

CS 229, Summer 2020

Problem Set #3 Solutions

YOUR NAME HERE (YOUR SUNET HERE)

Due Monday, August 10 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 Piazza forum, at <https://piazza.com/stanford/summer2020/cs229>. (3) This quarter, Summer 2020, students may submit in pairs. If you do so, make sure both names are attached to the Gradescope submission. However, students are not allowed to work with the same partner on more than one assignment. 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 10 at 11:59 pm. If you submit after Monday, August 10 at 11:59 pm, you will begin consuming your late days. If you wish to submit on time, submit before Monday, August 10 at 11:59 pm.

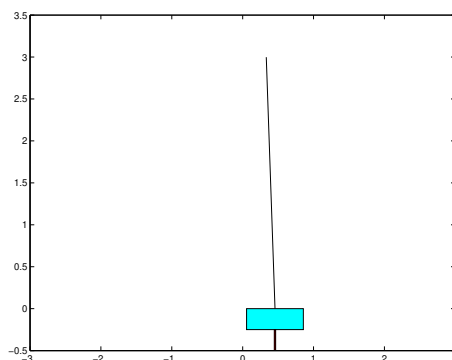
All students must submit an electronic PDF version of the written questions. We highly recommend typesetting your solutions via L^AT_EX, and we will award one bonus point for typeset submissions. All students must also submit a zip file of their source code to Gradescope, which should be created using 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/rlr/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.

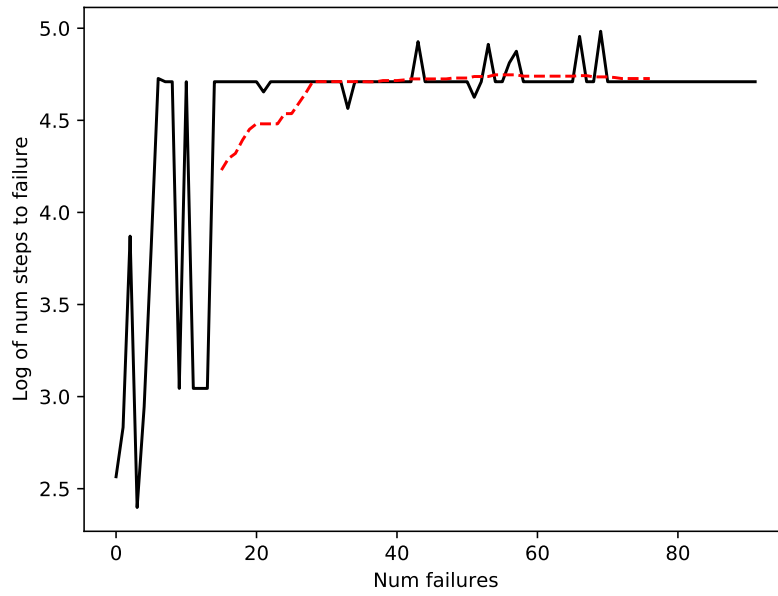
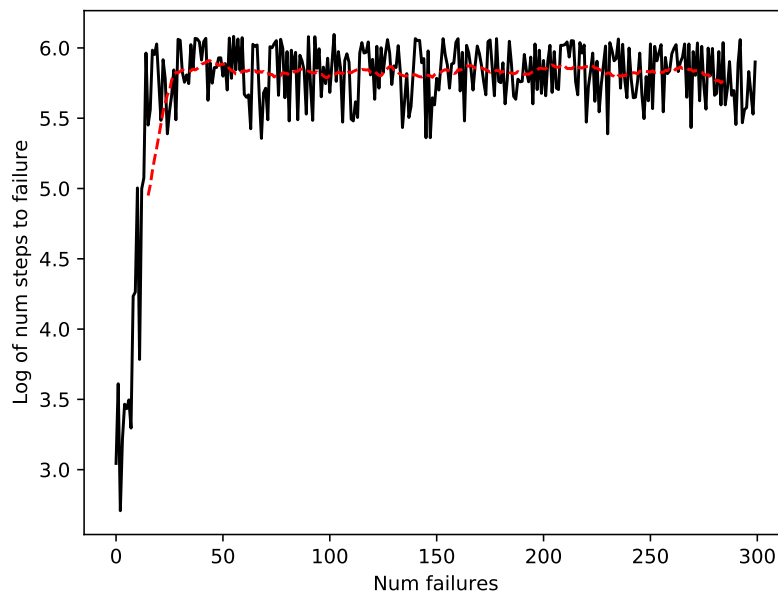
Finally, assume that the whole learning procedure has converged once several consecutive attempts (defined by the parameter `NO_LEARNING_THRESHOLD`) to solve Bellman's equation all converge in the first iteration. Intuitively, this indicates that the estimated model has stopped changing significantly.

