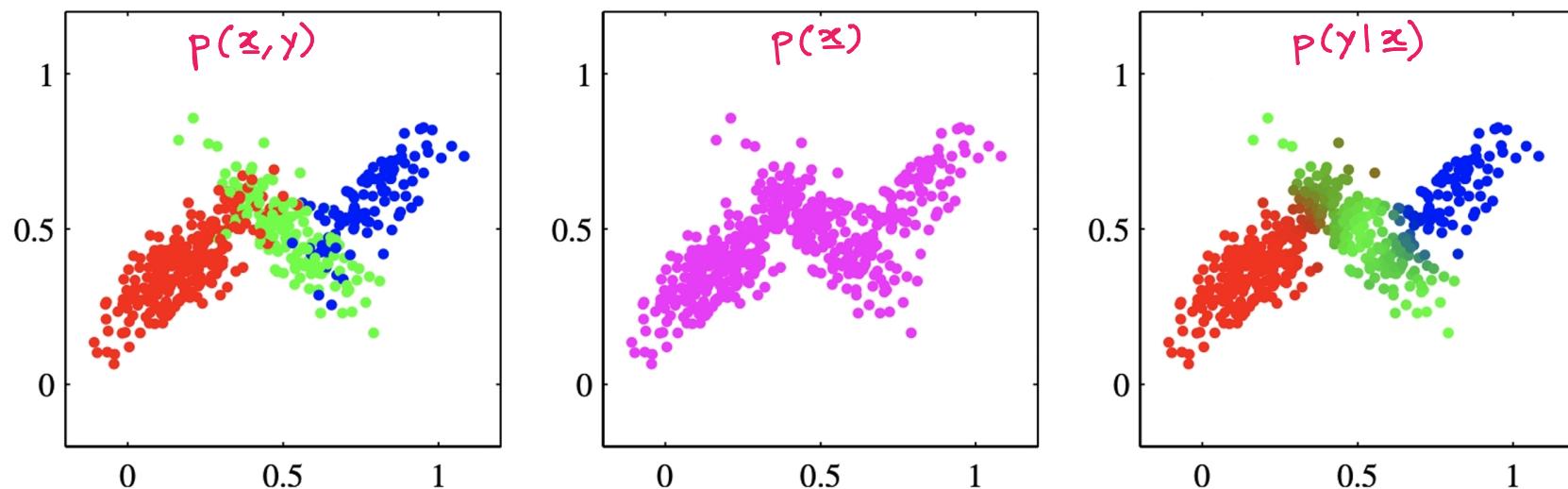


Recap of Generative Models

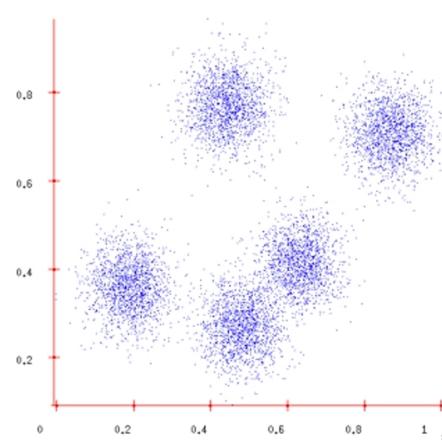
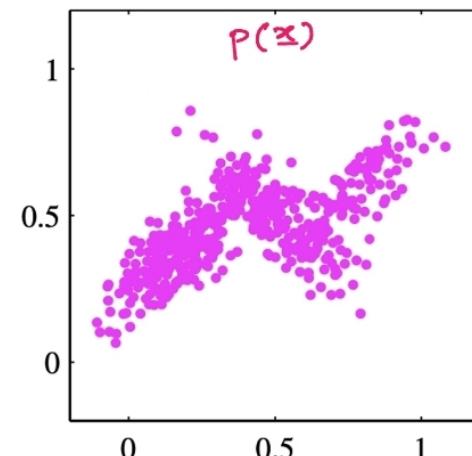
- In the last lecture, we introduced the concept of generative models that models the joint distribution $p(\mathbf{z}, \mathbf{y})$
- We used GMM to model the joint distribution $p(\mathbf{z}, \mathbf{y})$
- In supervised learning, the GMM was learned using both input-output data



E.g. Mixture of 3 Gaussians

Unsupervised Learning: Cluster Analysis

- In unsupervised learning, all the datapoints are unlabelled
 - Training data, $T = \{\underline{x}_i\}_{i=1}^N$
 - Objective: Build a model that can be used to find interesting patterns in data
In other words, build a model of the distribution $p(\underline{x})$
- Clustering → finding groups of similar \underline{x} values in the data space



GMM for clustering

- The GMM introduced in supervised learning is a joint model for \underline{x} & y

$$p(\underline{x}, y) = p(\underline{x}|y) p(y) = N(\underline{x} | \underline{\mu}_y, \underline{\Sigma}_y) \pi_y$$

- To obtain a model only for input \underline{x} , use marginalization

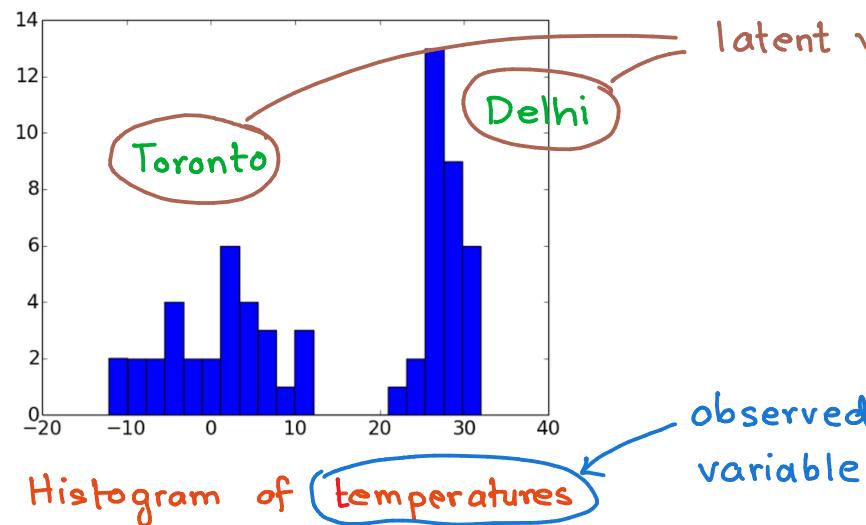
$$p(\underline{x}) = \sum_y p(\underline{x}, y)$$

