

CS 229, Summer 2022 Problem Set #3 Solutions

Musa Dildar Ahmed Cheema (`musadac`)

Due Monday, August 8 at 11:59 pm on Gradescope.

Notes: (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Ed forum, at <https://edstem.org/us/courses/23539>. (3) This quarter, Summer 2022, all homework assignments must be submitted individually. If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on the course website before starting work. (4) For the coding problems, you may not use any libraries except those defined in the provided `environment.yml` file. In particular, ML-specific libraries such as scikit-learn are not permitted. (5) To account for late days, the due date is Monday, August 8 at 11:59 pm. If you submit after Monday, August 8 at 11:59 pm, you will begin consuming your late days. If you wish to submit on time, submit before Monday, August 8 at 11:59 pm.

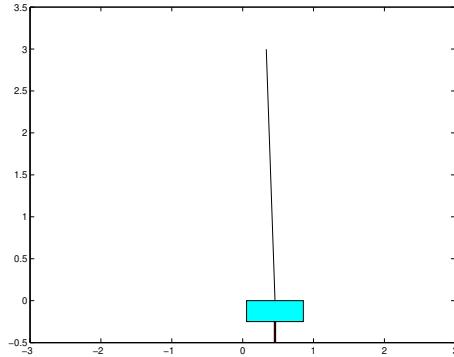
All students must submit an electronic PDF version of the written questions. We highly recommend typesetting your solutions via L^AT_EX. All students must also submit a zip file of their source code to Gradescope, which should be created using the `make_zip.py` script. You should make sure to (1) restrict yourself to only using libraries included in the `environment.yml` file, and (2) make sure your code runs without errors. Your submission may be evaluated by the auto-grader using a private test set, or used for verifying the outputs reported in the writeup.

1. [25 points] Reinforcement Learning: The inverted pendulum

In this problem, you will apply reinforcement learning to automatically design a policy for a difficult control task, without ever using any explicit knowledge of the dynamics of the underlying system.

The problem we will consider is the inverted pendulum or the pole-balancing problem.¹

Consider the figure shown. A thin pole is connected via a free hinge to a cart, which can move laterally on a smooth table surface. The controller is said to have failed if either the angle of the pole deviates by more than a certain amount from the vertical position (i.e., if the pole falls over), or if the cart's position goes out of bounds (i.e., if it falls off the end of the table). Our objective is to develop a controller to balance the pole with these constraints, by appropriately having the cart accelerate left and right.



We have written a simple simulator for this problem. The simulation proceeds in discrete time cycles (steps). The state of the cart and pole at any time is completely characterized by 4 parameters: the cart position x , the cart velocity \dot{x} , the angle of the pole θ measured as its deviation from the vertical position, and the angular velocity of the pole $\dot{\theta}$. Since it would be simpler to consider reinforcement learning in a discrete state space, we have approximated the state space by a discretization that maps a state vector $(x, \dot{x}, \theta, \dot{\theta})$ into a number from 0 to `NUM_STATES-1`. Your learning algorithm will need to deal only with this discretized representation of the states.

At every time step, the controller must choose one of two actions - push (accelerate) the cart right, or push the cart left. (To keep the problem simple, there is no *do-nothing* action.) These are represented as actions 0 and 1 respectively in the code. When the action choice is made, the simulator updates the state parameters according to the underlying dynamics, and provides a new discretized state.

We will assume that the reward $R(s)$ is a function of the current state only. When the pole angle goes beyond a certain limit or when the cart goes too far out, a negative reward is given, and the system is reinitialized randomly. At all other times, the reward is zero. Your program must learn to balance the pole using only the state transitions and rewards observed.

The files for this problem are in `src/cartpole/` directory. Most of the the code has already been written for you, and you need to make changes only to `cartpole.py` in the places specified. This file can be run to show a display and to plot a learning curve at the end. Read the comments at the top of the file for more details on the working of the simulation.

¹The dynamics are adapted from <http://www-anw.cs.umass.edu/rll/domains.html>

To solve the inverted pendulum problem, you will estimate a model (i.e., transition probabilities and rewards) for the underlying MDP, solve Bellman's equations for this estimated MDP to obtain a value function, and act greedily with respect to this value function.

Briefly, you will maintain a current model of the MDP and a current estimate of the value function. Initially, each state has estimated reward zero, and the estimated transition probabilities are uniform (equally likely to end up in any other state).

During the simulation, you must choose actions at each time step according to some current policy. As the program goes along taking actions, it will gather observations on transitions and rewards, which it can use to get a better estimate of the MDP model. Since it is inefficient to update the whole estimated MDP after every observation, we will store the state transitions and reward observations each time, and update the model and value function/policy only periodically. Thus, you must maintain counts of the total number of times the transition from state s_i to state s_j using action a has been observed (similarly for the rewards). Note that the rewards at any state are deterministic, but the state transitions are not because of the discretization of the state space (several different but close configurations may map onto the same discretized state).

Each time a failure occurs (such as if the pole falls over), you should re-estimate the transition probabilities and rewards as the average of the observed values (if any). Your program must then use value iteration to solve Bellman's equations on the estimated MDP, to get the value function and new optimal policy for the new model. For value iteration, use a convergence criterion that checks if the maximum absolute change in the value function on an iteration exceeds some specified tolerance.

Finally, assume that the whole learning procedure has converged once several consecutive attempts (defined by the parameter `NO_LEARNING_THRESHOLD`) to solve Bellman's equation all converge in the first iteration. Intuitively, this indicates that the estimated model has stopped changing significantly.

The code outline for this problem is already in `cartpole.py`, and you need to write code fragments only at the places specified in the file. There are several details (convergence criteria etc.) that are also explained inside the code. Use a discount factor of $\gamma = 0.995$.

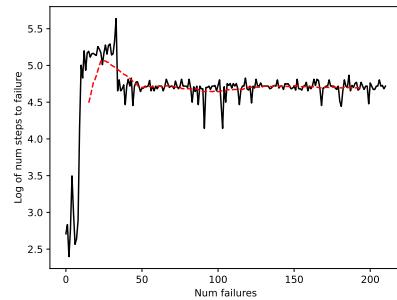
Implement the reinforcement learning algorithm as specified, and run it.

- How many trials (how many times did the pole fall over or the cart fall off) did it take before the algorithm converged? Hint: if your solution is correct, on the plot the red line indicating smoothed log num steps to failure should start to flatten out at about 60 iterations.
- Plot a learning curve showing the number of time-steps for which the pole was balanced on each trial. Python starter code already includes the code to plot. Include it in your submission.
- Find the line of code that says `np.random.seed`, and rerun the code with the seed set to 1, 2, and 3. What do you observe? What does this imply about the algorithm?

s Answer:

a) It takes 211 Failures before it converges.

b)



c) Seeds: 1 = 499 Failures

2 = 355 Failures

3 = 247 Failures

By looking at this very number of Failures and more, we can see that np and used in code generates random 1 and 0 which is set by this so this is will give different Failures if seed is not set.

2. [15 points] KL divergence and Maximum Likelihood

The Kullback-Leibler (KL) divergence is a measure of how much one probability distribution is different from a second one. It is a concept that originated in Information Theory, but has made its way into several other fields, including Statistics, Machine Learning, Information Geometry, and many more. In Machine Learning, the KL divergence plays a crucial role, connecting various concepts that might otherwise seem unrelated.

In this problem, we will introduce KL divergence over discrete distributions, practice some simple manipulations, and see its connection to Maximum Likelihood Estimation.

