

APPLIED MACHINE LEARNING

Pdf, GMM and E-M

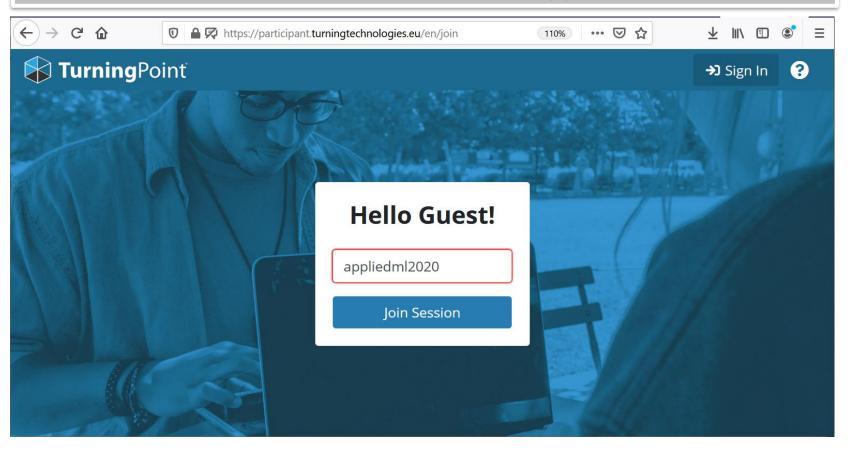
Interactive Lecture



Launch polling system

https://participant.turningtechnologies.eu/en/join

Acces as GUEST and enter the session id: appliedml2020



0.35 0.3 0.25

-0.2 -0.4

-0.6

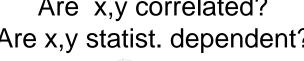
-0.8

-0.5



Statistical Independence and uncorrelatedness

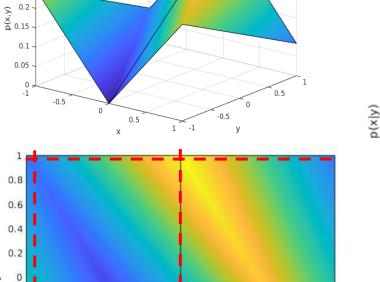
Are x,y correlated? Are x,y statist. dependent?



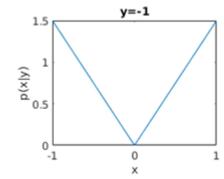
Uncorrelated: $E\{x, y\} = E\{x\} E\{y\}$

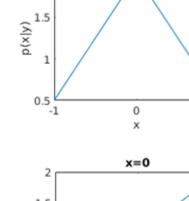
Statistical Ind.: p(x, y) = p(x) p(y)

$$E\{x, y\} = E\{x\} E\{y\} = 0$$

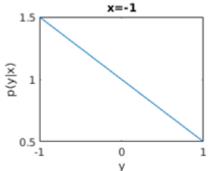


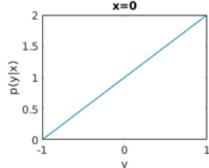
0.5





y=+1







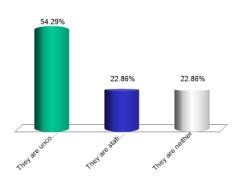
Statistical Independence and uncorrelatedness

Joint probabilities over two variables x_1, x_2

	x ₂ =-1	x ₂ =0	$x_2 = 1$	Total
$x_1 = -1$	3/12	0	3/12	1/2
$x_1=1$	1/12	4/12	1/12	1/2
Total	1/3	1/3	1/3	

Are x_1 and x_2 uncorrelated and statistically independent?

- A. They are uncorrelated
- B. They are statistically independent
- C. They are neither





Statistical Independence and uncorrelatedness

Joint probabilities over two variables x_1, x_2

	x ₂ =-1	x ₂ =0	$\mathbf{x_2} = 1$	Total
$x_1 = -1$	3/12	0	3/12	1/2
$x_1=1$	1/12	4/12	1/12	1/2
Total	1/3	1/3	1/3	

$$E\{x_{1}, x_{2}\} = E\{x_{1}\}E\{x_{2}\}$$

$$\sum_{i,j=1}^{3} (x_{1} = i, y = j) p((x_{1} = i), (y = j)) = \sum_{i=1}^{3} (x_{1} = i) p((x_{1} = i)) \sum_{i=1}^{3} (x_{2} = i) p((x_{2} = i))$$

$$\sum_{i=1}^{3} (x_{1} = i, y = j) p((x_{2} = i)) = \sum_{i=1}^{3} (x_{1} = i) p((x_{2} = i)) p((x_{2} = i))$$

Both sums are zero \Rightarrow x_1, x_2 : uncorrelated $p(x_1 = -1, x_2 = 1) = 3/12 = 0.25$

$$p(x_1 = -1) p(x_2 = 1) = 1/2*1/3 = 0.1667$$



Statistical Independence and uncorrelatedness

Independent



Uncorrelated

$$p(x_1, x_2) = p(x_1) p(x_2) \implies E\{x_1, x_2\} = E\{x_1\} E\{x_2\}$$

$$p(x_1, x_2) = p(x_1) p(x_2)$$
 $\not\leftarrow E\{x_1, x_2\} = E\{x_1\} E\{x_2\}$

Statistical independence ensures uncorrelatedness.

