

Assignment 2

Goal: 2: Sampling

Goal: Get familiar with sampling methods by implementing these methods and applying them to a given probability distribution.

In this assignment we are going to learn about sampling methods. The goal is to implement *Metropolis-Hastings (MH) algorithm and Simulated Annealing* (SA) algorithm and analyze their behavior. Importantly, we aim at noticing differences between these two methods.

Here, we are interested in **sampling** from a mixture of two Gaussians, namely:

$$p(\mathbf{x}) = 0.25 \cdot \mathcal{N}\left(\mu = \begin{bmatrix} 2 \\ 2 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right) + 0.75 \cdot \mathcal{N}\left(\mu = \begin{bmatrix} -3 \\ -3 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

In this assignemnt, you are asked to implement:

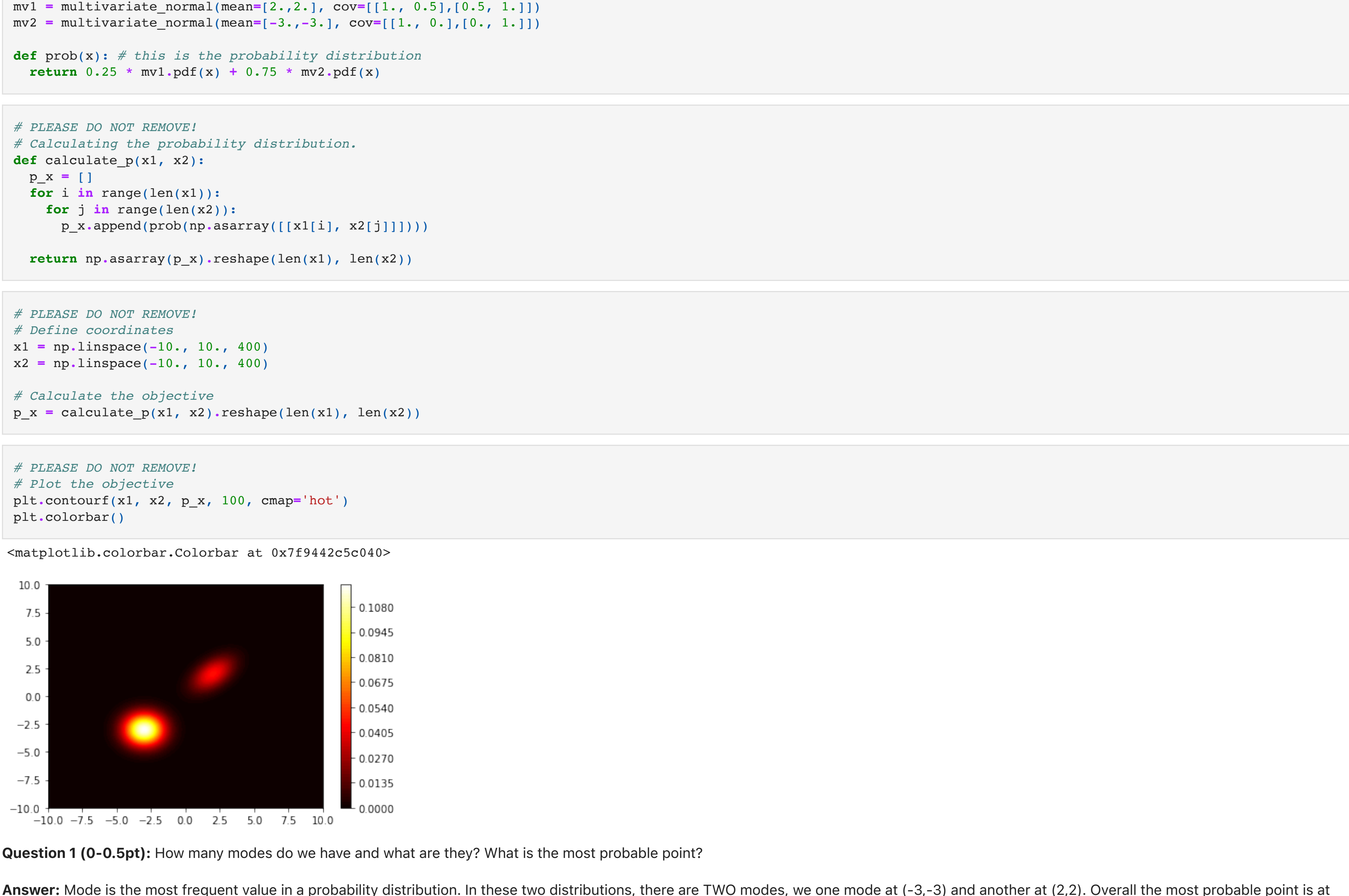
1. The Metropolis-Hastings (MH) algorithm.
2. The simulated annealing (SA) algorithm.

