

FALL17 10-701 Homework 4 Recitation 1

Ellen Vitercik Matthew Oresky Jessica Lin

GMM Question 1

Assume you have points that are generated by one of two possible Gaussian distributions. Which of the following are true?

- a) We know how to get a globally optimal solution by deriving the maximum likelihood estimate analytically.
- b) Using the EM algorithm to solve this problem, we assume that we know from which Gaussian each point originated.
- c) Once the EM algorithm has converged, we know for certain from which Gaussian each point originated.
- d) The EM algorithm for this problem guarantees that the likelihood of the data never decreases from one iteration to the next.

GMM Question 1

Assume you have points that are generated by one of two possible Gaussian distributions. Which of the following are true?

- a) We know how to get a globally optimal solution by deriving the maximum likelihood estimate analytically.
- b) Using the EM algorithm to solve this problem, we assume that we know from which Gaussian each point originated.
- c) Once the EM algorithm has converged, we know for certain from which Gaussian each point originated.
- d) The EM algorithm for this problem guarantees that the likelihood of the data never decreases from one iteration to the next.

Answer

Answer: (d). A - EM doesn't give the globally optimal solution. B - We can start out with one of the Gaussians being more likely for some points, but we don't know for sure. C - After convergence, we only know the probability values of belonging to a particular Gaussian.

GMM Question 2

Which of the following are true about the EM algorithm as applied to a Gaussian Mixture Model?

- a) The choice of initial values of parameters of the Gaussian affects the final estimates.
- b) The algorithm is guaranteed to converge.
- c) The algorithm is guaranteed to converge to a global maxima.
- d) The estimate of the parameters obtained at the end is the Maximum Likelihood Estimate.

GMM Question 2

Which of the following are true about the EM algorithm as applied to a Gaussian Mixture Model?

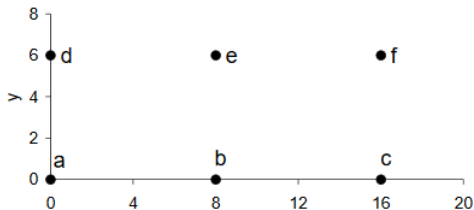
- a) The choice of initial values of parameters of the Gaussian affects the final estimates.
- b) The algorithm is guaranteed to converge.
- c) The algorithm is guaranteed to converge to a global maxima.
- d) The estimate of the parameters obtained at the end is the Maximum Likelihood Estimate.

Answer

A and B are true. C - EM doesn't give the globally optimal solution. D - We cannot solve GMM in closed form to get a clean maximum likelihood expression

K-means Question

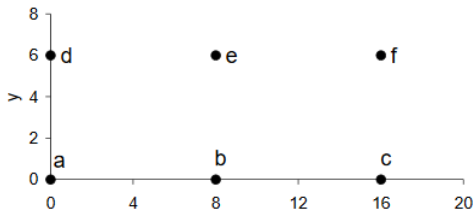
N.B.: Starting cluster centers can be any 3 of the 6 points given.



3-partition	Is it stable?	An example 3-starting configuration that can arrive at the 3-partition after 0 or more iterations of k -means (or write "none" if no such 3-starting configuration)	The number of unique starting configurations that can arrive at the 3-partition
$\{a, b, e\}, \{c, d\}, \{f\}$			
$\{a, b\}, \{d, e\}, \{c, f\}$			
$\{a, d\}, \{b, e\}, \{c, f\}$			
$\{a\}, \{d\}, \{b, c, e, f\}$			
$\{a, b\}, \{d\}, \{c, e, f\}$			
$\{a, b, d\}, \{c\}, \{e, f\}$			