The converse is not true.



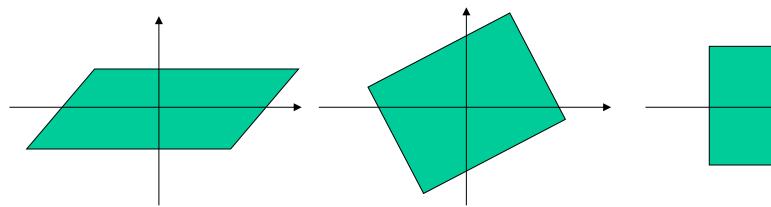
ICA: Preprocessing – Whitening & Independent Component Identification

Original Distribution

Uncorrelated distribution:

$$E\{x_1, x_2\} = E\{x_1\} E\{x_2\}$$

Statistically Indep. Distr.



Whitening preprocessing:

$$E\left\{ XX^{T}\right\} = I$$

After projection on independent components



Linear Correlation

Two variables x_1, x_2 are correlated if:

$$corr(x_1, x_2) = \frac{cov(x_1, x_2)}{var(x_1)var(x_2)} \neq 0$$

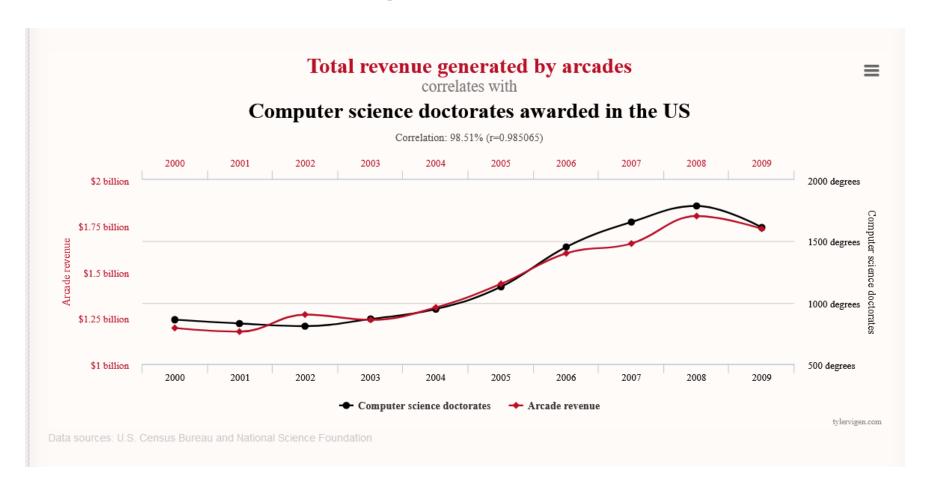
 $corr(x_1, x_2) >< 0$: positive / negative correlation

