

Clustering and mixture models

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Kai-Wei Chang, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Outline

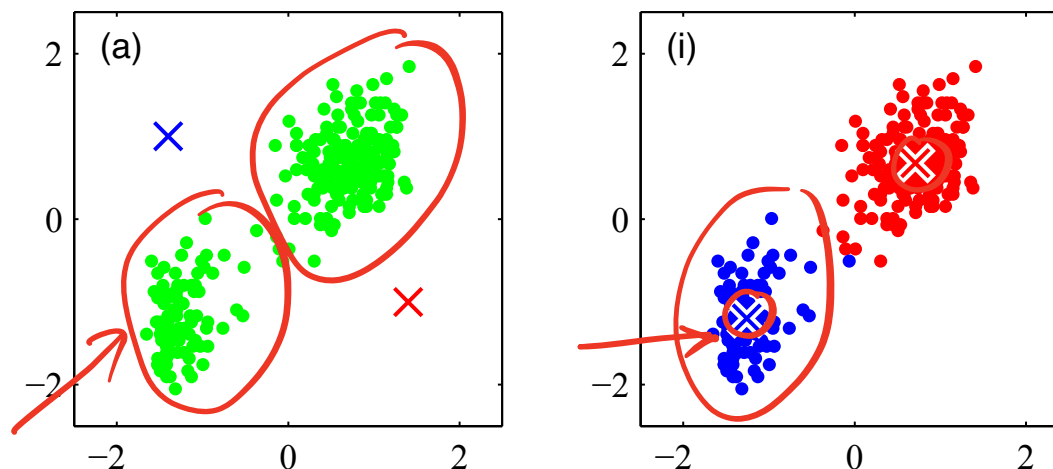
- 1 K-means
- 2 Gaussian mixture models
- 3 GMMs and Incomplete Data

Clustering

Setup Given $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$ and K , we want to output

- $\{\mu_k\}_{k=1}^K$: prototypes (or centroids) of clusters
- $A(\mathbf{x}_n) \in \{1, 2, \dots, K\}$: the cluster membership, i.e., the cluster ID assigned to \mathbf{x}_n

Toy Example Cluster data into two clusters.

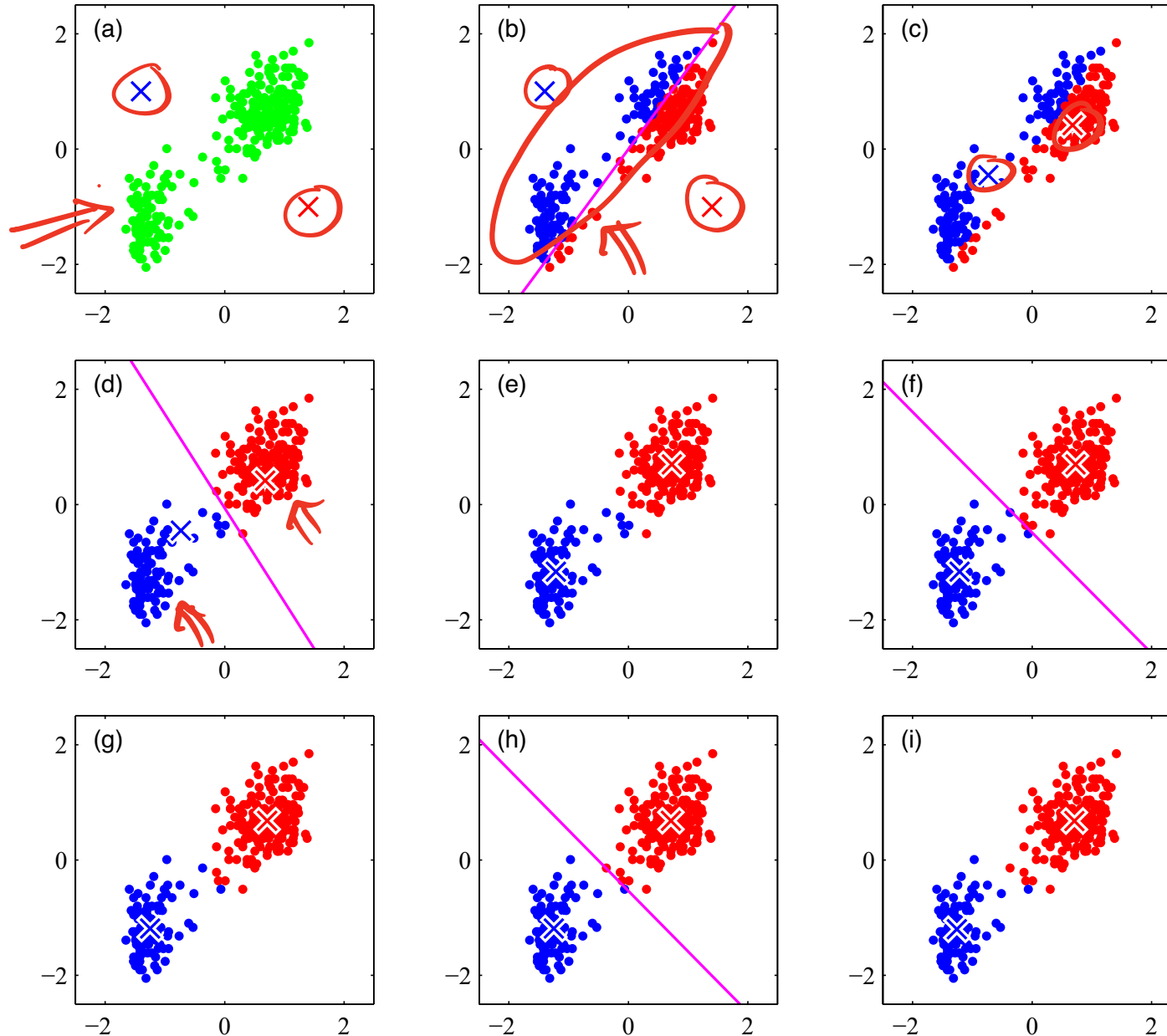


Applications

- Identify communities within social networks
- Find topics in news stories
- Group similar sequences into gene families

K-means example

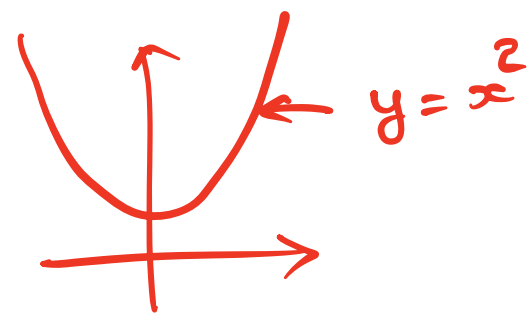
$K=2$



K-means clustering

Intuition Data points assigned to cluster k should be close to μ_k , the prototype.

K-means clustering



Intuition Data points assigned to cluster k should be close to μ_k , the prototype.