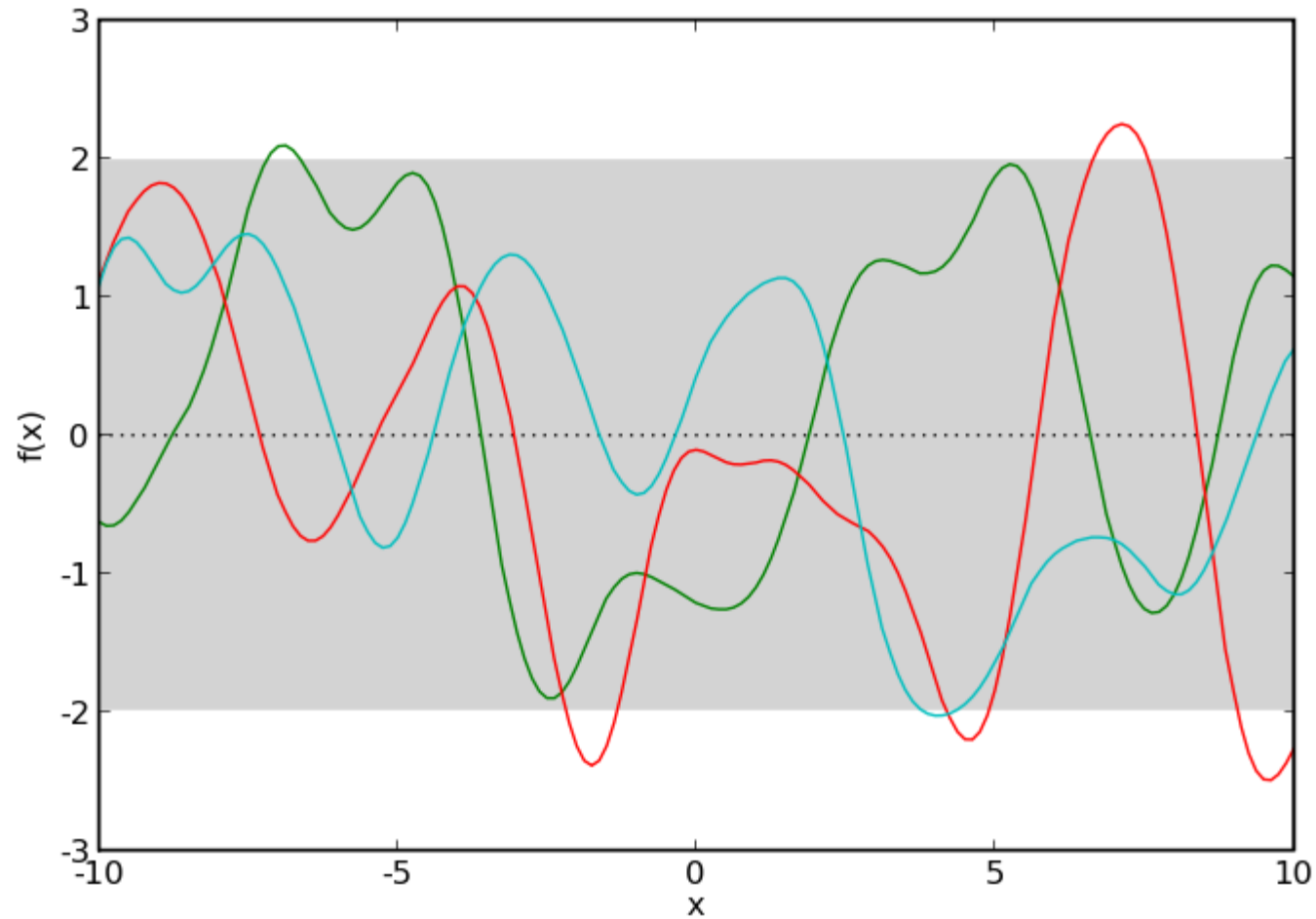
K-means Solution

N.B.: Starting cluster centers can be any 3 of the 6 points given.