$$|corr(x_1, x_2)| = 1$$
: perfectly correlated

 $|corr(x_1, x_2)| < 0.5$ weakly correlated

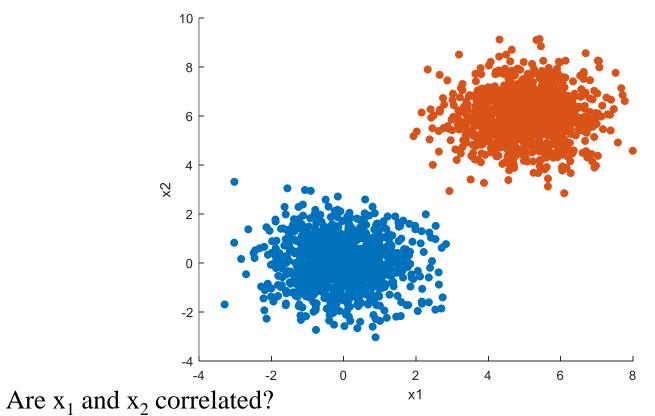


Real and Spurious Correlations

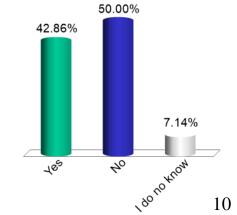






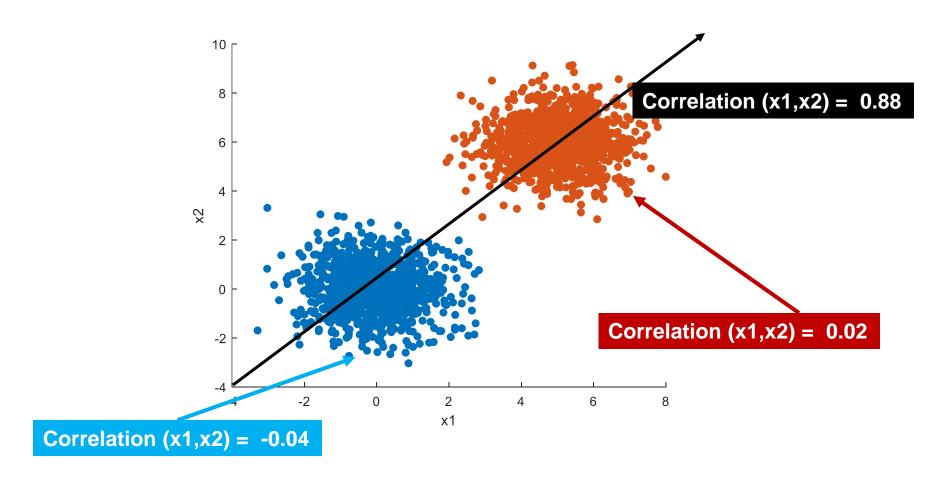


- A. Yes
- B. No
- C. I do no know





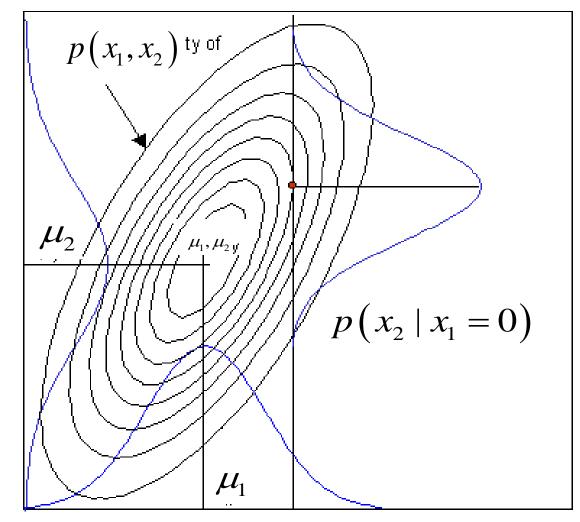
Correlations



Spurious correlations as we compare two groups of data that come from two different distributions

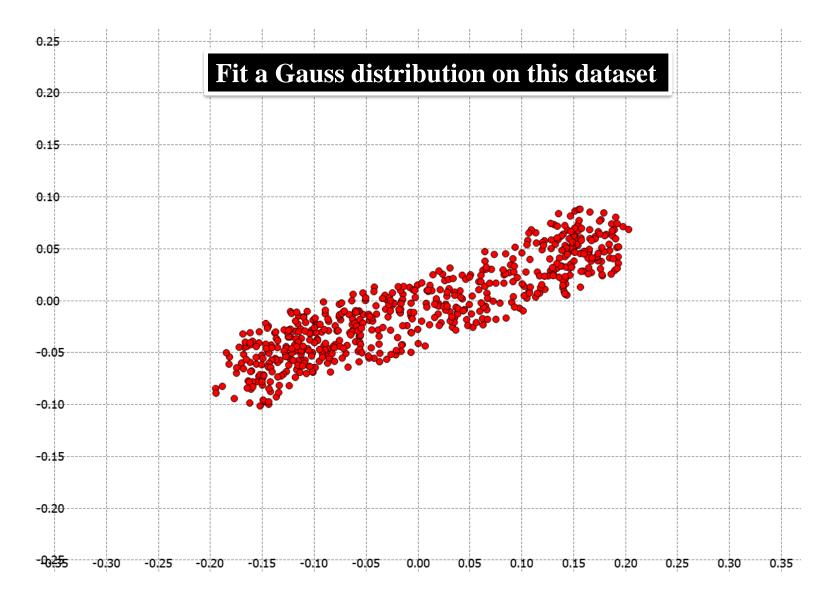


The conditional and marginal pdf of a multi-dimensional Gauss function are all Gauss functions!

