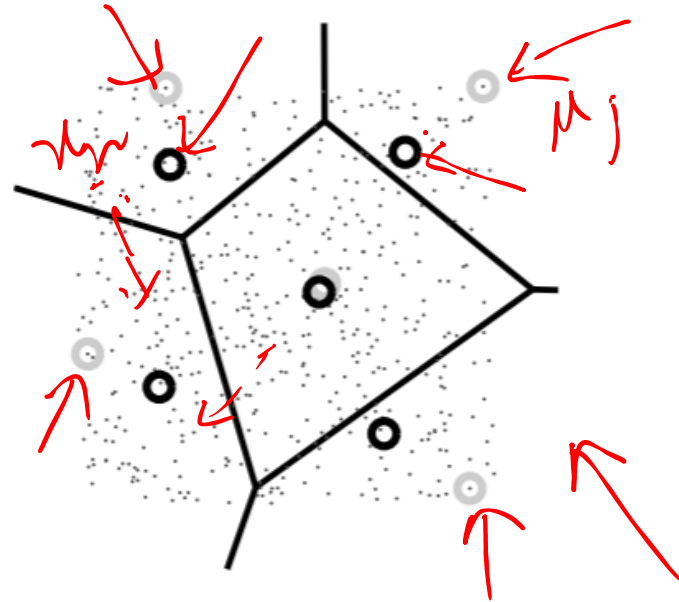
$\downarrow$

decreasing

$\longrightarrow$ local minima.
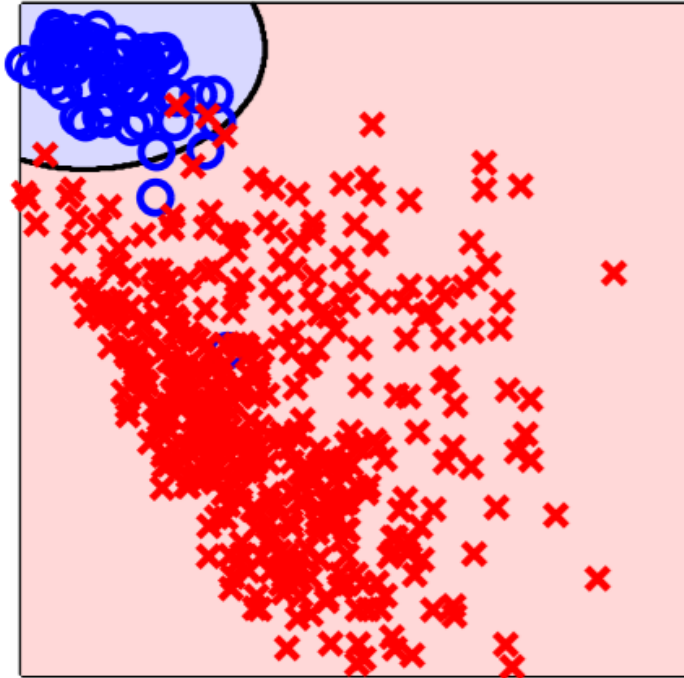
# Lloyd's Algorithm for $k$-Means Clustering

*Digits*

$$E_{\text{in}}(S_1, \ldots, S_k; \boldsymbol{\mu}_1, \ldots, \boldsymbol{\mu}_k) = \sum_{n=1}^{N} \| \mathbf{x}_n - \boldsymbol{\mu}(\mathbf{x}_n) \|^2$$

1: **Initialize** Pick well separated centers $\boldsymbol{\mu}_j$.

2: **Update** $S_j$ to be all points closest $\boldsymbol{\mu}_j$.

$$S_j \leftarrow \{ \mathbf{x}_n : \| \mathbf{x}_n - \boldsymbol{\mu}_j \| \leq \| \mathbf{x}_n - \boldsymbol{\mu}_\ell \| \text{ for } \ell = 1, \ldots, k \}.$$

3: **Update** $\boldsymbol{\mu}_j$ to the centroid of $S_j$.

$$\boldsymbol{\mu}_j \leftarrow \frac{1}{|S_j|} \sum_{\mathbf{x}_n \in S_j} \mathbf{x}_n$$

4: Repeat steps 2 and 3 until $E_{\text{in}}$ stops decreasing.