The *KL divergence* between two discrete-valued distributions $P(X), Q(X)$ over the outcome space \mathcal{X} is defined as follows²:

$$D_{KL}(P\|Q) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}$$

For notational convenience, we assume $P(x) > 0, \forall x$. (One other standard thing to do is to adopt the convention that “ $0 \log 0 = 0$.”) Sometimes, we also write the KL divergence more explicitly as $D_{KL}(P||Q) = D_{KL}(P(X)||Q(X))$.

Background on Information Theory

Before we dive deeper, we give a brief (optional) Information Theoretic background on KL divergence. While this introduction is not necessary to answer the assignment question, it may help you better understand and appreciate why we study KL divergence, and how Information Theory can be relevant to Machine Learning.

We start with the *entropy* $H(P)$ of a probability distribution $P(X)$, which is defined as

$$H(P) = - \sum_{x \in \mathcal{X}} P(x) \log P(x).$$

Intuitively, entropy measures how dispersed a probability distribution is. For example, a uniform distribution is considered to have very high entropy (i.e. a lot of uncertainty), whereas a distribution that assigns all its mass on a single point is considered to have zero entropy (i.e. no uncertainty). Notably, it can be shown that among continuous distributions over \mathbb{R} , the Gaussian distribution $\mathcal{N}(\mu, \sigma^2)$ has the highest entropy (highest uncertainty) among all possible distributions that have the given mean μ and variance σ^2 .

To further solidify our intuition, we present motivation from communication theory. Suppose we want to communicate from a source to a destination, and our messages are always (a sequence of) discrete symbols over space \mathcal{X} (for example, \mathcal{X} could be letters $\{a, b, \dots, z\}$). We want to construct an encoding scheme for our symbols in the form of sequences of binary bits that are transmitted over the channel. Further, suppose that in the long run the frequency of occurrence of symbols follow a probability distribution $P(X)$. This means, in the long run, the fraction of times the symbol x gets transmitted is $P(x)$.

A common desire is to construct an encoding scheme such that the average number of bits per symbol transmitted remains as small as possible. Intuitively, this means we want very frequent symbols to be assigned to a bit pattern having a small number of bits. Likewise, because we are

²If P and Q are densities for continuous-valued random variables, then the sum is replaced by an integral, and everything stated in this problem works fine as well. But for the sake of simplicity, in this problem we'll just work with this form of KL divergence for probability mass functions/discrete-valued distributions.

interested in reducing the average number of bits per symbol in the long term, it is tolerable for infrequent words to be assigned to bit patterns having a large number of bits, since their low frequency has little effect on the long term average. The encoding scheme can be as complex as we desire, for example, a single bit could possibly represent a long sequence of multiple symbols (if that specific pattern of symbols is very common). The entropy of a probability distribution $P(X)$ is its optimal bit rate, i.e., the lowest average bits per message that can possibly be achieved if the symbols $x \in \mathcal{X}$ occur according to $P(X)$. It does not specifically tell us *how* to construct that optimal encoding scheme. It only tells us that no encoding can possibly give us a lower long term bits per message than $H(P)$.

To see a concrete example, suppose our messages have a vocabulary of $K = 32$ symbols, and each symbol has an equal probability of transmission in the long term (i.e, uniform probability distribution). An encoding scheme that would work well for this scenario would be to have $\log_2 K$ bits per symbol, and assign each symbol some unique combination of the $\log_2 K$ bits. In fact, it turns out that this is the most efficient encoding one can come up with for the uniform distribution scenario.

It may have occurred to you by now that the long term average number of bits per message depends only on the frequency of occurrence of symbols. The encoding scheme of scenario A can in theory be reused in scenario B with a different set of symbols (assume equal vocabulary size for simplicity), with the same long term efficiency, as long as the symbols of scenario B follow the same probability distribution as the symbols of scenario A. It might also have occurred to you, that reusing the encoding scheme designed to be optimal for scenario A, for messages in scenario B having a *different probability* of symbols, will always be suboptimal for scenario B. To be clear, we do not need know *what* the specific optimal schemes are in either scenarios. As long as we know the distributions of their symbols, we can say that the optimal scheme designed for scenario A will be suboptimal for scenario B if the distributions are different.

Concretely, if we reuse the optimal scheme designed for a scenario having symbol distribution $Q(X)$, into a scenario that has symbol distribution $P(X)$, the long term average number of bits per symbol achieved is called the *cross entropy*, denoted by $H(P, Q)$:

$$H(P, Q) = - \sum_{x \in \mathcal{X}} P(x) \log Q(x).$$

To recap, the entropy $H(P)$ is the best possible long term average bits per message (optimal) that can be achieved under a symbol distribution $P(X)$ by using an encoding scheme (possibly unknown) specifically designed for $P(X)$. The cross entropy $H(P, Q)$ is the long term average bits per message (suboptimal) that results under a symbol distribution $P(X)$, by reusing an encoding scheme (possibly unknown) designed to be optimal for a scenario with symbol distribution $Q(X)$.

Now, KL divergence is the penalty we pay, as measured in average number of bits, for using the optimal scheme for $Q(X)$, under the scenario where symbols are actually distributed as $P(X)$. It is straightforward to see this

$$\begin{aligned} D_{KL}(P \| Q) &= \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)} \\ &= \sum_{x \in \mathcal{X}} P(x) \log P(x) - \sum_{x \in \mathcal{X}} P(x) \log Q(x) \\ &= H(P, Q) - H(P). \quad (\text{difference in average number of bits.}) \end{aligned}$$

If the cross entropy between P and Q is $H(P)$ (and hence $D_{KL}(P||Q) = 0$) then it necessarily means $P = Q$. In Machine Learning, it is a common task to find a distribution Q that is “close” to another distribution P . To achieve this, it is common to use $D_{KL}(Q||P)$ as the loss function to be optimized. As we will see in this question below, Maximum Likelihood Estimation, which is a commonly used optimization objective, turns out to be equivalent to minimizing the KL divergence between the training data (i.e. the empirical distribution over the data) and the model.

Now, we get back to showing some simple properties of KL divergence.

- (a) [5 points] **Nonnegativity.**

Prove the following:

$$\forall P, Q. \quad D_{KL}(P||Q) \geq 0$$

and

$$D_{KL}(P||Q) = 0 \quad \text{if and only if} \quad P = Q.$$

[Hint: You may use the following result, called **Jensen's inequality**. If f is a convex function, and X is a random variable, then $E[f(X)] \geq f(E[X])$. Moreover, if f is strictly convex (f is convex if its Hessian satisfies $H \geq 0$; it is *strictly* convex if $H > 0$; for instance $f(x) = -\log x$ is strictly convex), then $E[f(X)] = f(E[X])$ implies that $X = E[X]$ with probability 1; i.e., X is actually a constant.]

Answer:

$$\begin{aligned} D_{KL}(P||Q) &= -\sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= \sum_x P(x) \log \frac{P(x)}{E_p[\frac{Q(x)}{P(x)}]} \\ &= \log \sum_x P(x) \frac{Q(x)}{E_p[\frac{Q(x)}{P(x)}]} \\ &= \log \sum_x Q(x) \\ &= \log 1 \\ &= 0 \end{aligned}$$

As we have equality if and only if $\frac{Q(x)}{P(x)} = E_p[\frac{Q(x)}{P(x)}]$ and $P(x) = Q(x)$ Using Jensen's inequality.

So

$$D_{KL}(P||Q) = 0$$

- (b) [5 points] **Chain rule for KL divergence.**

The KL divergence between 2 conditional distributions $P(X|Y), Q(X|Y)$ is defined as follows:

$$D_{KL}(P(X|Y)||Q(X|Y)) = \sum_y P(y) \left(\sum_x P(x|y) \log \frac{P(x|y)}{Q(x|y)} \right)$$

This can be thought of as the expected KL divergence between the corresponding conditional distributions on x (that is, between $P(X|Y = y)$ and $Q(X|Y = y)$), where the expectation is taken over the random y .