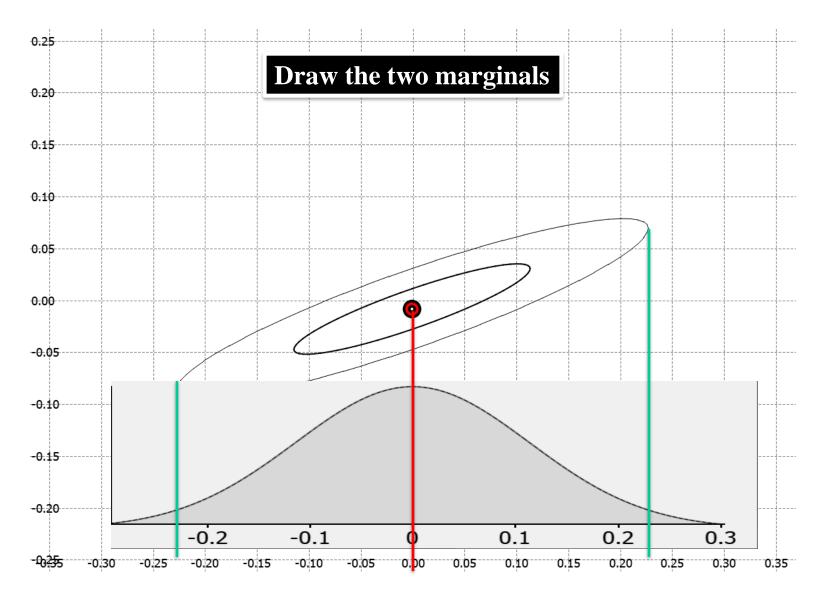


marginal density $p(x_2)$

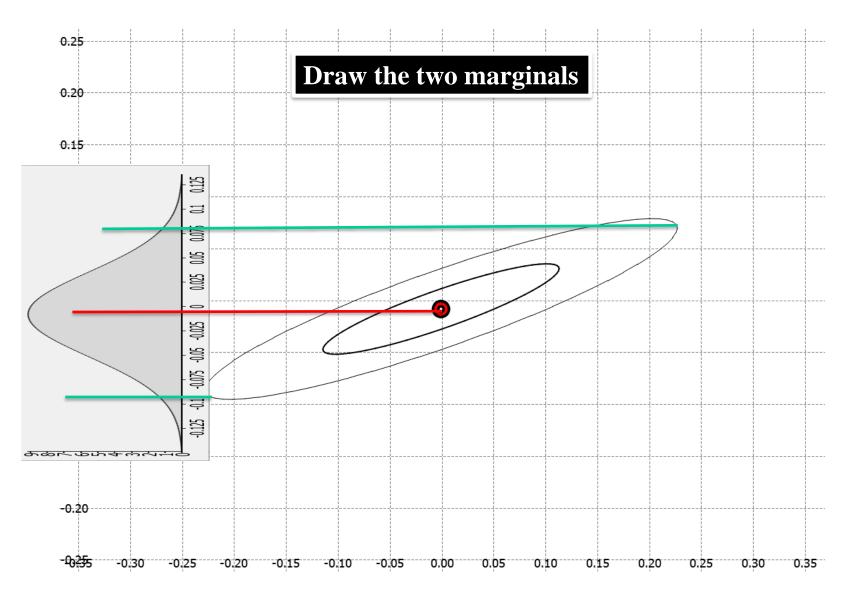










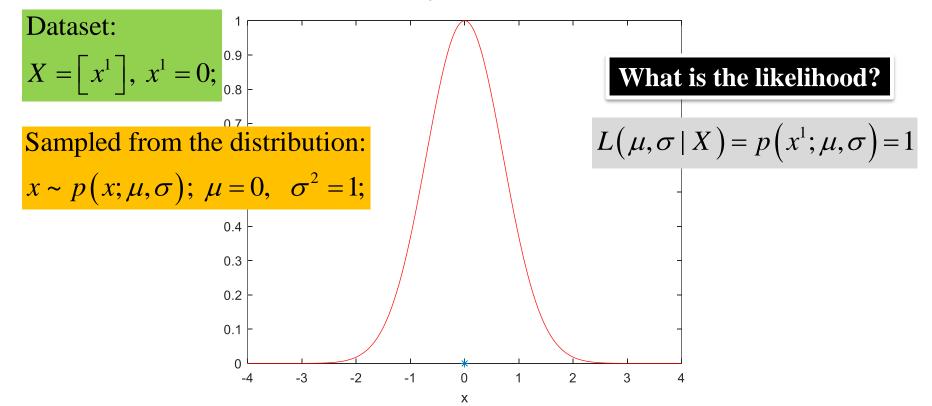




Likelihood

Likelihood for a single unnormalized Gauss function, given a set of M datapoints $\{x^i\}_{i=1}^M$

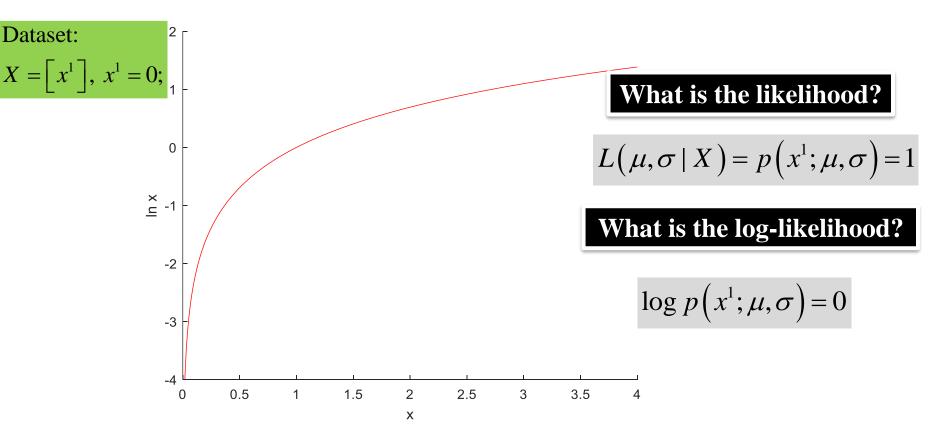
$$L(\mu, \sigma | X) = \prod_{i=1}^{M} p(x^{i}; \mu, \sigma) = e^{\left(-\frac{(x-\mu)^{2}}{2\sigma^{2}}\right)}$$



Here: we consider an un-normalized Gauss function, also called Radial Basis Functions (RBF).