Distortion measure (clustering objective function, cost function)

$$J(\{r_{nk}\}, \{\mu_k\}) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

where $\underline{r_{nk}} \in \{0, 1\}$ is an indicator variable

$$\underline{r_{nk}} = 1 \quad \text{if and only if} \quad \underline{A(\mathbf{x}_n)} = \underline{k}$$

r_{nk}

$$J(\{\underline{r_{nk}}\}, \{\underline{\mu_k}\}) = \sum_n \sum_k \underline{r_{nk}} \|\mathbf{x}_n - \underline{\mu_k}\|_2^2$$

$J(\{\underline{r_{nk}}\}, \{\underline{\mu_k}\})$

K-means clustering

K-means objective

$$\operatorname{argmin}_{\{r_{nk}\}, \{\boldsymbol{\mu}_k\}} J(\{r_{nk}\}, \{\boldsymbol{\mu}_k\}) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|_2^2$$

where $r_{nk} \in \{0, 1\}$ is an indicator variable

$$r_{nk} = 1 \quad \text{if and only if} \quad A(\mathbf{x}_n) = k$$

- Is a non-convex objective function.
- Minimizing the K-means objective function is NP-hard.

Lloyd's algorithm for minimizing the K-means objective

Often simply called the **K-means algorithm**

Minimize cost function alternative optimization between $\{r_{nk}\}$ and $\{\mu_k\}$

- **Step 0** Initialize $\{\mu_k\}$ to some values

Lloyd's algorithm for minimizing the K-means objective

Often simply called the K-means algorithm

Minimize cost function alternative optimization between $\{r_{nk}\}$ and $\{\mu_k\}$

- **Step 0** Initialize $\{\mu_k\}$ to some values
- **Step 1** Assume the current value of $\{\mu_k\}$ fixed, minimize J over $\{\underline{r_{nk}}\}$, which leads to the following cluster assignment rule

$$\underline{r_{nk}} = \begin{cases} \underline{1} & \text{if } k = \arg \min_j \underline{\|x_n - \mu_j\|_2^2} \\ 0 & \text{otherwise} \end{cases}$$

Lloyd's algorithm for minimizing the K-means objective

Often simply called the K-means algorithm

Minimize cost function alternative optimization between $\{r_{nk}\}$ and $\{\mu_k\}$

$$J(\{r_{nk}\}, \{\mu_k\}) = \sum_n \sum_k r_{nk} \|x_n - \mu_k\|_2^2$$

$\nabla_{\mu_k} J = 0$

- **Step 0** Initialize $\{\mu_k\}$ to some values
- **Step 1** Assume the current value of $\{\mu_k\}$ fixed, minimize J over $\{r_{nk}\}$, which leads to the following cluster assignment rule

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|x_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

- **Step 2** Assume the current value of $\{r_{nk}\}$ fixed, minimize J over $\{\mu_k\}$, which leads to the following rule to update the prototypes of the clusters

$$\mu_k = \frac{\sum_n r_{nk} x_n}{\sum_n r_{nk}} = \# \text{ pts assigned to } k$$

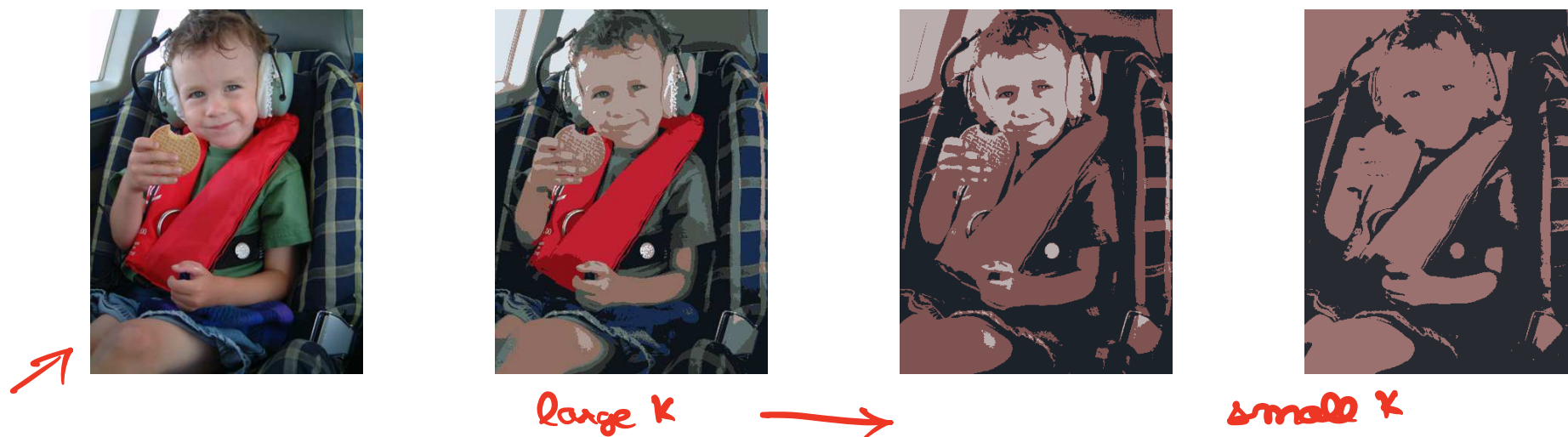
- **Step 3** Stop if the objective function J stays the same or return to Step 1

Remarks

- Prototype $\underline{\mu_k}$ is the mean of data points assigned to the cluster k , hence 'K-means'
- The procedure reduces \underline{J} in both Step 1 and Step 2 and thus makes improvements on each iteration

Application: vector quantization

- Replace data point with associated prototype $\underline{\mu_k}$
- In other words, compress the data points into i) a codebook of all the prototypes; ii) a list of indices to the codebook for the data points
- Lossy compression, especially for small K



Clustering pixels and vector quantizing them. From left to right: Original image, quantized with large K , medium K , and a small K . Details are missing due to the higher compression (smaller K).

Properties of the K-means algorithm

$$\underline{J}(\underline{\{x_{nk}\}}, \underline{\{\mu_k\}}) = \sum_n \sum_k \underbrace{r_{nk}}_{r_{nk} \in \{0,1\}} \underbrace{\|x_n - \mu_k\|_2^2}_{\geq 0} \geq 0$$

- Does the K-means algorithm converge (*i.e.*, terminate)?
 - ▶ Yes.
- How long does it take to converge ?
 - ▶ In the worst case, exponential in the number of data points.
 - ▶ In practice, very quick.

Properties of the K-means algorithm

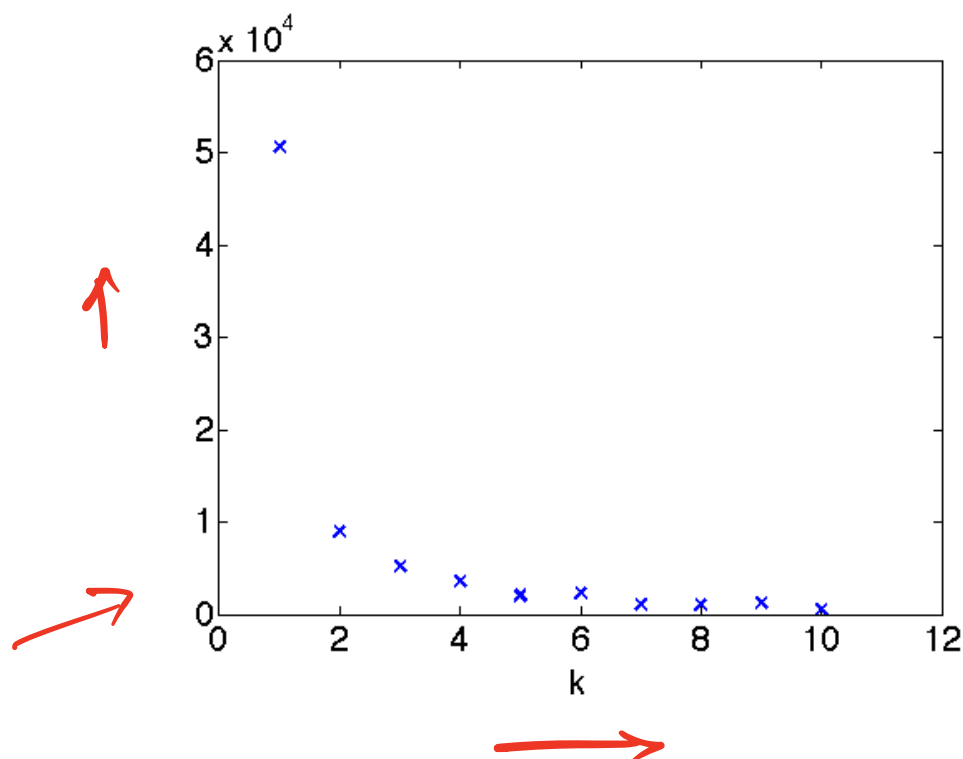
How good is the K-means solution?

