

Unit 4 Unsupervised Learning (2

Course > weeks)

> Homework 5 > 3. EM Algorithm

Audit Access Expires May 11, 2020

You lose all access to this course, including your progress, on May 11, 2020.

3. EM Algorithm

Consider the following mixture of two Gaussians:

$$p\left(x; heta
ight) = \pi_{1}\mathcal{N}\left(x;\mu_{1},\sigma_{1}^{2}
ight) + \pi_{2}\mathcal{N}\left(x;\mu_{2},\sigma_{2}^{2}
ight)$$

This mixture has parameters $\theta = \{\pi_1, \pi_2, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2\}$. They correspond to the mixing proportions, means, and variances of each Gaussian. We initialize θ as $\theta_0 = \{0.5, 0.5, 6, 7, 1, 4\}$.

We have a dataset $\mathcal D$ with the following samples of x: $x^{(0)}=-1$, $x^{(1)}=0$, $x^{(2)}=4$, $x^{(3)}=5$, $x^{(4)}=6$.

We want to set our parameters θ such that the data log-likelihood $l\left(\mathcal{D};\theta\right)$ is maximized:

$$\operatorname*{argmax}_{ heta} \sum_{i=0}^{4} \log p\left(x^{(i)}; heta\right).$$

Recall that we can do this with the EM algorithm. The algorithm optimizes a lower bound on the log-likelihood, thus iteratively pushing the data likelihood upwards. The iterative algorithm is specified by two steps applied successively:

1. E-step: infer component assignments from current $heta_0= heta$ (complete the data)

$$p(y = k \mid x^{(i)}) := p(y = k \mid x^{(i)}; \theta_0), \text{ for } k = 1, 2, \text{ and } i = 0, \dots, 4.$$

2. M-step: maximize the expected log-likelihood

$$ilde{l}\left(D; heta
ight) := \sum_{i} \sum_{k} p\left(y = k \mid x^{(i)}
ight) \log rac{p\left(x^{(i)}, y = k; heta
ight)}{p\left(y = k \mid x^{(i)}
ight)}$$

with respect to heta while keeping $p\left(y=k\mid x^{(i)}
ight)$ fixed.

To see why this optimizes a lower bound, consider the following inequality:

$$\begin{split} \log p\left(x;\theta\right) &= \log \sum_{y} p\left(x,y;\theta\right) \\ &= \log \sum_{y} q\left(y|x\right) \frac{p\left(x,y;\theta\right)}{q\left(y|x\right)} \\ &= \log \mathbb{E}_{y \sim q\left(y|x\right)} \left[\frac{p\left(x,y;\theta\right)}{q\left(y|x\right)} \right] \\ &\geq \mathbb{E}_{y \sim q\left(y|x\right)} \left[\log \frac{p\left(x,y;\theta\right)}{q\left(y|x\right)} \right] \\ &= \sum_{y} q\left(y|x\right) \log \frac{p\left(x,y;\theta\right)}{q\left(y|x\right)} \end{split}$$

where the inequality comes from **Jensen's inequality** . EM makes this bound tight for the current setting of θ by setting q(y|x) to be $p(y|x;\theta_0)$.

Note: If you have taken 6.431x Probability–The Science of Uncertainty, you could review the video in Unit 8: Limit Theorems and Classical Statistics, Additional Theoretical Material, 2. Jensen's Inequality.

Likelihood Function

1 point possible (graded)

What is the log-likelihood of the data $l(\mathcal{D}; \theta)$ given the initial setting of θ ? Please round to the nearest tenth.

Note: You will want to write a script to calculate this, using the natural log (np.log) and np.float64 data types.

Answer: -24.5

Solution:

The likelihood can be written as:

$$egin{align} P\left(\mathcal{D}; heta
ight) &= \prod_{i=0}^4 p\left(x; heta
ight) \ &= \prod_{i=0}^4 \pi_1 \mathcal{N}\left(x^{(i)};\mu_1,\sigma_1^2
ight) + \pi_2 \mathcal{N}\left(x^{(i)};\mu_2,\sigma_2^2
ight) \end{aligned}$$

Taking the log gives:

$$l\left(\mathcal{D}; heta
ight) = \sum_{i=0}^{4} \log\left(\pi_{1}\mathcal{N}\left(x^{(i)};\mu_{1},\sigma_{1}^{2}
ight) + \pi_{2}\mathcal{N}\left(x^{(i)};\mu_{2},\sigma_{2}^{2}
ight)
ight)$$

We then evaluate each Gaussian using the standard formulation:

$$\mathcal{N}\left(x;\mu,\sigma^{2}
ight)=rac{1}{\sqrt{2\pi\sigma^{2}}}e^{-rac{\left(x-\mu
ight)^{2}}{2\sigma^{2}}}$$

Submit

You have used 0 of 3 attempts

1 Answers are displayed within the problem

E-Step

1 point possible (graded)

What is the formula for p $(y=k\mid x,\theta)$? Write in terms of π_k , π_1 , π_2 , N_k , N_1 , and N_2 (where $N_k=\mathcal{N}\ (x\mid \mu_k,\sigma_k^2)$).