After implementing both methods, please run experimnts and compare both methods. Please follow all instructions.

1. Understanding the probability distribution

Please run the code below and visualize the probability distribution. Please try to understand this distribution, what the modes are (you can do it by inspecting the plot). What are possible problems here?

If any code line is unclear to you, please read on that in numpy or matplotlib docs.



Question 1 (0-0.5pt): How many modes do we have and what are they? What is the most probable point?

Answer: Mode is the most frequent value in a probability distribution. In these two distributions, there are TWO modes, we one mode at (-3,-3) and another at (2,2). Overall the most probable point is at (-3,-3) where the probability is the maximum in the whole distribution.

2. The Metropolis-Hastings algorithm

First, you are asked to implement the Metropolis-Hastings (MH) algorithm. Please take a look at the class below and fill in the missing parts.

NOTE: Please pay attention to the inputs and outputs of each function.

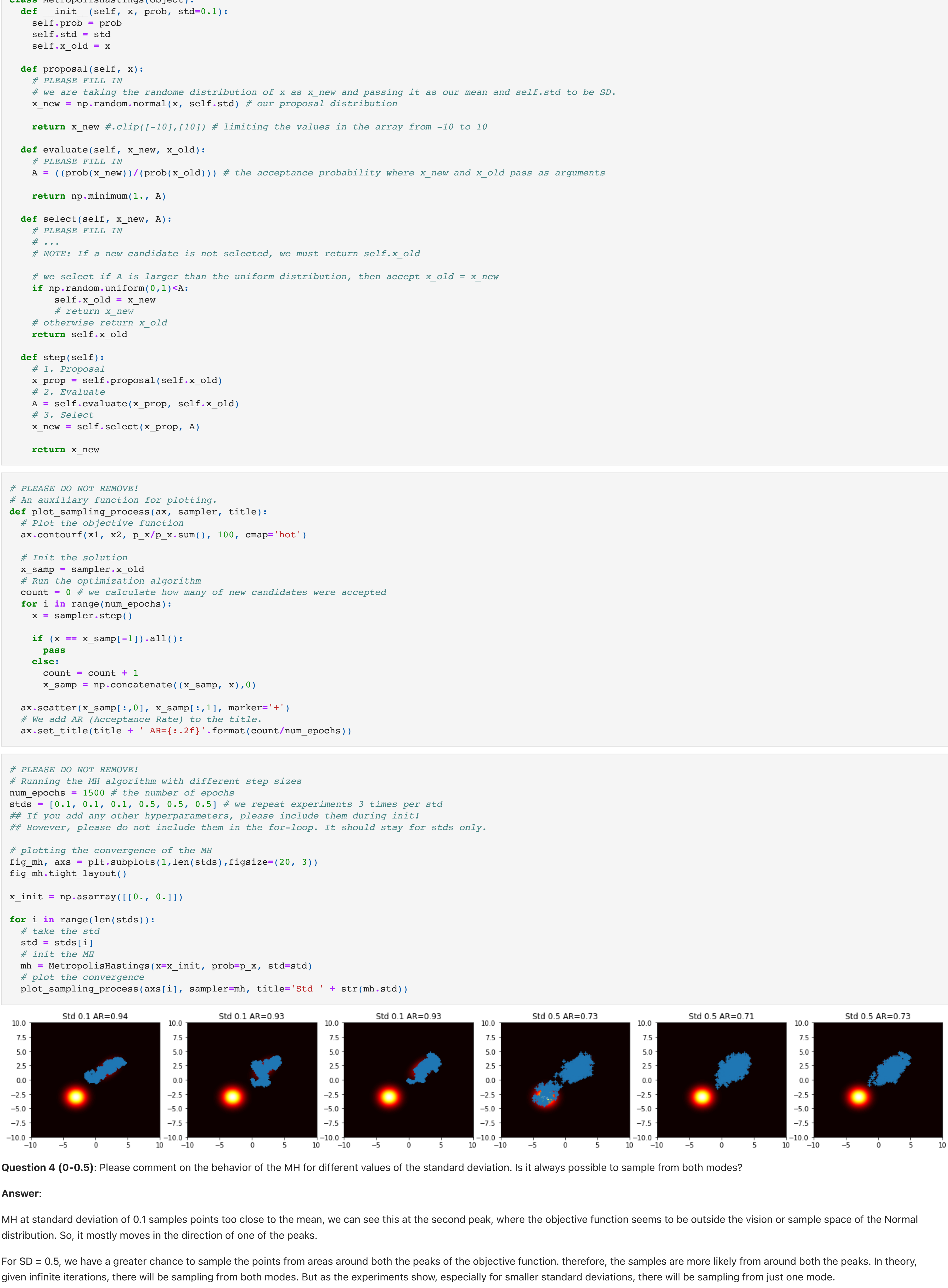
NOTE: To implement the MH algorithm, you need to specify the proposal distribution.