3-partition	Stable?	An example 3-starting configuration that can arrive at the 3-partition after 0 or more iterations of k -means (or write “none” if no such 3-starting configuration exists)	# of unique 3-starting configurations that arrive at the 3-partition
$\{a, b, e\}, \{c, d\}, \{f\}$	N	none	0
$\{a, b\}, \{d, e\}, \{c, f\}$	Y	$\{b, c, e\}$	4
$\{a, d\}, \{b, e\}, \{c, f\}$	Y	$\{a, b, c\}$	8
$\{a\}, \{d\}, \{b, c, e, f\}$	Y	$\{a, b, d\}$	2
$\{a, b\}, \{d\}, \{c, e, f\}$	Y	none	0
$\{a, b, d\}, \{c\}, \{e, f\}$	Y	$\{a, c, f\}$	1

A *Gaussian process* is a distribution over functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$.



A *Gaussian process* is a distribution over functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$.

It has a mean function:

$$m(\boldsymbol{x}) = \mathbb{E}_{f \sim D}[f(\boldsymbol{x})]$$

A *Gaussian process* is a distribution over functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$.

It has a mean function:

$$m(\mathbf{x}) = \mathbb{E}_{f \sim D}[f(\mathbf{x})]$$

and a covariance function:

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{f \sim D}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

You need to **two things** to define a Gaussian process.

You need to choose the form of:

1. The mean function $m(\mathbf{x})$
2. The covariance function $k(\mathbf{x}, \mathbf{x}')$



How do I choose these functions?

You've chosen the mean and the covariance functions.

You get a training set of labeled points:

$$\{(\boldsymbol{x}_1, f_1), (\boldsymbol{x}_2, f_2)\}$$

You've chosen the mean and the covariance functions.

You get a training set of labeled points:

$$\{(\mathbf{x}_1, f_1), (\mathbf{x}_2, f_2)\}$$



How do I label a new point \mathbf{x} ?

Let $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

By definition of a Gaussian process,

$$\begin{bmatrix} f_1 \\ f_2 \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}_1) \\ m(\mathbf{x}_2) \\ m(\mathbf{x}) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & k(\mathbf{x}_1, \mathbf{x}_3) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & k(\mathbf{x}_2, \mathbf{x}_3) \\ k(\mathbf{x}_3, \mathbf{x}_1) & k(\mathbf{x}_3, \mathbf{x}_2) & k(\mathbf{x}_3, \mathbf{x}_3) \end{bmatrix} \right)$$



How do I label a new point \mathbf{x} ?

Question.

You decide $m(\mathbf{x}) = \mathbf{0}$ and $k(x, x') = \min\{x, x'\}$.

You see your training set:

$$\begin{aligned}\mathbf{x}_1 &= 1, f_1 = 2 \\ \mathbf{x}_2 &= 2, f_2 = 5\end{aligned}$$

You see a new datapoint $\mathbf{x} = 0.5$. Again, let $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

Write down the joint distribution for:

$$\begin{bmatrix} f_1 \\ f_2 \\ f_* \end{bmatrix}$$

Answer.

You decide $m(\mathbf{x}) = \mathbf{0}$ and $k(x, x') = \min\{x, x'\}$.

You see your training set:

$$\begin{aligned}\mathbf{x}_1 &= 1, f_1 = 2 \\ \mathbf{x}_2 &= 2, f_2 = 5\end{aligned}$$

You see a new datapoint $\mathbf{x} = 0.5$. Again, let $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

Write down the joint distribution for:

$$\begin{bmatrix} f_1 \\ f_2 \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 & 0.5 \\ 1 & 2 & 0.5 \\ 0.5 & 0.5 & 0.5 \end{bmatrix} \right)$$

Let \mathbf{x} be an unlabeled example and $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

$$\text{Let } X = \begin{bmatrix} - & \mathbf{x}_1 & - \\ \vdots & \ddots & \vdots \\ - & \mathbf{x}_m & - \end{bmatrix} \text{ and } f = \begin{bmatrix} f_1 \\ \vdots \\ f_m \end{bmatrix}.$$

$$\text{Then } \begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}_1) \\ \vdots \\ m(\mathbf{x}_m) \\ m(\mathbf{x}) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_m, \mathbf{x}_1) & k(\mathbf{x}, \mathbf{x}_1) \\ \vdots & \ddots & \vdots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_m) & \cdots & k(\mathbf{x}_m, \mathbf{x}_m) & k(\mathbf{x}, \mathbf{x}_m) \\ k(\mathbf{x}_1, \mathbf{x}) & \cdots & k(\mathbf{x}, \mathbf{x}_m) & k(\mathbf{x}, \mathbf{x}) \end{bmatrix} \right)$$

