

6.036: Final Exam, Spring 2019

Solutions

- This is a closed book exam. Two sheets of paper (8 1/2 in. by 11 in.) of notes, front and back, are permitted. Calculators are not permitted.
- The problems are not necessarily in any order of difficulty.
- Record all your answers in the places provided. If you run out of room for an answer, continue on a blank page and mark it clearly.
- If a question seems vague or under-specified to you, make an assumption, write it down, and solve the problem given your assumption.
- If you absolutely *have* to ask a question, come to the front.
- **Write your name on every page.**

Name: _____ Athena ID: _____
(username)

Question:	1	2	3	4	5	6	7	8	9	Total
Points:	8	8	12	10	16	14	14	8	10	100
Score:										

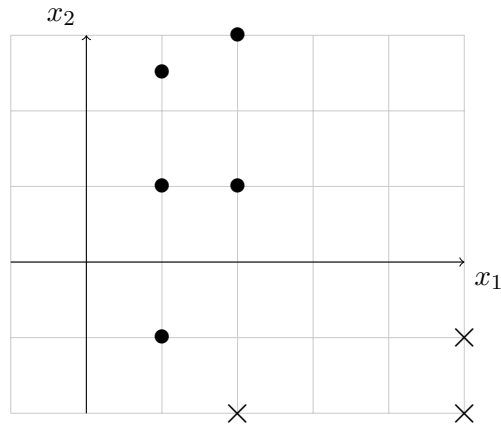
Name: _____

Nearest Neighbors

1. (8 points) Consider the following 2D dataset:

x	y
(1, -1)	+1
(1, 1)	+1
(1, 2.5)	+1
(2, -2)	-1
(2, 1)	+1
(2, 3)	+1
(5, -1)	-1
(5, -2)	-1

The dataset is plotted below, with positively labeled points as solid points (\bullet) and negatively labeled points as X marks (\times):



Break ties in distance by choosing the point with smaller x_1 coordinate, and if still tied, by smaller x_2 coordinate.

- (a) Compute the leave-one-out cross validation accuracy (i.e., average 8-fold cross validation accuracy) of the 1-nearest-neighbor learning algorithm on this dataset.

Solution: 6/8. When left out of the training set, the point at (1,-1) will be misclassified during testing; similarly for the point at (2,-2).

- (b) Compute the leave-one-out cross validation accuracy of the 3-nearest-neighbor learning algorithm on this dataset.

Name: _____

Solution: $7/8$. Now only the point at $(2,-2)$ will be misclassified during testing, when left out of the training set.

- (c) In the case of the 1-nearest-neighbor learning algorithm, is it possible to strictly increase the leave-one-out cross validation accuracy on this dataset by changing the label of a single point in the original dataset? If so, give such a point.

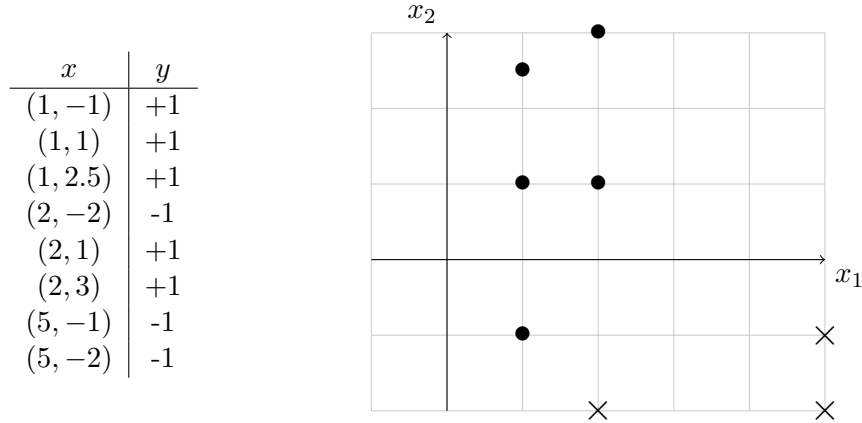
Solution: Yes. Change either point at $(2, -2)$ to $+1$, or point at $(1, -1)$ to -1 .

- (d) How about in the case of the 3-nearest neighbor algorithm? If so, give such a point.

Solution: No, not possible. If we try to change the point at $(2, -2)$ to $+1$, then that point will be correctly predicted during cross-validation as $+1$. Unfortunately, with that change the two points at $(5,-1)$ and $(5,-2)$ will now be misclassified, making our cross-validation accuracy worse.