- Converges to a local minimum.
- The solution depends on the initialization.
- In practice, run many times with different initializations and pick the best.
- K-means++ is a neat approximation algorithm that has theoretical guarantees on the final value of the objective.
 - ▶ Still no guarantee that you will reach the global minimum
 - ▶ You are guaranteed to get reasonably close (approximation guarantee on the final value).

Other practical issues

Choosing K

- Increasing K will always decrease the optimal value of the K-means objective.
 - ▶ Analogous to overfitting in supervised learning.
- Information criteria that effectively regularize more complex models.



K-medoids

- K-means is sensitive to outliers.
- In some applications we want the prototypes to be one of the points.
- Leads to K-medoids.

K-medoids



- **Step 0** Initialize $\{\underline{\mu_k}\}$ by randomly selecting K of the N points
- **Step 1** Assume the current value of $\{\mu_k\}$ fixed, assign points to clusters:

$$\underline{r_{nk}} = \begin{cases} 1 & \text{if } k = \arg \min_j \underline{\|\mathbf{x}_n - \mu_j\|_2^2} \\ 0 & \text{otherwise} \end{cases}$$

- **Step 2** Assume the current value of $\{r_{nk}\}$ fixed, update the prototype of cluster k . In K-medoids, the prototype for a cluster is the data point that is closest to all other data points in the cluster

$$\underline{k^*} = \arg \min_{m:r_{mk}=1} \sum_n \underline{r_{nk}} \underline{\|\mathbf{x}_n - \mathbf{x}_m\|_2^2}$$
$$\mu_k = \mathbf{x}_{k^*}$$

- **Step 3** Stop if the objective function J stays the same or return to Step 1

Outline

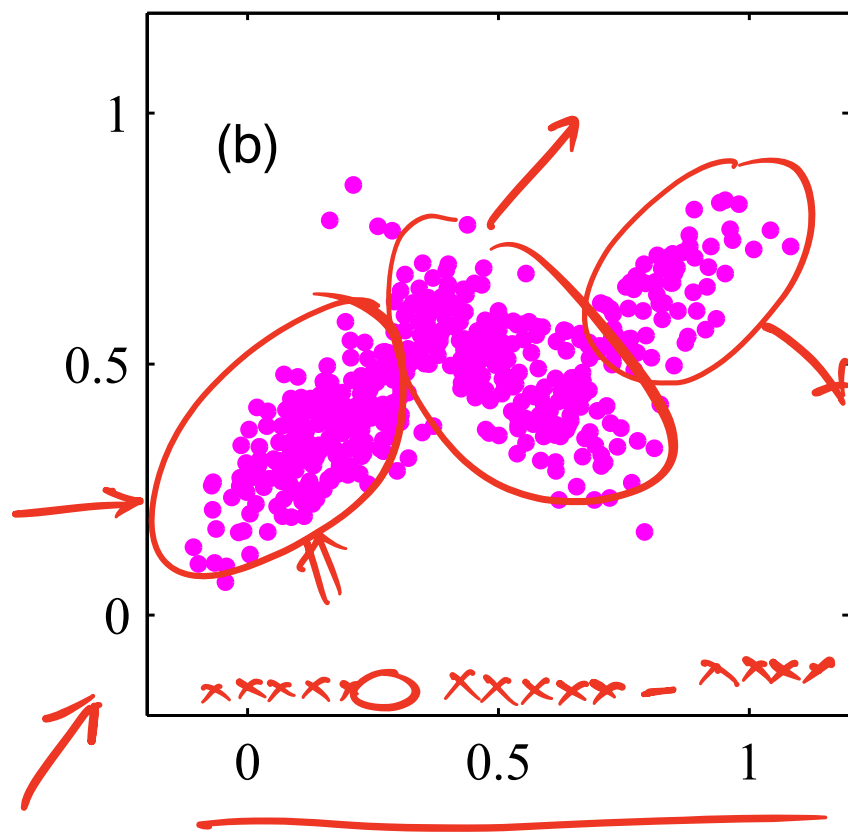
- 1 K-means
- 2 Gaussian mixture models
- 3 GMMs and Incomplete Data

RSS Linear regression
↓
Probabilistic interpretation

Probabilistic interpretation of clustering?

We can impose a probabilistic interpretation of our intuition that points stay close to their cluster centers

- How can we model $p(\mathbf{x})$ to reflect this?



- Data points seem to form 3 clusters

$$p(\mathbf{x}) \sim \mathcal{N}(\mu, \Sigma)$$

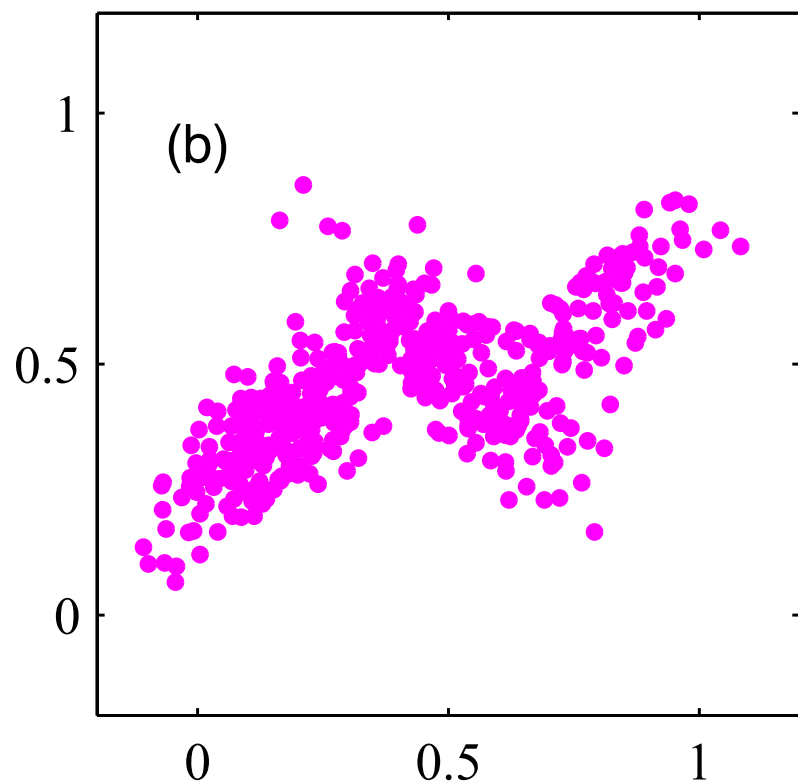
$$x \sim \mathcal{N}(\mu, \sigma^2)$$

A hand-drawn diagram of a normal distribution curve. The curve is drawn with a red line. Below the curve, a horizontal axis is shown with several red 'x' marks. A red circle labeled with the Greek letter mu is placed below the axis, with an arrow pointing to the peak of the curve.