Prove the following chain rule for KL divergence:

$$D_{KL}(P(X, Y) \| Q(X, Y)) = D_{KL}(P(X) \| Q(X)) + D_{KL}(P(Y|X) \| Q(Y|X)).$$

Answer:

$$\begin{aligned} D_{KL}(P(X, Y) \| Q(X, Y)) &= \sum_{x,y} P(x, y) \log \frac{P(x, y)}{Q(x, y)} \\ &= \sum_{x,y} P(x, y) \log \frac{P(x)P(y|x)}{Q(x)Q(y|x)} \\ &= \sum_{x,y} P(x, y) \log \frac{P(x)}{Q(x)} + P(x, y) \log \frac{P(y|x)}{Q(y|x)} \\ &= \sum_x P(x, y) \log \frac{P(x)}{Q(x)} + \sum_x P(x) \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \\ &= D_{KL}(P(X) \| Q(X)) + D_{KL}(P(Y|X) \| Q(Y|X)) \end{aligned}$$

(c) [5 points] **KL and maximum likelihood.**

Consider a density estimation problem, and suppose we are given a training set $\{x^{(i)}; i = 1, \dots, n\}$. Let the empirical distribution be $\hat{P}(x) = \frac{1}{n} \sum_{i=1}^n 1\{x^{(i)} = x\}$. (\hat{P} is just the uniform distribution over the training set; i.e., sampling from the empirical distribution is the same as picking a random example from the training set.)

Suppose we have some family of distributions P_θ parameterized by θ . (If you like, think of $P_\theta(x)$ as an alternative notation for $P(x; \theta)$.) Prove that finding the maximum likelihood estimate for the parameter θ is equivalent to finding P_θ with minimal KL divergence from \hat{P} . I.e. prove:

$$\arg \min_{\theta} D_{KL}(\hat{P} \| P_\theta) = \arg \max_{\theta} \sum_{i=1}^n \log P_\theta(x^{(i)})$$

Remark. Consider the relationship between parts (b-c) and multi-variate Bernoulli Naive Bayes parameter estimation. In the Naive Bayes model we assumed P_θ is of the following form: $P_\theta(x, y) = p(y) \prod_{i=1}^d p(x_i|y)$. By the chain rule for KL divergence, we therefore have:

$$D_{KL}(\hat{P} \| P_\theta) = D_{KL}(\hat{P}(y) \| p(y)) + \sum_{i=1}^d D_{KL}(\hat{P}(x_i|y) \| p(x_i|y)).$$

This shows that finding the maximum likelihood/minimum KL-divergence estimate of the parameters decomposes into $2n + 1$ independent optimization problems: One for the class priors $p(y)$, and one for each of the conditional distributions $p(x_i|y)$ for each feature x_i given each of the two possible labels for y . Specifically, finding the maximum likelihood estimates for each of these problems individually results in also maximizing the likelihood of the joint distribution. (If you know what Bayesian networks are, a similar remark applies to parameter estimation for them.)

Answer:

$$\begin{aligned} \operatorname{argmin}_{\theta} D_{KL}(\hat{P} \| P_{\theta}) &= \operatorname{argmin}_{\theta} \sum_x \hat{P}(x) \log \hat{P}(x) - \hat{P}(x) \log P_{\theta}(x) \\ &= \operatorname{argmin}_{\theta} \sum_x -\hat{P}(x) \log P_{\theta}(x) \\ &= \operatorname{argmax}_{\theta} \sum_x \hat{P}(x) \log P_{\theta}(x) \\ &= \operatorname{argmax}_{\theta} \frac{1}{n} \sum_x \sum_{i=1}^n \{x^i = x\} \log P_{\theta}(x) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \log P_{\theta}(x^i) \end{aligned}$$

3. [20 points] K-means for compression

In this problem, we will apply the K-means algorithm to lossy image compression, by reducing the number of colors used in an image.

We will be using the files `src/k_means/peppers-small.tiff` and `src/k_means/peppers-large.tiff`.

The `peppers-large.tiff` file contains a 512x512 image of peppers represented in 24-bit color. This means that, for each of the 262144 pixels in the image, there are three 8-bit numbers (each ranging from 0 to 255) that represent the red, green, and blue intensity values for that pixel. The straightforward representation of this image therefore takes about $262144 \times 3 = 786432$ bytes (a byte being 8 bits). To compress the image, we will use K-means to reduce the image to $k = 16$ colors. More specifically, each pixel in the image is considered a point in the three-dimensional (r, g, b) -space. To compress the image, we will cluster these points in color-space into 16 clusters, and replace each pixel with the closest cluster centroid.

Follow the instructions below. Be warned that some of these operations can take a while (several minutes even on a fast computer)!

- (a) [15 points] **[Coding Problem] K-Means Compression Implementation.** First let us *look* at our data. From the `src/k_means/` directory, open an interactive Python prompt, and type

```
from matplotlib.image import imread; import matplotlib.pyplot as plt;
```

and run `A = imread('peppers-large.tiff')`. Now, `A` is a “three dimensional matrix,” and `A[:, :, 0]`, `A[:, :, 1]` and `A[:, :, 2]` are 512x512 arrays that respectively contain the red, green, and blue values for each pixel. Enter `plt.imshow(A); plt.show()` to display the image.

Since the large image has 262,144 pixels and would take a while to cluster, we will instead run vector quantization on a smaller image. Repeat (a) with `peppers-small.tiff`.

Next we will implement image compression in the file `src/k_means/k_means.py` which has some starter code. Treating each pixel’s (r, g, b) values as an element of \mathbb{R}^3 , implement K-means with 16 clusters on the pixel data from this smaller image, iterating (preferably) to convergence, but in no case for less than 30 iterations. For initialization, set each cluster centroid to the (r, g, b) -values of a randomly chosen pixel in the image.

Take the image of `peppers-large.tiff`, and replace each pixel’s (r, g, b) values with the value of the closest cluster centroid from the set of centroids computed with `peppers-small.tiff`. Visually compare it to the original image to verify that your implementation is reasonable. **Include in your write-up a copy of this compressed image alongside the original image.**

Answer:



(b) [5 points] **Compression Factor.**

If we represent the image with these reduced (16) colors, by (approximately) what factor have we compressed the image?

Answer:

Original Image = 24 bits per pixel

For 16 You should have 4 bits per pixel $\log_2 16$

Compression Factor $\frac{24}{4} = 6$

So 6 is the compression factor for this image

4. [35 points] Semi-supervised EM

Expectation Maximization (EM) is a classical algorithm for unsupervised learning (*i.e.*, learning with hidden or latent variables). In this problem we will explore one of the ways in which EM algorithm can be adapted to the semi-supervised setting, where we have some labeled examples along with unlabeled examples.

