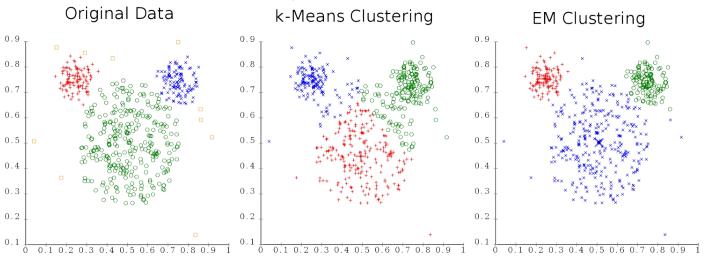# Lecture 17 – The Expectation-Maximization Algorithm

In some settings, the k-means algorithm does not give desirable results:

## Different cluster analysis results on "mouse" data set:



The issue is that k-means does not account for the size of the different clusters.

Alternative approach: assume $\underline{x}_i$'s are drawn at random from a mixture of Gaussians distribution and cluster using the Expectation Maximization (EM) algorithm.

A Gaussian distribution is a probability distribution that characterizes how likely different values of a random variable are.

$$f(\underline{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \Sigma^{-1}(\underline{x}-\underline{\mu})\right\} \quad \text{(shorthand } \mathcal{N}(\underline{\mu}, \Sigma))$$
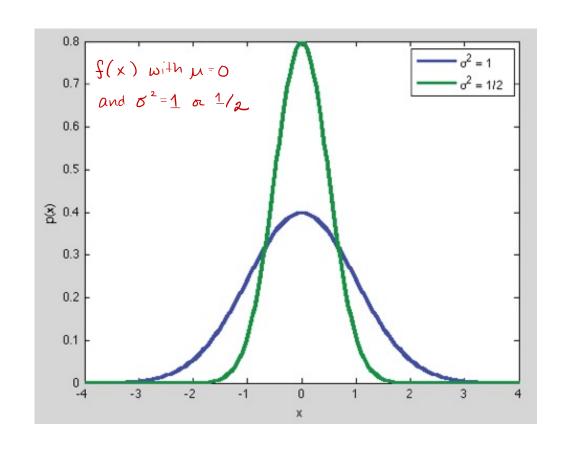
$|\Sigma| =$ matrix determinant
$=$ product of singular values

$\underline{\mu} =$ mean $= \mathbb{E}\,\underline{x} =$ expected (average) value of $\underline{x}$

$\Sigma =$ covariance $= \mathbb{E}\left[(\underline{x}-\underline{\mu})(\underline{x}-\underline{\mu})^T\right] \iff \Sigma_{ij} = \mathbb{E}\left[(x_i - \mu_i)(x_j - \mu_j)\right]$

Ex 1: $p = 1$

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

$$\left(\text{i.e. } \Sigma = \sigma^2\right)$$

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$



$f(x)$ with $\mu = 0$ and $\sigma^2 = 1$ or $1/2$

Legend: $\sigma^2 = 1$, $\sigma^2 = 1/2$

Ex.  p = 2

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

$\Sigma$

$\begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$

$f(\underline{x}) = \gamma \iff \dfrac{(x_1 - \mu_1)^2}{\sigma^2} + \dfrac{(x_2 - \mu_2)^2}{\sigma^2} = \dfrac{\|\underline{x} - \underline{\mu}\|_2^2}{\sigma^2} = \gamma'$

Contour plot. each circle is set of all $\underline{x}$
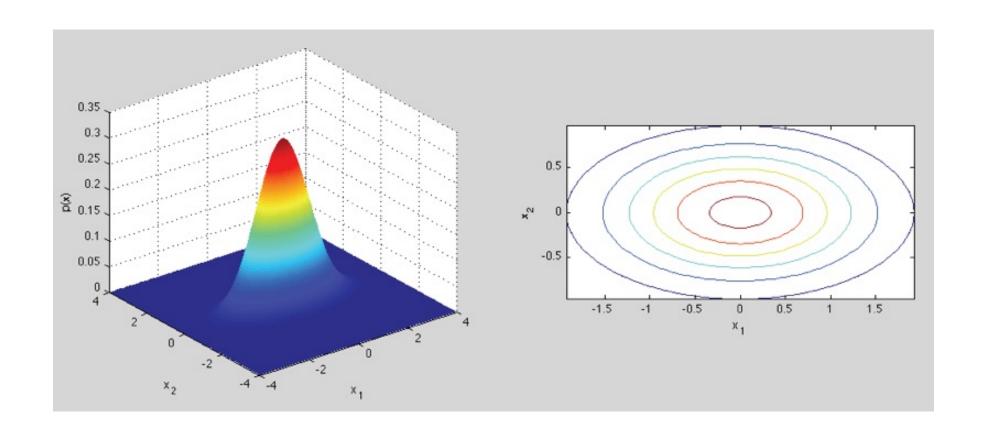
such that $f(\underline{x}) = \gamma$ for some $\gamma$