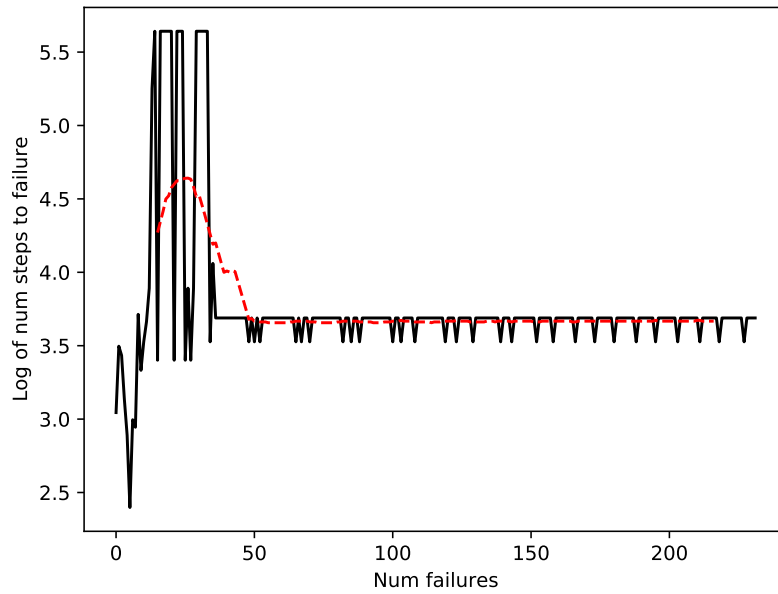
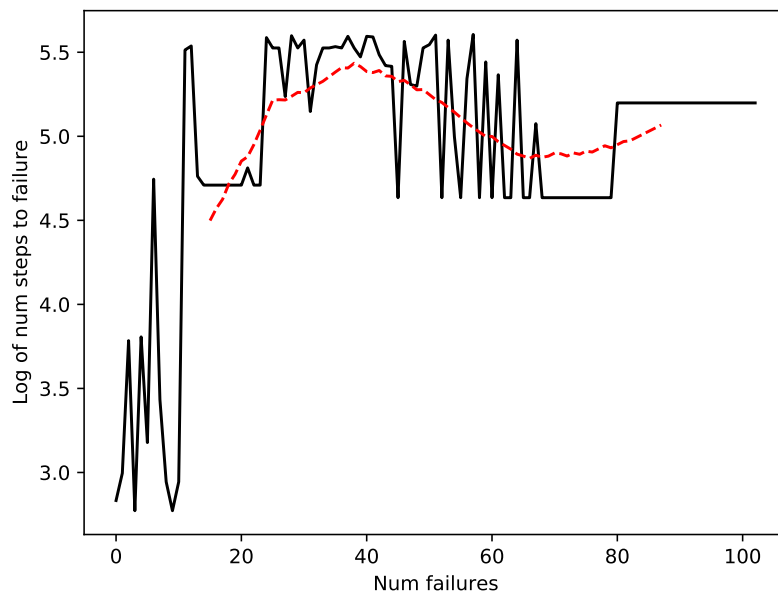
The code outline for this problem is already in `cartpole.py`, and you need to write code fragments only at the places specified in the file. There are several details (convergence criteria etc.) that are also explained inside the code. Use a discount factor of $\gamma = 0.995$.

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?

Answer: We get our plot of falling times for seeds 0,1,2,3

Figure 1: `np.random.seed(0)`Figure 2: `np.random.seed(1)`

Figure 3: `np.random.seed(2)`Figure 4: `np.random.seed(3)`

We have pretty different convergence behaviours for different random seeds. Max num trials survived is an order of magnitude different between e.g. seeds 1 and 2. Seed 1 had the best performance, although it took a while to converge fully (while maintaining that high performance) - I observed its value function difference spiking at multiple occasions which prevented the no learning threshold of 20 from triggering a few times.

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:**

$$\begin{aligned} D_{KL}(P||Q) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= E_{X \sim P} \left[-\log \frac{Q(X)}{P(X)} \right] \\ \text{using Jensen's inequality with } f &= -\log, \quad \geq -\log E_{X \sim P} \left[\frac{Q(X)}{P(X)} \right] \\ &= -\log \left(\sum_x Q(x) \right) \\ \text{probability sums to 1 so} \quad &= 0 \end{aligned}$$

- (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:

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

(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:

$$\begin{aligned}
 \operatorname{argmin}_{\theta} D_{KL}(\hat{P}||P_{\theta}) &= \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_{\theta}(x)} \\
 &= \operatorname{argmax}_{\theta} \sum_x \hat{P}(x) \log P_{\theta}(x) \\
 &= \operatorname{argmax}_{\theta} \sum_x \left[\sum_{i=1}^n 1_{x^{(i)}=x} \log P_{\theta}(x) \right] \\
 &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \left[\sum_x 1_{x^{(i)}=x} \log P_{\theta}(x) \right] \\
 (x^{(i)} = x \text{ for precisely one } x \in \chi) &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \log P_{\theta}(x) \\
 &= MLE_{\theta}
 \end{aligned}$$

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 512×512 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 512×512 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:

Original large image



Updated large image



Compressed pepper looks kinda as you'd imagine :D.