Probabilistic interpretation of clustering?

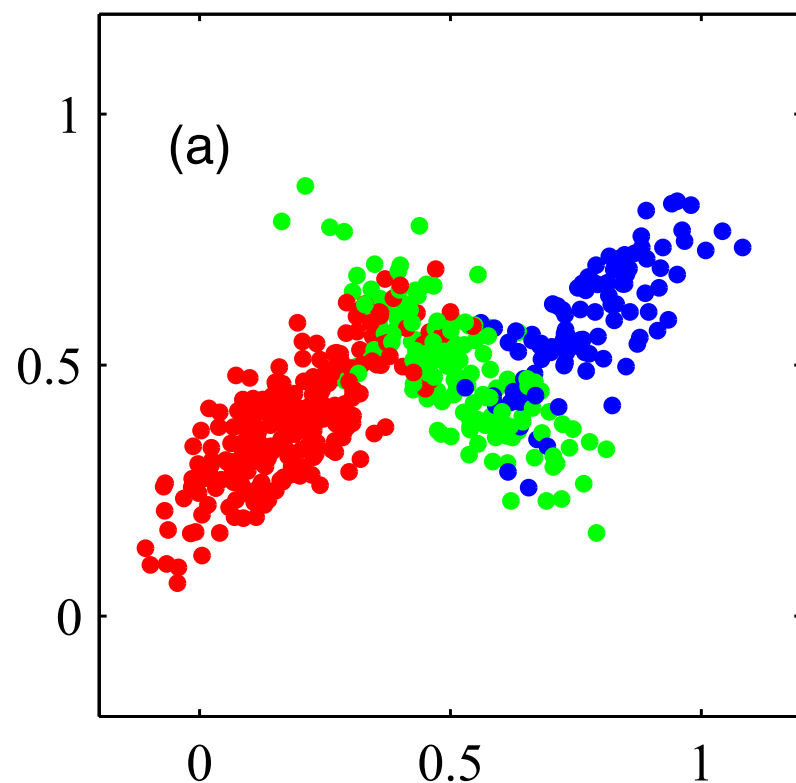
We can impose a probabilistic interpretation of our intuition that points stay close to their cluster centers

- How can we model $p(\mathbf{x})$ to reflect this?



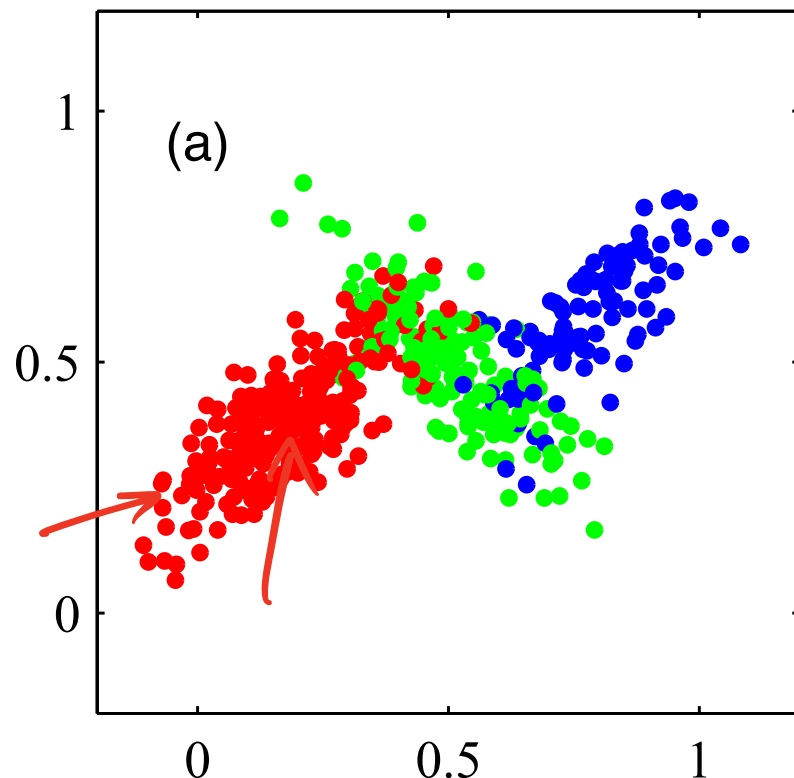
- Data points seem to form 3 clusters
- We cannot model $p(\mathbf{x})$ with simple and known distributions
- The data is not a Gaussian as we have 3 distinct concentrated regions

Gaussian mixture models: intuition



- We can model *each* region with a distinct distribution
- Common to use Gaussians, i.e., Gaussian mixture models (GMMs) or mixture of Gaussians (MoGs).

Gaussian mixture models: intuition



- We can model *each* region with a distinct distribution
- Common to use Gaussians, i.e., Gaussian mixture models (GMMs) or mixture of Gaussians (MoGs).
- We don't know *cluster assignments* (label) or *parameters* of Gaussians or *mixture components*!
- We need to learn them all from our *unlabeled* data

$$\mathcal{D} = \{\underline{\underline{x_n}}\}_{n=1}^N$$

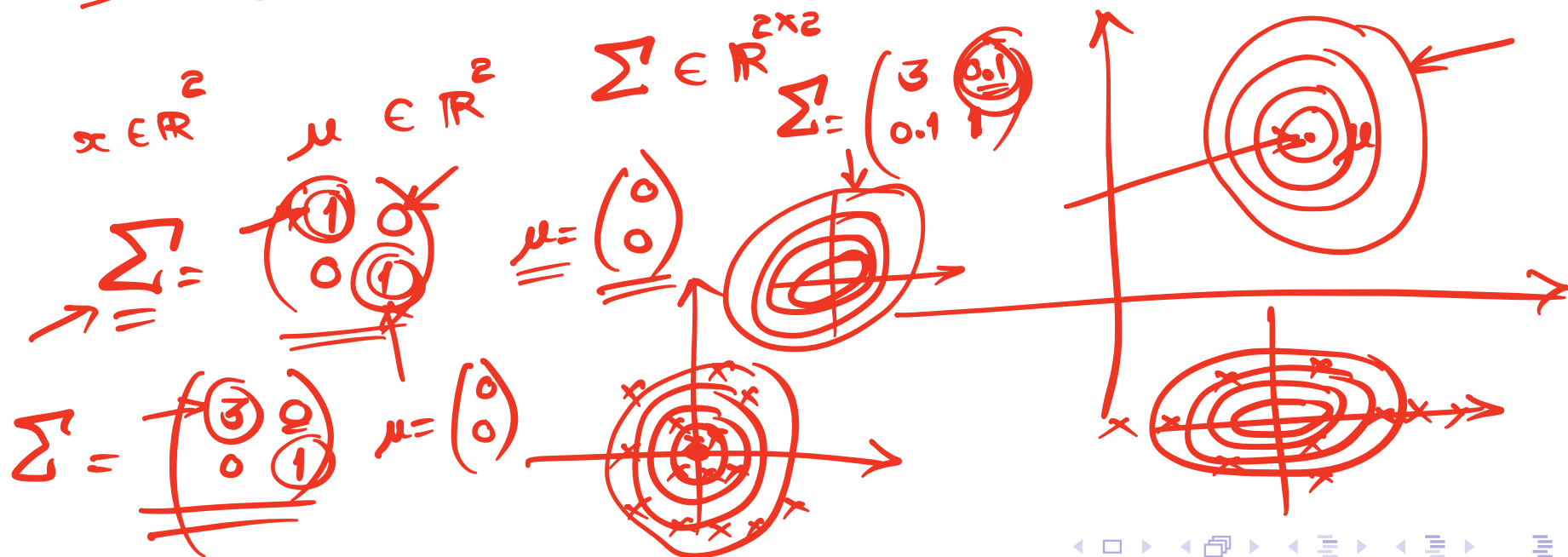
Gaussian mixture models: formal definition