= set of $\underline{x}$ such that $\|\underline{x} - \underline{\mu}\|^2 = \gamma'$

Ex. $p = 2$.

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$\underline{x} \sim N(\underline{\mu}, \Sigma)$$

Here density contours are ellipses whose axes align with the coordinate axes. Note:

$$f(\underline{x}) = \gamma \iff \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} = \gamma'$$
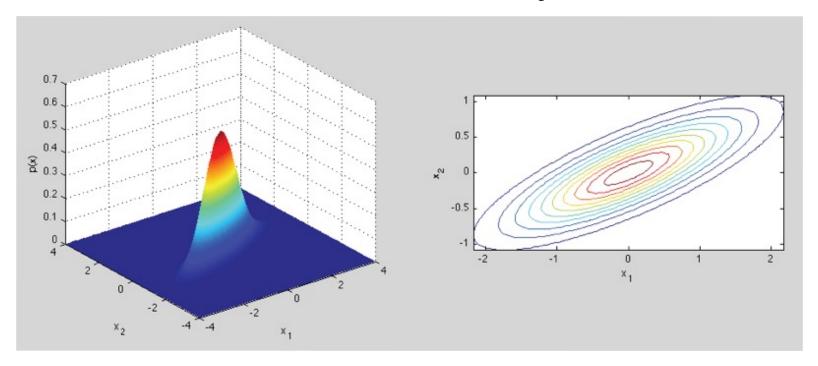
Ex. $\Sigma$ is an arbitrary positive-definite matrix.

Let $\Sigma = V \Lambda V^T$ (eigenvalue decomposition)

Let $\underline{x}' = V^T \underline{x}$ and $\underline{\mu}' = V^T \underline{\mu}$. $\Leftarrow$ $\underline{x}' = \underline{x}$ in rotated coordinate system defined by colums of $V$.

Then contour = set of all $\underline{x}$ s.t. $(\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) = \gamma$ for some $\gamma$

$$(\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) = (\underline{x} - \underline{\mu})^T V \Lambda^{-1} V^T (\underline{x} - \underline{\mu}) = (\underline{x}' - \underline{\mu}')^T \Lambda^{-1} (\underline{x}' - \underline{\mu}') = \frac{(x_1' - \mu_1')^2}{\lambda_1} + \frac{(x_2' - \mu_2')^2}{\lambda_2}$$

= ellipse in rotated coordinate system, where $V$ defines rotation.



The SVD of the covariance matrix $\Sigma$ tells us how the distribution is "oriented" and spread out in different directions

# Maximum likelihood estimation

Given $\underline{x}_i \sim N(\mu, \Sigma)$ for $i = 1, 2, \ldots, n$, we may wish to estimate $\underline{\mu}$ and $\Sigma$.

Maximum likelihood strategy: choose $\hat{\underline{\mu}}$ and $\hat{\Sigma}$ to maximize the likelihood of the $\underline{x}_i$'s

$$(\hat{\underline{\mu}}, \hat{\Sigma}) = \underset{\mu, \Sigma}{\text{argmax}} \prod_{i=1}^{n} f(\underline{x}_i \; ; \mu, \Sigma)$$

$$= \underset{\mu, \Sigma}{\text{argmax}} \log \left( \prod_{i=1}^{n} f(\underline{x}_i \; , \underline{\mu}, \Sigma) \right)$$

$$= \underset{\mu, \Sigma}{\text{argmin}} -\log \left( \prod_{i=1}^{n} f(\underline{x}_i \; ; \mu, \Sigma) \right)$$

$$= \underset{\mu, \Sigma}{\text{argmin}} \sum_{i=1}^{n} -\log f(\underline{x}_i \; ; \mu, \Sigma)$$

$$= \underset{\mu, \Sigma}{\text{argmin}} \sum_{i=1}^{n} \frac{1}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\underline{x}_i - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu})$$

Compute gradients, set to zero $\Rightarrow$

$$\hat{\underline{\mu}} = \frac{1}{n} \sum_{i=1}^{n} \underline{x}_i \quad , \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (\underline{x}_i - \hat{\underline{\mu}})(\underline{x}_i - \hat{\underline{\mu}})^T$$

Assume there are K Gaussians (each will correspond to a different cluster) with

means $\underline{\mu}_k$ and covariances $\Sigma_k$ for $k = 1, 2, \ldots, K$.