(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: We went from 256^3 possible colours to just 16 colours, so in effect our compression ratio is $256^3/16 = 2^{24}/2^4 = 2^{20}$. I.e. we went from needing 24 bits to represent any colour to needing just 4.

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:

Simplest way to think about the setup is that we have $a(\theta) = l_{\text{unsup}}(\theta)$, $b(\theta) = l_{\text{sup}}(\theta)$ and $\tilde{a}(Q, \theta) = \text{ELBO}(Q, \theta)$. We have from previous theory by Jensen's inequality that for any Q , we have $a(\theta) \leq \tilde{a}(Q, \theta)$, but that by argmaxing over Q we have $Q_i(z_i) = p(z_i|x_i; \theta)$ s.t. $a(\theta) = \tilde{a}(Q^*, \theta)$.

Thus at the E-step we set Q to max out this lower bound, $a(\theta^{(t)}) = \tilde{a}(Q^{(t)}, \theta)$, and then at the M-step we argmax over θ the sum of $\tilde{a}(Q, \theta) + b(\theta)$ to get $\theta^{(t+1)}$. Hence at the E-step we raise \tilde{a} such that $a = \tilde{a}$ and so $a + b = \tilde{a} + b$. And then on the M-step we overall raise $\tilde{a} + b$. This means that after the M-step, our new $a(\theta^{(t+1)}) + b(\theta^{(t+1)}) \geq a(\theta^{(t)}) + b(\theta^{(t)})$.

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: semi supervised EM algorithm for Gaussian mixture model is straightforward. We have log likelihood $l_{\text{semisup}}(\theta) = l_{\text{unsup}}(\theta) + \alpha l_{\text{sup}}(\theta)$ which we wish to maximise, where here $\theta = \phi, \mu, \Sigma$. for the E-step we compute $Q_i(z_i) = p(z_i|x_i; \phi, \mu, \Sigma)$ for the unsupervised points $x_i \in \{x_1, \dots, x_n\}$.

We parametrise the distribution $Q_i(z_i) \in \mathbb{R}^k$ for k choices of Gaussian distribution by writing

$$\begin{aligned}
 Q_i(z_i)_j &= w_j^{(i)} = p(z_i = j | x_i; \phi, \mu, \Sigma) \\
 &= p(x | z) \frac{p(z)}{p(x)} \\
 &= p(x_i | z_i = j) \frac{p(z_i = j)}{p(x_i)} \\
 &= N(\mu_j, \Sigma_j) \frac{\phi_j}{\sum_{l=1}^k p(x_i | z_i = l) p(z_i = l)} \\
 &= \frac{N(x_i; \mu_j, \Sigma_j) \phi_j}{\sum_{l=1}^k N(x_i; \mu_l, \Sigma_l) \phi_l}
 \end{aligned}$$

- (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.

Answer: For the semi supervised M-step we maximise over $\theta = \phi, \mu, \Sigma$ while keeping Q fixed the expression $ELBO(Q, \theta) + l_{\text{sup}}(\theta) = \sum_{i=1}^n \sum_{j=1}^k Q_i(z_i)_j \log \frac{p(x_i, z_i; \phi, \mu, \Sigma)}{Q_i(z_i)_j} + \alpha \sum_{i=1}^{\tilde{n}} \log p(x_i, z_i; \phi, \mu, \Sigma)$. So having pre-computed in E-step the fixed $Q_i(z_i)_j = w_j^i$, we can differentiate in turn for each of ϕ, μ, Σ . Rewriting the whole expression to be maximised as

$$\sum_{i=1}^n \sum_{j=1}^k w_j^i \log \frac{\phi_j N(x_i; \mu_j, \Sigma_j)}{w_j^i} + \alpha \sum_{i=1}^{\tilde{n}} \log (\phi_{\tilde{z}_i} N(\tilde{x}_i | \tilde{z}_i; \mu, \Sigma))$$

mu

We get

$$\begin{aligned}
 \frac{\partial}{\partial \mu_l} &= \sum_{i=1}^n w_l^i \frac{\partial}{\partial \mu_l} \log N(x_i; \mu_j, \Sigma_j) + \alpha \sum_{i=1}^{\tilde{n}} \frac{\partial}{\partial \mu_l} \log N(\tilde{x}_i | \tilde{z}_i; \mu, \Sigma) \\
 &= \sum_{i=1}^n w_l^i \Sigma_l^{-1} (x_i - \mu_l) + \alpha \sum_{\tilde{z}_i=l} \Sigma_l^{-1} (\tilde{x}_i - \mu_l)
 \end{aligned}$$

