

# CS 229, Fall 2021

## Problem Set #4 Solutions

YOUR NAME HERE (YOUR SUNET HERE)

---

**Due Saturday, Nov 20 at 11:59pm 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/14428/discussion/>. (3) This quarter, Fall 2021, students may submit in pairs. If you do so, make sure both names are attached to the Gradescope submission. 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 Saturday, Nov 20 at 11:59pm. If you submit after Saturday, Nov 20 at 11:59pm, you will begin consuming your late days. If you wish to submit on time, submit before Saturday, Nov 20 at 11:59pm.

All students must submit an electronic PDF version of the written questions. We highly recommend typesetting your solutions via L<sup>A</sup>T<sub>E</sub>X, 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. [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  $j^{\text{th}}$  row of  $W$ . Note that this implies that the  $j^{\text{th}}$  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:**

## (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:**

## (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  $W$  ( $5 \times 5$ ) 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  $\alpha$  (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:**

## 2. [15 points] Markov decision processes

Consider an MDP with finite state and action spaces, and discount factor  $\gamma < 1$ . Let  $B$  be the Bellman update operator with  $V$  a vector of values for each state. I.e., if  $V' = B(V)$ , then

$$V'(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V(s').$$

- (a) [10 points] Prove that, for any two finite-valued vectors  $V_1, V_2$ , it holds true that

$$\|B(V_1) - B(V_2)\|_\infty \leq \gamma \|V_1 - V_2\|_\infty.$$

where

$$\|V\|_\infty = \max_{s \in S} |V(s)|.$$

(This shows that the Bellman update operator is a “ $\gamma$ -contraction in the max-norm.”)

**Remark:** The result you proved in part(a) implies that value iteration converges geometrically (i.e. exponentially) to the optimal value function  $V^*$ .

**Answer:**

- (b) [5 points] We say that  $V$  is a **fixed point** of  $B$  if  $B(V) = V$ . Using the fact that the Bellman update operator is a  $\gamma$ -contraction in the max-norm, prove that  $B$  has at most one fixed point—i.e., that there is at most one solution to the Bellman equations. You may assume that  $B$  has at least one fixed point.

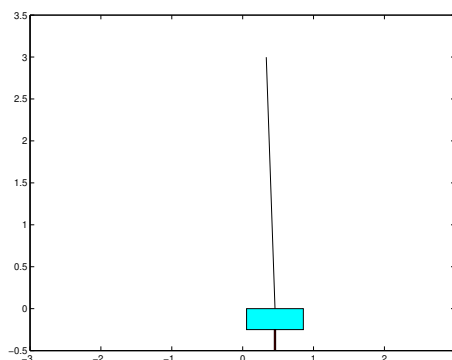
**Answer:**

### 3. [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.<sup>1</sup>

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  $\theta$  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.

<sup>1</sup>The 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:**

**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 :)**