

Clustering: Models and Algorithms

Shikui Tu

Shanghai Jiao Tong University

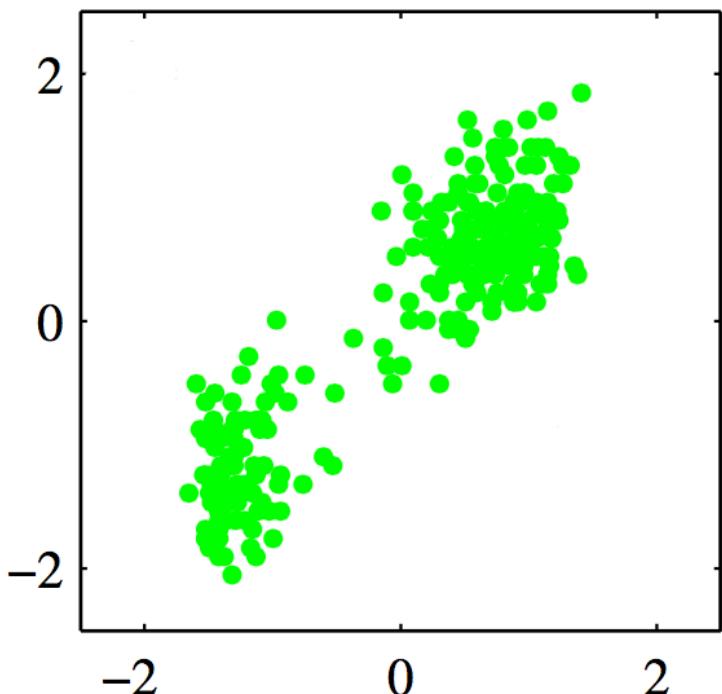
2021-03-09

Outline

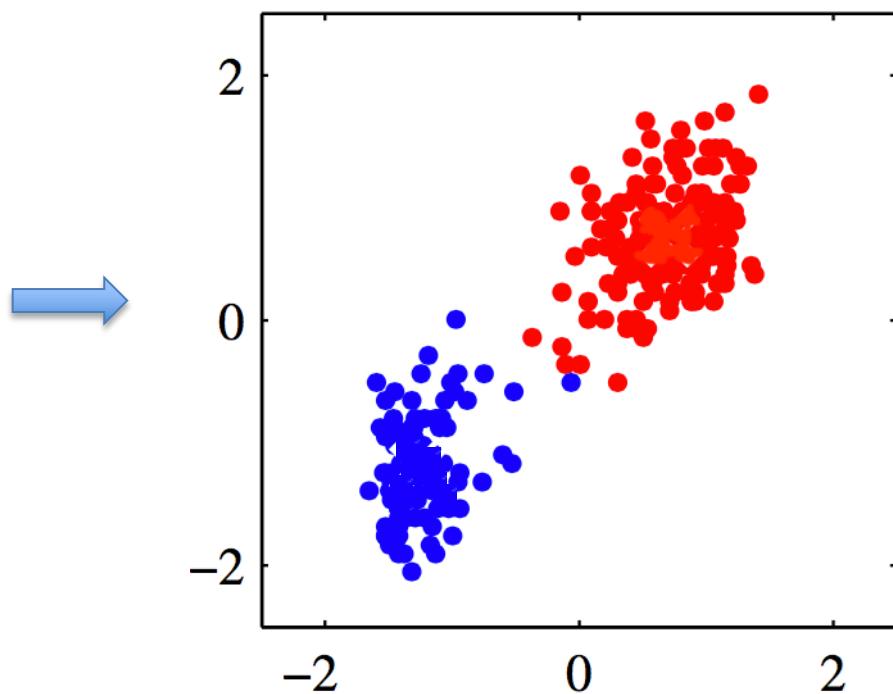
- A brief review of K-mean clustering
- Gaussian Mixture Models (GMM)
- Expectation-Maximization (EM) for maximum likelihood

Clustering the data

We have the following data:



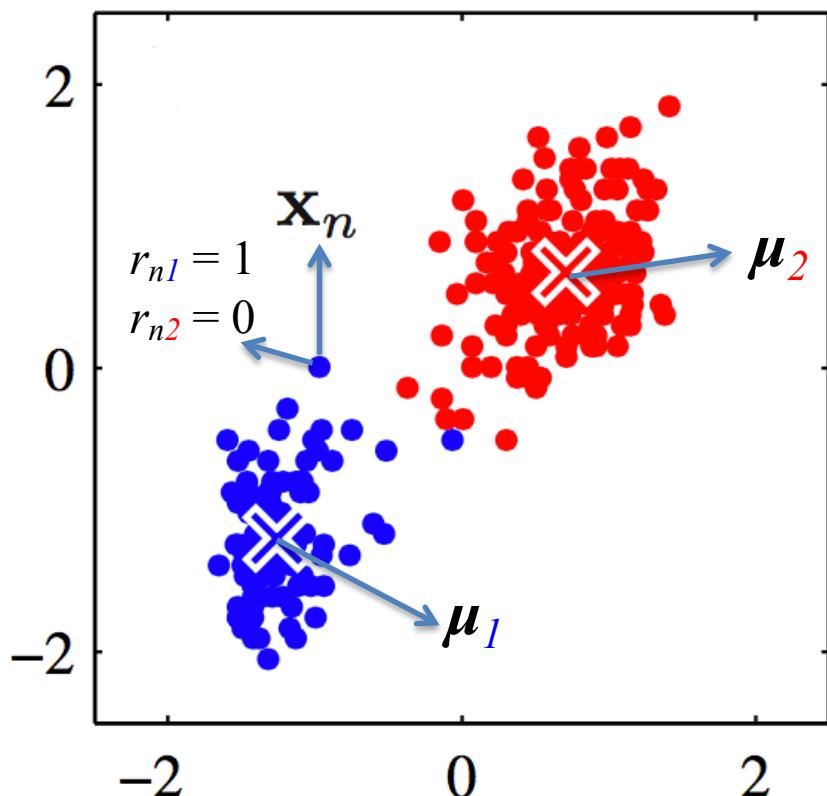
We want to cluster the data into two clusters (red and blue)



How?

Minimize the sum of square distances J

$$\text{minimize} \quad J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$



$r_{nk} = 1$ if and only if data point \mathbf{x}_n is assigned to cluster k ;
otherwise $r_{nk} = 0$.

$k = 1, 2$; $K = 2$ clusters

$n = 1, \dots, N$;
 N : the total number of points.

We need to calculate $\{r_{nk}\}$ and $\{\boldsymbol{\mu}_k\}$.