In the standard unsupervised setting, we have $n \in \mathbb{N}$ unlabeled examples $\{x^{(1)}, \dots, x^{(n)}\}$. We wish to learn the parameters of $p(x, z; \theta)$ from the data, but $z^{(i)}$'s are not observed. The classical EM algorithm is designed for this very purpose, where we maximize the intractable $p(x; \theta)$ indirectly by iteratively performing the E-step and M-step, each time maximizing a tractable lower bound of $p(x; \theta)$. Our objective can be concretely written as:

$$\begin{aligned}\ell_{\text{unsup}}(\theta) &= \sum_{i=1}^n \log p(x^{(i)}; \theta) \\ &= \sum_{i=1}^n \log \sum_{z^{(i)}} p(x^{(i)}, z^{(i)}; \theta)\end{aligned}$$

Now, we will attempt to construct an extension of EM to the semi-supervised setting. Let us suppose we have an *additional* $\tilde{n} \in \mathbb{N}$ labeled examples $\{(\tilde{x}^{(1)}, \tilde{z}^{(1)}), \dots, (\tilde{x}^{(\tilde{n})}, \tilde{z}^{(\tilde{n})})\}$ where both x and z are observed. We want to simultaneously maximize the marginal likelihood of the parameters using the unlabeled examples, and full likelihood of the parameters using the labeled examples, by optimizing their weighted sum (with some hyperparameter α). More concretely, our semi-supervised objective $\ell_{\text{semi-sup}}(\theta)$ can be written as:

$$\begin{aligned}\ell_{\text{sup}}(\theta) &= \sum_{i=1}^{\tilde{n}} \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \\ \ell_{\text{semi-sup}}(\theta) &= \ell_{\text{unsup}}(\theta) + \alpha \ell_{\text{sup}}(\theta)\end{aligned}$$

We can derive the EM steps for the semi-supervised setting using the same approach and steps as before. You are *strongly encouraged* to show to yourself (no need to include in the write-up) that we end up with:

E-step (semi-supervised)

For each $i \in \{1, \dots, n\}$, set

$$Q_i^{(t)}(z^{(i)}) := p(z^{(i)} | x^{(i)}; \theta^{(t)})$$

M-step (semi-supervised)

$$\theta^{(t+1)} := \arg \max_{\theta} \left[\sum_{i=1}^n \left(\sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{Q_i^{(t)}(z^{(i)})} \right) + \alpha \left(\sum_{i=1}^{\tilde{n}} \log p(\tilde{x}^{(i)}, \tilde{z}^{(i)}; \theta) \right) \right]$$

- (a) [5 points] **Convergence.** First we will show that this algorithm eventually converges. In order to prove this, it is sufficient to show that our semi-supervised objective $\ell_{\text{semi-sup}}(\theta)$

monotonically increases with each iteration of E and M step. Specifically, let $\theta^{(t)}$ be the parameters obtained at the end of t EM-steps. Show that $\ell_{\text{semi-sup}}(\theta^{(t+1)}) \geq \ell_{\text{semi-sup}}(\theta^{(t)})$.

Answer:

$$l(\theta^{(t+1)}) = \alpha l_{\text{sup}}(\theta^{(t+1)}) + l_{\text{unsup}}(\theta^{(t+1)})$$

Jensen Equality

$$\begin{aligned} &\geq \alpha l_{\text{sup}}(\theta^{(t+1)}) + \sum_{i=1}^n \sum_{z^i} Q_i^t z^i \log \frac{p(x^i, z^i; \theta^{t+1})}{Q_i^t z^i} \\ &\geq \alpha l_{\text{sup}}(\theta^{(t)}) + \sum_{i=1}^n \sum_{z^i} Q_i^t z^i \log \frac{p(x^i, z^i; \theta^{t+1})}{p(z^i | x^i; \theta^t)} \\ &\geq \alpha l_{\text{sup}}(\theta^{(t)}) + l_{\text{unsup}}(\theta^{(t)}) \\ &\geq l(\theta^{(t)}) \end{aligned}$$

Semi-supervised GMM

Now we will revisit the Gaussian Mixture Model (GMM), to apply our semi-supervised EM algorithm. Let us consider a scenario where data is generated from $k \in \mathbb{N}$ Gaussian distributions, with unknown means $\mu_j \in \mathbb{R}^d$ and covariances $\Sigma_j \in \mathbb{S}_+^d$ where $j \in \{1, \dots, k\}$. We have n data points $x^{(i)} \in \mathbb{R}^d, i \in \{1, \dots, n\}$, and each data point has a corresponding latent (hidden/unknown) variable $z^{(i)} \in \{1, \dots, k\}$ indicating which distribution $x^{(i)}$ belongs to. Specifically, $z^{(i)} \sim \text{Multinomial}(\phi)$, such that $\sum_{j=1}^k \phi_j = 1$ and $\phi_j \geq 0$ for all j , and $x^{(i)} | z^{(i)} \sim \mathcal{N}(\mu_{z^{(i)}}, \Sigma_{z^{(i)}})$ i.i.d. So, μ , Σ , and ϕ are the model parameters.

We also have additional \tilde{n} data points $\tilde{x}^{(i)} \in \mathbb{R}^d, i \in \{1, \dots, \tilde{n}\}$, and an associated *observed* variable $\tilde{z}^{(i)} \in \{1, \dots, k\}$ indicating the distribution $\tilde{x}^{(i)}$ belongs to. Note that $\tilde{z}^{(i)}$ are known constants (in contrast to $z^{(i)}$ which are unknown *random* variables). As before, we assume $\tilde{x}^{(i)} | \tilde{z}^{(i)} \sim \mathcal{N}(\mu_{\tilde{z}^{(i)}}, \Sigma_{\tilde{z}^{(i)}})$ i.i.d.

In summary we have $n + \tilde{n}$ examples, of which n are unlabeled data points x 's with unobserved z 's, and \tilde{n} are labeled data points $\tilde{x}^{(i)}$ with corresponding observed labels $\tilde{z}^{(i)}$. The traditional EM algorithm is designed to take only the n unlabeled examples as input, and learn the model parameters μ , Σ , and ϕ .

Our task now will be to apply the semi-supervised EM algorithm to GMMs in order to also leverage the additional \tilde{n} labeled examples, and come up with semi-supervised E-step and M-step update rules specific to GMMs. Whenever required, you can cite the lecture notes for derivations and steps.

- (b) [5 points] **Semi-supervised E-Step.** Clearly state which are all the latent variables that need to be re-estimated in the E-step. Derive the E-step to re-estimate all the stated latent variables. Your final E-step expression must only involve x, z, μ, Σ, ϕ and universal constants.

Answer:

The latent variables are all the $z(i)$,

From the lecture notes:

$$\begin{aligned}
w_j^i &= Q_i(z^i = j) = \frac{p(x^i | z^i = j; \mu, \Sigma) p(z^i = j)}{\sum_{l=1}^k p(x^i | z^i = l; \mu, \Sigma) p(z^i = l)} \\
&= \frac{\frac{1}{\sqrt{(2\pi)^d |\Sigma_j|}} \exp(-\frac{1}{2}(x^i - \mu_j)^T \Sigma_j^{-1} (x^i - \mu_j)) \phi_j}{\sum_{l=1}^k \frac{1}{\sqrt{(2\pi)^d |\Sigma_l|}} \exp(-\frac{1}{2}(x^i - \mu_l)^T \Sigma_l^{-1} (x^i - \mu_l)) \phi_j} \\
&= \frac{|\Sigma_j|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x^i - \mu_j)^T \Sigma_j^{-1} (x^i - \mu_j)) \phi_j}{\sum_{l=1}^k |\Sigma_l|^{-\frac{1}{2}} \exp(-\frac{1}{2}(x^i - \mu_l)^T \Sigma_l^{-1} (x^i - \mu_l)) \phi_j}
\end{aligned}$$

- (c) [10 points] **Semi-supervised M-Step.** Clearly state which are all the parameters that need to be re-estimated in the M-step. Derive the M-step to re-estimate all the stated parameters. Specifically, derive closed form expressions for the parameter update rules for $\mu^{(t+1)}$, $\Sigma^{(t+1)}$ and $\phi^{(t+1)}$ based on the semi-supervised objective.