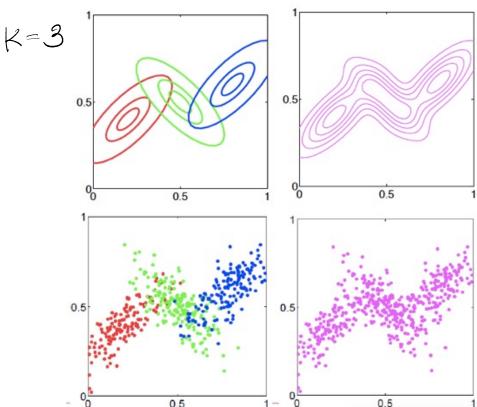We model the observed $\underline{x}_i$'s drawn from a mixture of these Gaussians as follows.

- Choose one of the K Gaussians — ie. the $k^{th}$ Gaussian is chosen with probability $\pi_k$,

  where $\sum\limits_{k=1}^{K} \pi_k = 1$; denote $k_i$

- draw $\underline{x}_i \sim N(\underline{\mu}_{k_i}, \Sigma_{k_i})$

$$\Rightarrow f(\underline{x}_i) = \sum_{k=1}^{K} \pi_k \, f(\underline{x}; \mu_k, \Sigma_k)$$

$K = 3$



Given $\underline{x}_i$'s, we want to cluster them **without**

knowing $\underline{\mu}_k$'s or $\Sigma_k$'s or $\pi_k$'s.

if we knew cluster membership $k_i$ for each $\underline{x}_i$,

then we could use maximum likelihood estimation to

compute $\mu_k$'s, $\Sigma_k$'s, and $\pi_k$'s. Without $k_i$'s,

maximum likelihood estimation is hard.

# Expectation - Maximization Algorithm

- Initialize means $\hat{\mu}_k$, covariances $\hat{\Sigma}_k$, and mixture weights $\hat{\pi}_k$ for $k = 1, 2, \ldots, K$

- E-step: compute $p_k(\underline{x}) = $ Probability that $\underline{x}_i$ was drawn from $k^{th}$ Gaussian given value of $\underline{x}_i$

$$= Pr(k_i = k \mid \underline{x}_i) = \frac{Pr(k_i = k) \, f(\underline{x} \mid k_i = k)}{f(\underline{x})} \qquad \text{(Bayes rule)}$$

$$\hat{p}_k(\underline{x}_i) = \frac{\hat{\pi}_k \, f(\underline{x}_i \mid \hat{\underline{\mu}}_k, \hat{\Sigma}_k)}{\sum_{j=1}^{K} \hat{\pi}_j \, f(\underline{x}_i \mid \hat{\underline{\mu}}_j, \hat{\Sigma}_j)} \qquad \begin{array}{l} \text{for } k = 1, 2, \ldots, K \\ i = 1, 2, \ldots, n \end{array}$$

- M-step: using $\hat{p}_k(\underline{x}_i)$'s, update estimates of $\hat{\underline{\mu}}_k, \hat{\Sigma}_k, \hat{\pi}_k$:

$$\hat{\underline{\mu}}_k = \frac{\sum_{i=1}^{n} \hat{p}_k(\underline{x}_i) \, \underline{x}_i}{\sum_{i=1}^{n} \hat{p}_k(\underline{x}_i)}$$

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^{n} \hat{p}_k(\underline{x}_i) (\underline{x}_i - \hat{\underline{\mu}}_k)(\underline{x}_i - \hat{\underline{\mu}}_k)^T}{\sum_{j=1}^{n} \hat{p}_k(\underline{x}_j)}$$

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^{n} \hat{p}_k(\underline{x}_i)$$

- if not converged, return to E-step.

E-step

M-step

L = 1

E-step + M-step

L = 2

L = 5

L = 20