A Gaussian mixture model has the following density function for x

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

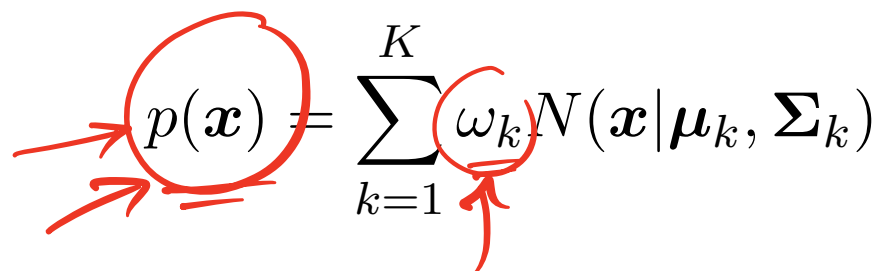
$\underline{x} \in \mathbb{R}^D$

- K : the number of Gaussians — they are called (mixture) components
- $\underline{\mu}_k$ and $\underline{\Sigma}_k$: mean and covariance matrix of the k -th component



Gaussian mixture models: formal definition

A Gaussian mixture model has the following density function for \mathbf{x}

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


- K : the number of Gaussians — they are called (mixture) components
- $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$: mean and covariance matrix of the k -th component
- ω_k : mixture weights — they represent how much each component contributes to the final distribution. It satisfies two properties:

$$\omega_k > 0$$

$$\sum_k \omega_k = 1$$

Gaussian mixture models: formal definition

A Gaussian mixture model has the following density function for \mathbf{x}

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

- K : the number of Gaussians — they are called (mixture) components
- $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$: mean and covariance matrix of the k -th component
- ω_k : mixture weights – they represent how much each component contributes to the final distribution. It satisfies two properties:

$$\forall k, \omega_k > 0, \quad \text{and} \quad \sum_k \omega_k = 1$$

The properties ensure $p(\mathbf{x})$ is a properly normalized probability density function.

GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$\rightarrow p(\underline{x}, z) = p(z)p(x|z)$$

where z is a discrete random variable taking values between 1 and K .

$$\underline{p(x, z)} = \underline{p(z)} \underset{\substack{\uparrow \\ \text{Condition}}}{\underline{p(x|z)}}$$

probability

GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$\Rightarrow p(\mathbf{x}, z) = p(z)p(\mathbf{x}|z)$$

where z is a discrete random variable taking values between 1 and K .

Denote

$$\omega_k = p(z = k)$$

Now, assume the conditional distributions are Gaussian distributions

$$p(\mathbf{x}|z = k) = N(\mathbf{x}|\mu_k, \Sigma_k)$$

$$\begin{aligned} p(\mathbf{x}) &= p(\mathbf{x}, z=1) + p(\mathbf{x}, z=2) + \dots + p(\mathbf{x}, z=K) \\ &= \sum_k p(\mathbf{x}, z=k) = \sum_k p(z=k) p(\mathbf{x}|z=k) \end{aligned}$$

GMM as the marginal distribution of a joint distribution

Consider the following joint distribution

$$p(\underline{\mathbf{x}}, \underline{z}) = p(z)p(\mathbf{x}|z)$$

where z is a discrete random variable taking values between 1 and K .
Denote

$$\omega_k = p(z = k)$$

Now, assume the conditional distributions are Gaussian distributions

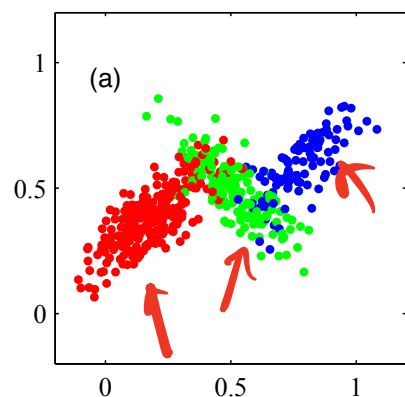
$$p(\mathbf{x}|z = k) = N(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Then, the marginal distribution of \mathbf{x} is

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

Namely, the Gaussian mixture model

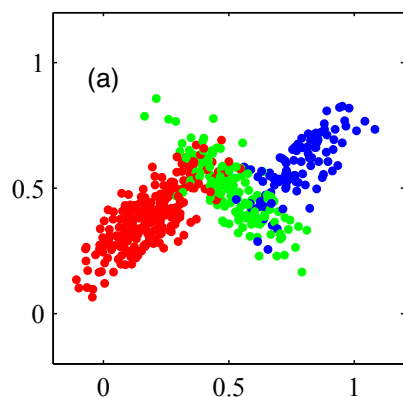
GMMs: example



The conditional distribution between \mathbf{x} and z (representing color) are

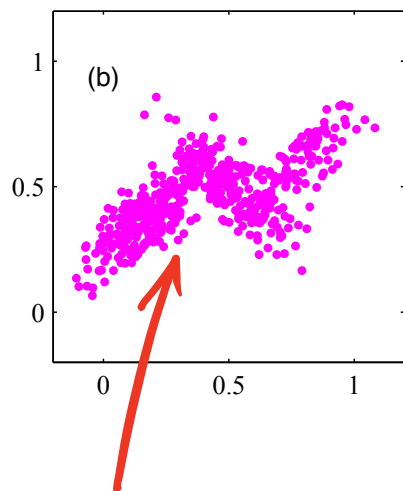
$$\left. \begin{aligned} p(\mathbf{x}|z = red) &= N(\mathbf{x}|\mu_1, \Sigma_1) \\ p(\mathbf{x}|z = blue) &= N(\mathbf{x}|\mu_2, \Sigma_2) \\ p(\mathbf{x}|z = green) &= N(\mathbf{x}|\mu_3, \Sigma_3) \end{aligned} \right\}$$

GMMs: example



The conditional distribution between \mathbf{x} and z (representing color) are

$$\left\{ \begin{aligned} \Rightarrow p(\mathbf{x}|z = \text{red}) &= N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ \Rightarrow p(\mathbf{x}|z = \text{blue}) &= N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \\ \Rightarrow p(\mathbf{x}|z = \text{green}) &= N(\mathbf{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3) \end{aligned} \right.$$



The marginal distribution is thus

$$\underline{p(\mathbf{x})} = \underline{p(\text{red})} \underline{N(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)} + \underline{p(\text{blue})} \underline{N(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} + \underline{p(\text{green})} \underline{N(\mathbf{x}|\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)}$$

Parameter estimation for Gaussian mixture models

The parameters in GMMs are $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$. To estimate, consider the simple (and unrealistic) case first.