If we know r_{n1} , r_{n2} for all $n=1,\dots,N$

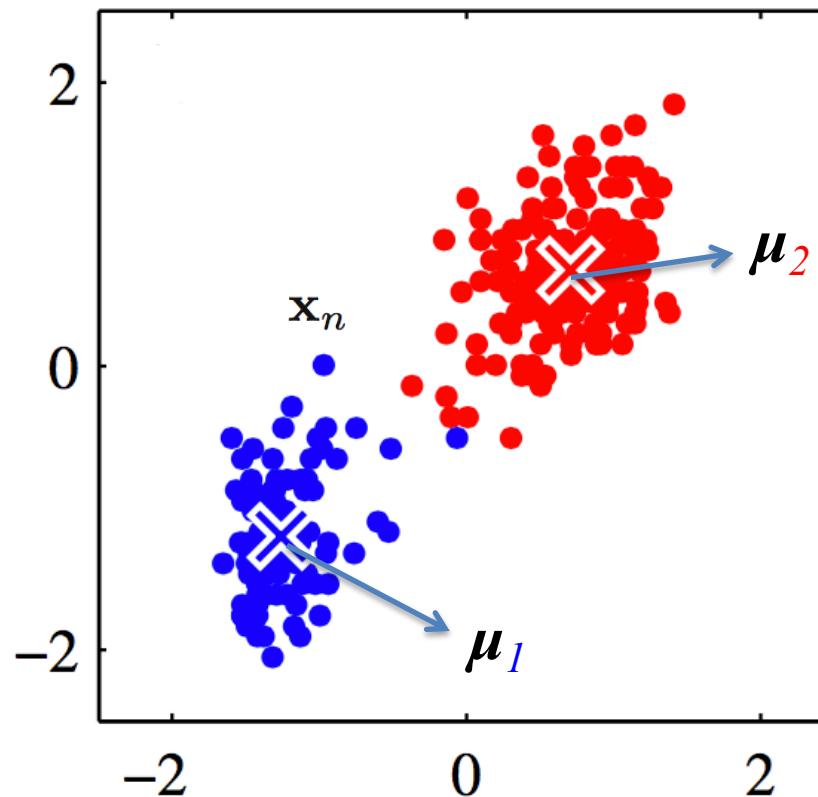
Since the points have been assigned to cluster 1 or cluster 2, we calculate

μ_1 = mean of the points in cluster 1

μ_2 = mean of the points in cluster 2

Or formally

$$\mu_k = \frac{\sum_n r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}$$



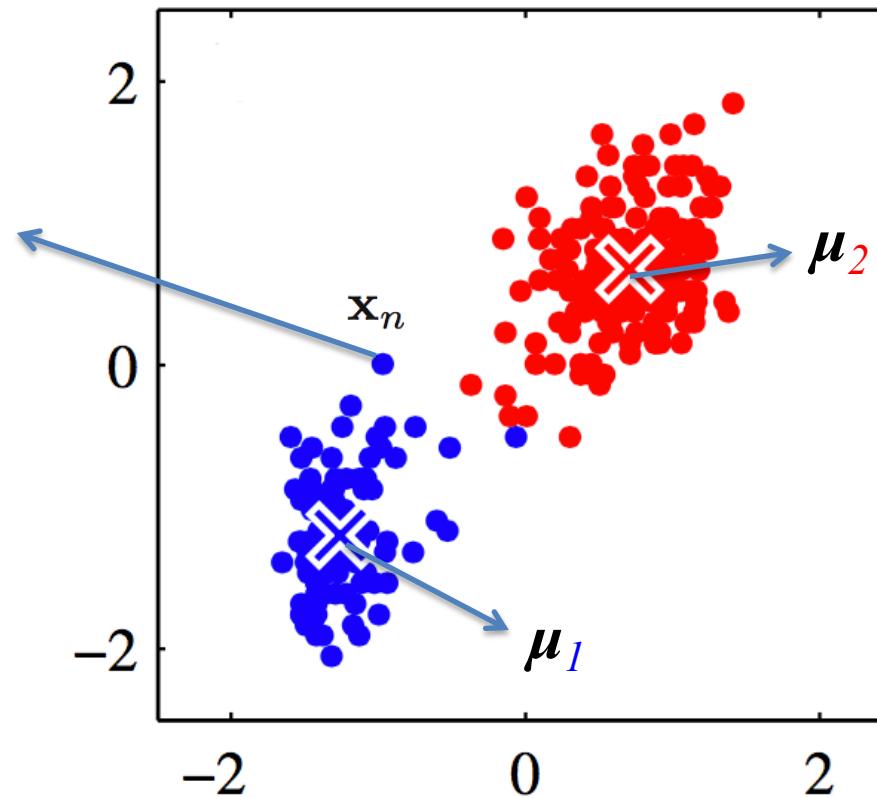
We call it the **M Step.**

If we know μ_1, μ_2

We should assign point \mathbf{x}_n to cluster 1, because

$$\|\mathbf{x}_n - \mu_1\|^2 < \|\mathbf{x}_n - \mu_2\|^2$$

Then, $r_{n1} = 1$
 $r_{n2} = 0$

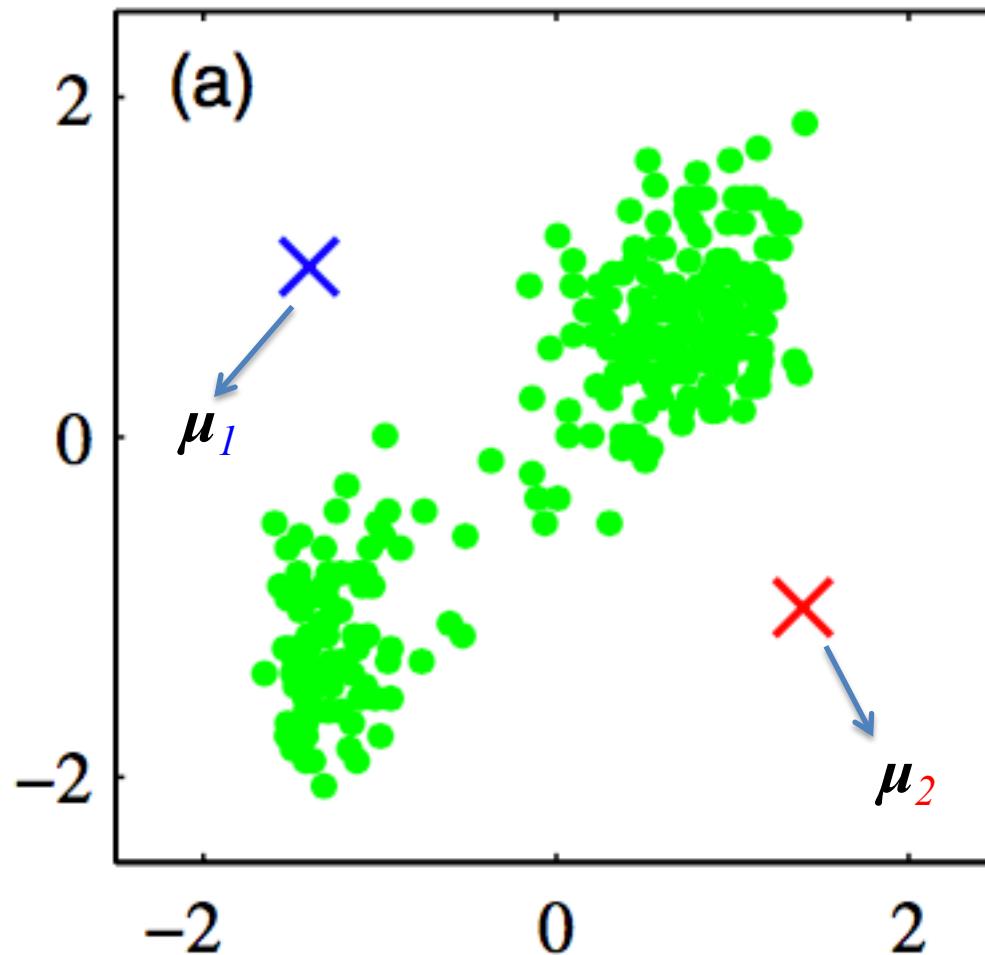


Or formally

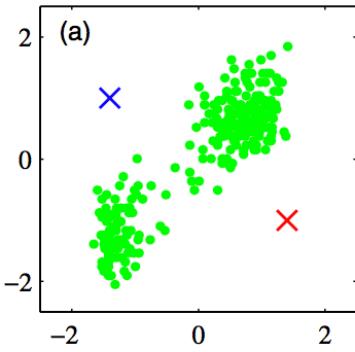
$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