(c) In Order to Simplify derivation, we denote

$$w_j^{(i)} = \phi_i^t (z^i = j)$$

and

$$\tilde{w}_j^{(i)} = \begin{cases} \infty & z^i = j \\ 0 & \text{otherwise} \end{cases}$$

We further denote $S = \sum$

writing m-step. in terms of S and maximize it becomes

$$\phi^{t+1}, \mu^{t+1}, S^{t+1} = \underset{\phi, \mu, S}{\operatorname{argmax}} \sum_{i=1}^n \sum_{j=1}^k \phi_j^t \log p(z^i, z_j; \phi, \mu, S) + \sum_{i=1}^n \log p(x^i, z_i; \phi, \mu, S)$$

$$= \underset{\phi, \mu, S}{\operatorname{argmax}} \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \frac{|S_j|^{1/2}}{(2\pi)^{d/2}} \exp \left(-\frac{1}{2} (x^i - \mu_j)^T S_j^{-1} (x^i - \mu_j) \right) \phi_j^t +$$

$$\sum_{i=1}^n \sum_{j=1}^k \tilde{w}_j^{(i)} \log \frac{|S_j|^{1/2}}{(2\pi)^{d/2}} \exp \left(-\frac{1}{2} (x^i - \mu_j)^T S_j^{-1} (x^i - \mu_j) \right) \phi_j^t$$

maximize ϕ_j :

$$\mathcal{L}(\phi, \beta) = C + \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \phi_j + \sum_{i=1}^n \sum_{j=1}^k \tilde{w}_j^{(i)} \log \phi_j + \beta \left(\sum_{j=1}^k \phi_j - 1 \right)$$

$$\nabla_{\phi_j} \mathcal{L}(\phi, \beta) = \sum_{i=1}^n w_j^{(i)} / \phi_j + \sum_{i=1}^n \tilde{w}_j^{(i)} / \phi_j + \beta = 0$$

$$\phi_j = \underbrace{\sum_{i=1}^n w_j^{(i)}}_{-\beta} + \underbrace{\sum_{i=1}^n \tilde{w}_j^{(i)}}_{\beta}$$

$$\nabla_{\beta} \mathcal{L}(\phi, \beta) = \sum_{j=1}^k \phi_j - 1 = 0$$

$$\sum_{j=1}^k \underbrace{\sum_{i=1}^n w_j^{(i)}}_{-\beta} + \underbrace{\sum_{i=1}^n \tilde{w}_j^{(i)}}_{\beta} = 1 = -\beta$$

Answer:

$$\begin{aligned}
 \phi_j^{(t+1)} &= \sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \\
 &\quad \overbrace{\left(\sum_{j=1}^K \left(\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \right) \right)} \\
 &= \frac{\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i}{n + \alpha \tilde{n}}
 \end{aligned}$$

Now we drive μ_j

$$\begin{aligned}
 &= -\nabla_{\mu_j} \left(C + \sum_{i=1}^n w_j^i (\tilde{x}^i - \mu_j)^T S_j (x^i - \mu_j) + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (\tilde{x}^i - \mu_j)^T S_j (\tilde{x}^i - \mu_j) \right) \\
 &= -\nabla_{\mu_j} \left(\sum_{i=1}^n w_j^i (-2 \tilde{x}^i S_j + \mu_j^T S_j \mu_j) + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (-S_j^T + \mu_j^T) \right) \\
 &= 2S \left(\sum_{i=1}^n w_j^i x^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \tilde{x}^i \right) - 2S \left(\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \right) \mu_j \\
 &\overbrace{\mu_j^{(t+1)} = \frac{\sum_{i=1}^n w_j^i x^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \tilde{x}^i}{\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i}}
 \end{aligned}$$

Now we drive S_j using s_j :

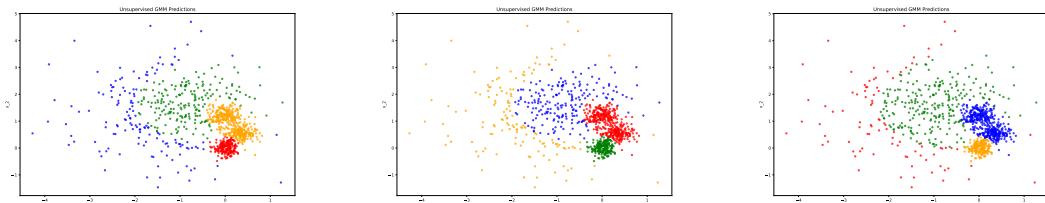
$$\begin{aligned}
 &= \nabla_{S_j} \left(C + \sum_{i=1}^n w_j^i (\log |S_j| - (x^i - \mu_j)^T S_j (x^i - \mu_j)) + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (\log |\tilde{S}_j| - (\tilde{x}^i - \mu_j)^T S_j (\tilde{x}^i - \mu_j)) \right) \\
 &= \sum_{i=1}^n w_j^i (S_j^{-1} - (x^i - \mu_j)(x^i - \mu_j)^T) + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (S_j^{-1} - (\tilde{x}^i - \mu_j)(\tilde{x}^i - \mu_j)^T) \\
 &= \left(\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i \right) S_j^{-1} - \left(\sum_{i=1}^n w_j^i (x^i - \mu_j)(x^i - \mu_j)^T + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (\tilde{x}^i - \mu_j)(\tilde{x}^i - \mu_j)^T \right) \\
 &\overbrace{s_j^{(t+1)} = \frac{\sum_{i=1}^n w_j^i (x^i - \mu_j)(x^i - \mu_j)^T + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i (\tilde{x}^i - \mu_j)(\tilde{x}^i - \mu_j)^T}{\sum_{i=1}^n w_j^i + \sum_{i=1}^{\tilde{n}} \tilde{w}_j^i}}
 \end{aligned}$$

- (d) [5 points] **Classical (Unsupervised) EM Implementation.** For this sub-question, we are only going to consider the n unlabelled examples. Follow the instructions in `src/semi_supervised_em/gmm.py` to implement the traditional EM algorithm, and run it on the unlabelled data-set until convergence.

Run three trials and use the provided plotting function to construct a scatter plot of the resulting assignments to clusters (one plot for each trial). Your plot should indicate cluster assignments with colors they got assigned to (*i.e.*, the cluster which had the highest probability in the final E-step).

Submit the three plots obtained above in your write-up.

Answer:

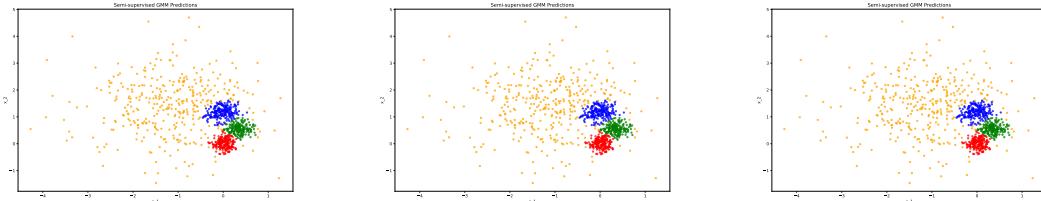


- (e) [7 points] **Semi-supervised EM Implementation.** Now we will consider both the labelled and unlabelled examples (a total of $n + \tilde{n}$), with 5 labelled examples per cluster. We have provided starter code for splitting the dataset into matrices \mathbf{x} and $\mathbf{x}_{\text{tilde}}$ of unlabelled and labelled examples respectively. Add to your code in `src/semi_supervised_em/gmm.py` to implement the modified EM algorithm, and run it on the dataset until convergence.