We have labels z If we assume z is observed for every x , then our estimation problem is easier to solve. Our training data is augmented:

$$\mathcal{D}' = \{\underline{x}_n, z_n\}_{n=1}^N$$

z_n denotes the component where x_n comes from. \mathcal{D}' is the *complete* data and \mathcal{D} the *incomplete* data. How can we learn our parameters?

Parameter estimation for Gaussian mixture models

The parameters in GMMs are $\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$. To estimate, consider the simple (and unrealistic) case first.

We have labels z If we assume z is observed for every x , then our estimation problem is easier to solve. Our training data is augmented:

$$\mathcal{D}' = \{\mathbf{x}_n, z_n\}_{n=1}^N$$

z_n denotes the component where \mathbf{x}_n comes from. \mathcal{D}' is the *complete* data and \mathcal{D} the *incomplete* data. How can we learn our parameters?

Given \mathcal{D}' , the maximum likelihood estimation of the θ is given by

$$\theta = \arg \max \log P(\mathcal{D}') = \sum_n \log p(\mathbf{x}_n, z_n)$$

Parameter estimation for GMMs: complete data

The complete likelihood is decomposable

$$\sum_n \log p(\underline{x}_n, z_n) = \sum_n \log p(z_n) p(\underline{x}_n | z_n) = \sum_k \sum_{n: z_n = k} \log p(z_n) p(\underline{x}_n | z_n)$$

where we have grouped data by its values z_n . Let us introduce a binary variable $\gamma_{nk} \in \{0, 1\}$ to indicate whether $z_n = k$. We then have

Parameter estimation for GMMs: complete data

The *complete* likelihood is decomposable

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_n \log p(z_n) p(\mathbf{x}_n | z_n) = \sum_k \sum_{n: z_n=k} \log p(z_n) p(\mathbf{x}_n | z_n)$$

where we have grouped data by its values z_n . Let us introduce a binary variable $\gamma_{nk} \in \{0, 1\}$ to indicate whether $z_n = k$. We then have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log p(z = k) p(\mathbf{x}_n | z = k)$$

We use a “dummy” variable z to denote all the possible values cluster assignment values for \mathbf{x}_n

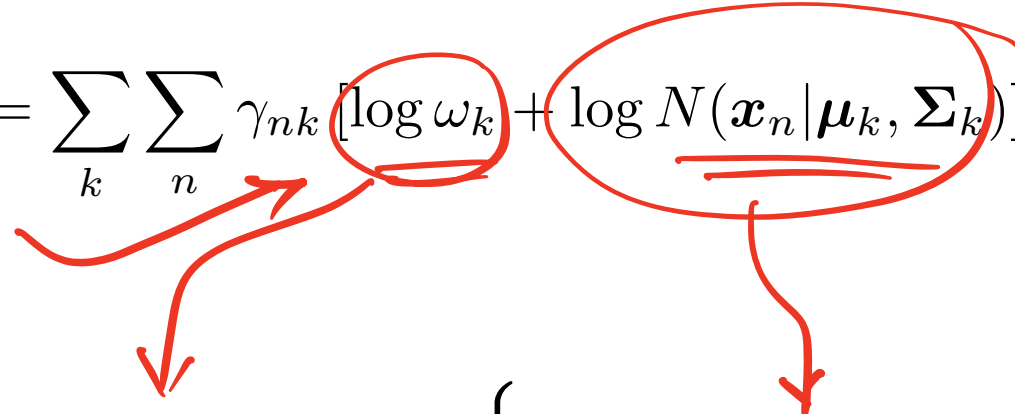
$$\sum_k \sum_n \gamma_{nk} \log p(z=k) p(\mathbf{x}_n | z=k)$$

\mathcal{D}' specifies this value in the complete data setting

$$\log(\omega_k \mathcal{N}(\mathbf{x}_n; \mu_k, \Sigma_k)) = (\log(\omega_k) + \log \mathcal{N}(\mathbf{x}_n; \mu_k, \Sigma_k))$$

Parameter estimation for GMMs: complete data

We now have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(\mathbf{x}_n | \underline{\underline{\mu_k}}, \underline{\underline{\Sigma_k}})]$$


Regrouping, we have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \underline{\underline{\omega_k}} + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \underline{\underline{\mu_k}}, \underline{\underline{\Sigma_k}}) \right\}$$

Parameter estimation for GMMs: complete data

We now have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

Regrouping, we have

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

The term inside the braces depends on k -th component's parameters. It can be shown that the MLE is:

$\gamma_{nk} \in \{0,1\}$

$$\omega_k = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad \boldsymbol{\mu}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$

$N =$

$$\boldsymbol{\Sigma}_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

What's the intuition?

Intuition

Since γ_{nk} is binary, the previous solution is nothing but

- For ω_k : count the number of data points whose $\underline{z_n}$ is k and divide by the total number of data points (note that $\sum_k \sum_n \gamma_{nk} = N$)
- For μ_k : get all the data points whose z_n is k , compute their mean
- For Σ_k : get all the data points whose z_n is k , compute their covariance matrix

This intuition is going to help us to develop an algorithm for estimating θ when we do not know z_n (incomplete data).

Parameter estimation for GMMs: incomplete data

$$p(z_n = k | \underline{x_n})$$

When z_n is not given, we can guess it via the posterior probability

$$p(z_n = k | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | z_n = k) p(z_n = k)}{p(\mathbf{x}_n)} = \frac{p(\mathbf{x}_n | z_n = k) p(z_n = k)}{\sum_{k'=1}^K p(\mathbf{x}_n | z_n = k') p(z_n = k')}$$

$$\underbrace{p(y|x)}_{\substack{\uparrow \\ \text{Random} \\ \text{variable}}} = \frac{\underbrace{p(x|y)}_{\substack{\uparrow \\ \text{Random} \\ \text{variable}}} \underbrace{p(y)}}{\underbrace{p(x)}} \quad (\text{Bayes theorem})$$

$$\begin{aligned} p(z_n = k | \mathbf{x}_n) &= \frac{p(\underline{x_n} | \underline{z_n = k}) \underline{p(z_n = k)}}{\underline{p(\mathbf{x}_n)}} \\ &= \mathcal{N}(\mathbf{x}_n; \mu_k, \Sigma_k) \omega_k \end{aligned}$$

Parameter estimation for GMMs: incomplete data

When z_n is not given, we can guess it via the posterior probability

$$p(z_n = k | \mathbf{x}_n) = \frac{p(\mathbf{x}_n | z_n = k)p(z_n = k)}{p(\mathbf{x}_n)} = \frac{p(\mathbf{x}_n | z_n = k)p(z_n = k)}{\sum_{k'=1}^K p(\mathbf{x}_n | z_n = k')p(z_n = k')}$$

To compute the posterior probability, we need to know the parameters θ !

Let's pretend we know the value of the parameters so we can compute the posterior probability.

How is that going to help us?