We call it the **E Step**

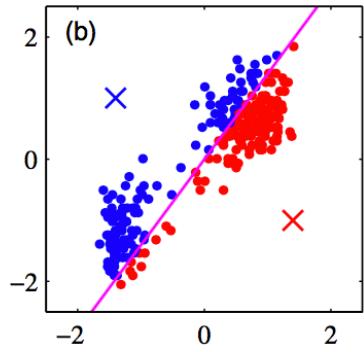
Initialization



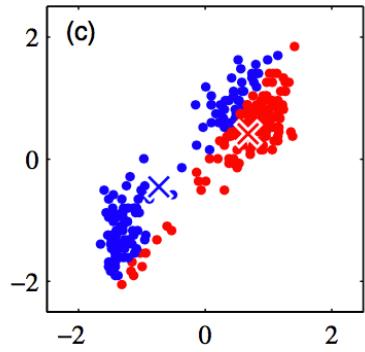
Initialization



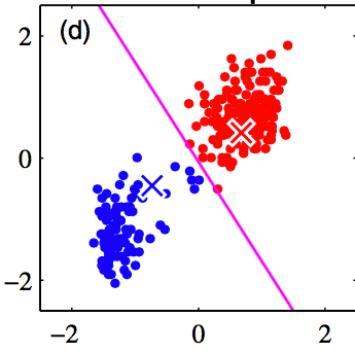
E-Step



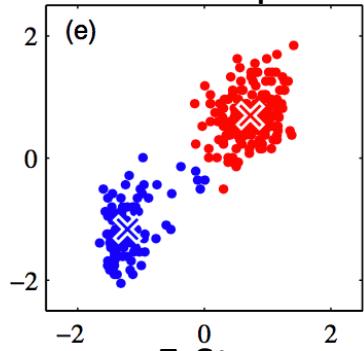
M-Step



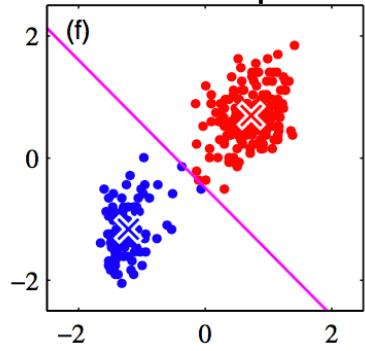
E-Step



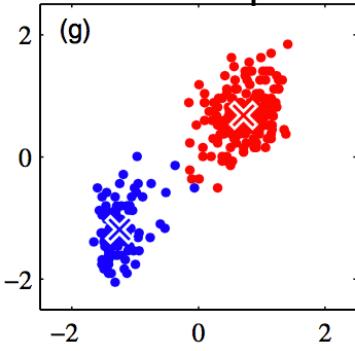
M-Step



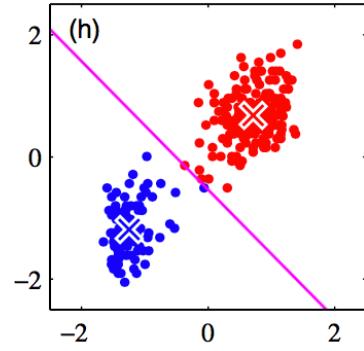
E-Step



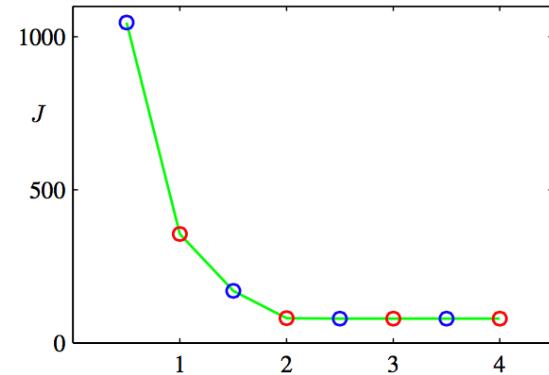
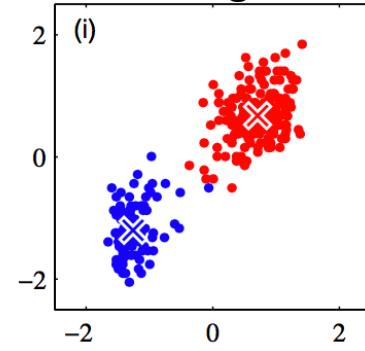
M-Step



E-Step



Convergence



If J does not change, or $\{\mu_1, \mu_2\}$ do not change, then the algorithm converges.

K均值法小结

- 初始化均值点 μ_1, \dots, μ_k
- 迭代如下
 - 把每个数据点按照就近原则分配给相应的 μ_i
 - 把 μ_i 更新为所分配的数据点的均值
- 迭代停止，如果聚类分配不变

Initialize $\mathbf{m}_i, i = 1, \dots, k$, for example, to k random \mathbf{x}^t

Repeat

For all $\mathbf{x}^t \in \mathcal{X}$

$$b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

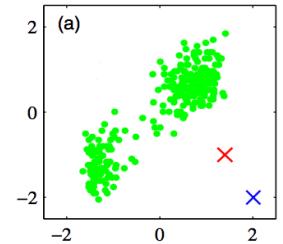
For all $\mathbf{m}_i, i = 1, \dots, k$

$$\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$$

Until \mathbf{m}_i converge

Issues for K-mean algorithm

- Does it find the global optimum of J ?
 - No, the nearest local optimum, depending on initialization
- If Euclidean distance is not good for some data, do we have other choices?
- Can we assign each data point to the clusters probabilistically?
- If K (the total number of clusters) is unknown, can we estimate it from the data?