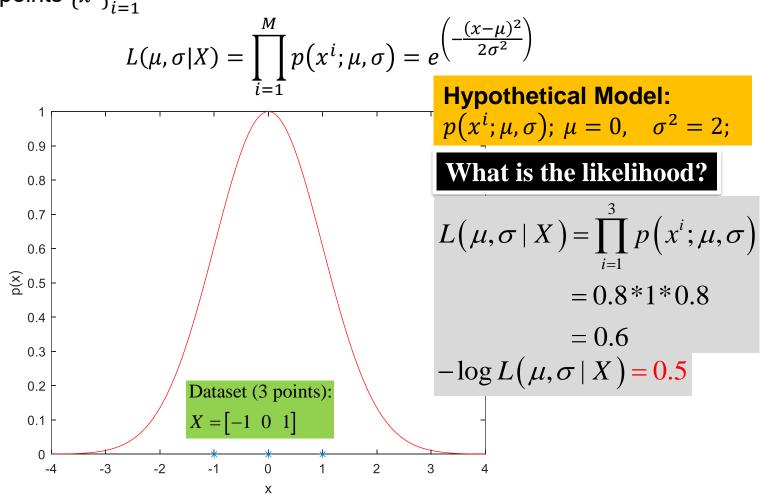
Likelihood





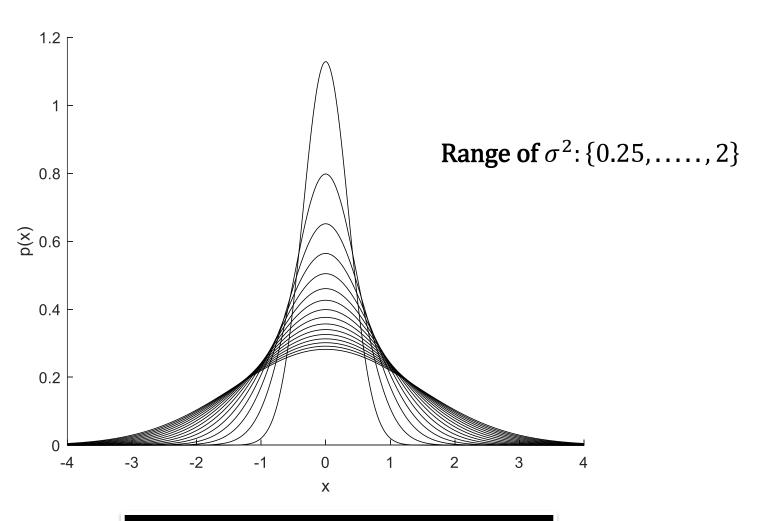
Likelihood

Likelihood for a single unnormalized Gauss function, given a set of M datapoints $\{x^i\}_{i=1}^M$





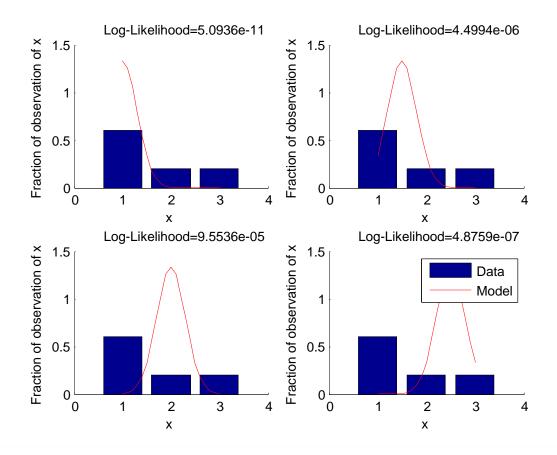
Gauss pdfs



Pdf for the normalized distributions



Likelihood Function for data not gauss distributed

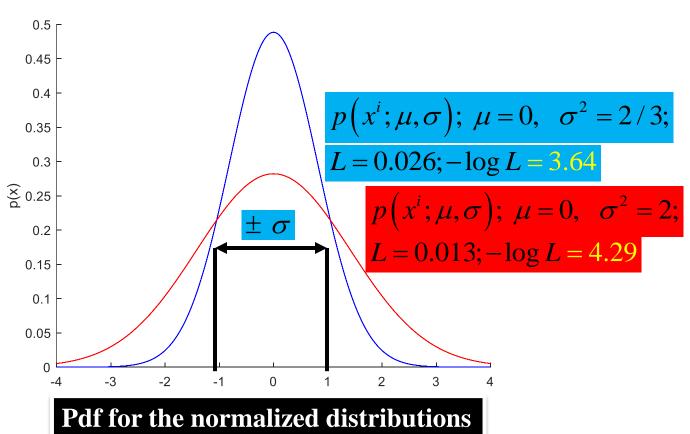


Log-Likelihood for a series of Gauss functions applied to datasets with pdfs that do not follow a Gauss distribution. The Likelihood increases as the fit is closer to the real mean of the data, even if this may appear as a poorer fit.