Let \mathbf{x} be an unlabeled example and $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

$$\text{Let } k(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_m, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_m) & \cdots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix} \text{ and } k(X, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}_m) \end{bmatrix}$$

$$\text{Then } \begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}_1) \\ \vdots \\ m(\mathbf{x}_m) \\ m(\mathbf{x}) \end{bmatrix}, \begin{bmatrix} k(X, X) & k(X, \mathbf{x}) \\ k(X, \mathbf{x})^T & k(\mathbf{x}, \mathbf{x}) \end{bmatrix} \right)$$

Let \mathbf{x} be an unlabeled example and $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

$$\text{Let } k(X, X) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_m, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_m) & \cdots & k(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix} \text{ and } k(X, \mathbf{x}) = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}_m) \end{bmatrix}$$

$$\begin{aligned} & f_* \mid X, \mathbf{f}, \mathbf{x} \\ & \sim \mathcal{N} \left(\begin{array}{c} m(\mathbf{x}) + k(X, \mathbf{x})^T k(X, X)^{-1} (\mathbf{f} - m(X)), \\ k(\mathbf{x}, \mathbf{x}) - k(X, \mathbf{x})^T k(X, X)^{-1} k(X, \mathbf{x}) \end{array} \right) \end{aligned}$$

Question.

You decide $m(\mathbf{x}) = \mathbf{0}$ and $k(x, x') = \min\{x, x'\}$.

You see your training set:

$$\mathbf{x}_1 = 1, f_1 = 2$$

$$\mathbf{x}_2 = 2, f_2 = 5$$

You see a new datapoint $\mathbf{x} = 0.5$. Again, let $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

Write down the joint distribution for:

$$f_* \mid X, \mathbf{f}, \mathbf{x}$$

$$\text{Hint: } \begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

Answer.

You decide $m(\mathbf{x}) = \mathbf{0}$ and $k(x, x') = \min\{x, x'\}$.

You see your training set:

$$\mathbf{x}_1 = 1, f_1 = 2$$

$$\mathbf{x}_2 = 2, f_2 = 5$$

You see a new datapoint $\mathbf{x} = 0.5$. Again, let $f_* = f(\mathbf{x})$, where f is sampled from our Gaussian process.

Write down the joint distribution for:

$$f_* \mid X, \mathbf{f}, \mathbf{x} \sim \mathcal{N}(1, 0.25)$$

You need to **two things** to define a Gaussian process.

You need to choose the form of:

1. The mean function $m(\mathbf{x})$
2. The covariance function $k(\mathbf{x}, \mathbf{x}')$



How do I choose these functions?