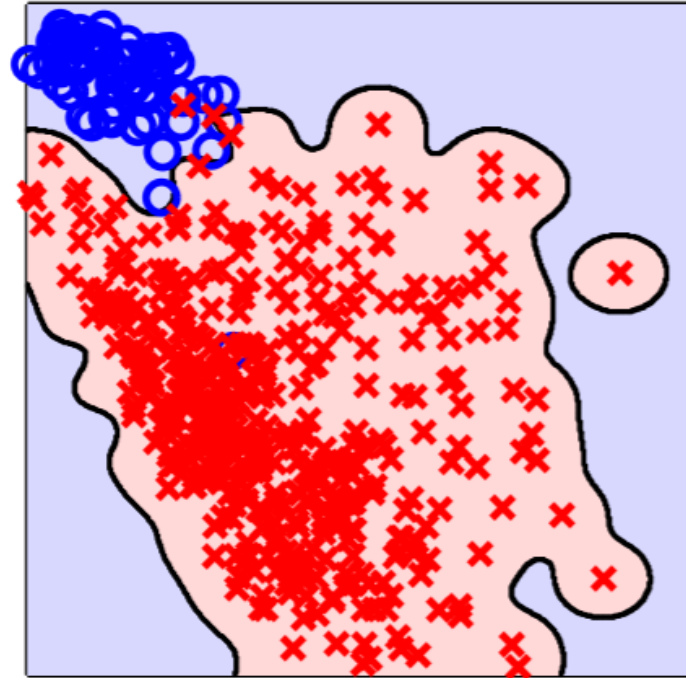$\boldsymbol{\mu}_j$

# Application to $k$-RBF-Network

## 10-center RBF-network

## 300-center RBF-network



Choosing $k$ - knowledge of problem (10 digits) or CV.

1 VS Not 1

Overfitting
↓
Regularization

# Probability Density Estimation

$$P(\mathbf{x})$$

$P(\mathbf{x})$ measures how likely it is to generate inputs *similar* to $\mathbf{x}$.
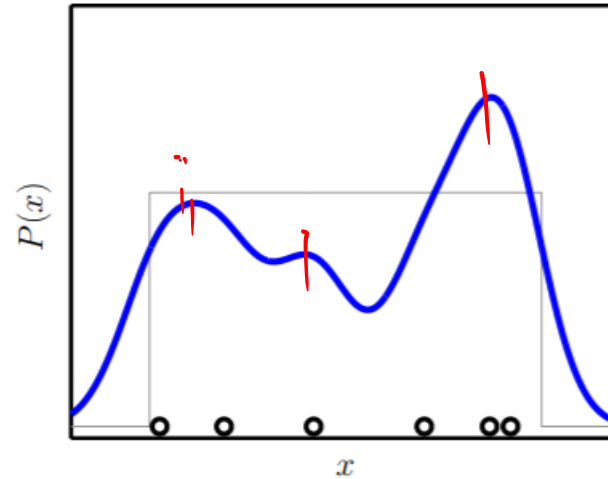
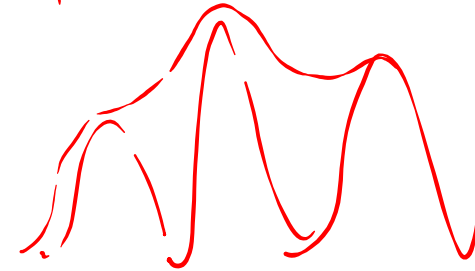Estimating $P(\mathbf{x})$ results in a 'softer/finer' representation than clustering

Clusters are regions of high probability.

Data $\rightarrow P(x)$

Test

Similarity

# Parzen Windows – RBF density estimation

Basic idea: put a bump of 'size' (volume) $\frac{1}{N}$ on each data point.



$$\hat{P}(\mathbf{x}) = \frac{1}{Nr^d} \sum_{i=1}^{N} \phi \left( \frac{\|\mathbf{x} - \mathbf{x}_i\|}{r} \right)$$

$$\phi(z) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}z^2}$$

Normalized
Prob. density

$\mathbf{x}_i$

scale

$N$

$\frac{1}{N}$

# Digits Data

Non-parametric RBF

### RBF Density Estimate

high prob.