Create a plot for each trial, as done in the previous sub-question.

Submit the three plots obtained above in your write-up.

Answer:



- (f) [3 points] **Comparison of Unsupervised and Semi-supervised EM.** Briefly describe the differences you saw in unsupervised *vs.* semi-supervised EM for each of the following:
- Number of iterations taken to converge.
 - Stability (*i.e.*, how much did assignments change with different random initializations?)
 - Overall quality of assignments.

Note: The dataset was sampled from a mixture of three low-variance Gaussian distributions, and a fourth, high-variance Gaussian distribution. This should be useful in determining the overall quality of the assignments that were found by the two algorithms.

Answer:

Semi-supervised EM outperforms classical EM in all three categories, even with only 5 labeled instances per cluster.

- Unsupervised EM needs about 158 iterations to converge, whereas semi-supervised EM does so substantially faster. About 22 are needed for semi-supervised EM.
- The stability of semi-supervised is substantially higher: In contrast to unsupervised learning, semi-supervised EM assignments are nearly identical.
- The assignment is generally of higher quality when semi-supervised: Unsupervised EM has several different kinds of mistakes failing to discover there is just a single relatively high-variance Gaussian distribution in the mixture. The underlying distribution is almost accurately revealed by semi-supervised EM.

5. [10 points] PCA

In class, we showed that PCA finds the “variance maximizing” directions onto which to project the data. In this problem, we find another interpretation of PCA.

Suppose we are given a set of points $\{x^{(1)}, \dots, x^{(n)}\}$. Let us assume that we have as usual preprocessed the data to have zero-mean and unit variance in each coordinate. For a given unit-length vector u , let $f_u(x)$ be the projection of point x onto the direction given by u . I.e., if $\mathcal{V} = \{\alpha u : \alpha \in \mathbb{R}\}$, then

$$f_u(x) = \arg \min_{v \in \mathcal{V}} \|x - v\|^2.$$

Show that the unit-length vector u that minimizes the mean squared error between projected points and original points corresponds to the first principal component for the data. I.e., show that

$$\arg \min_{u: u^T u=1} \sum_{i=1}^n \|x^{(i)} - f_u(x^{(i)})\|_2^2.$$

gives the first principal component.

Remark. If we are asked to find a k -dimensional subspace onto which to project the data so as to minimize the sum of squares distance between the original data and their projections, then we should choose the k -dimensional subspace spanned by the first k principal components of the data. This problem shows that this result holds for the case of $k = 1$.

Answer:

$$\begin{aligned} \text{argmin}_{u: u^T u=1} \sum_{i=1}^m \|x^i - f_u(x^i)\|_2^2 &= \text{argmin}_{u: u^T u=1} \sum_{i=1}^m \|x^i - u^T x^i u\|_2^2 \\ &= \text{argmin}_{u: u^T u=1} \sum_{i=1}^m ((x^i - u^T x^i u)^T (x^i - u^T x^i u)) \\ &= \text{argmin}_{u: u^T u=1} \sum_{i=1}^m (x^{i^T} x^i - 2(u^T x^i)^2 + u^T u (u^T x^i)^2) \\ &= \text{argmin}_{u: u^T u=1} \sum_{i=1}^m (-2(u^T x^i)^2 + (u^T x^i)^2) \\ &= \text{argmin}_{u: u^T u=1} \sum_{i=1}^m -(u^T x^i)^2 \\ &= \text{argmax}_{u: u^T u=1} u^T (\sum_{i=1}^m x^i x^{i^T}) u \end{aligned}$$

Corresponding to Optimization Problem

6. [20 points] Independent components analysis

While studying Independent Component Analysis (ICA) in class, we made an informal argument about why Gaussian distributed sources will not work. We also mentioned that any other distribution (except Gaussian) for the sources will work for ICA, and hence used the logistic distribution instead. In this problem, we will go deeper into understanding why Gaussian distributed sources are a problem. We will also derive ICA with the Laplace distribution, and apply it to the cocktail party problem.

Reintroducing notation, let $s \in \mathbb{R}^d$ be source data that is generated from d independent sources. Let $x \in \mathbb{R}^d$ be observed data such that $x = As$, where $A \in \mathbb{R}^{d \times d}$ is called the *mixing matrix*. We assume A is invertible, and $W = A^{-1}$ is called the *unmixing matrix*. So, $s = Wx$. The goal of ICA is to estimate W . Similar to the notes, we denote w_j^T to be the j^{th} row of W . Note that this implies that the j^{th} source can be reconstructed with w_j and x , since $s_j = w_j^T x$. We are given a training set $\{x^{(1)}, \dots, x^{(n)}\}$ for the following sub-questions. Let us denote the entire training set by the design matrix $X \in \mathbb{R}^{n \times d}$ where each example corresponds to a row in the matrix.

(a) [5 points] Gaussian source

For this sub-question, we assume sources are distributed according to a standard normal distribution, i.e $s_j \sim \mathcal{N}(0, 1)$, $j = \{1, \dots, d\}$. The log-likelihood of our unmixing matrix, as described in the notes, is

$$\ell(W) = \sum_{i=1}^n \left(\log |W| + \sum_{j=1}^d \log g'(w_j^T x^{(i)}) \right),$$

where g is the cumulative distribution function, and g' is the probability density function of the source distribution (in this sub-question it is a standard normal distribution). Whereas in the notes we derive an update rule to train W iteratively, for the cause of Gaussian distributed sources, we can analytically reason about the resulting W .

Try to derive a closed form expression for W in terms of X when g is the standard normal CDF. Deduce the relation between W and X in the simplest terms, and highlight the ambiguity (in terms of rotational invariance) in computing W .

Answer:

$$l(W) = n \log |W| + \sum_{i=1}^n \sum_{j=1}^d \log \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(w_j^T x^{(i)})^2)$$

$$= n \log |W| + C - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^d (w_j^T x^{(i)})^2$$

$$\nabla_w l(W) = nW^{-T} - \frac{1}{2} \sum_{i=1}^n \begin{bmatrix} \nabla w_1^T (w_1^T x^{(i)})^2 \\ \nabla w_2^T (w_2^T x^{(i)})^2 \\ \vdots \\ \nabla w_d^T (w_d^T x^{(i)})^2 \end{bmatrix}$$

$$= nW^{-T} - \frac{1}{2} \sum_{i=1}^n \begin{bmatrix} 2(w_1^T x^{(i)}) x^{iT} \\ 2(w_2^T x^{(i)}) x^{iT} \\ \vdots \\ 2(w_d^T x^{(i)}) x^{iT} \end{bmatrix}$$

$$\begin{aligned}
&= nW^{-T} - \sum_{i=1}^n \begin{bmatrix} w_1^T(x^i x^{iT}) \\ w_2^T(x^i x^{iT}) \\ \dots \\ w_d^T(x^i x^{iT}) \end{bmatrix} \\
&= nW^{-T} - W \sum_{i=1}^n x^i x^{iT} \\
&= nW^{-T} - WX^T X \\
nW^{-1}W^{-T} &= X^T X \\
W^T W &= (\frac{1}{n} X^T X)^{-1}
\end{aligned}$$

This suggests that any $W = WU$, where U is some full rank unitary matrix, is likewise a valid solution if some W is a valid solution to the relation stated above. The main cause of this uncertainty is the quadratic form that the Gaussian log-likelihood assumes.