Question 2 (0-0.5pt): What is the proposal distribution, $q(\mathbf{x}_{new}|\mathbf{x}_{old})$, in your case?

Answer: using gaussian distribution where the data is normally distributed, illustrating the mean to be \mathbf{x}_{old} and in this case, $q(\mathbf{x}_{new}|\mathbf{x}_{old}) = 1/\sigma\sqrt{2\pi} * e^{*-((\mathbf{x}_{new}-\mathbf{x}_{old})*2)/(2\sigma*2)}$ Since Normal distribution is symmetric the value of the $q(\mathbf{x}_{new}|\mathbf{x}_{old})$ is same as $q(\mathbf{x}_{old}|\mathbf{x}_{new})$

Question 3 (0-0.5pt): Is your proposal a proper proposal distribution? (That is, it must fulfill irreducibility and aperiodicity, see Andrieu et al., "An Introduction to MCMC for Machine Learning".)

Answer: The Gaussian distribution is positive in the range of $-\infty$ to $+\infty$, allowing the Probability to reach any point on a real line from any other point. The rea line is always positive, even if the values become infinitely small, thus making this proposl distribution irreducible.



Question 4 (0-0.5): Please comment on the behavior of the MH for different values of the standard deviation. Is it always possible to sample from both modes?

Answer:

MH at standard deviation of 0.1 samples points too close to the mean, we can see this at the second peak, where the objective function seems to be outside the vision or sample space of the Normal distribution. So, it mostly moves in the direction of one of the peaks.

For SD = 0.5, we have a greater chance to sample the points from areas around both the peaks of the objective function. therefore, the samples are more likely from around both the peaks. In theory, given infinite iterations, there will be sampling from both modes. But as the experiments show, especially for smaller standard deviations, there will be sampling from just one mode.

Question 5 (0-0.5): Please comment on the acceptance ratio (AR) for std=0.1 and std=0.5. How can you influence the difference?

Answer: The suggested samples from proposal distribution with std=0.1 are closer to the mean, so the value of f(xnew) is similar to f(xold). Therefore, the acceptance probability f(xnew)/f(xold) is closer to 1, making most of the samples suggested have a high probability of getting accepted.

However, the suggested samples from proposal distribution with std=0.5, illustrates samples which are further from the mean point. Depcting that ther are high chances that the acceptance probability is not always closer to 1.Thus, more samples from this distribution are rejected.

Consequently, proposal distribution with std=0.1 has an acceptance rate of around 0.93 in this case, while distribution with std=0.5 has around 0.73 acceptance rate.