Setting this to 0 we get

$$\mu_l = \frac{\sum_{i=1}^n w_l^i x_i + \alpha \sum_{\tilde{z}_i=l} \tilde{x}_i}{\sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \}}$$

phi

Similarly for ϕ we have

$$\frac{\partial}{\partial \phi_l} = \sum_{i=1}^n w_l^i (1/\phi_l) + \alpha \sum_{\tilde{z}_i=l} (1/\phi_l) = (1/\phi_l) \sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \}$$

with constraint that $\sum_l \phi_l = 1$, so we apply the Lagrangian multipliers method:

$$\mathcal{L}(\phi, \lambda) = \sum_{l=1}^k \sum_{i=1}^n w_l^i \log \phi_l + \alpha \# \{ \tilde{z}_i = l \} \phi_l + \lambda \left(\sum_l \phi_l - 1 \right)$$

$$\frac{\partial \mathcal{L}}{\partial \phi_l} = (1/\phi_l) \sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \} + \lambda = 0$$

this applies for each $l = 1, \dots, k$, with λ constant, so we surmise that

$$\phi_l = C \sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \} \quad \text{for some constant } C, \text{ again constant across each } l$$

$$\text{so applying constraint } \sum_l \phi_l = 1 \text{ we get } \phi_l = \frac{\sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \}}{n + \alpha \tilde{n}}$$

Sigma

I cannot be bothered to write this all out again, you get the pattern, so we end up with

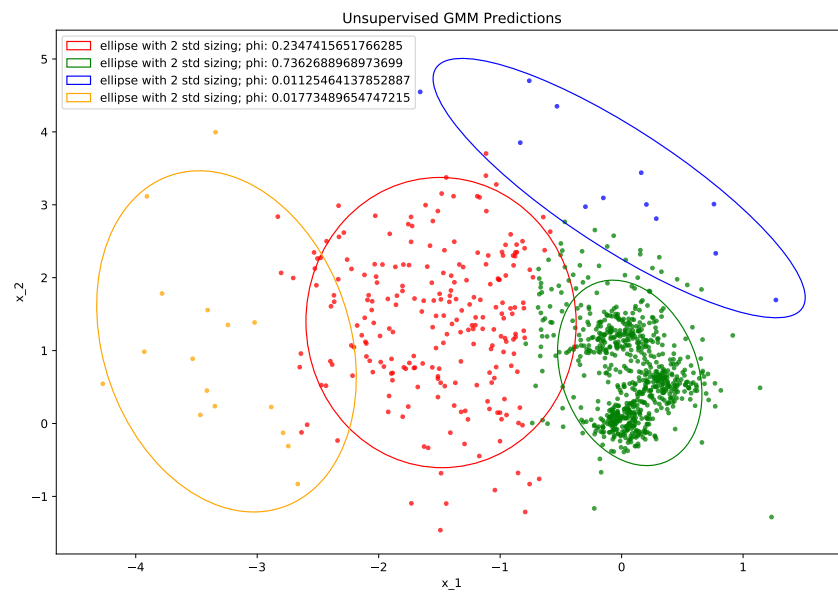
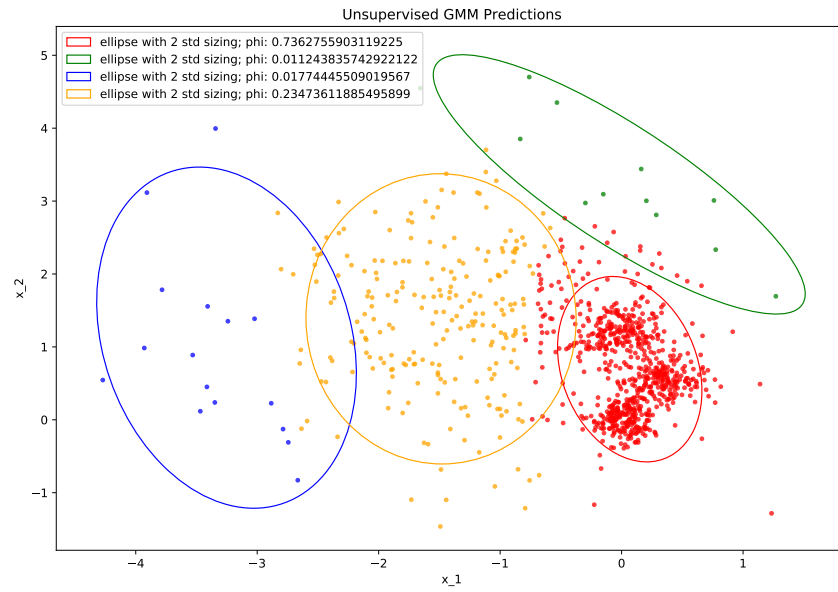
$$\Sigma_l = \frac{\sum_{i=1}^n w_l^i (x_i - \mu_l)(x_i - \mu_l)^T + \alpha \sum_{\tilde{z}_i=l} (\tilde{x}_i - \mu_l)(\tilde{x}_i - \mu_l)^T}{\sum_{i=1}^n w_l^i + \alpha \# \{ \tilde{z}_i = l \}}$$