- (b) [10 points] **Laplace source.**

For this sub-question, we assume sources are distributed according to a standard Laplace distribution, i.e $s_i \sim \mathcal{L}(0, 1)$. The Laplace distribution $\mathcal{L}(0, 1)$ has PDF $f_{\mathcal{L}}(s) = \frac{1}{2} \exp(-|s|)$. With this assumption, derive the update rule for a single example in the form

$$W := W + \alpha(\dots).$$

Answer:

Starting with the PDF of the Laplace distribution, $g(s) = 1/2\exp(-|s|)$, is the first step. When we incorporate this into the likelihood given a lone example x^i , we obtain

$$l(W) = \log|W| + \sum_{j=1}^d \log \exp(-|w_j^T x^i|)$$

$$= \log|W| - \sum_{j=1}^d |w_j^T x^i|$$

$$\nabla_w l(W) = W^{-T} - \begin{bmatrix} \nabla w_1 |w_1^T x^i| \\ \nabla w_2 |w_2^T x^i| \\ \dots \\ \nabla w_d |w_d^T x^i| \end{bmatrix}$$

$$= W^{-T} - \begin{bmatrix} sign(w_1^T x^i) x^{iT} \\ sign(w_2^T x^i) x^{iT} \\ \dots \\ sign(w_d^T x^i) x^{iT} \end{bmatrix}$$

$$= W^{-T} - \begin{bmatrix} sign(w_1^T x^i) \\ sign(w_2^T x^i) \\ \dots \\ sign(w_d^T x^i) \end{bmatrix} x^{iT}$$

The Updated Rule is

$$W := W + \alpha(W^{-T} - \begin{bmatrix} sign(w_1^T x^i) \\ sign(w_2^T x^i) \\ \dots \\ sign(w_d^T x^i) \end{bmatrix} x^{iT})$$

(c) [5 points] **Cocktail Party Problem**

For this question you will implement the Bell and Sejnowski ICA algorithm, but assuming a Laplace source (as derived in part-b), instead of the Logistic distribution covered in class. The file `src/ica/mix.dat` contains the input data which consists of a matrix with 5 columns, with each column corresponding to one of the mixed signals x_i . The code for this question can be found in `src/ica/ica.py`.

Implement the `update_W` and `unmix` functions in `src/ica/ica.py`.

You can then run `ica.py` in order to split the mixed audio into its components. The mixed audio tracks are written to `mixed_i.wav` in the output folder. The split audio tracks are written to `split_i.wav` in the output folder.

To make sure your code is correct, you should listen to the resulting unmixed sources. (Some overlap or noise in the sources may be present, but the different sources should be pretty clearly separated.)

Submit the full unmixing matrix W (5×5) that you obtained, by including the `W.txt` the code outputs along with your code.

If your implementation is correct, your output `split_0.wav` should sound similar to the file `correct_split_0.wav` included with the source code.

Note: In our implementation, we **anneal** the learning rate α (slowly decreased it over time) to speed up learning. In addition to using the variable learning rate to speed up convergence, one thing that we also do is choose a random permutation of the training data, and running stochastic gradient ascent visiting the training data in that order (each of the specified learning rates was then used for one full pass through the data).

Answer:

Matrix Obtained by Laplace

52.83492974	16.79598806	19.9411949	-10.19841036	-20.8977174
-9.9368057	-0.97879563	-4.68186342	8.0430365	1.79099473
8.31143332	-7.47699382	19.31554724	15.17460858	-14.32640472
-14.66729873	-26.64481368	2.44071692	21.38223128	-8.42094492
-0.26917605	18.37373974	9.31200636	9.10275731	30.59390495

7. [0 points / ungraded] KL Divergence, Fisher Information, and the Natural Gradient

As seen before, the Kullback-Leibler divergence between two distributions is an asymmetric measure of how different two distributions are. Consider two distributions over the same space given by densities $p(x)$ and $q(x)$. The KL divergence between two continuous distributions, q and p is defined as,

$$\begin{aligned} D_{KL}(p||q) &= \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx \\ &= \int_{-\infty}^{\infty} p(x) \log p(x) dx - \int_{-\infty}^{\infty} p(x) \log q(x) dx \\ &= \mathbb{E}_{x \sim p(x)} [\log p(x)] - \mathbb{E}_{x \sim p(x)} [\log q(x)]. \end{aligned}$$

A nice property of KL divergence is that it invariant to parametrization. This means, KL divergence evaluates to the same value no matter how we parametrize the distributions P and Q . For e.g, if P and Q are in the exponential family, the KL divergence between them is the same whether we are using natural parameters, or canonical parameters, or any arbitrary reparametrization.

Now we consider the problem of fitting model parameters using gradient descent (or stochastic gradient descent). As seen previously, fitting model parameters using Maximum Likelihood is equivalent to minimizing the KL divergence between the data and the model. While KL divergence is invariant to parametrization, the gradient w.r.t the model parameters (i.e, direction of steepest descent) is *not invariant to parametrization*. To see its implication, suppose we are at a particular value of parameters (either randomly initialized, or mid-way through the optimization process). The value of the parameters correspond to some probability distribution (and in case of regression, a conditional probability distribution). If we follow the direction of steepest descent from the current parameter, take a small step along that direction to a new parameter, we end up with a new distribution corresponding to the new parameters. The non-invariance to reparametrization means, a step of fixed size in the parameter space could end up in a distribution that could either be extremely far away in D_{KL} from the previous distribution, or on the other hand not move very much at all w.r.t D_{KL} from the previous distributions.

This is where the *natural gradient* comes into picture. It is best introduced in contrast with the usual gradient descent. In the usual gradient descent, we *first choose the direction* in the *parameter space* by calculating the gradient of the MLE objective w.r.t the parameters, and then move a magnitude of step size (where size is measured in the *parameter space*) along that direction. Whereas in natural gradient, we *first choose a divergence* amount by which we would like to move, in the D_{KL} sense. This effectively gives us a perimeter (of some arbitrary shape) around the current parameter, such that all points on this perimeter correspond to distributions which are at an equal D_{KL} -distance away from the current parameter. Among the set of all distributions on this perimeter, we move to the distribution that maximizes the objective the most (i.e minimize D_{KL} between data and itself the most). This approach makes the optimization process invariant to parametrization. That means, even if we chose a new arbitrary reparametrization, the natural gradient ensures that by starting from a particular distribution, we always descend down the same sequence of distributions towards the optimum. Whereas the usual gradient will choose a path that is specific to the choice of parametrization.

In the rest of this problem, we will construct and derive the natural gradient update rule. For that, we will break down the process into smaller sub-problems, and give you hints to answer

them. Along the way, we will encounter important statistical concepts such as the *score function* and *Fisher Information* (which play a prominent role in Statistical Theory as well). Finally, we will see how this new natural gradient based optimization is actually equivalent to Newton's method for Generalized Linear Models.

Let the distribution of a random variable Y parameterized by $\theta \in \mathbb{R}^d$ be $p(y; \theta)$.

(a) [3 points] **Score function**

The score function associated with $p(y; \theta)$ is defined as $\nabla_{\theta} \log p(y; \theta)$, which signifies the sensitivity of the likelihood function with respect to the parameters. Note that the score function is actually a vector since it's the gradient of a scalar quantity with respect to the vector θ .