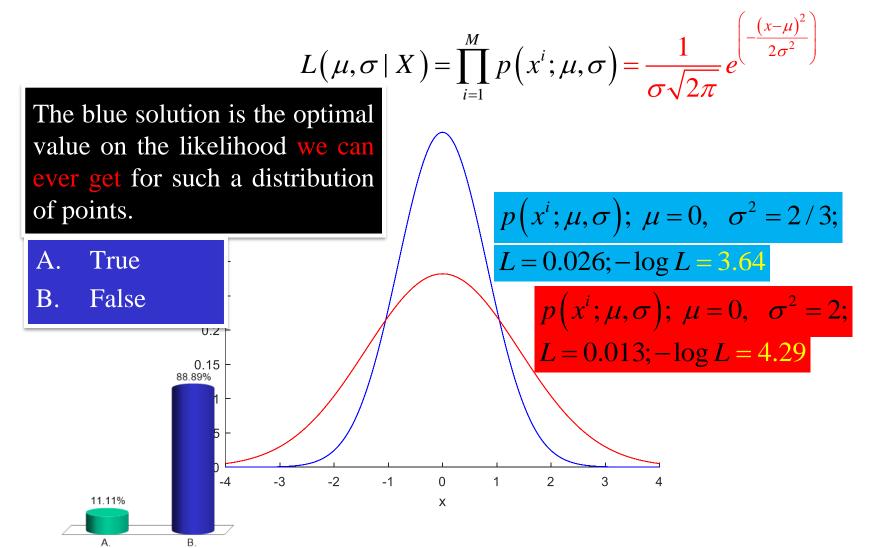
Likelihood for a single Gauss function, given a set of M datapoints $\{x^i\}_{i=1}^M$

$$L(\mu, \sigma \mid X) = \prod_{i=1}^{M} p(x^{i}; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{\left(-\frac{(x-\mu)^{2}}{2\sigma^{2}}\right)}$$



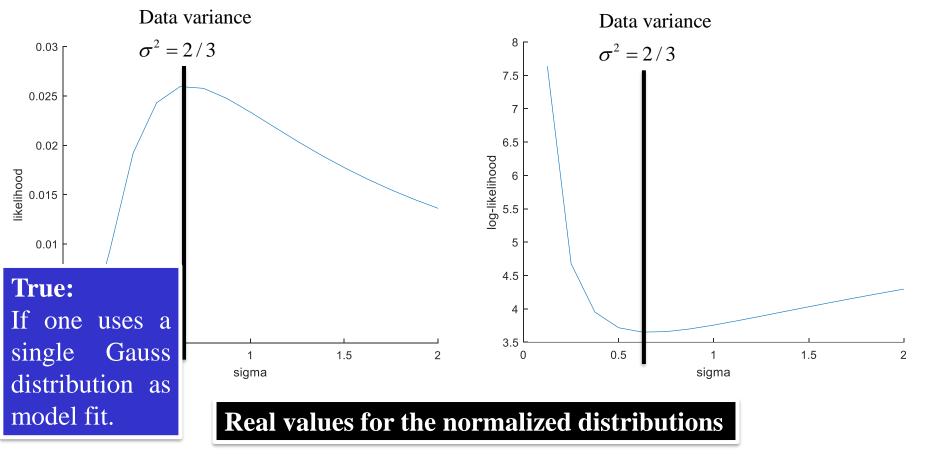


Likelihood for a single Gauss function, given a set of M datapoints $\left\{x^i\right\}_{i=1}^{M}$





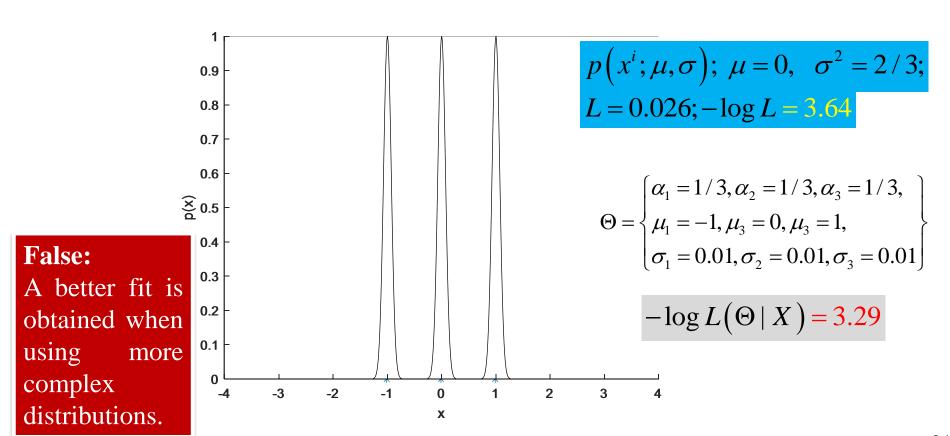
The maximum of the likelihood, minimum of -log-likelihood, is obtained for a distribution with same variance as that of the data.