Squared exponential kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$$

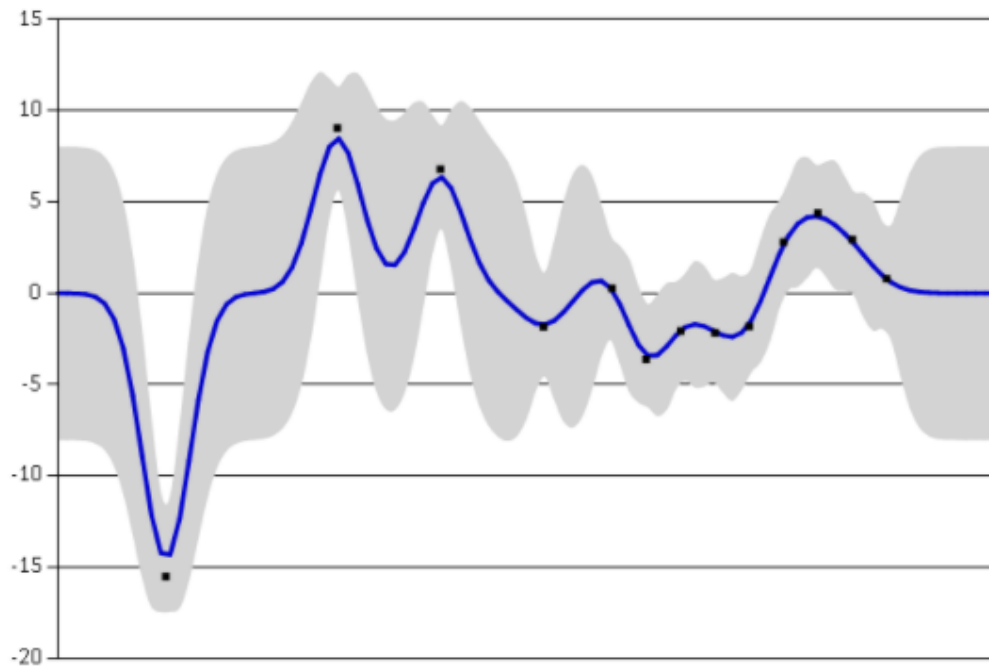
This kernel is **infinitely differentiable**.

It is appropriate for modelling **very smooth** functions.

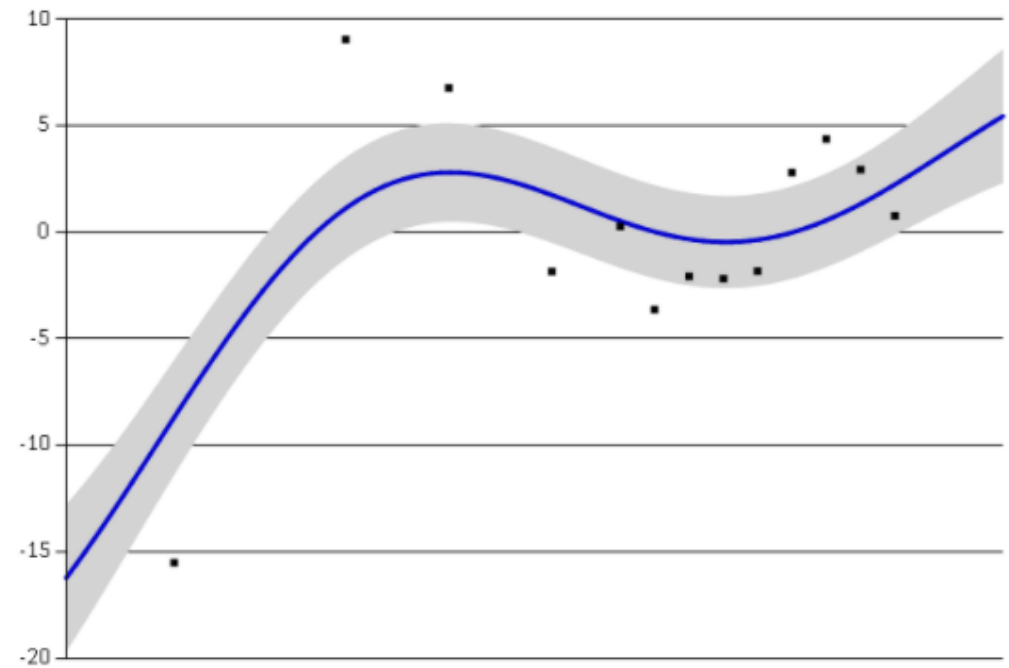
Question.

What happens when ℓ is small? What about when it's large?

As $\ell \rightarrow 0$, $k(\mathbf{x}, \mathbf{x}') \rightarrow 0$. The values look more and more independent.



As $\ell \rightarrow \infty$, $k(\mathbf{x}, \mathbf{x}') \rightarrow 1$. The values look more and more dependent.



Periodic kernel

$$k(x, x') = \exp\left(-\frac{2 \sin(\pi|x - x'|/p)^2}{\ell^2}\right)$$