summation because y is a discrete-valued class

- Since the output y is not observed, we will treat it as a LATENT variable

latent r.v. \leftrightarrow exists in the model but not observed in data

Example:



latent variables, if you are only given the temperature values and the city

observed variable

Maximum Likelihood with latent variables

- How should we choose the parameters $\underline{\Theta} = \{\pi_m, \underline{\mu}_m, \underline{\Sigma}_m\}_{m=1}^M$
- We will use Maximum likelihood principle \rightarrow Choose parameters that maximize the likelihood of observed data
- Note that, we don't observe the cluster labels y , we only get the input \underline{x}

- Given data $T = \{\underline{x}_i\}_{i=1}^N$, choose parameters to maximize:

$$\hat{\underline{\Theta}} = \arg \max_{\underline{\Theta}} \ln p(\{\underline{x}_i\}_{i=1}^N | \underline{\Theta})$$

$$= \arg \max_{\underline{\Theta}} \sum_{i=1}^N \ln p(\underline{x}_i | \underline{\Theta})$$

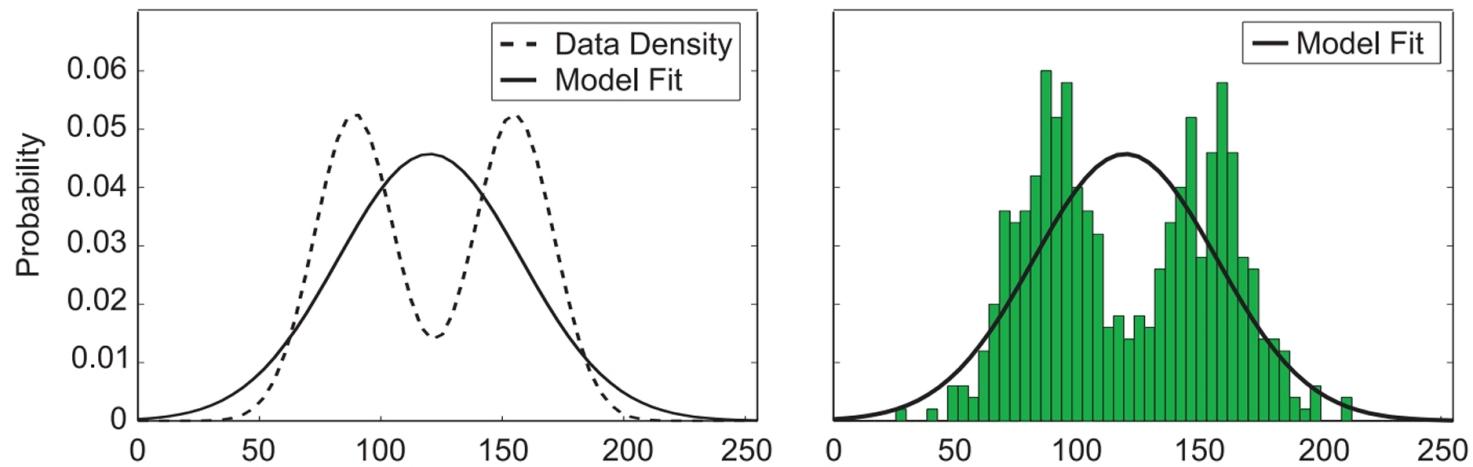
[Assuming independence
of inputs]

- We can find $p(\underline{x})$ by marginalizing out y :

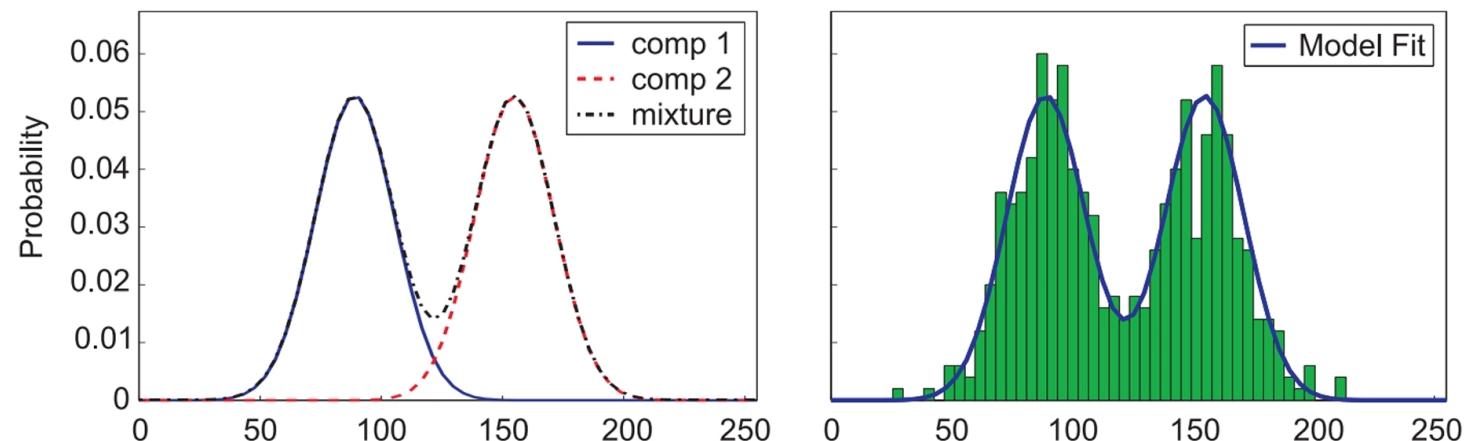
$$p(\underline{x}) = \sum_{m=1}^M p(\underline{x}, y=m) = \sum_{m=1}^M p(y=m) p(\underline{x}|y=m)$$