Recall that $E_{y \sim p(y)}[g(y)] = \int_{-\infty}^{\infty} p(y)g(y)dy$. Using this fact, show that the expected value of the score is 0, i.e.

$$E_{y \sim p(y; \theta)}[\nabla_{\theta'} \log p(y; \theta')|_{\theta'=\theta}] = 0$$

Answer:

(b) [2 points] **Fisher Information**

Let us now introduce a quantity known as the Fisher information. It is defined as the covariance matrix of the score function,

$$\mathcal{I}(\theta) = \text{Cov}_{y \sim p(y; \theta)}[\nabla_{\theta'} \log p(y; \theta')|_{\theta'=\theta}]$$

Intuitively, the Fisher information represents the amount of information that a random variable Y carries about a parameter θ of interest. When the parameter of interest is a vector (as in our case, since $\theta \in \mathbb{R}^d$), this information becomes a matrix. Show that the Fisher information can equivalently be given by

$$\mathcal{I}(\theta) = \mathbb{E}_{y \sim p(y; \theta)}[\nabla_{\theta'} \log p(y; \theta') \nabla_{\theta'} \log p(y; \theta')^\top|_{\theta'=\theta}]$$

Note that the Fisher Information is a function of the parameter. The parameter of the Fisher information is both a) the parameter value at which the score function is evaluated, and b) the parameter of the distribution with respect to which the expectation and variance is calculated.

Answer:

(c) [5 points] **Fisher Information (alternate form)**

It turns out that the Fisher information can not only be defined as the covariance of the score function, but in most situations it can also be represented as the expected negative Hessian of the log-likelihood.

Show that $\mathbb{E}_{y \sim p(y; \theta)}[-\nabla_{\theta'}^2 \log p(y; \theta')|_{\theta'=\theta}] = \mathcal{I}(\theta)$.

Remark. The Hessian represents the curvature of a function at a point. This shows that the expected curvature of the log-likelihood function is also equal to the Fisher information matrix. If the curvature of the log-likelihood at a parameter is very steep (i.e. Fisher information is very high), this generally means you need fewer data samples to estimate that parameter well (assuming data was generated from the distribution with those parameters), and vice versa. The Fisher information matrix associated with a statistical model parameterized by θ is extremely important in determining how a model behaves as a function of the number of training set examples.

Answer:

(d) [5 points] **Approximating D_{KL} with Fisher Information**

As we explained at the start of this problem, we are interested in the set of all distributions that are at a small fixed D_{KL} distance away from the current distribution. In order to calculate D_{KL} between $p(y; \theta)$ and $p(y; \theta + d)$, where $d \in \mathbb{R}^d$ is a small magnitude “delta” vector, we approximate it using the Fisher Information at θ . Eventually d will be the natural gradient update we will add to θ . To approximate the KL-divergence with Fisher Information, we will start with the Taylor Series expansion of D_{KL} and see that the Fisher Information pops up in the expansion.

$$\text{Show that } D_{KL}(p_\theta || p_{\theta+d}) \approx \frac{1}{2} d^T \mathcal{I}(\theta) d.$$

Hint: Start with the Taylor Series expansion of $D_{KL}(p_\theta || p_{\tilde{\theta}})$ where θ is a constant and $\tilde{\theta}$ is a variable. Later set $\tilde{\theta} = \theta + d$. Recall that the Taylor Series allows us to approximate a scalar function $f(\tilde{\theta})$ near θ by:

$$f(\tilde{\theta}) \approx f(\theta) + (\tilde{\theta} - \theta)^T \nabla_{\theta'} f(\theta')|_{\theta'=\theta} + \frac{1}{2} (\tilde{\theta} - \theta)^T (\nabla_{\theta'}^2 f(\theta')|_{\theta'=\theta}) (\tilde{\theta} - \theta)$$

Answer:

(e) [8 points] **Natural Gradient**

Now we move on to calculating the natural gradient. Recall that we want to maximize the log-likelihood by moving only by a fixed D_{KL} distance from the current position. In the previous sub-question we came up with a way to approximate D_{KL} distance with Fisher Information. Now we will set up the constrained optimization problem that will yield the natural gradient update d . Let the log-likelihood objective be $\ell(\theta) = \log p(y; \theta)$. Let the D_{KL} distance we want to move by, be some small positive constant c . The natural gradient update d^* is

$$d^* = \arg \max_d \ell(\theta + d) \quad \text{subject to} \quad D_{KL}(p_\theta || p_{\theta+d}) = c \quad (1)$$

First we note that we can use Taylor approximation on $\ell(\theta+d) \approx \ell(\theta) + d^T \nabla_{\theta'} \ell(\theta')|_{\theta'=\theta}$. Also note that we calculated the Taylor approximation $D_{KL}(p_\theta || p_{\theta+d})$ in the previous subproblem. We shall substitute both these approximations into the above constrained optimization problem.

In order to solve this constrained optimization problem, we employ the *method of Lagrange multipliers*. If you are familiar with Lagrange multipliers, you can proceed directly to solve for d^* . If you are not familiar with Lagrange multipliers, here is a simplified introduction. (You may also refer to a slightly more comprehensive introduction in the Convex Optimization section notes, but for the purposes of this problem, the simplified introduction provided here should suffice).

Consider the following constrained optimization problem

$$d^* = \arg \max_d f(d) \quad \text{subject to} \quad g(d) = c$$

The function f is the objective function and g is the constraint. We instead optimize the *Lagrangian* $\mathcal{L}(d, \lambda)$, which is defined as

$$\mathcal{L}(d, \lambda) = f(d) - \lambda[g(d) - c]$$

with respect to both d and λ . Here $\lambda \in \mathbb{R}_+$ is called the Lagrange multiplier. In order to optimize the above, we construct the following system of equations:

$$\begin{aligned}\nabla_d \mathcal{L}(d, \lambda) &= 0, & (a) \\ \nabla_\lambda \mathcal{L}(d, \lambda) &= 0. & (b)\end{aligned}$$

So we have two equations (a and b above) with two unknowns (d and λ), which can be sometimes be solved analytically (in our case, we can).

The following steps guide you through solving the constrained optimization problem:

- Construct the Lagrangian for the constrained optimization problem (1) with the Taylor approximations substituted in for both the objective and the constraint.
- Then construct the system of linear equations (like (a) and (b)) from the Lagrangian you obtained.
- From (a), come up with an expression for d that *involves* λ .

At this stage we have already found the “direction” of the natural gradient d , since λ is only a positive scaling constant. For most practical purposes, the solution we obtain here is sufficient. This is because we almost always include a learning rate hyperparameter in our optimization algorithms, or perform some kind of a line search for algorithmic stability. This can make the exact calculation of λ less critical. Let’s call this expression \tilde{d} (involving λ) as the *unscaled natural gradient*. Clearly state what is \tilde{d} as a function of λ .

The remaining steps are to figure out the value of the scaling constant λ along the direction of d , for completeness.

- Plug in that expression for d into (b). Now we have an equation that has λ but not d . Come up with an expression for λ that does *not include* d .
- Plug that expression for λ (without d) back into (a). Now we have an equation that has d but not λ . Come up with an expression for d that does *not include* λ .

The expression for d obtained this way will be the desired natural gradient update d^* . Clearly state and highlight your final expression for d^* . This expression cannot include λ .

Answer:

(f) [2 points] **Relation to Newton’s Method**

After going through all these steps to calculate the natural gradient, you might wonder if this is something used in practice. We will now see that the familiar Newton’s method that we studied earlier, when applied to Generalized Linear Models, is equivalent to natural gradient on Generalized Linear Models. While the two methods (Newton’s method and natural gradient) agree on GLMs, in general they need not be equivalent.

Show that the direction of update of Newton’s method, and the direction of natural gradient, are exactly the same for Generalized Linear Models. You may want to recall and cite the results you derived in problem set 1 question 4 (Convexity of GLMs). For the natural gradient, it is sufficient to use \tilde{d} , the unscaled natural gradient.

Answer:

If you got here and finished all the above problems, you are done with the final PSet of CS 229! We know these assignments are not easy, so well done :)