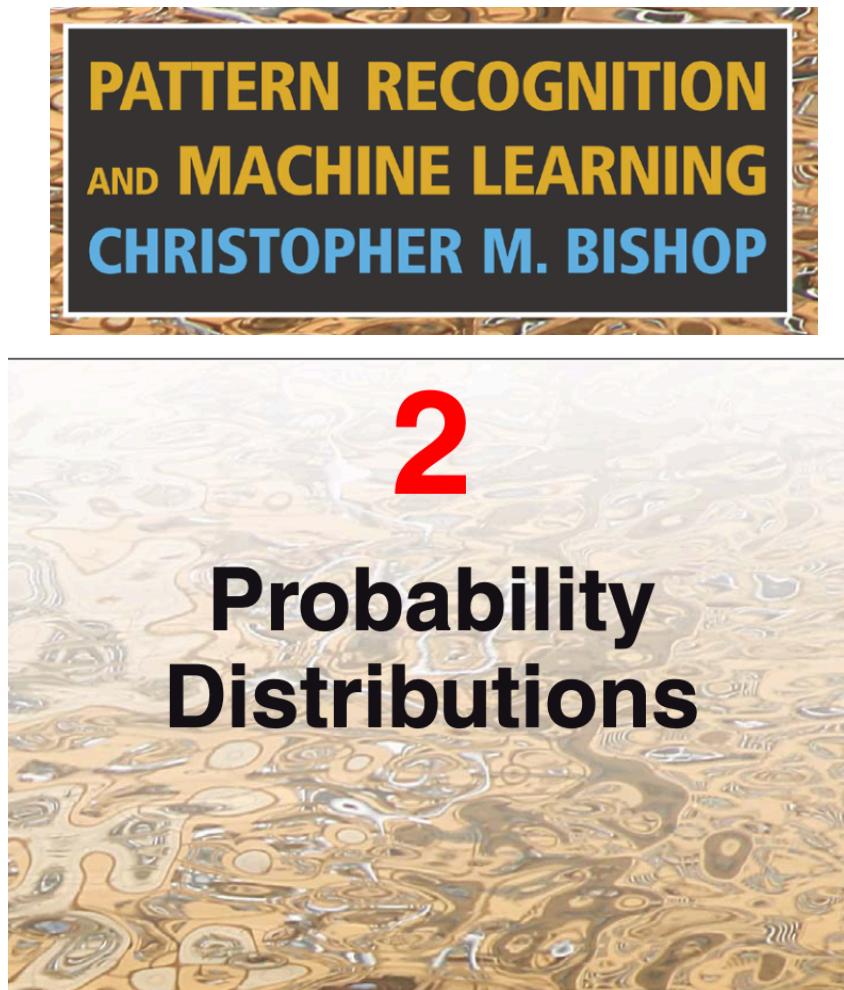


Outline

- Clustering
 - K-mean clustering, hierarchical clustering
- Adaptive learning (online learning)
 - CL, FSCL, RPCL
- Gaussian Mixture Models (GMM)
- Expectation-Maximization (EM) for maximum likelihood

Probability distributions

- Read Bishop's PRML book, Chapter 2



Binary variables

$$x \in \{0, 1\}$$

$$0 \leq \mu \leq 1$$

- Bernoulli distribution

The probability of $x = 1$

$$p(x = 1|\mu) = \mu$$

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}$$

$$p(x = 0|\mu) = 1 - \mu$$

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu)$$

Binary variables

$$x \in \{0, 1\}$$

$$0 \leq \mu \leq 1$$

- Bernoulli distribution

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

Now suppose we have a data set $\mathcal{D} = \{x_1, \dots, x_N\}$ of observed values of x . We can construct the likelihood function, which is a function of μ , on the assumption that the observations are drawn independently from $p(x|\mu)$, so that

$$p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}. \quad (2.5)$$

$$\ln p(\mathcal{D}|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}.$$

Multinomial variables

1-of- K scheme

- Example $\mathbf{x} = (0, 0, 1, 0, 0, 0)^T \quad K = 6 \quad x_3 = 1$
- Multinomial distribution $\sum_{k=1}^K x_k = 1$

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T \quad \mu_k \geq 0$$
$$\sum_k \mu_k = 1$$

$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = (\mu_1, \dots, \mu_M)^T = \boldsymbol{\mu}$$

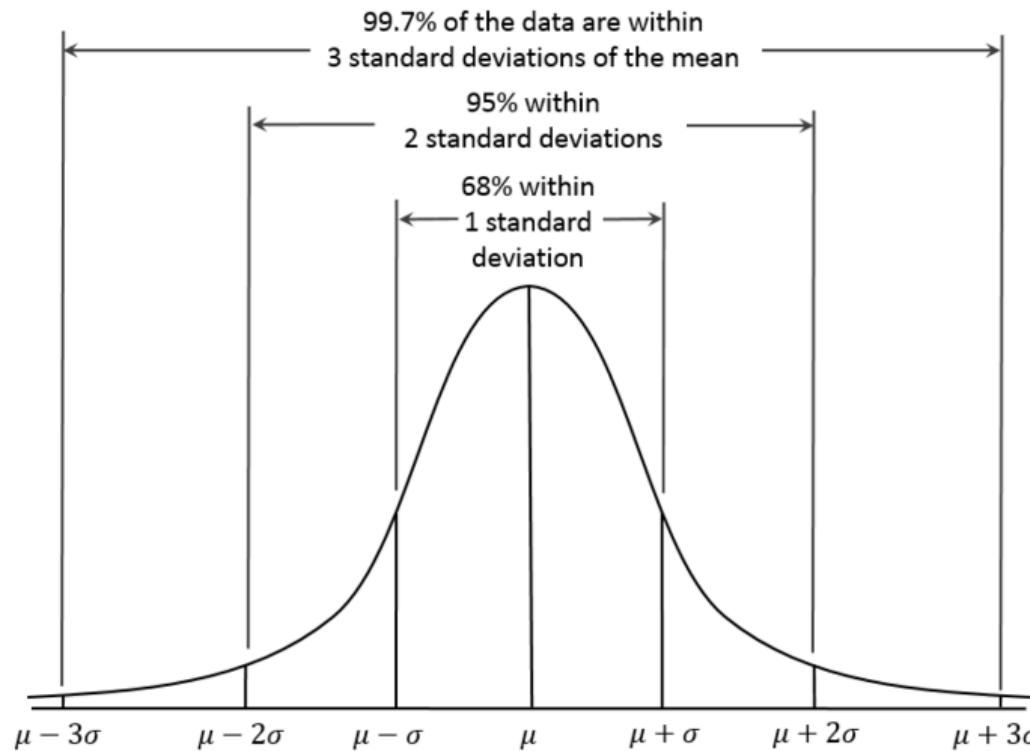
$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k} \quad m_k = \sum_n x_{nk}$$

$$\text{Mult}(m_1, m_2, \dots, m_K | \boldsymbol{\mu}, N) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k}$$