Visualizing a mixture of Gaussians (1D - Gaussians)

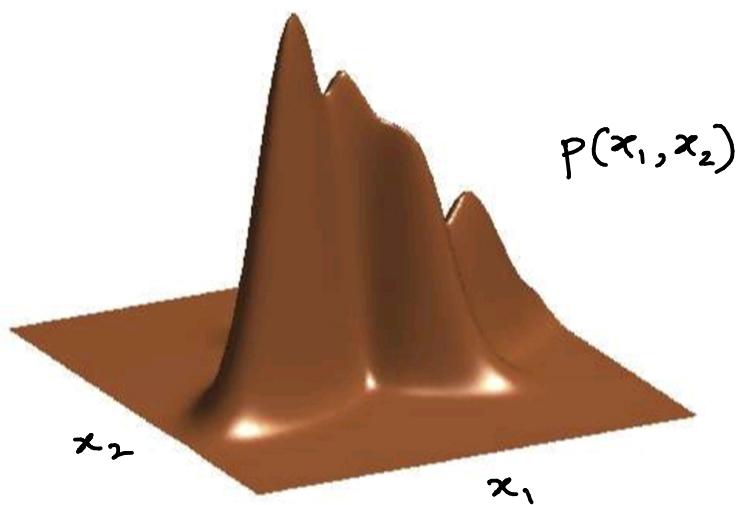
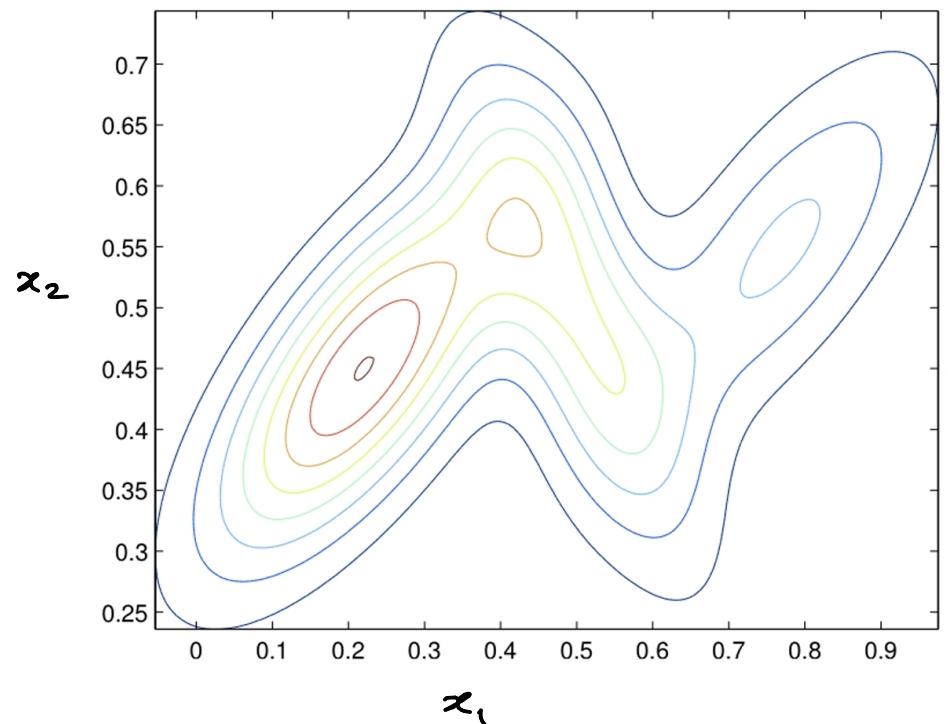
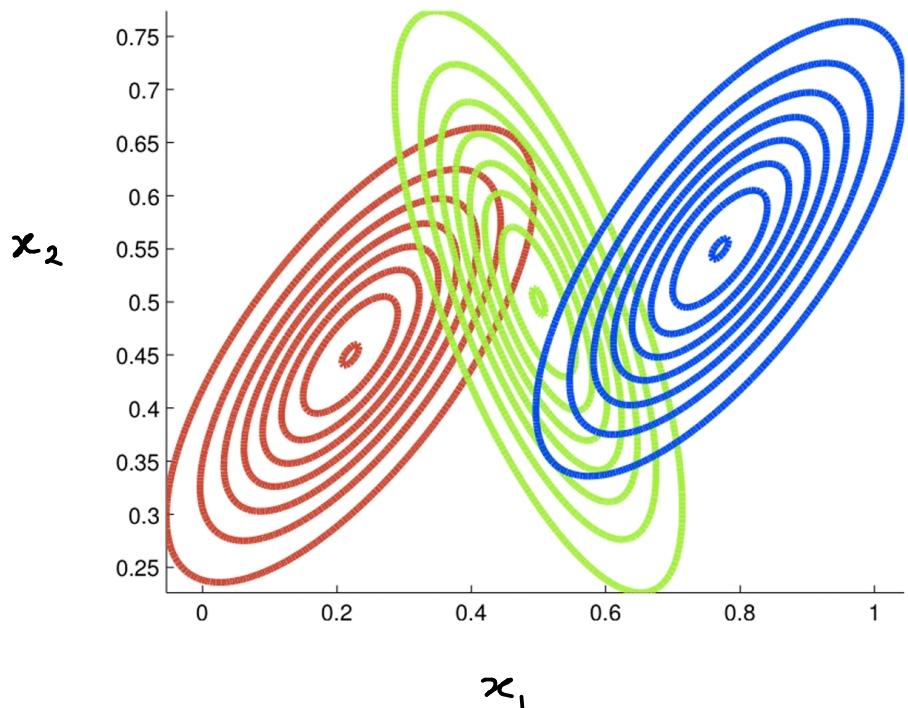
- If you fit a single Gaussian to bimodal data:



- If you fit a GMM with $M=2$ (or 2 Gaussians)



Visualizing a mixture of Gaussians (2D - Gaussians)



Fitting GMMs for unsupervised learning

- Notation: $\underline{\Theta} = \{\pi_m, \underline{M}_m, \underline{\Sigma}_m\}_{m=1}^M$, $\underline{x} = \{x_i\}_{i=1}^N$, $\underline{y} = \{y_i\}_{i=1}^N$
- The maximum likelihood objective:
$$\operatorname{argmax}_{\underline{\Theta}} \ln p(\underline{x} | \underline{\Theta}) = \operatorname{argmax}_{\underline{\Theta}} \sum_{i=1}^N \ln \left(\sum_{m=1}^M \pi_m N(x_i | \underline{M}_m, \underline{\Sigma}_m) \right)$$
- In general, there is no closed-form solution of the above optimization unlike the supervised learning case
- Use numerical iterative optimization
 - Gradient descent will face challenges:
 - * Non-convex
 - * Need to enforce non-negativity constraint on π_m and PSD constraint on Σ_m
 - * Derivatives of Σ_m are expensive

A Different Approach for optimization

- Estimating (or inferring) the latent variable \underline{y} with known $\underline{\Theta}$
- Estimating the GMM parameters $\underline{\Theta}$ with known latent variable \underline{y}

Alternating steps

Estimating (or inferring) the latent variable \underline{y} with known $\underline{\Theta}$

- If we knew the parameters $\underline{\Theta} = \{\pi_m, \underline{m}_m, \Sigma_m\}_{m=1}^M$, we could infer which component a data point \underline{x}_i probably belongs to by inferring y_i (i.e Gaussian component)

• This is just finding: $P(y_i = m | \underline{x}_i)$

$$P(y_i = m | \underline{x}_i) = \frac{P(y_i = m) \pi_m P(\underline{x}_i | y_i = m)}{\sum_{j=1}^M P(y_i = j) P(\underline{x}_i | y_i = j)}$$

A Different Approach for optimization

- Estimating (or inferring) the latent variable \underline{y} with known $\underline{\Theta}$
- Estimating the GMM parameters $\underline{\Theta}$ with known latent variable \underline{y}

Alternating steps

Estimating the GMM parameters $\underline{\Theta}$ with known latent variable \underline{y}

- If somehow we knew the latent variable \underline{y} , we could simply maximize the likelihood of the joint distribution:

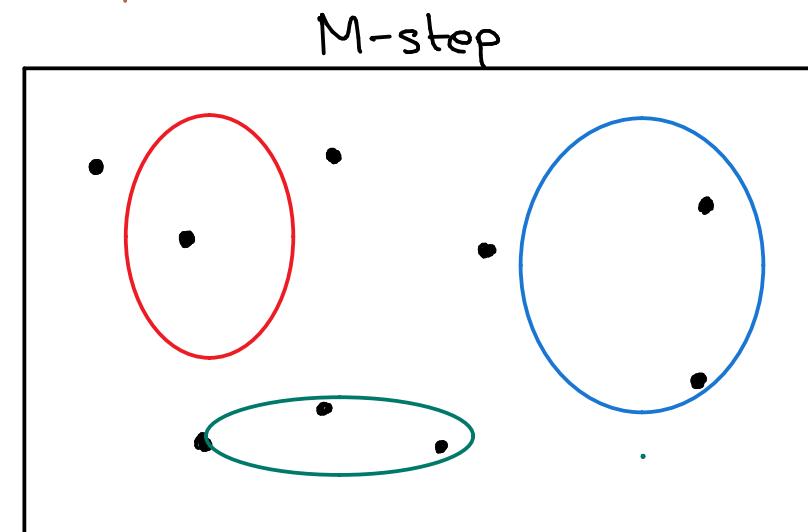
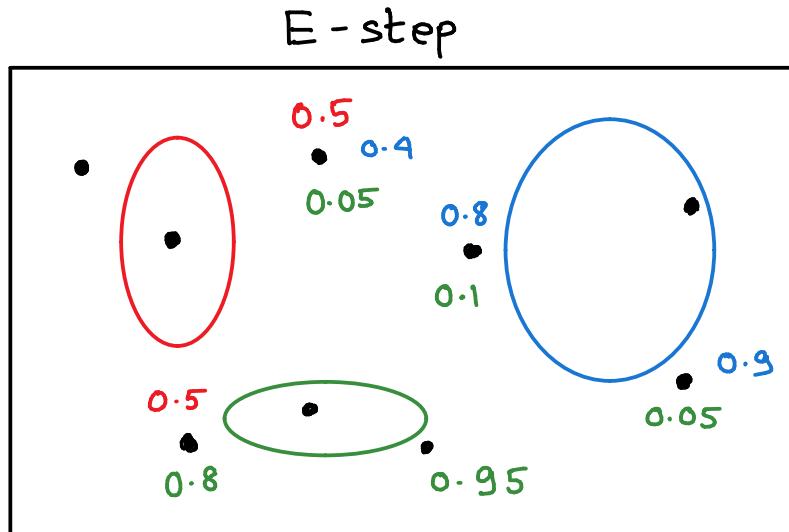
$$\ln p(\{\underline{x}_i, y_i\}_{i=1}^N | \underline{\Theta}) = \sum_{i=1}^N \sum_{m=1}^M \mathbb{I}\{y_i = m\} \left\{ \ln N(\underline{x}_i | \underline{\mu}_m, \underline{\Sigma}_m) + \ln p(y_i | \underline{\Theta}) \right\}$$