$y$     $x$ (Intensity)

symmetry

### Density Contours

$y$

$x$

# Gaussian Mixture Model (GMM)

$k = 4$

$M_1 \ M_2 \ M_3 \ M_4$

Shape : Covariance Matrix

$\Sigma_1 \ \Sigma_2 \ \Sigma_3 \ \Sigma_4$

Bump weights

$\omega_1 \ \omega_2 \ \omega_3 \ \omega_4$

gaussian density.

$$\hat{P}(x) = \sum_{k=1}^{K} \omega_k \, N(x, M_k, \Sigma_k)$$

GMM

$$N\left(x, \overset{\frown}{\underset{\text{mean}}{\textcircled{$\mu$}}}, \underset{\substack{\text{covariance} \\ \text{matrix}}}{\Sigma}\right)$$

$$x, \ x_2 \text{ --- } x_n$$

$$N(x, \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \, e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

$$\frac{1}{M} \sum x_i = E[x] = \mu$$

$$\frac{1}{M} \sum x_i x_i^T \Longrightarrow E[x x^T] = \mu \mu^T + \Sigma$$

$$\hat{P}(x) = \sum_{k=1}^{K} \omega_k \, N(x, \mu_k, \Sigma_k), \quad \omega_k > 0$$

$$\sum_{k=1}^{K} \omega_k = 1$$

$$\left( \omega_j, \; \mu_k, \; \Sigma_j \right)$$

weights   centers  shapes.

Pick in such a way $\longrightarrow$ max. the prob. to generate the data $x_1 \, x_2 \cdots x_L$

Maximum Likelihood.

# Expectation Maximization

1) we knew which $x_n$ belongs to which bump

$\longrightarrow \mu_j \propto \sum_j \propto w_j$

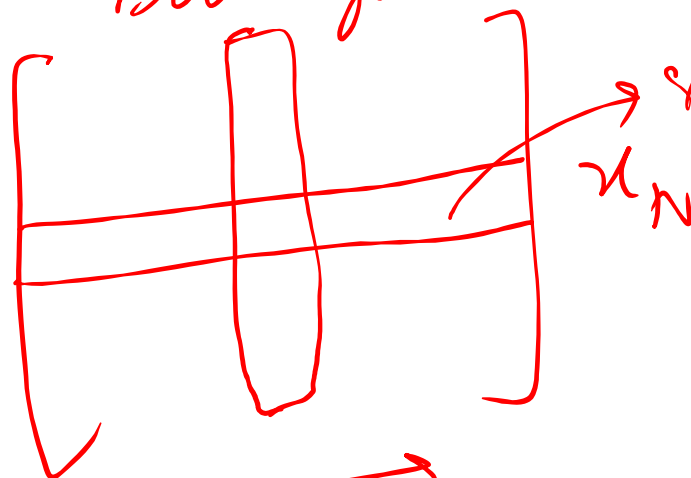2) If we knew $\{ w_j \ \mu_j \ \Sigma_j \} \longrightarrow$ which $x_n$ belongs to $j$

# Algo

1) **Initialization:**

membership variable $r_{nj}$ = The fraction of Data point $x_n$ that belongs to bump $j$

$$r = N$$

$$\text{sum} = 1 \quad r_{nj} \geq 0$$

$$x_N$$

$$\sum_j r_{nj} = 1$$

$$\leftarrow k \rightarrow$$

Assume we know $\{r_{nj}\}$

$$N_j = \sum_{n=1}^{N} r_{nj}$$

$$N = \sum_j N_j$$

Volume $\omega_j = \dfrac{N_j}{N}$

$$\mu_j = \dfrac{\sum_{n=1}^{N} r_{nj} x_n}{N_j}$$

$$\Sigma_j = \dfrac{1}{N_j} \sum_{n=1}^{N} r_{nj} x_n x_n^T - \mu_j \mu_j^T$$

② Update $r_{nj} = \dfrac{P[x_n \text{ came from bump } j]}{\sum_j P[x_n \text{ came from bump } j]}$