Answer: (pi_k * N_k) / (pi_1 * N_1 + pi_2 * N_2)

STANDARD NOTATION

Solution:

Following Bayes Rule we have:

$$p\left(y\mid x
ight) = rac{p\left(y
ight)p\left(x\mid y
ight)}{\sum_{y'}p\left(y'
ight)p\left(x|y'
ight)}$$

For this problem, this equates to:

$$p\left(y=k\mid x; heta
ight)=rac{\pi_{k}\mathcal{N}\left(x;\mu_{y},\sigma_{y}^{2}
ight)}{\sum_{i=1}^{2}\pi_{i}\mathcal{N}\left(x;\mu_{i},\sigma_{i}^{2}
ight)}$$

Submit

You have used 0 of 3 attempts

• Answers are displayed within the problem

E-Step Weights

5 points possible (graded)

For each of the given data points say which Gaussian (1 or 2) they are given more weight towards in the first E-step using the given setting of θ_0 . This is, answer 2 if $p\left(y=2\mid x,\theta_0\right)>p\left(y=1\mid x,\theta_0\right)$ and 1 otherwise.

 $x^{(0)}$: Answer: 2

 $x^{(1)}$:

 $x^{(2)}:$ Answer: 2

 $x^{(3)}:$ Answer: 1

 $x^{(4)}:$ Answer: 1

Solution:

Note that x will more likely be assigned to Gaussian 2 (y=2) instead of Gaussian 1 (y=1) when the following is true:

$$\begin{split} \frac{P\left(y=2|x^{(i)},\theta_{0}\right)}{P\left(y=1|x^{(i)},\theta_{0}\right)} &> 1 \\ \Leftrightarrow \frac{P\left(x^{(i)}|y=2\right)P\left(y=2\right)}{P\left(x^{(i)}|y=1\right)P\left(y=1\right)} &> 1 \\ \Leftrightarrow \frac{\frac{1}{\sqrt{(2\pi\sigma_{2}^{2})}} exp\{-\frac{1}{2}(x-\mu_{2})^{2}/\sigma_{2}^{2}\}}{\frac{1}{\sqrt{(2\pi\times4)}} exp\{-\frac{1}{2}(x-\mu_{1})^{2}/\sigma_{1}^{2}\}} &> 1 \\ \Leftrightarrow \frac{\frac{1}{\sqrt{(2\pi\times4)}} exp\{-\frac{1}{2}(x-7)^{2}/4\}}{\frac{1}{\sqrt{(2\pi\times1)}} exp\{-\frac{1}{2}(x-6)^{2}\}} &> 1 \\ \Leftrightarrow \frac{1}{2} exp\{-\frac{1}{2}((x-7)^{2}/4-(x-6)^{2})\} &> 1 \\ \Leftrightarrow \frac{1}{2} exp\{\frac{1}{8}(x-5)\left(3x-19\right)\} &> 1 \\ \Leftrightarrow \log\left(\frac{1}{2}\right) + \frac{1}{8}(x-5)\left(3x-19\right) &> 0 \end{split}$$

The x-intercepts of this parabola are $x_1 \approx 4.1525, x_2 \approx 7.1809$. Thus, we can see that all points $x \in [4.15, 7.18]$ have higher probability under class y=1, and all other points have higher probability under y=2. Thus, $x^{(0)}$, $x^{(1)}$, and $x^{(2)}$ are more likely (but not entirely) assigned to Gaussian 2, and the rest of the points $(x^{(3)}, x^{(4)})$ are more likely (but not entirely) assigned to Gaussian 1.

Submit

You have used 0 of 3 attempts

• Answers are displayed within the problem

M-Step

3. EM Algorithm | Homework 5 | 6.86x Coursewar...

https://courses.edx.org/courses/course-v1:MITx+...

3 points possible (graded)

Fixing $p(y = k \mid x, \theta_0)$, we want to update θ such that our lower bound is maximized.

What is the optimal $\hat{\mu}_k$? Answer in terms of $x^{(1)}$, $x^{(2)}$, and γ_{k1} , γ_{k2} , which are defined to be $\gamma_{ki}=p\,(y=k\mid x^{(i)};\theta_0)$

(For ease of input, use subscripts instead superscripts, i.e. type x_i for $x^{(i)}$. Type $_{ t gamma_ki}$ for γ_{ki} .)