- Optimized parameters (from closed-form):

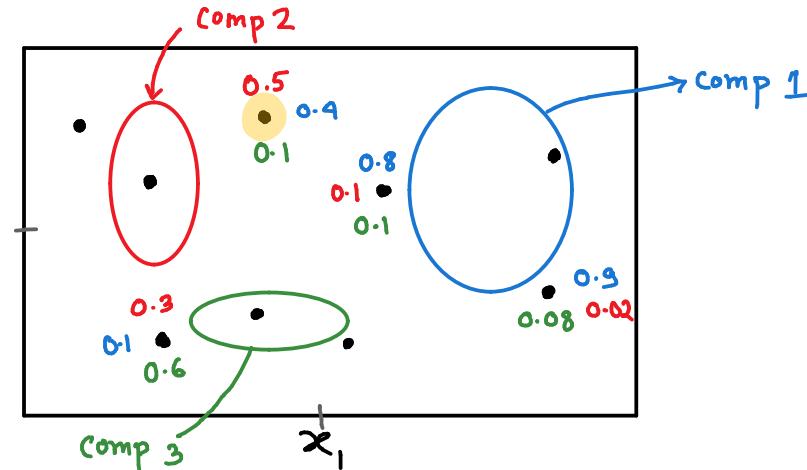
$$\hat{\pi}_m = \frac{n_m}{N}, \quad \hat{\mu}_m = \frac{1}{n_m} \sum_{i: y_i = m} \underline{x}_i, \quad \hat{\Sigma}_m = \frac{1}{n_m} \sum_{i: y_i = m} (\underline{x}_i - \hat{\mu}_m)(\underline{x}_i - \hat{\mu}_m)^T$$

Expectation - Maximization (EM) algorithm

- We will use the EM algorithm that alternates between the two steps:
 - Expectation step (E-step): Compute the probability over \underline{y} given we know our current model parameters Θ
i.e. how much do we think each Gaussian generates each datapoint
 - Maximization step (M-step): Assuming that the complete data was generated this way, tune the parameters of each Gaussian to maximize the likelihood of the complete data



E-step



- Assign the responsibility $w_m^{(i)}$ of component m for each data point i using probability

$$w_m^{(i)} = p(y_i = m \mid \underline{x}_i, \underline{\theta})$$

Example: Suppose we observe $\underline{x}_i = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -0.1 \\ 0.5 \end{bmatrix}$

$$p(y=1) = 0.4 \quad \pi_1$$

$$p(\underline{x}_i | y=1) = \mathcal{N}\left(\begin{bmatrix} -0.1 \\ 0.5 \end{bmatrix}; \begin{bmatrix} 0.8 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.5 & 0 \\ 0 & 0.6 \end{bmatrix}\right) \quad \underline{\mu}_1, \quad \underline{\Sigma}_1$$

$$p(y=2) = 0.35 \quad \pi_2$$

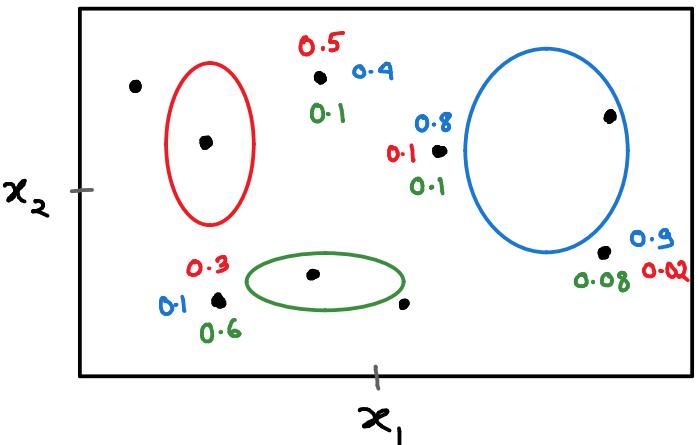
$$p(\underline{x}_i | y=2) = \mathcal{N}\left(\begin{bmatrix} -0.1 \\ 0.5 \end{bmatrix}; \begin{bmatrix} -0.5 \\ 0.2 \end{bmatrix}, \begin{bmatrix} 0.3 & 0 \\ 0 & 0.4 \end{bmatrix}\right) \quad \underline{\mu}_2, \quad \underline{\Sigma}_2$$

$$p(y=3) = 0.25 \quad \pi_3$$

$$p(\underline{x}_i | y=3) = \mathcal{N}\left(\begin{bmatrix} -0.1 \\ 0.5 \end{bmatrix}; \begin{bmatrix} -0.1 \\ -0.3 \end{bmatrix}, \begin{bmatrix} 0.7 & 0 \\ 0 & 0.2 \end{bmatrix}\right) \quad \underline{\mu}_3, \quad \underline{\Sigma}_3$$

$$\omega_1^{(i)} = p(y=1 | \underline{x}_i) = \frac{p(\underline{x}_i | y=1) p(y=1)}{\sum_{m=1}^3 p(\underline{x}_i | y=m) p(y=m)} = 0.4$$

E-step

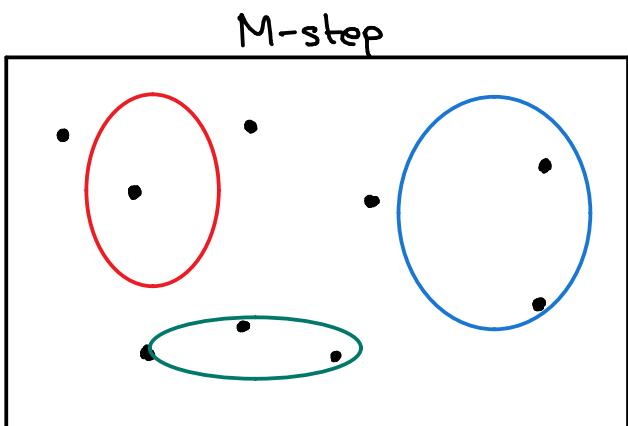


- Assign the **responsibility** $w_m^{(i)}$ of component m for each data point i using probability

$$w_m^{(i)} = p(y_i = m \mid \underline{x}_i, \underline{\theta})$$

M-step

- Apply maximum likelihood updates, where the parameters of each Gaussian component is fit with a weighted dataset