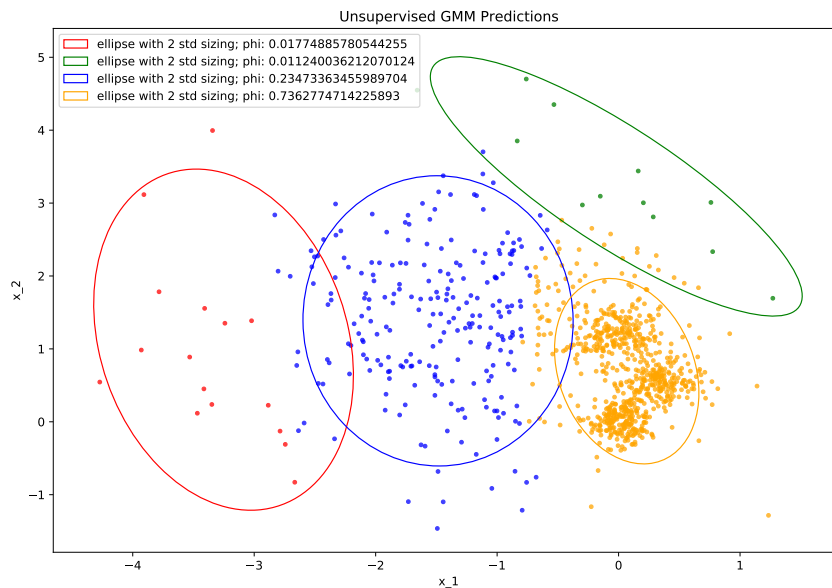
- (d) [5 points] **Classical (Unsupervised) EM Implementation.** For this sub-question, we are only going to consider the n unlabelled examples. Follow the instructions in `src/semi_supervised_em/gmm.py` to implement the traditional EM algorithm, and run it on the unlabelled data-set until convergence.

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:





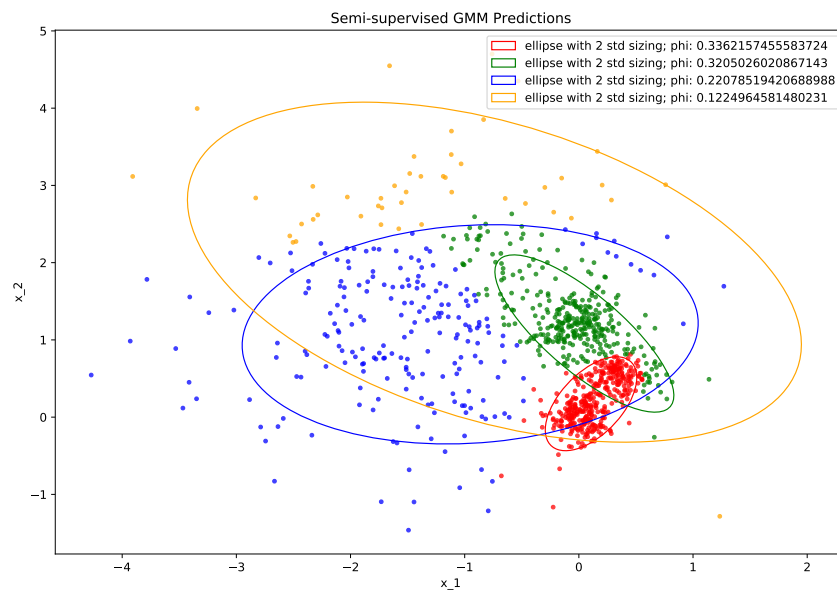
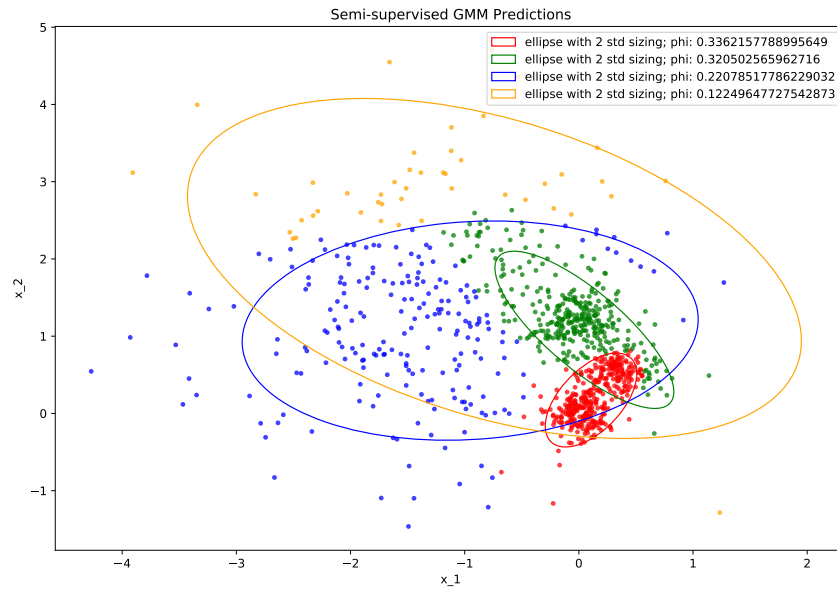
Our 3 trials took 101, 57, 95 iterations to converge respectively, however they've all converged to a very similar 4 clusters. I struggled a bit to get the ellipses to work but displayed are the mean centered 2 standard deviation width / height ellipses for each estimated cluster gaussian covariance. This is not quite the same thing as a 95% confidence interval but it's a useful verification of what's occurring