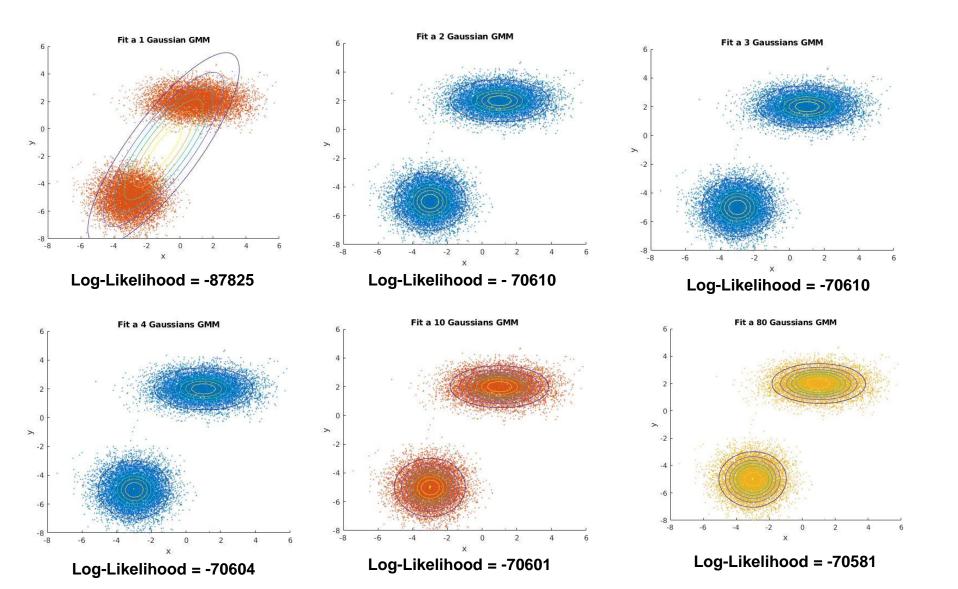
Likelihood for a single Gauss function, given a set of M datapoints $\left\{x^i\right\}_{i=1}^M$

$$L(\Theta \mid X) = \prod_{i=1}^{M} \sum_{k=1}^{3} \alpha_{k} p(x^{i}; \mu_{k}, \sigma_{k})$$