$$= \dfrac{\omega_j N(x_n \mu_j \Sigma_j)}{\sum_j \omega_j N(x_n \mu_j \Sigma_j)}$$
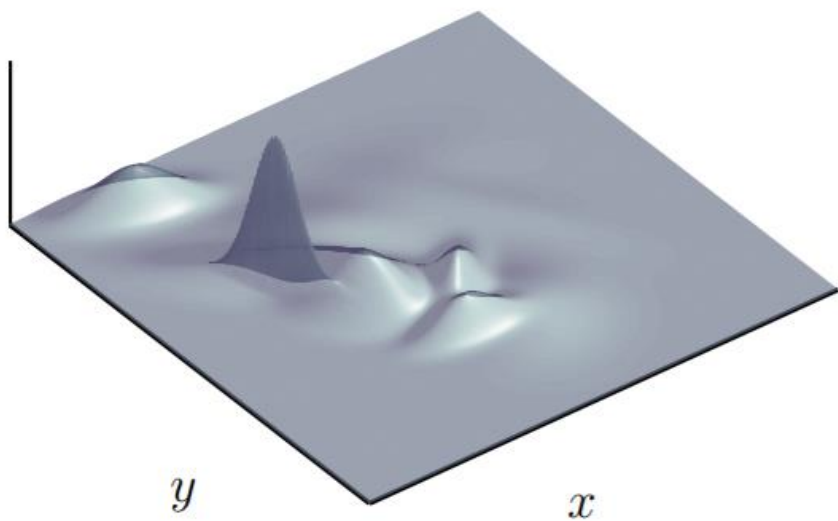
③ Repeat Until convergence

# E-M Algorithm

**E-M Algorithm for GMMs:**
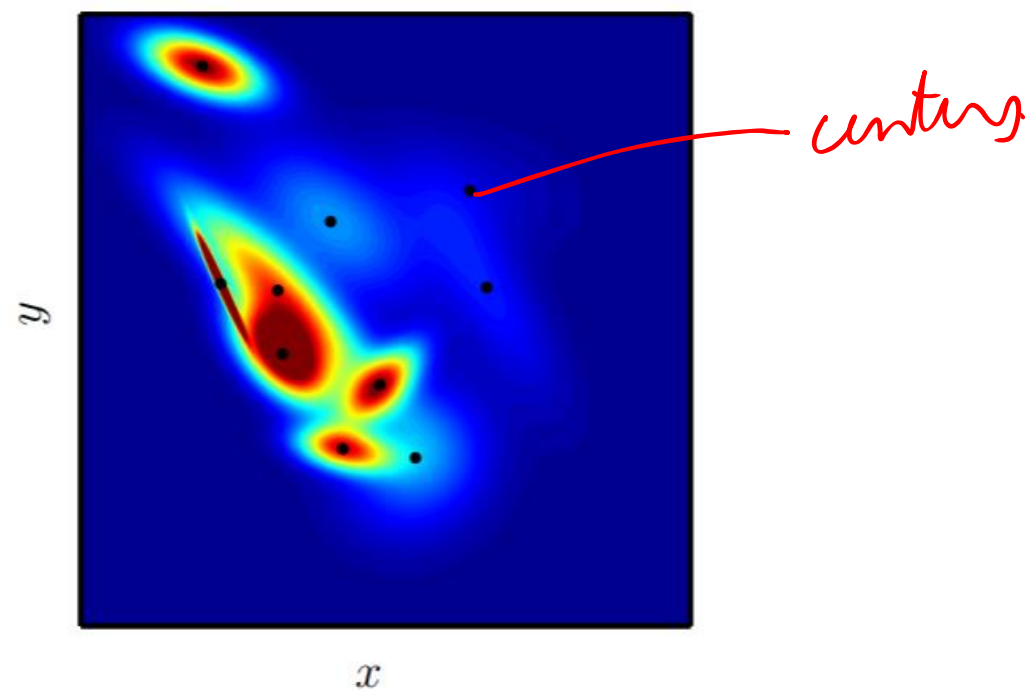
1: Start with estimates for the bump membership $\gamma_{nj}$.

2: Estimate $w_j, \boldsymbol{\mu}_j, \Sigma_j$ given the bump memberships.

3: Update the bump memberships given $w_j, \boldsymbol{\mu}_j, \Sigma_j$;

4: Iterate to step 2 until convergence.

# GMM on Digits Data

10-center GMM

Density Contours



centers

# Thanks!