Decision Trees

2. (8 points) Consider the following 2D dataset (same as that used in the previous question). Positively labeled points are solid points (\bullet) and negatively labeled points are X marks (\times).



We will construct a tree using the algorithm discussed in the lecture notes, i.e., a greedy algorithm that recursively minimizes weighted average entropy. Recall that the weighted average entropy of a split into subsets A and B is:

$$(\text{fraction of points in } A) \cdot H(R_{j,s}^A) + (\text{fraction of points in } B) \cdot H(R_{j,s}^B)$$

where the entropy $H(R_m)$ of data in a region R_m is given by

$$H(R_m) = - \sum_k \hat{P}_{mk} \log_2 \hat{P}_{mk}.$$

The \hat{P}_{mk} is the *empirical probability*, which is in this case the fraction of items in region m that are of class k .

Some facts that might be useful to you:

$$H(0) = 0$$

$$H(3/5) = 0.97$$

$$H(3/8) = 0.95$$

$$H(3/4) = 0.81$$

$$H(5/6) = 0.65$$

$$H(1) = 0$$

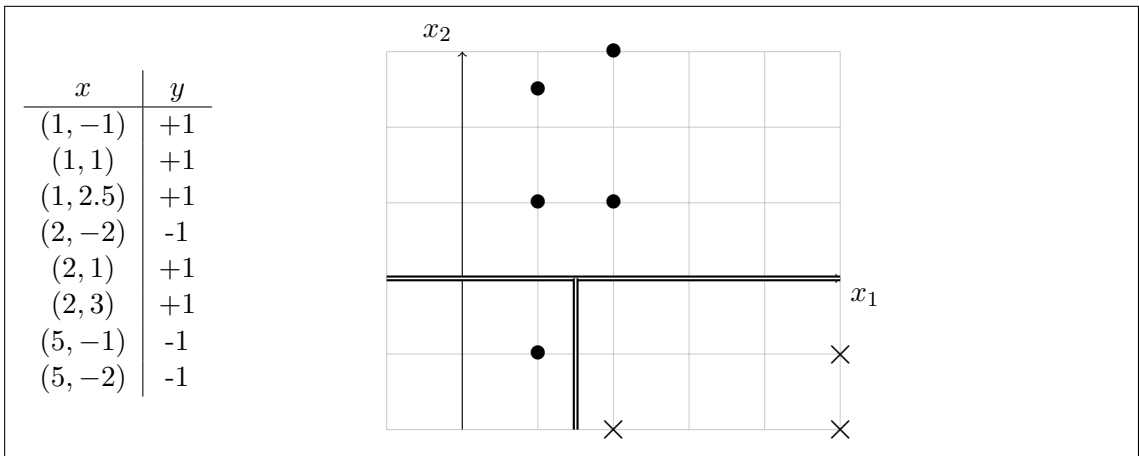
Name: _____

- (a) Draw the decision tree that would be constructed by our tree algorithm for this dataset. Clearly label the test in each node, which case (yes or no) each branch corresponds to, and the prediction that will be made at each leaf. Assume there is no pruning and that the algorithm runs until each leaf has only members of a single class.

Solution:
 $x_2 < 0$
Yes branch:
 $x_1 < 1.5$
 Yes branch: +1
 No branch: -1
No branch: +1

Note: a first (top level) split on $x_1 < 3.5$ looks appealing, but has slightly worse weighted average entropy than the top level split on $x_2 < 0$. Compare $(4/8)H(3/4) \approx 0.405$ with $(6/8)H(5/6) \approx 0.488$.

- (b) Draw the decision tree boundaries represented by your decision tree on the data plot figure below.



- (c) What class does your decision tree (above) predict for the new point: (1, -2)?

Solution: +1

Name: _____

- (d) Decision trees built using our greedy algorithm are a good choice of classifiers for images.
 T **F** Explain.

Solution: While possible to construct decision trees to classify images (and some advanced versions of tree algorithms succeed in doing so), image data tends to have content whose positions vary in the image that are better handled using convolutional neural networks, rather than by the fixed splits along single dimensions selected by our greedy tree algorithm.

- (e) For decision trees built using our greedy algorithm, standardizing feature values is important.
 T **F** Explain.

Solution: Trees split on individual variables and split values, and these splits are simply scaled or shifted by standardization but otherwise do not change, so standardizing features is not necessary. Interpretation of trees is also made more difficult with standardized data.

- (f) A disadvantage of using decision trees for classification is that they can only be used to classify data having two classes.
 T **F** Explain.

Solution: Decision trees extend nicely to multiple class data.