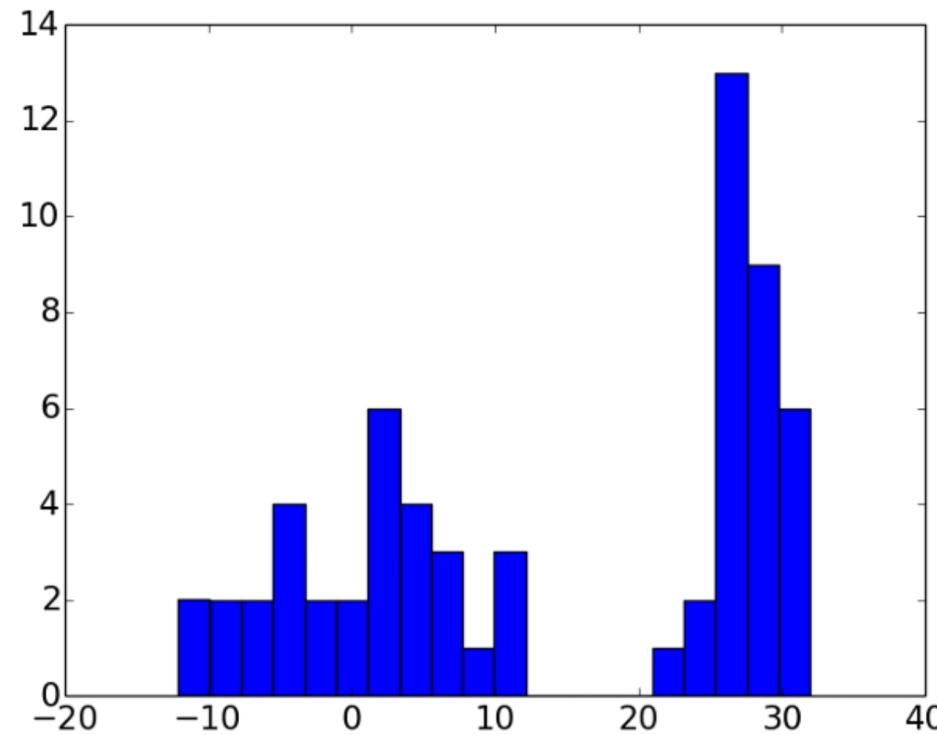


$$\hat{\pi}_m = \frac{1}{N} \sum_{i=1}^N w_m^{(i)}$$

$$\hat{\mu}_m = \frac{1}{\sum_{i=1}^N w_m^{(i)}} \sum_{i=1}^N w_m^{(i)} \underline{x}_i$$

$$\hat{\Sigma}_m = \frac{1}{\sum_{i=1}^N w_m^{(i)}} \sum_{i=1}^N w_m^{(i)} (\underline{x}_i - \underline{\mu}_m)(\underline{x}_i - \underline{\mu}_m)^T$$

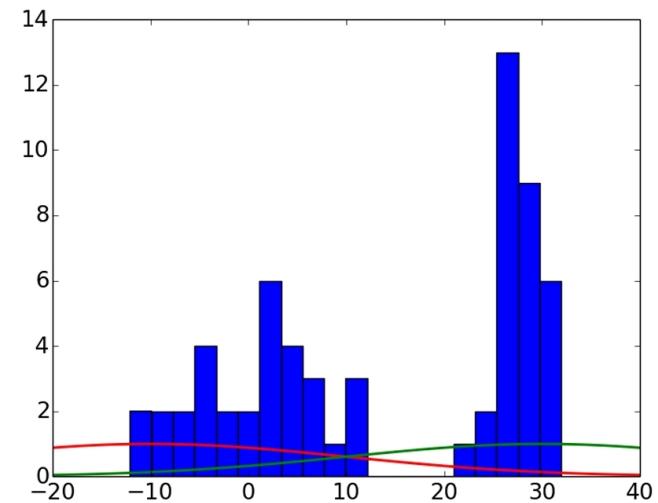
Example : Suppose you have recorded a bunch of temperatures in March for Toronto and Delhi , but forgot which was which, and they are all jumbled up together



Lets try to separate them using a mixture of Gaussians with EM

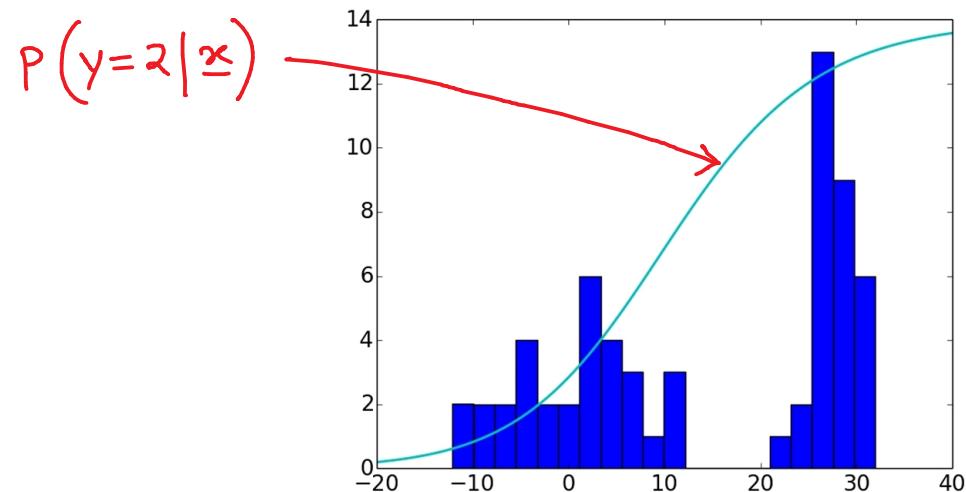
- Choose the number of Gaussians, $M=2$
- Randomly initialize the parameters Θ

$$\Theta = \begin{bmatrix} \pi_1 \\ \pi_2 \\ \mu_1 \\ \mu_2 \\ \sigma_1 \\ \sigma_2 \end{bmatrix}$$