Estimation with soft γ_{nk}

We define γ_{nk} = $p(z_n = k | \mathbf{x}_n)$

Estimation with soft γ_{nk}

We define $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$

- Recall that γ_{nk} was binary.
- Now it's a “soft” assignment of \mathbf{x}_n to k -th component
- Each \mathbf{x}_n is assigned to a component fractionally according to $p(z_n = k | \mathbf{x}_n)$

Parameter estimation for GMMs: incomplete data

With the soft assignment γ_{nk} plugged into the complete data log likelihood, we now have:

$$\sum_k \sum_n \gamma_{nk} [\log \omega_k + \log N(\mathbf{x}_n | \underline{\boldsymbol{\mu}_k}, \underline{\boldsymbol{\Sigma}_k})]$$

Regrouping, we have

$$\sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Parameter estimation for GMMs: incomplete data

With the soft assignment γ_{nk} plugged into the complete data log likelihood, we now have:

$$\sum_k \sum_n \underline{\gamma_{nk}} [\log \omega_k + \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]$$

Regrouping, we have

$$\sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

We now get the same expression for the MLE as before!

$$\underline{\omega_k} = \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, \quad \underline{\boldsymbol{\mu}_k} = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n$$
$$\underline{\boldsymbol{\Sigma}_k} = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

But remember, we're 'cheating' by using $\boldsymbol{\theta}$ to compute γ_{nk} !

Iterative procedure

$$\theta = \underbrace{(\omega_k, \mu_k, \Sigma_k)_{k=1}^K}$$

We can alternate between estimating γ_{nk} and using the estimated γ_{nk} to compute the parameters (same idea as with K -means!)

- Step 0: initialize θ with some values (random or otherwise)
- Step 1: compute γ_{nk} using the current θ
- Step 2: update θ using the just computed γ_{nk}
- Step 3: go back to Step 1

Questions:

- Is this procedure reasonable, i.e., are we optimizing a sensible criteria?
- Will this procedure converge?

The answers lie in the *EM algorithm* — a powerful procedure for model estimation with unknown data.

Outline

- 1 K-means
- 2 Gaussian mixture models
- 3 GMMs and Incomplete Data**

Parameter estimation for GMMs: complete data

GMM Parameters

$$\theta = \{\omega_k, \mu_k, \Sigma_k\}_{k=1}^K$$

Complete Data: We (unrealistically) assume z is observed for every x ,

$$\mathcal{D}' = \{\underline{x}_n, \underline{z}_n\}_{n=1}^N$$

MLE: Maximize the complete likelihood

$$\theta = \arg \max \log P(\mathcal{D}') = \sum_n \log \underline{p(\underline{x}_n, \underline{z}_n)}$$

Leads to closed-form solution!

Parameter estimation for GMMs: Incomplete data

GMM Parameters

$$\boldsymbol{\theta} = \{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$$

Incomplete Data

Our data contains observed and unobserved random variables, and hence is incomplete

- Observed: $\mathcal{D} = \{\mathbf{x}_n\}$
- Unobserved (hidden): $\{\mathbf{z}_n\}$

Goal Obtain the maximum likelihood estimate of $\boldsymbol{\theta}$:

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log P(\mathcal{D}) = \arg \max_{\boldsymbol{\theta}} \sum_n \log p(\mathbf{x}_n | \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_n \log \sum_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta})\end{aligned}$$

The objective function $\ell(\boldsymbol{\theta})$ is called the *incomplete* log-likelihood.

Issue with Incomplete log-likelihood

No simple way to optimize the incomplete log-likelihood

Expectation-Maximization (EM) algorithm provides a strategy for iteratively optimizing this function

Two steps as they apply to GMM:

- E-step: 'guess' values of the z_n using existing values of θ
- M-step: solve for new values of θ given imputed values for z_n
(maximize complete likelihood!)

E-step: Soft cluster assignments

We define γ_{nk} as $p(z_n = k | \mathbf{x}_n, \boldsymbol{\theta})$

- This is the posterior distribution of z_n given \mathbf{x}_n and $\boldsymbol{\theta}$
- Recall that in complete data setting γ_{nk} was binary
- Now it's a “soft” assignment of \mathbf{x}_n to k -th component, with \mathbf{x}_n assigned to each component with some probability

E-step: Soft cluster assignments

$$0 \leq \gamma_{nk} \leq 1$$

We define γ_{nk} as $p(z_n = k | \mathbf{x}_n, \boldsymbol{\theta})$

- This is the posterior distribution of z_n given \mathbf{x}_n and $\boldsymbol{\theta}$
- Recall that in complete data setting γ_{nk} was binary
- Now it's a “soft” assignment of \mathbf{x}_n to k -th component, with \mathbf{x}_n assigned to each component with some probability

Given $\boldsymbol{\theta} = \{\omega_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$, we can compute γ_{nk} using **Bayes theorem**:

$$\begin{aligned}\gamma_{nk} &= p(z_n = k | \mathbf{x}_n) \\ &= \frac{p(\mathbf{x}_n | z_n = k) p(z_n = k)}{p(\mathbf{x}_n)} \\ &= \frac{p(\mathbf{x}_n | z_n = k) p(z_n = k)}{\sum_{k'=1}^K p(\mathbf{x}_n | z_n = k') p(z_n = k')} = \frac{\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \omega_k}{\sum_{k'=1}^K \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'}) \omega_{k'}}\end{aligned}$$

M-step: Maximize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_n \log p(\underline{\mathbf{x}_n}, \underline{z_n}) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Previously γ_{nk} was binary, but now we define $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$ (E-step)

M-step: Maximize complete likelihood

Recall definition of complete likelihood from earlier:

$$\sum_n \log p(\mathbf{x}_n, z_n) = \sum_k \sum_n \gamma_{nk} \log \omega_k + \sum_k \left\{ \sum_n \gamma_{nk} \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right\}$$

Previously γ_{nk} was binary, but now we define $\gamma_{nk} = p(z_n = k | \mathbf{x}_n)$ (E-step)

We get the same simple expression for the MLE as before!

$$\left\{ \begin{aligned} \omega_k &= \frac{\sum_n \gamma_{nk}}{\sum_k \sum_n \gamma_{nk}}, & \boldsymbol{\mu}_k &= \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} \mathbf{x}_n \\ \boldsymbol{\Sigma}_k &= \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \end{aligned} \right.$$

Intuition: Each point now contributes some fractional component to each of the parameters, with weights determined by γ_{nk}