Answer: $(gamma_k1 * x_1 + gamma_k2 * x_2) / (gamma_k1 + gamma_k2)$

What is the optimal $\hat{\sigma}_k^2$? Answer in terms of $x^{(1)}$, $x^{(2)}$, γ_{k1} and γ_{k2} , which are defined as above to be $\gamma_{ki}=p\left(y=k\mid x^{(i)};\theta_0\right)$, and $\hat{\mu}_k$.

(Type hatmu_k for $\hat{\mu}_k$. As above, for ease of input, use subscripts instead superscripts, i.e. type x_i for $x^{(i)}$. Type gamma_ki for γ_{ki} .)

Answer: $(gamma_k1 * (x_1 - hatmu_k)^2 + gamma_k2 * (x_2 - hatmu_k)^2) / (gamma_k1 + gamma_k2)$

What is the optimal $\hat{\pi}_k$? Answer in terms of γ_{k1} and γ_{k2} , which are defined as above to be $\gamma_{ki}=p$ $(y=k\mid x^{(i)}; heta_0)$,

(As above, type gamma_ki for γ_{ki} .)

Note: that you must account for the constraint that $\pi_1 + \pi_2 = 1$ where $\pi_1, \pi_2 \geq 0$.

Note: If you know that some aspect of your formula equals an exact constant, simplify and use this number, i.e. $\gamma_{11}+\gamma_{21}=1$.

Answer: (gamma_k1 + gamma_k2) / 2

STANDARD NOTATION

Solution:

The function we are optimizing is now:

$$\sum_{i}\sum_{k}\gamma_{ki}\log\left(\pi_{k}\mathcal{N}\left(x^{(i)};\mu_{k},\sigma_{k}^{2}
ight)
ight)$$

Taking $\frac{\partial}{\partial \mu_k}$ and setting to 0 gives:

3. EM Algorithm | Homework 5 | 6.86x Coursewar...

https://courses.edx.org/courses/course-v1:MITx+...

$$\begin{split} \frac{\partial}{\partial \mu_k} \sum_i \sum_k \gamma_{ki} \log \left(\pi_k \mathcal{N} \left(x^{(i)}; \mu_k, \sigma_k^2 \right) \right) &= \sum_i \gamma_{ki} \frac{\partial}{\partial \mu_k} \log \left(\pi_k \mathcal{N} \left(x^{(i)}; \mu_k, \sigma_k^2 \right) \right) \\ &= \sum_i \gamma_{ki} \frac{\partial}{\partial \mu_k} \left(\log \left(\frac{1}{\sqrt{2\pi\sigma_k^2}} \right) - \frac{\left(x^{(i)} - \mu_k \right)^2}{2\sigma_k^2} \right) \\ &= \sum_i \gamma_{ki} \frac{x^{(i)} - \mu_k}{\sigma_k^2} = 0 \end{split}$$

Separating out μ_k gives:

$$\mu_k = rac{\sum_i \gamma_{ki} x^{(i)}}{\sum_i \gamma_{ki}}$$

We can interpret this as a weighted average of the data points, normalized by the "total mass" assigned to Gaussian k. The weight is the probability that point $x^{(i)}$ "belongs" to Gaussian k.

Solving for σ_k^2 is similar:

$$\begin{split} \frac{\partial}{\partial \sigma_k^2} \sum_i \sum_k \gamma_{ki} \log \left(\pi_k \mathcal{N} \left(x^{(i)}; \mu_k, \sigma_k^2 \right) \right) &= \sum_i \gamma_{ki} \frac{\partial}{\partial \sigma_k^2} \log \left(\pi_k \mathcal{N} \left(x^{(i)}; \mu_k, \sigma_k^2 \right) \right) \\ &= \sum_i \gamma_{ki} \frac{\partial}{\partial \sigma_k^2} (\log \left(\frac{1}{\sqrt{2\pi\sigma_k^2}} \right) - \frac{\left(x^{(i)} - \mu_k \right)^2}{2\sigma_k^2}) \\ &= \sum_i \gamma_{ki} \left(-\frac{1}{2\sigma_k^2} + \frac{\left(x^{(i)} - \mu_k \right)^2}{2\sigma_k^4} \right) = 0 \end{split}$$

Separating out σ_k^2 gives:

$$\sigma_k^2 = rac{\sum_i \gamma_{ki} {(x^{(i)} - \mu_k)}^2}{\sum_i \gamma_{ki}}$$