Tic-Tac-Toe Revised

3. (12 points) Tic-tac-toe is a paper-and-pencil game for two players, X and O, who take turns marking the spaces in a 3×3 grid. The player who succeeds in placing three of their marks sequentially in a horizontal, vertical, or diagonal row wins the game. The following example game is won by the first player, X:



In this question, we'll consider a "solitaire" version of tic-tac-toe, in which we assume:

- We are the X player;
- The O player is a fixed (but possibly stochastic) algorithm;
- The initial state of the board is empty, and X has the first move;
- We can select any of the nine squares on our turn;
- We don't know the strategy of the O player or the reward function used by O.

We place an X in an empty square, then an O appears in some other square, and then it's our turn to play again. We receive a +1 reward for getting three X's in a row, reward -1 if there are three O's in a row, and reward 0 otherwise. If we select a square that already has an X or an O in it, nothing changes and it's still our turn.

- (a) We can model this problem as a Markov decision process in several different ways. Here are some possible choices for the state space.
- Jody suggests letting the state space be all possible 3×3 grids in which each square contains one of the following: a space, an O, and an X.
 - Dana suggests using all possible 3×3 grids in which each square contains one of the three options (a space, an O, and an X), and there is an equal number of O's and X's.
 - Chris suggests using all 3×3 tic-tac-toe game grids which appear in games where the players both employ optimal strategies.

- i. Is Dana's suggestion better or worse for tabular Q learning than Jody's? Explain your answer.

Solution: Dana's is better, because some of Jody's states might not be part of any plausible games. Jody's approach covers a much larger state space, including states that cannot arise given the rules of the game.

- ii. Is Chris' suggestion better or worse for tabular Q learning than Jody's? Explain your answer.

Name: _____

Solution: Worse, since we do not know if O plays optimally. We might not cover all possible states. Also, Chris' suggestion may be infeasible, if we do not know the optimal strategies for both players.

(b) Many states of the game are effectively the same due to symmetry.

i. Draw a pair of such states which are the same due to symmetry:

Solution: Horizontal, vertical, two different diagonal symmetries with respect to the line passing through the center; rotations through 90, 180, 270 degrees.

ii. Jordan suggests using a state-space that includes one state that stands for each set of board games that are equivalent due to symmetry. Would this be better or worse for learning than Jody's representation? Explain your answer.

Solution: Better. Jordan's state space representation has fewer states, and should facilitate faster learning.

(c) What is the action space of the MDP with Dana's state space definition?

Solution: Selection of one of the 9 squares.

(d) You get to sit and watch an expert player (who always makes optimal moves) play this game for a long time, and you observe the sequence of state-action pairs that occur in many games. Which of the following machine-learning problem formulations is most appropriate, for you to learn how to play the game? For the item you select, provide the specified additional information (where not "none").

1. supervised regression (describe the loss function)
2. supervised classification (describe the loss function)
3. reinforcement learning of a policy (none)
4. reinforcement learning of a value function (none)

Explain your answer.

Solution: supervised classification (loss function). You learn the mapping from input to output (e.g., the position on the grid, where you need to make the next move). The loss function could be the negative log likelihood between the expert's move and your predicted move.

(e) You get to interact with an implementation of this game for many game instances, selecting your actions, observing the results and rewards. Which of the following machine-learning problem formulations is most appropriate, for you to learn how to play the game? For the item you select, provide the specified additional information (where not "none").

1. supervised regression (describe the loss function)

Name: _____

2. supervised classification (describe the loss function)
3. reinforcement learning of a policy (none)
4. reinforcement learning of a value function (none)

Explain your answer.

Solution: Reinforcement learning of a policy (none).

- (f) Barney wants to solve a tic-tac-toe problem that is exactly the same as the above game (i.e., three in a row/column/diagonal wins), except that it is played on a 100 x 100 grid.
- i. Is it better for Barney to use tabular Q learning or neural-net Q learning? Explain.

Solution: Neural-net Q-learning. A table would be too large.

- ii. Suppose Barney were to use neural-net Q learning; would it help for him to start with a convolutional layer? If your answer is yes, describe four 3 x 3 convolutional filters that would be particularly helpful for this problem.

Solution: Yes. A 3x3 filter that detects vertical, horizontal, or diagonal lines can be very useful in detecting local solution (both for X and O).

- (g) Suppose you apply Q-learning to the 3x3 tic-tac-toe problem, and your actions always select an unfilled square. Bert suggests that it is okay to let the discount factor be 1. Is that true? Explain why or why not.

Solution: Yes. The game has a finite number of steps.