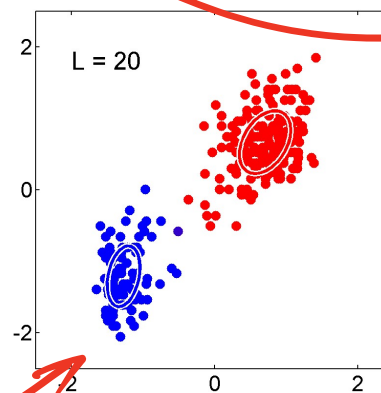
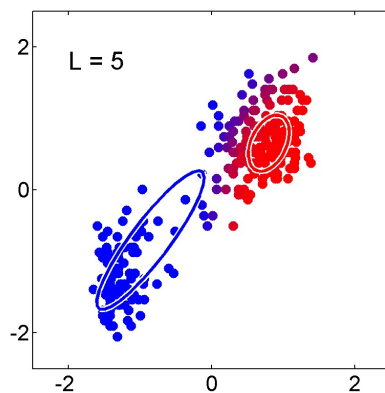
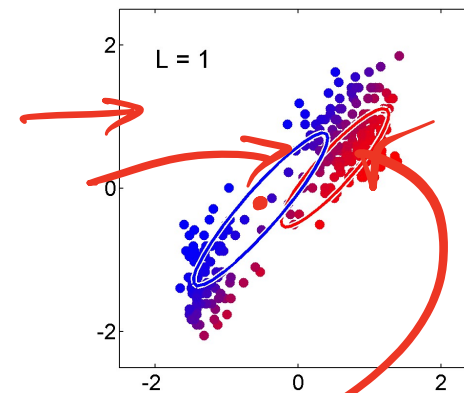
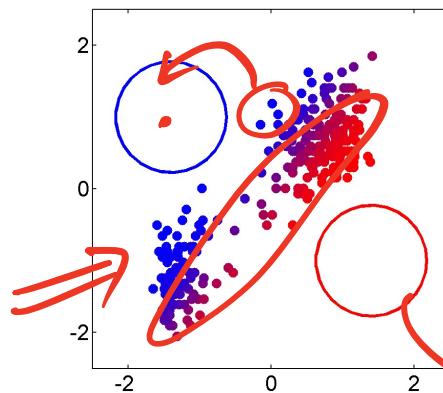
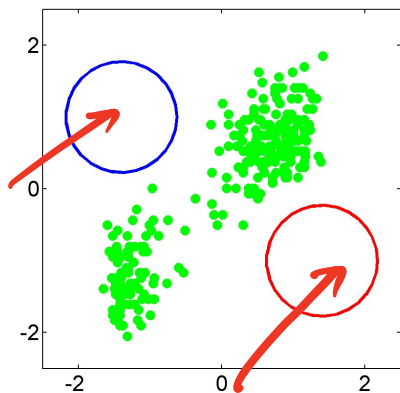
EM procedure for GMM

Alternate between estimating γ_{nk} and estimating θ

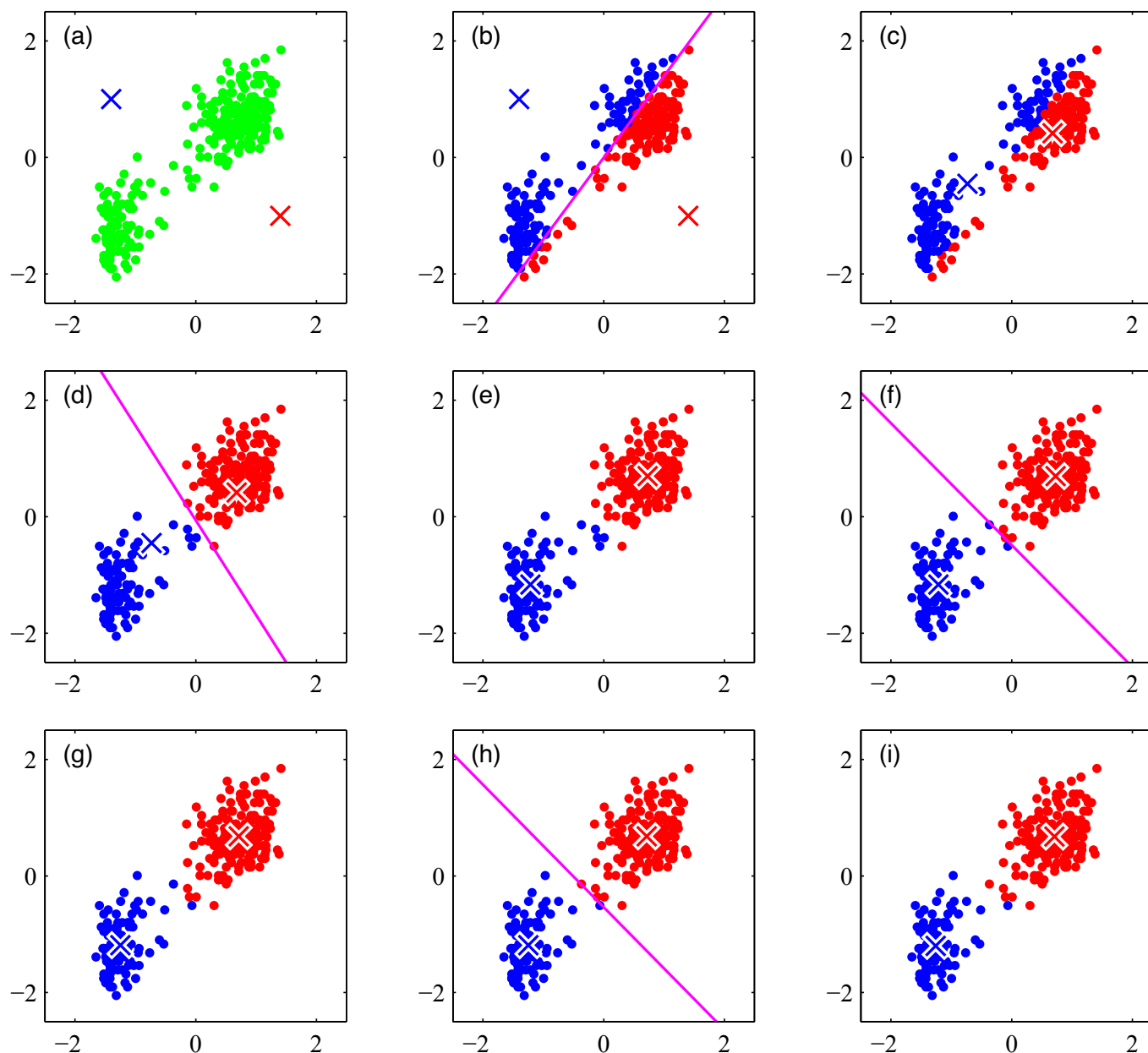
- Initialize θ with some values (random or otherwise)
- Repeat
 - ▶ E-Step: Compute γ_{nk} using the current θ
 - ▶ M-Step: Update θ using the γ_{nk} we just computed
- Until Convergence

Example of GMM

(!)



Compare to K-means example



EM procedure for GMM

Questions to be answered next

- How does GMM relate to K -means?
- Is this procedure reasonable, i.e., are we optimizing a sensible criterion?
- Will this procedure converge?

GMMs and K-means

GMMs provide probabilistic interpretation for K-means

GMMs reduce to K-means under the following assumptions (in which case EM for GMM parameter estimation simplifies to K-means):

- Assume all mixture weights ω_k are equal
- Assume all Gaussians have $\sigma^2 \mathbf{I}$ covariance matrices
- Further assume $\sigma \rightarrow 0$, so we only need to estimate $\underline{\mu_k}$, i.e., means
- GMMs are more general model.



K-means is often called “hard” GMM or GMMs is called “soft” K-means

The posterior γ_{nk} provides a probabilistic assignment for x_n to cluster k

EM algorithm

- The estimates of the parameter θ in each iteration increase the likelihood.
- EM algorithm converges but only to a local optimum.

Summary

Clustering

- Group similar instances
- K-means
 - ▶ Minimize a cost function that measures the sum of squared distances from the cluster prototypes.
 - ▶ Iterative algorithm for minimizing the cost function.
- Variants: K-medoids
- Probabilistic interpretation of K-means: Gaussian Mixture Model
- Can define a number of mixture models for other kinds of data.
- Probabilistic interpretation: GMMs
 - ▶ Generalization of K-means
 - ▶ Estimation using an iterative EM algorithm.