Gaussian distribution

also known as the **normal distribution**

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$



Multivariate Gaussian

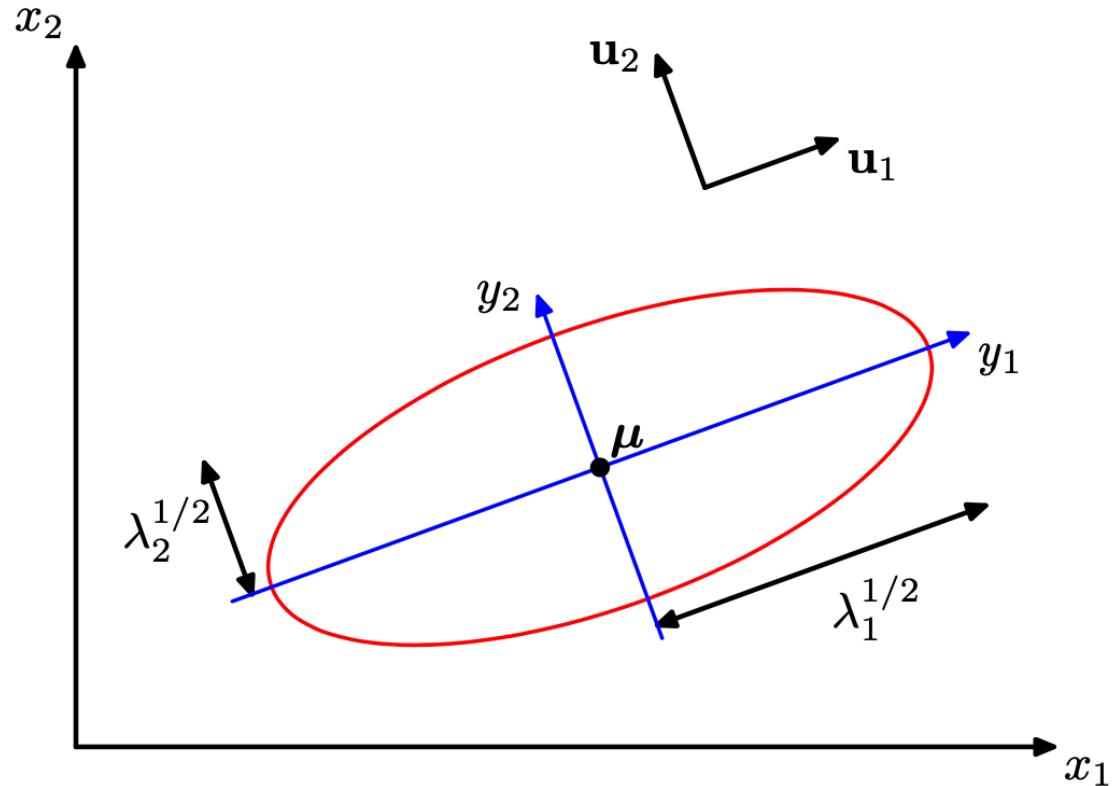
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

The red curve shows the elliptical surface of constant probability density for a Gaussian in a two-dimensional space $\mathbf{x} = (x_1, x_2)$ on which the density is $\exp(-1/2)$ of its value at $\mathbf{x} = \boldsymbol{\mu}$. The major axes of the ellipse are defined by the eigenvectors \mathbf{u}_i of the covariance matrix, with corresponding eigenvalues λ_i .

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$$

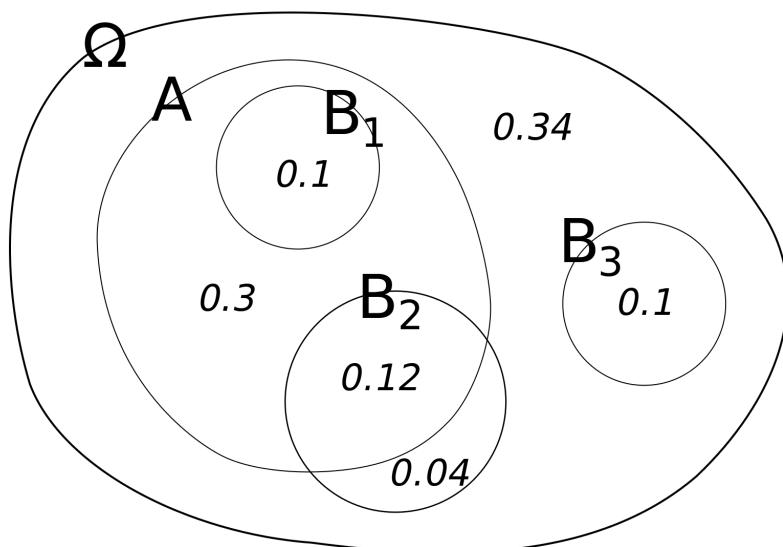
$$\text{cov}[\mathbf{x}] = \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] = \boldsymbol{\Sigma}.$$



Conditional probability

From wiki

$$P(A | B) = \frac{P(A \cap B)}{P(B)},$$



The **unconditional** probability

$$P(A) = 0.30 + 0.10 + 0.12 = 0.52.$$

The **conditional** probability

$$P(A|B_1) = 1,$$

$$P(A|B_2) = 0.12 \div (0.12 + 0.04) = 0.75,$$

$$P(A|B_3) = 0.$$

Bayes Theorem

From wiki

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

$$p(X) = \sum_Y p(X|Y)p(Y)$$

Number of occurrences	Beard: No beard:		sum
	B	\bar{B}	
Astigmatic: A	2	3	5
Not astigmatic: \bar{A}	6	9	15
sum	8	12	20

	B	\bar{B}	
A	2	3	A
\bar{A}	6	9	\bar{A}

$$P(B, \text{ given } A) \cdot P(A) = P(B|A) \cdot P(A)$$

$$\frac{2}{2+3} \cdot \frac{2+3}{2+3+6+9} = \frac{2}{2+3+6+9}$$

	B	\bar{B}	
A	2	3	A
\bar{A}	6	9	\bar{A}

$$P(A, \text{ given } B) \cdot P(B) = P(A|B) \cdot P(B)$$

$$\frac{2}{2+6} \cdot \frac{2+6}{2+3+6+9} = \frac{2}{2+3+6+9}$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$\therefore P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