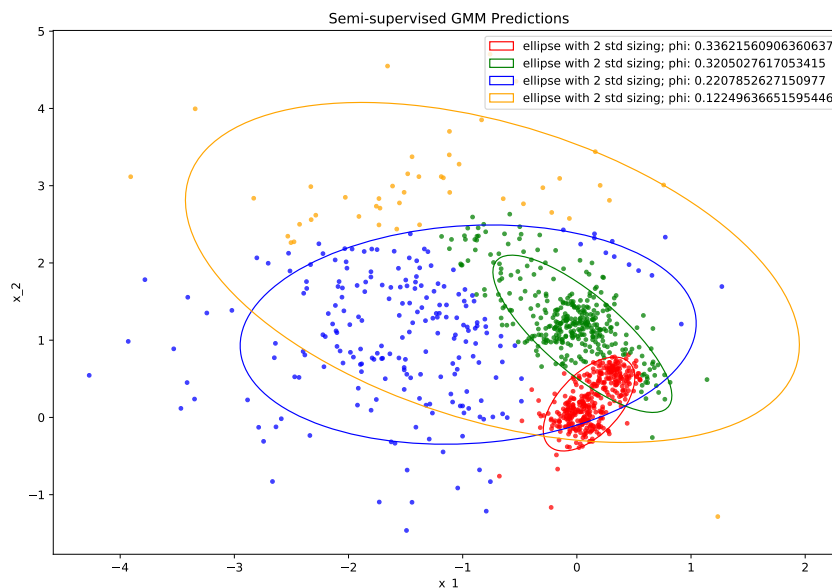
- (e) [7 points] **Semi-supervised EM Implementation.** Now we will consider both the labelled and unlabelled examples (a total of $n + \tilde{n}$), with 5 labelled examples per cluster. We have provided starter code for splitting the dataset into matrices \mathbf{x} and $\mathbf{\tilde{x}}$ of unlabelled 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:





We have very different gaussians here than in the unsupervised case.