Finally we solve for π_k while including a lagrange multiplier for the constraint that $\sum_k \pi_k = 1$.

$$egin{aligned} rac{\partial}{\partial \pi_k} \sum_i \sum_k \gamma_{ki} \log \left(\pi_k \mathcal{N}\left(x^{(i)}; \mu_k, \sigma_k^2
ight)
ight) + \lambda \left(\sum_k \pi_k - 1
ight) \ &= \sum_i \gamma_{ki} rac{\partial}{\partial \pi_k} \log \left(\pi_k
ight) + rac{\partial}{\partial \pi_k} \lambda \left(\sum_k \pi_k - 1
ight) \ &= rac{\sum_i \gamma_{ki}}{\pi_k} + \lambda = 0 \end{aligned}$$

Giving
$$\pi_k = -rac{\sum_i \gamma_{ki}}{\lambda}.$$

Solving for λ gives:

$$rac{\partial}{\partial \lambda} \sum_{i} \sum_{k} \gamma_{ki} \log \left(\pi_{k} \mathcal{N} \left(x^{(i)}; \mu_{k}, \sigma_{k}^{2}
ight)
ight) + \lambda \left(\sum_{k} \pi_{k} - 1
ight) = \sum_{k} \pi_{k} - 1 = 0$$

Combining the two gives:

$$\lambda = -\sum_i \sum_k \gamma_{ki}$$

which we recognize as N, the total number of points. Thus $\hat{\pi}_k$ is $rac{\sum_i \gamma_{ki}}{N}$.

Submit

You have used 0 of 3 attempts

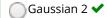
1 Answers are displayed within the problem

Training 1

1 point possible (graded)

In the first M-step, which Gaussian will shift to the left more (relatively)?

Gaussian 1



Solution:

Intuitively, Gaussian 2 is influenced most by the points $x^{(0)}$, $x^{(1)}$, and so it will move to the left. Gaussian 1 will be more influenced by the points at $x^{(2)}$, $x^{(3)}$ and $x^{(4)}$ and so it will not move very much to the left. If we computed the actual values, we would see that the updated means for the two Gaussians are approximately $\mu_1=5.1317$ and $\mu_2=1.4710$.

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem

Training 2

1 point possible (graded)

In the first M-step, which Gaussian's variance will increase more (relatively)?

Gaussian 1

Gaussian 2 🗸

Solution:

Intuitively, the variance of Gaussian 2 spreads out to cover points $x^{(0)}$ and $x^{(1)}$ which it is most influenced by. The 3 points which most influence Gaussian 1 are concentrated around its mean, we would not expect the variance to increase. Numerically, σ_1 decreases to approximately 0.7846 while σ_2 increases to 2.6395.

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem

Training 3

1 point possible (graded)

After convergence, which variance will be larger?





Solution:

Gaussian 1 will be centered around the cluster of 3 points on the right, while Gaussian 2 will be centered around the 2 points on the left. Gaussian 1 will have larger variance because of the larger spread of the right cluster.

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem

Discussion

Topic: Unit 4 Unsupervised Learning (2 weeks): Homework 5 / 3. EM Algorithm

Hide Discussion

Add a Post

Sho	ow all posts by recent a	ictivity
Ą	[staff] EM Algorithm at the top of the page – 7 questions Pinned	9
∀	MLE or MAP? Since we are using Bayes to get the estimator, aren't we using MAP (Maximum A Posteriori) estimation? it just happens that MLE, in this case,	6
Q	Likelihood Functionstupid question incomming	8
2	Implementation in code It was supposed that we should solved the numeric parts by hand or using a programming language? Because at least for me the calculation	2

3. EM Algorithm | Homework 5 | 6.86x Coursewar...

https://courses.edx.org/courses/course-v1:MITx+...

?	<u>Training 3 correct answer review.</u> Lgot every single answer right with the petty excel sheet I made. Attempted training 3 quest first and got it wrong and in confusion marked tr	4
2	Hint By the way, you can "cheat" this problem by doing the derivation the same way it was taught in the lecture.	1
?	<u>Likelihood Function</u> <u>I somehow get wrong answers although I implemented the code as specified in the formulas, any suggestions by my fellow students? I get d</u>	12
?	Likelihood lower bound	2
?	What is E - Step Question is asking?	4
2	M-step part 3 [STAFF]	7
∀	Training 2 and 3: I need hints. Training 2: My variances decrease, not increase. Are my calculations wrong? Training 3: How can I predict what will happen after convergenc	6
2	norm how to input norm	2
∀	Don't we need norm In second question of M step, don't we need the norm of vector? Why there is no expalanation about using the norm?	2

© All Rights Reserved