The Rules of Probability

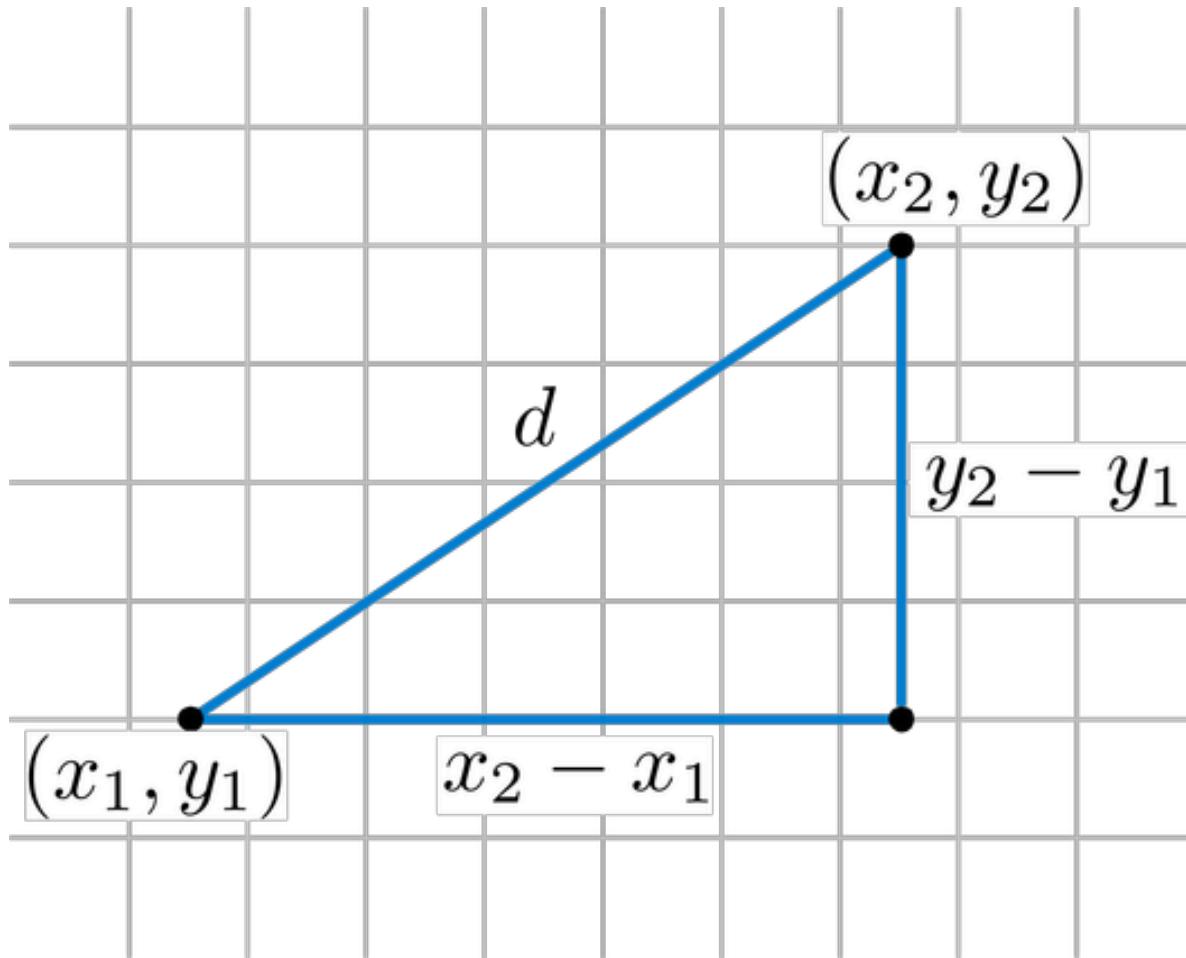
sum rule

$$p(X) = \sum_Y p(X, Y)$$

product rule

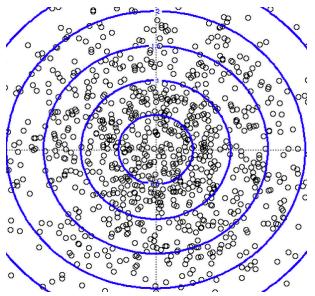
$$p(X, Y) = p(Y|X)p(X).$$

Euclidean Distance

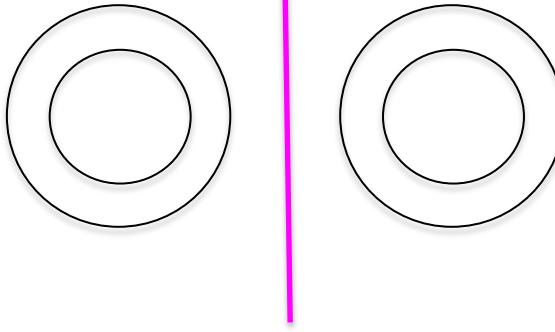


Euclidian distance may not be a good measure for some data

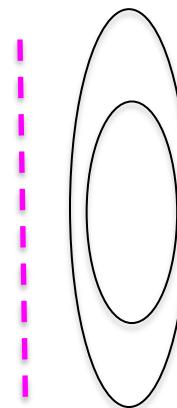
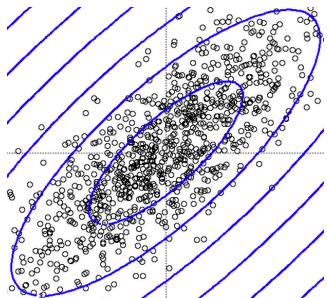
Euclidean distance



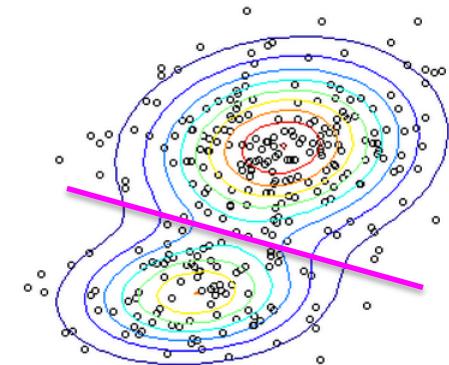
Equal distance line



Mahalanobis distance



In general



Distances at different directions could be different!

$$d(\vec{x}, \vec{y}) = \sqrt{(\vec{x} - \vec{y})^T \Sigma^{-1} (\vec{x} - \vec{y})}$$

Σ is the covariance matrix

More Distance Measures

Table 1 Gene expression similarity measures

Manhattan distance (city-block distance, L1 norm)	$d_{fg} = \sum_c e_{fc} - e_{gc} $
Euclidean distance (L2 norm)	$d_{fg} = \sqrt{\sum_c (e_{fc} - e_{gc})^2}$
Mahalanobis distance	$d_{fg} = (e_f - e_g)^\top \Sigma^{-1} (e_f - e_g)$, where Σ is the (full or within-cluster) covariance matrix of the data
Pearson correlation (centered correlation)	$d_{fg} = 1 - r_{fg}$, with $r_{fg} = \frac{\sum_c (e_{fc} - \bar{e}_f)(e_{gc} - \bar{e}_g)}{\sqrt{\sum_c (e_{fc} - \bar{e}_f)^2} \sqrt{\sum_c (e_{gc} - \bar{e}_g)^2}}$
Uncentered correlation (angular separation, cosine angle)	$d_{fg} = 1 - r_{fg}$, with $r_{fg} = \frac{\sum_c e_{fc} e_{gc}}{\sqrt{\sum_c e_{fc}^2} \sqrt{\sum_c e_{gc}^2}}$
Spellman rank correlation	As Pearson correlation, but replace e_{gc} with the rank of e_{gc} within the expression values of gene g across all conditions $c = 1 \dots C$
Absolute or squared correlation	$d_{fg} = 1 - r_{fg} $ or $d_{fg} = 1 - r_{fg}^2$
d_{fg} , distance between expression patterns for genes f and g . e_{gc} , expression level of gene g under condition c .	

From distance to probability

distance

$$\|x - \mu\|^2$$

likely



$$\exp\{-\lambda \|x - \mu\|^2\}$$

“The closer, the more likely.”