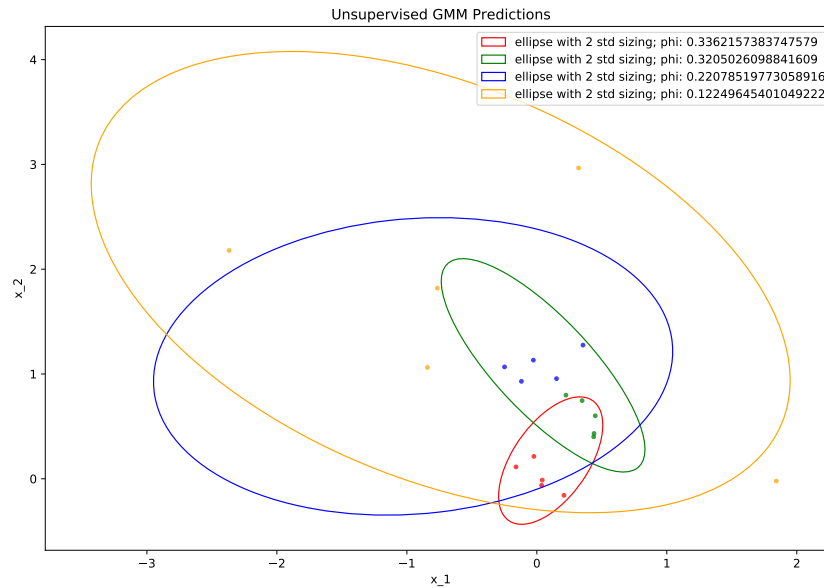
- (f) [3 points] **Comparison of Unsupervised and Semi-supervised EM.** Briefly describe the differences you saw in unsupervised *vs.* semi-supervised EM for each of the following:
- Number of iterations taken to converge.
 - Stability (*i.e.*, how much did assignments change with different random initializations?)
 - 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: It's odd that the semi supervised case seems to fit the datapoints more intuitively, despite the fact that the unsupervised case must somehow be getting better overall log likelihoods for the non labelled points. Another thing to note is that the unsupervised case produces pretty separate gaussians (they're not really overlayed at all unlike the semi supervised case).

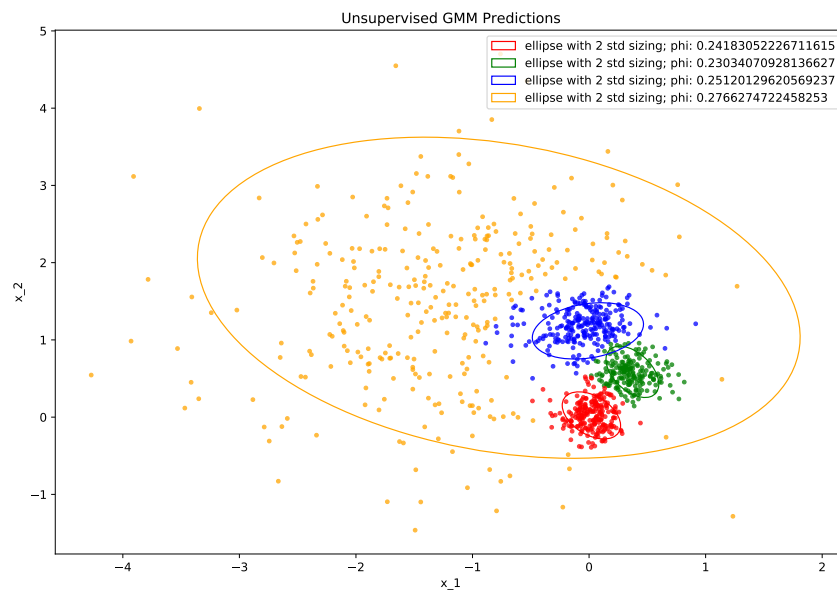
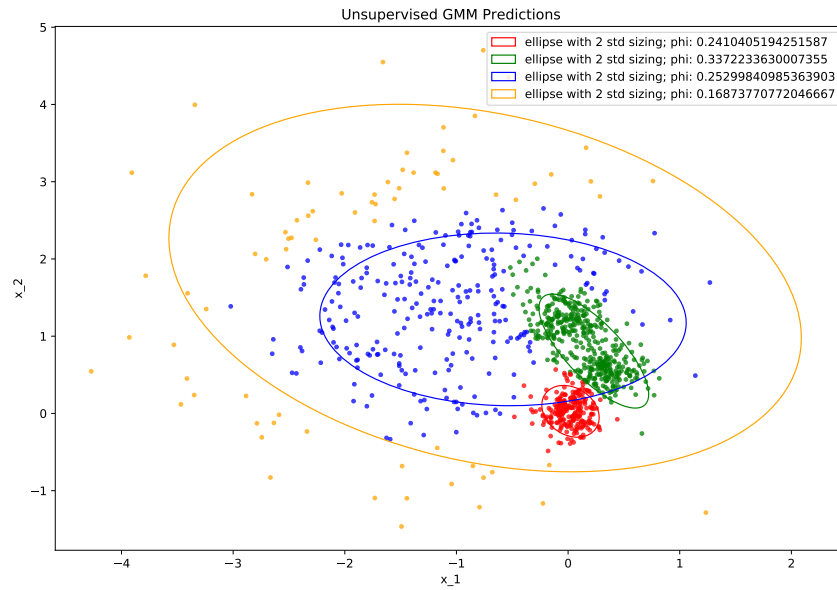
The number of iterations to convergence seems similar between unsupervised and semi supervised - or at least the variation within each class between trials is more noticeable anyway.