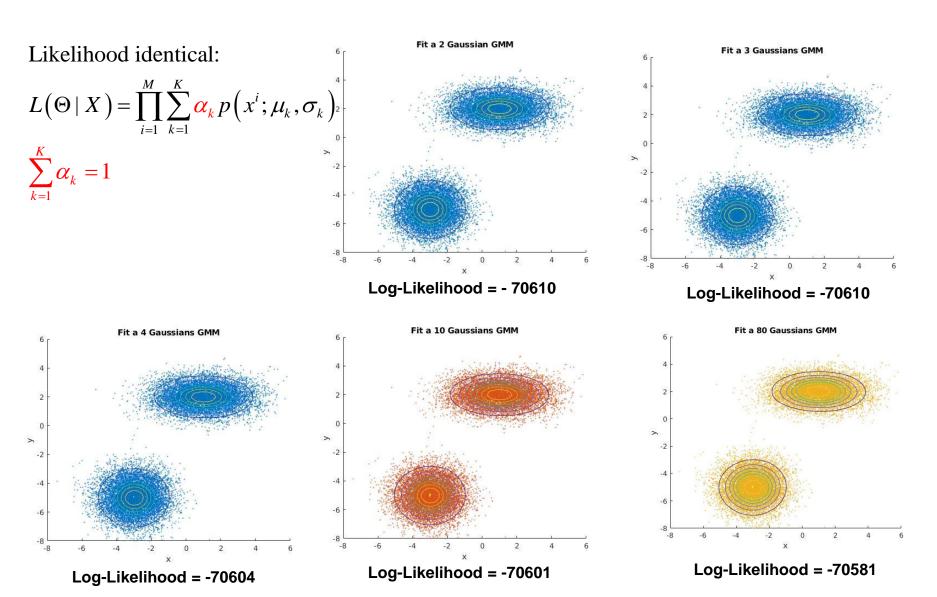


Maximum likelihood with Mixture of Gaussians



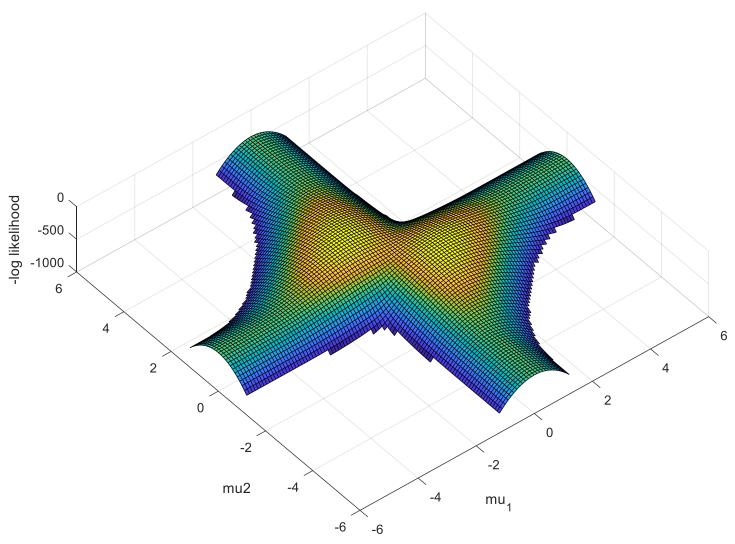


Maximum likelihood with Mixture of Gaussians



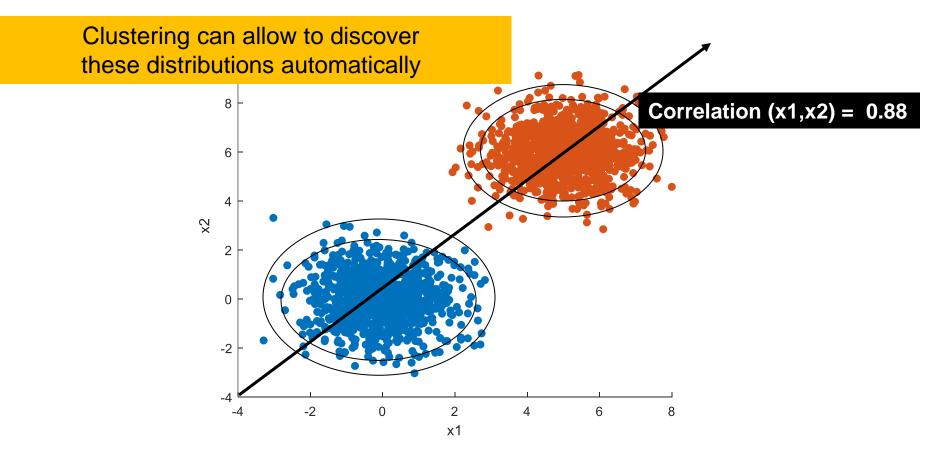


Non-convexity of the likelihood





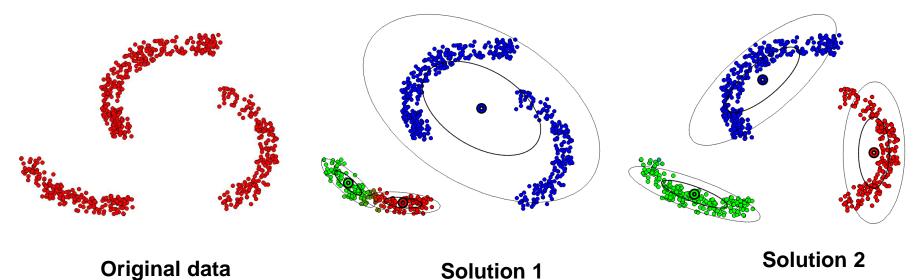
Correlations



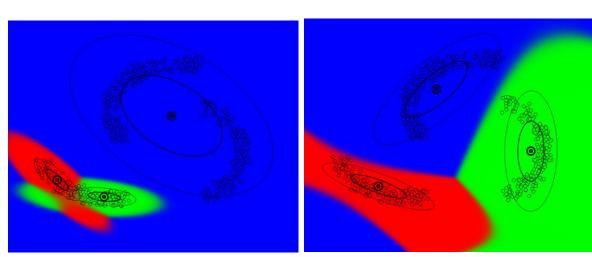
Spurious correlations as we compare two groups of data that come from two different distributions



Clustering with Gaussian Mixture Models



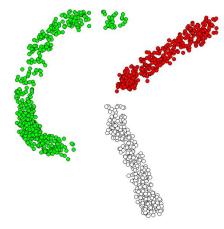
Compute boundary across clusters by comparing likelihood of each cluster, i.e. of each Gauss function.



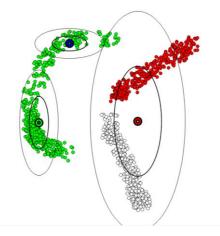
Boundaries

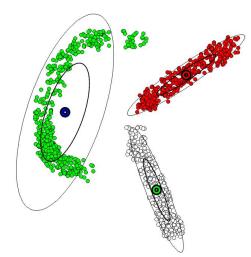


Which Model?



Original data





Solution with full matrices

Diagonal matrices:

$$\Sigma^k = egin{bmatrix} \sigma_1^k & 0 \ 0 & \sigma_2^k \end{bmatrix}$$

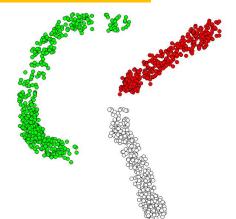
Full matrices:

$$\Sigma^k = egin{bmatrix} oldsymbol{\sigma}_1^k & oldsymbol{\sigma}_{12}^k \ oldsymbol{\sigma}_{21}^k & oldsymbol{\sigma}_2^k \end{bmatrix}$$

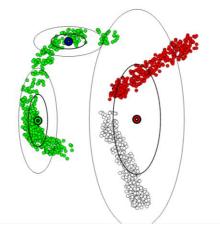


Metrics to choose model

1180 sample points



Original data



Solution with full matrices

Parameters: 3*(2+2)=12

Likelihood: 1999,

AIC: - 2912 BIC: -2972

Solution with diagonal matrices

Parameters: 3*(2+3)=15

Likelihood: 2719,

AIC: - 5419 BIC: -5373

Even if it requires more parameters, the gain on likelihood is important
→ Optimal solution

Which of the two solutions would get the best values on AIC or BIC?

 $AIC = -2\ln(L) + 2B$

 $BIC = -2\ln(L) + \ln(M)B$

B=# Parameters