Sum or integral to
be one

Probability

$$\frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

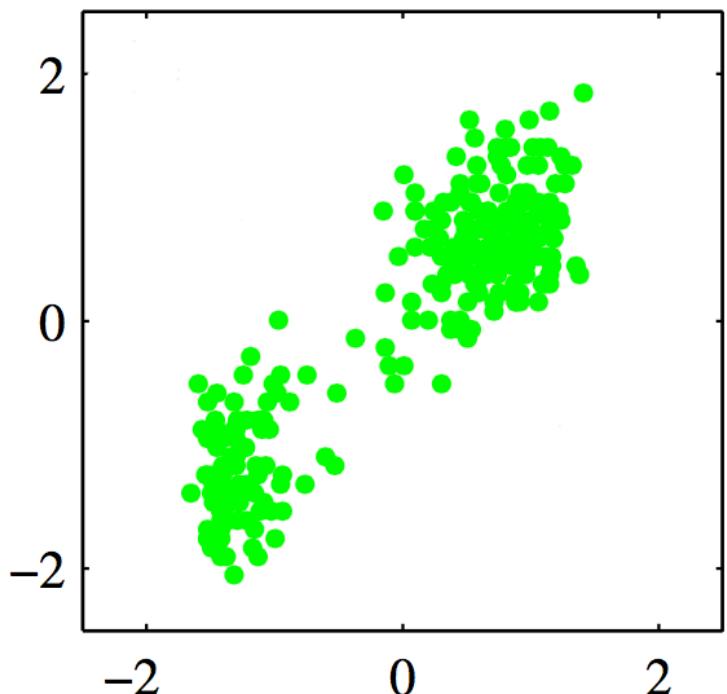
It is **more powerful** to consider everything in probability framework!

Gaussian distribution with the Mahalanobis distance

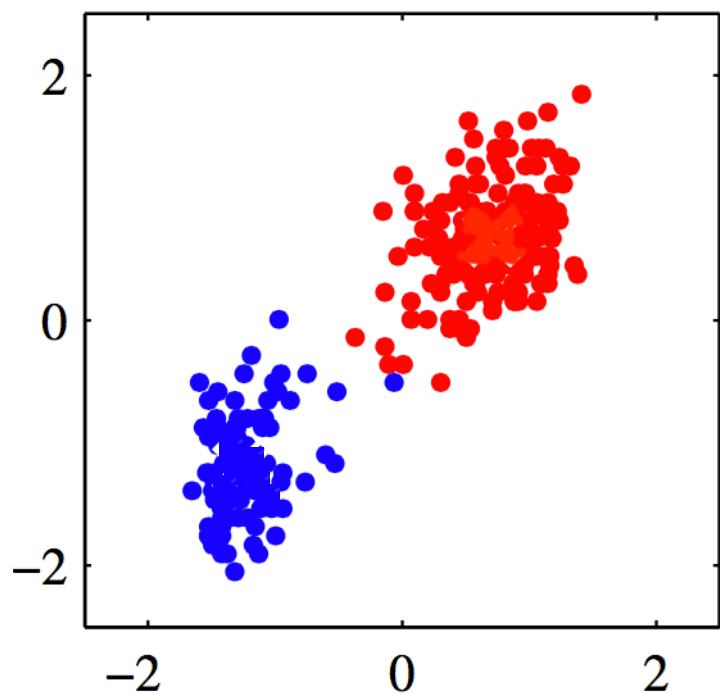
$$D_M(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}.$$

Review the clustering problem again

We have the following data:

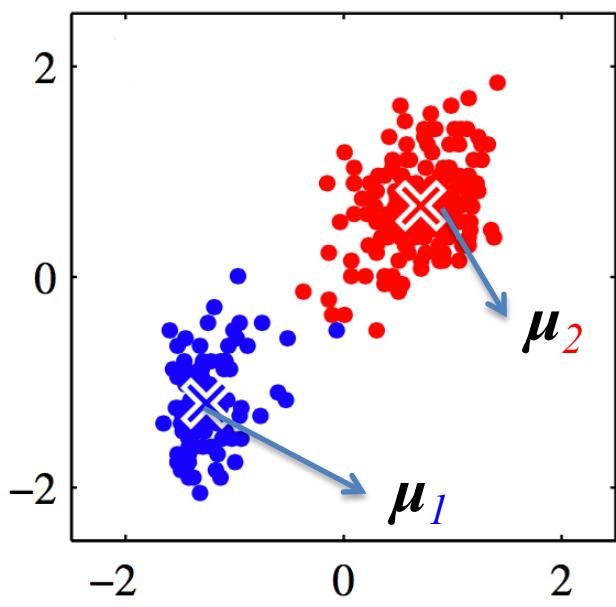


We want to cluster the data into two clusters (**red** and **blue**)

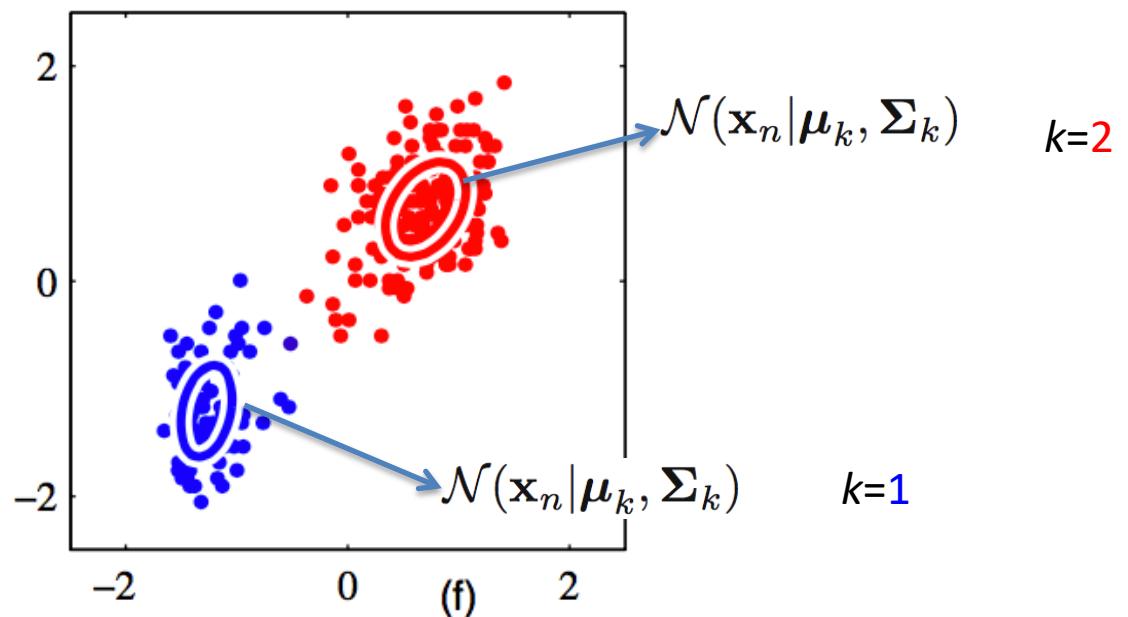


Instead if using $\{\mu_1, \mu_2\}$, each cluster is represented as a Gaussian distribution

K-means

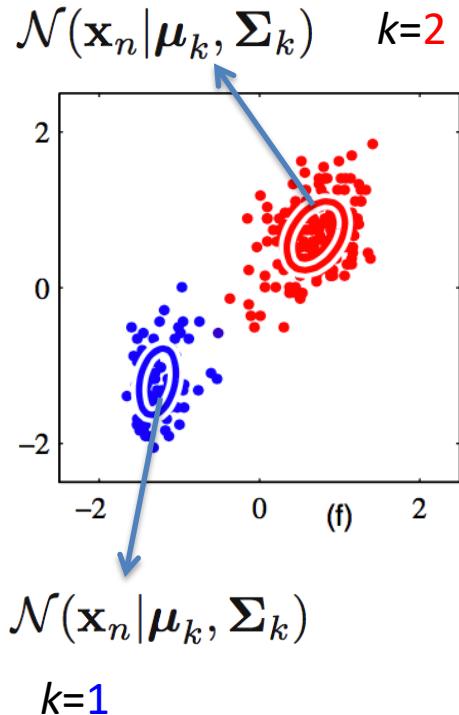


Gaussian Mixture Model (GMM)