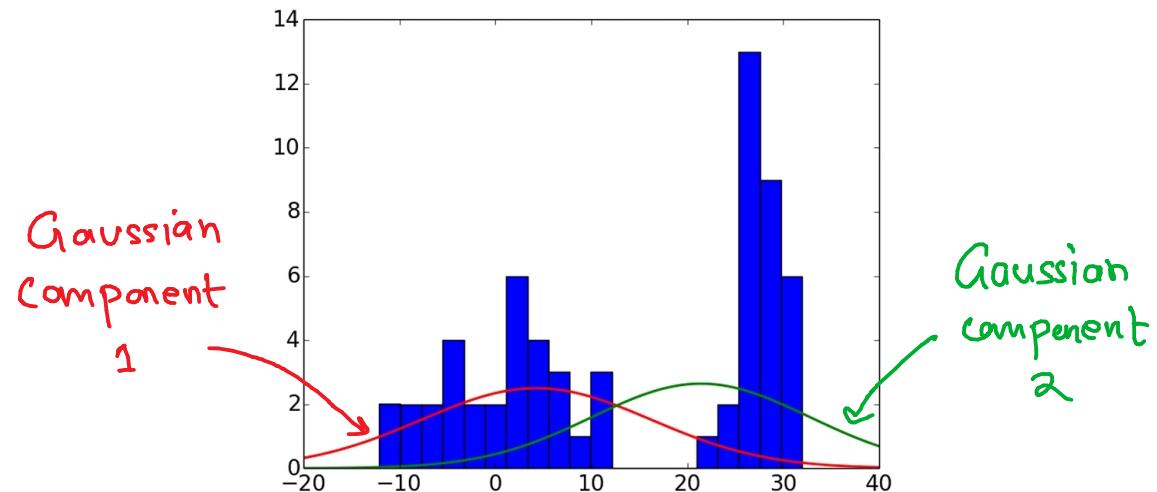


- Iteration 1: E-step and M-step

E-step

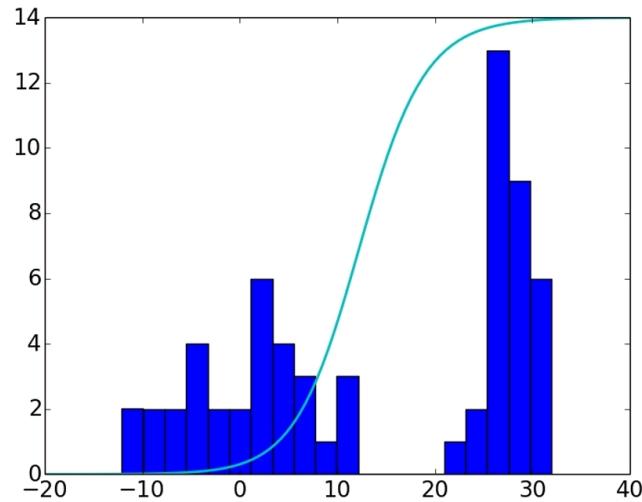


M-step

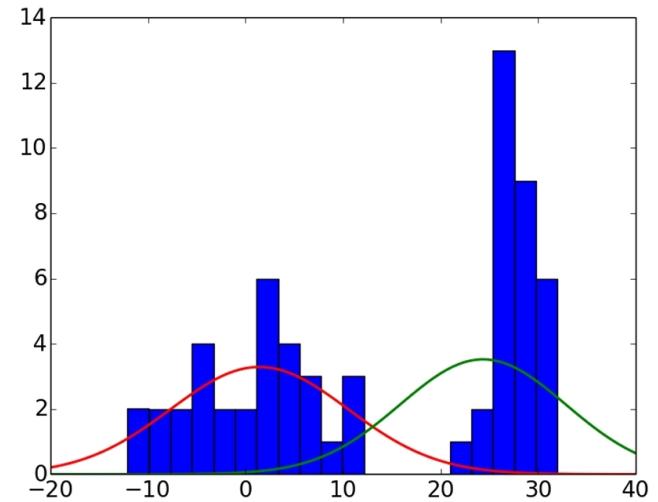


- Iteration 2

E-step

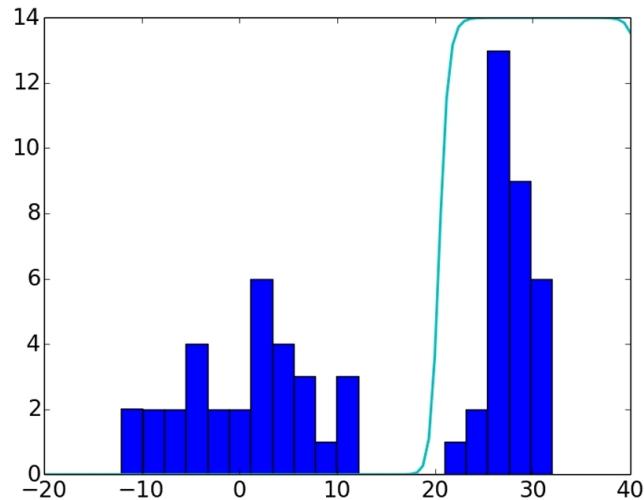


M-step

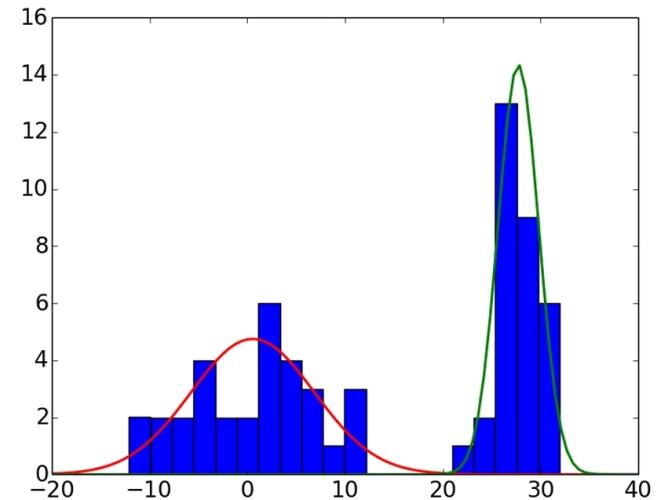


- Iteration 10

E-step

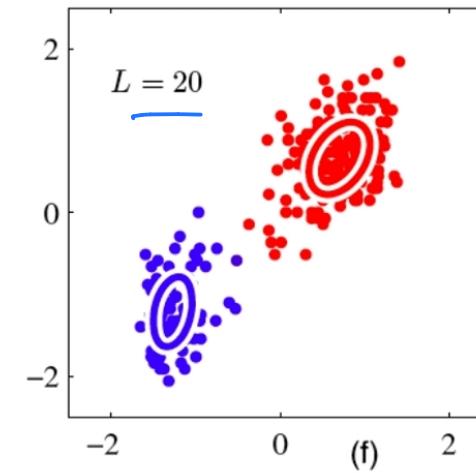
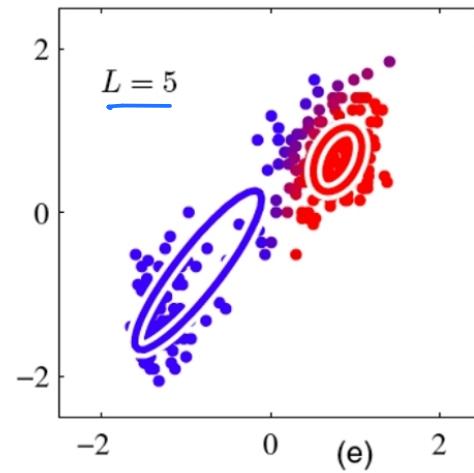
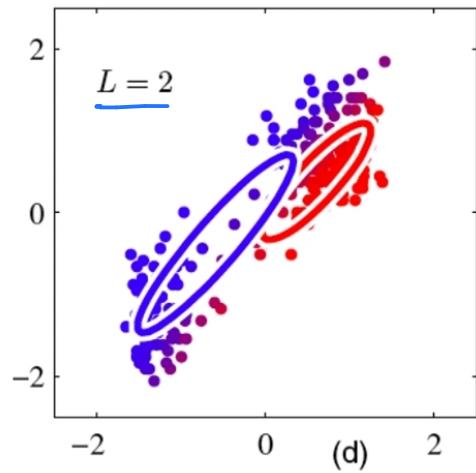
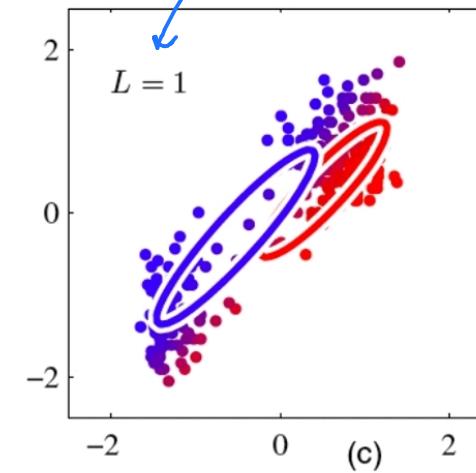
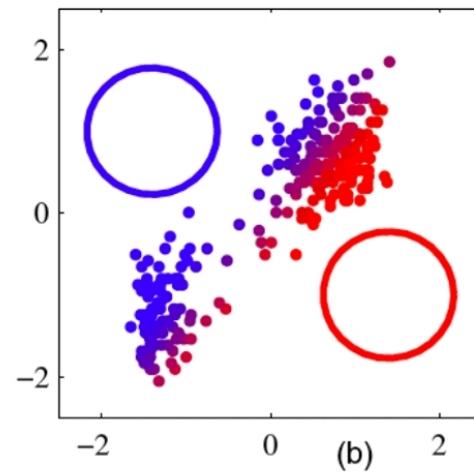
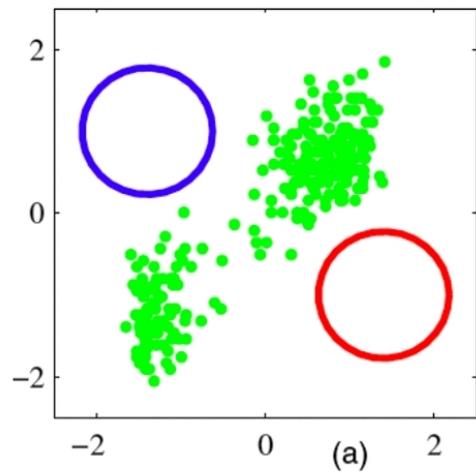


M-step



EM for multivariate Gaussians (2D)

Initialization



Few points to be remembered

- The number of clusters (i.e. the # of Gaussian components M) has to be specified before running the algorithm
 - M is a hyperparameter in GMM for clustering
- EM is a local optimizer and hence is only guaranteed to converge to a local optima / saddle point
 - Poor initialization can result in convergence to poor local optimum
 - Use random initializations