For both algorithms between trials a near identical point classification and predicted gaussians is produced.



Plotting the semi supervised predicted gaussian distributions against the 20 labelled points, and noting the note on the true distribution being 3 low variance gaussians and one high variance. So we see that the semi supervised prediction was better in capturing this than the unsupervised, but still not perfect - it only hits 2 out of the 3 low variance gaussians. We see that it mistakenly has the blue labelled points' ellipse blown up much bigger than the low variance. We somehow ended up with the red ellipse cannibalising some of the unlabelled green points, the green ellipse totally cannibalising the unlabelled blue points, and the blue ellipse left to blow up massively, still covering okishly its' labelled blue points.

The problem above is the small number of supervised points, only 20 compared to 980 unsupervised. At least we have a good spread between clusters. We can try to fix this by adjusting the weight α , increasing it to try to force the blue gaussian smaller to truly capture its cluster. This works - plotted are for $\alpha = 40$ and $\alpha = 80$ - we see the blue gaussian shrinking and then subsequently pretty well fitting its true cluster



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: We had the variance derivation of PCA for $k = 1$ being seeking to maximise the variance of the projections, i.e. maximising $\sum_{i=1}^n (u^T x_i)^2$ (where we use the fact that the vectors x_i have been mean zeroed already to ensure that the mean of $u^T x_i$ is 0, so we are indeed capturing the sample variance). So then we're equivalently maximising $u^T (\sum_{i=1}^n x_i x_i^T) u$ which shows us that the optimal unit vector u is the eigenvector of the largest eigenvalue of the sample covariance matrix.

In the minimal projection error setup we want to instead minimise over unit vector u the value of $\sum_{i=1}^n \|x_i - f_u(x_i)\|^2$, however $f_u(x_i) = (u^T x_i)u$, so we're minimising over

$$\sum_{i=1}^n \|x_i - (u^T x_i)u\|^2 = \sum_{i=1}^n (x_i - (u^T x_i)u)^T (x_i - (u^T x_i)u) = \sum_{i=1}^n x_i^T x_i - (u^T x_i)^2$$

where the last line comes through since $u^T u = 1$. Hence equivalent to the previous variance maximisation.

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: Since $s_j \sim \mathcal{N}(0, 1)$ we have that $s_j = w_j^T x_i$ and $g'(s_j) = (2\pi)^{-1/2} e^{-\frac{1}{2}s_j^2}$. So then we get

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

So then taking ∇_W of the log likelihood to maximise it, noting that $\nabla_{W_{ab}} \sum_{j=1}^d (w_j^T x_i)^2 = 2(w_a^T x^{(i)})x_b^{(i)}$ we get

$$\begin{aligned} \nabla_W l(W) &= \sum_{i=1}^n \left(\frac{1}{|W|} \text{adj}(W)^T - (W x_i) x_i^T \right) \\ &= nW^{-T} - W \sum_{i=1}^n x_i x_i^T \\ &= nW^{-T} - W X^T X \end{aligned}$$

Setting this to zero we have an optimal matrix W satisfies the equation $nI = W^T W X^T X$, which we see is invariant under transformations that don't change $W^T W$, such as $\tilde{W} = PW$ for some orthogonal matrix P , then $\tilde{W}^T \tilde{W} = W^T P^T P W = W^T W$.

(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: With pdf $g'(s_j) = \frac{1}{2} \exp(-|s_j|)$ we get

$$l(W) = \sum_{i=1}^n \left(\log |W| + \sum_{j=1}^d -\log 2 - |w_j^T x_i| \right)$$

So then taking ∇_W of the log likelihood to maximise it, noting that $\nabla_{W_{ab}} \sum_{j=1}^d |w_j^T x_i| = \text{sign}(w_a^T x^{(i)}) x_b^{(i)}$ we get

$$\nabla_W l(W) = \sum_{i=1}^n \left(\frac{1}{|W|} \text{adj}(W)^T - \text{sign}(W x_i) x_i^T \right)$$

giving us stochastic update rule

$$W := W + \alpha(W^{-T} - \text{sign}(W x_i) x_i^T)$$

(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:

Our unmixing matrix W obtained through gradient descent is shown below, where its rows w_j follow $s_j = w_j^T x$.

$$\begin{pmatrix} 52.833 & 16.795 & 19.941 & -10.198 & -20.897 \\ -9.933 & -0.979 & -4.680 & 8.044 & 1.790 \\ 8.311 & -7.477 & 19.315 & 15.174 & -14.326 \\ -14.667 & -26.645 & 2.441 & 21.382 & -8.421 \\ -0.269 & 18.374 & 9.312 & 9.103 & 30.594 \end{pmatrix}$$

It's magical that this actually works so well!