$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

Gaussian Mixture Model (GMM)



We use $z_k = 1$ to indicate a point \mathbf{x} belongs to cluster k

$$\mathbf{z} = (z_1, \dots, z_K) \quad z_k \in \{0, 1\} \quad \sum_k z_k = 1$$

Assume the points in the same cluster follow a **Gaussian distribution**

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

A mixing weight for each cluster:

$$p(z_k = 1) = \pi_k \quad 0 \leq \pi_k \leq 1 \quad \sum_{k=1}^K \pi_k = 1$$

prior probability of point belonging to a cluster

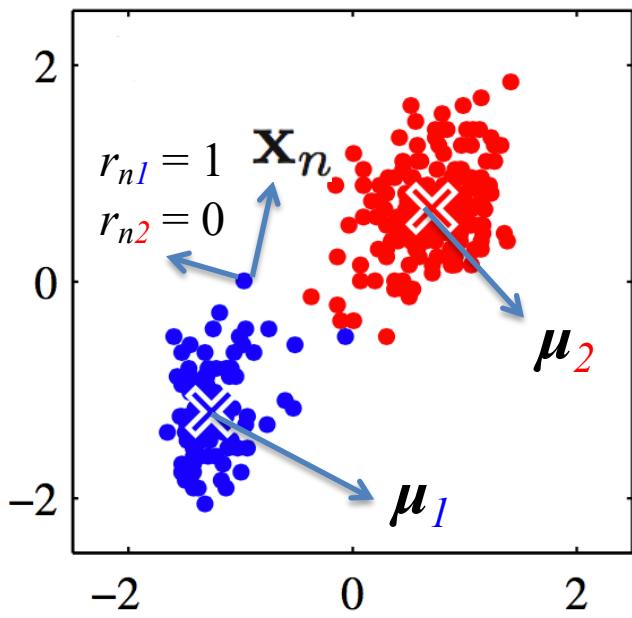
So, we get a distribution for the data point \mathbf{x} :

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

From minimizing sum of square distances to finding maximum likelihood

minimize

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|\mathbf{x}_n - \boldsymbol{\mu}_k\|^2$$



maximize likelihood

$$p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})$$

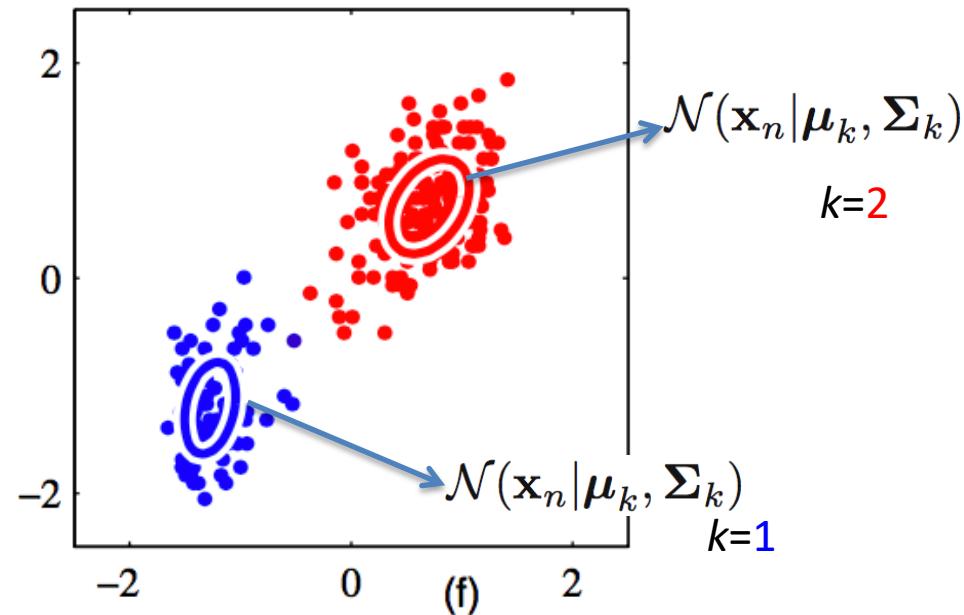


$$X = \{x_1, \dots, x_N\}$$

$$\pi = \{\pi_1, \dots, \pi_K\}$$

$$\boldsymbol{\mu} = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}$$

$$\boldsymbol{\Sigma} = \{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\}$$



Remember: The closer the distance, the more likely the probability.

Maximum likelihood

Given a data set $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ in which the observations $\{\mathbf{x}_n\}$ are assumed to be drawn independently from a multivariate Gaussian distribution, we can estimate the parameters of the distribution by maximum likelihood. The log likelihood function is given by

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

Maximizing the log-likelihood function:

$$\frac{\partial}{\partial \boldsymbol{\mu}} \ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) = 0 \quad \longrightarrow \quad \boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$\text{Similarly we get} \quad \boldsymbol{\Sigma}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})(\mathbf{x}_n - \boldsymbol{\mu}_{\text{ML}})^T$$

$\boldsymbol{\mu}_{\text{ML}}$ and $\boldsymbol{\Sigma}_{\text{ML}}$ are the maximum likelihood estimates of the mean and the co-variance matrix.

Thank you!

Matrix-cook-book

$$\begin{aligned}\partial \mathbf{A} &= 0 & (\mathbf{A} \text{ is a constant}) \\ \partial(\alpha \mathbf{X}) &= \alpha \partial \mathbf{X} \\ \partial(\mathbf{X} + \mathbf{Y}) &= \partial \mathbf{X} + \partial \mathbf{Y} \\ \partial(\text{Tr}(\mathbf{X})) &= \text{Tr}(\partial \mathbf{X}) \\ \partial(\mathbf{X} \mathbf{Y}) &= (\partial \mathbf{X}) \mathbf{Y} + \mathbf{X} (\partial \mathbf{Y}) \\ \partial(\mathbf{X} \circ \mathbf{Y}) &= (\partial \mathbf{X}) \circ \mathbf{Y} + \mathbf{X} \circ (\partial \mathbf{Y}) \\ \partial(\mathbf{X} \otimes \mathbf{Y}) &= (\partial \mathbf{X}) \otimes \mathbf{Y} + \mathbf{X} \otimes (\partial \mathbf{Y}) \\ \partial(\mathbf{X}^{-1}) &= -\mathbf{X}^{-1} (\partial \mathbf{X}) \mathbf{X}^{-1} \\ \partial(\det(\mathbf{X})) &= \det(\mathbf{X}) \text{Tr}(\mathbf{X}^{-1} \partial \mathbf{X}) \\ \partial(\ln(\det(\mathbf{X}))) &= \text{Tr}(\mathbf{X}^{-1} \partial \mathbf{X}) \\ \partial \mathbf{X}^T &= (\partial \mathbf{X})^T \\ \partial \mathbf{X}^H &= (\partial \mathbf{X})^H\end{aligned}$$