3. The simulated annealing (SA) sampling

In the second part of this assignment, you are asked to implement the Simulated Annealing (SA) algorithm with cooling scheme as discussed during the lecture.

Question 6 (0-0.5pt): Please explain what is the difference between MH and SA?

Answer: Both algorithms are stochastic, generating new points to move to at random. Where they differ is in their acceptance/rejection criterion. Both algorithms move to a new random point with a certain probability, which is based on the difference (or the ratio) of the current and new proposed point in the search space.

Simulated annealing is a meta-heuristic algorithm used for optimization, that is finding the minimum/maximum of a function with the application of temperature. It moves to a new point depending if the temperature is high or low. If it is high then it will move at random to a new point, but if the temperature is low the algorithm will be more deterministic and select the next point to move to.

Metropolis-Hastings is an algorithm used for exploring a function (finding possible values/samples). the algorithm calculates the ration between \mathbf{x}_{new} and \mathbf{x}_{old} . If the ratio is > 1 then the algorithms determines to move to the respective point. If the ratio is less than 1, then it will determine a point between 0-1 and evaluate if the new point is greater than or less than the ratio. If it is greater the new point will be accepted, otherwise it will be rejected.

Question 7 (0-0.5pt): Why is SA sometimes more preferable than MH? Which of these two methods would you use for optimization (not sampling)?

Answer: Simulated annealing is a more preferable because the algorithm is a meta-herustic algorithm which is used for optimization. It is similar to metropolis hastings but with the application of temperature. SA is better, because when it determines which point to go to, the algorithm becomes more deterministic, choosing the next best point, in contrast to MH it chooses a random point and determines if it is the best point to go to.



Question 8 (0-0.5pt) How does the standard deviation influence the SA?

Answer:

By correlating a rate of cooling in terms of a distance from equilibrium state—the higher the standard deviation (the system is far away from the equilibrium), the higher the rate of cooling. We can see that as the standard deviation increases to 0.5 the points cool and accumulate around the mean point, illustrating that as the temperature keeps cooling down, the results slowly become worse as SA has to take decisive steps instead of randomly selecting point.

Standard Deviation at 0.1 limited the search space for proposing next points. On one hand it is slow to move, but on the other it allowed the SA to sample more points from the given probability distribution. It has a consistently higher acceptance rate as compared to the SA with SD as 0.5, and it sampled the objective function more accurately.

Question 9 (0-0.5pt) How does the initial temperature, T_0 , influence the SA?

Answer: PLEASE FILL IN

Question 10 (0-0.5pt) How does the constant C influence the SA?

Answer: PLEASE FILL IN

Question 11 (0-0.5pt) Which setting (i.e., std, T_0 , C) did perform the best in terms of sampling?

Answer:

Looking at the consistency of sampling, we can see the best sampling to be one one of the best sampling was obtained by the following parameters: std=0.1, T0=0.1, C=10

Question 12 (0-0.5pt) How do different values of the hyperparameters (i.e., std, T_0 , C) influence the acceptance ratio (AR)? Why?

Answer: PLEASE FILL IN

4. Final remarks: MH vs. SA

Eventually, please answer the following last questions that will allow you to conclude the assignment and draw conclusions.

Question 13 (0-0.5-1pt): Which of the two algorithms did perform better? Why?

Answer: Both algorithms performed well, but MH performed better. We can see in MH that it often circled around the mean point and sampled close to the target pointsm but with the right parameters SA was more likely to find the global maxima.

MH focussed around the objective function, thus the proposals were nearby relevant values. Thus, the acceptance rate was high and more accurate sampling was present. MH didn't explore much and at lower SD=0.1, it didn't move from one of the modes.

SA "explored" much more points. It was more successful in reaching global maxima but needed the right hyperparameters for that. It was bad at sampling initially but with lowered temperatures, the sampling improved.

Question 14 (0-0.5-1pt): Which of the two algorithms is easier to use? Why?

Answer: MH with Gaussian distribution was easy to apply, because in contrast to SA we needed hyperparametr tuning, setting up and calculating T. Although SA is better for real-life optimization problems, it is trickier to set-up.

MH is easy for accurate sampling.