CNN Backpropagation

4. (10 points) Conne von Lucien has many pictures from her trip to Flatland and wants to determine which ones have her in the image. All of the pictures are arrays of size 4×1 , with array values of either 0 or 1. Conne looks like the vector $[1, 0, 1]$ in one dimension, so if a picture contains the pattern $[1, 0, 1]$ anywhere inside it, it should be classified as a positive example, otherwise as a negative example.

Fortunately, you learned about CNNs and have helped Conne by designing the following network architecture with three layers:

1. A convolutional layer with one filter W that is size 3×1 , and stride 1, and a single bias w_0 (where the output pixel corresponds to the input pixel that the filter is centered on). Input values of 0 should be assumed beyond the boundaries of the input.
2. A max-pooling layer P with size 2×1 and stride 2.
3. A fully connected layer $\sigma(\cdot)$ with a single output unit having a sigmoidal activation function.

- (a) What is the shape of the output of each layer?

Solution: 4×1

Solution: 2×1

Solution: 1×1 , scalar

- (b) What loss function is most appropriate here, especially if you want your neural network package to be useful with few modifications, to other Flatland visitors (who may appear as longer vectors)?

- A. NLL loss**
- B. Hinge loss
- C. Quadratic loss

Name: _____

- (c) We can express the loss function as $L(\sigma(P), y)$ where P is the output from the max pooling layer of the CNN and y is the true label for the input. Given $\frac{dL}{dP}$, derive the update rule for w_1 if the filter is composed of $W = [w_1, w_2, w_3]^T$ with bias w_0 , and step size is η .

Solution:

Consider Z to be the outputs of layer 1, $Z = [z_1, z_2, z_3, z_4]^T$.

$$z_1 = w_1 \cdot 0 + w_2 x_1 + w_3 x_2 + w_0$$

$$z_2 = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

$$z_3 = w_1 x_2 + w_2 x_3 + w_3 x_4 + w_0$$

$$z_4 = w_1 x_3 + w_2 x_4 + w_3 \cdot 0 + w_0$$

$$P = [p_1, p_2]^T$$

$$p_1 = \max(z_1, z_2)$$

$$\frac{dp_1}{dw_1} = 0 \text{ if } z_1 > z_2 \text{ else } x_1$$

$$p_2 = \max(z_3, z_4)$$

$$\frac{dp_2}{dw_1} = x_2 \text{ if } z_3 > z_4 \text{ else } x_3$$

$$\frac{dP}{dw_1} = \left[\frac{dp_1}{dw_1}, \frac{dp_2}{dw_1} \right]^T$$

$$w_1 := w_1 - \eta \frac{dL}{dP} \frac{dP}{dw_1}$$

- (d) Given $\frac{dL}{dP}$, provide the update rule for w_0 , the bias to the filter.

Solution:

$$\frac{dP}{dw_0} = [1, 1]^T$$

$$w_0 := w_0 - \eta \frac{dL}{dP} \frac{dP}{dw_0}$$

Name: _____

- (e) Conne decides to use the neural network code as written by a 6.036 student for the 6.036 homework (and that actually was a correct implementation) to train her CNN using SGD. The `sgd` procedure may be called multiple times from elsewhere (e.g., to implement multiple epochs of SGD). Conne thinks she has a better `sgd` python procedure than that given in the package; her code is:

```
1 def sgd(nn, X, Y, iters=100, lrate=0.005):
2     D, N = X.shape
3     sum_loss = 0
4     for k in range(iters):
5         Xt = X[:, k:k+1]
6         Yt = Y[:, k:k+1]
7         Ypred = nn.forward(Xt)
8         sum_loss += nn.loss.forward(Ypred, Yt)
9         err = nn.loss.backward()
10        nn.backward(err)
11        nn.sgd_step(lrate)
```

Here, `nn` is an instance of the `Sequential` class implementing the CNN. She knows from the unit tests that the `nn` routines function properly. In particular, `nn.forward` properly computes the predicted outputs `Ypred` from input data `Xt`, `nn.loss.forward` also properly computes the forward loss, `nn.loss.backward` properly computes the backward loss, `nn.backward` properly computes the backward gradients, and `nn.sgd_step` properly applies an SGD update step with the specified learning rate `lrate`. And the N sets of dimension D input data `X`, and labels `Y` are known to be correct.

However, Conne's procedure consistently gives poor results (and occasionally throws errors), compared with the 6.036 student's correct SGD routine, when run with identical arguments.

Why? Specify the line(s) which have errors, and describe how the code should be improved to do as well as the correct implementation of the 6.036 student:

Solution: Lines 5 and 6

Solution:

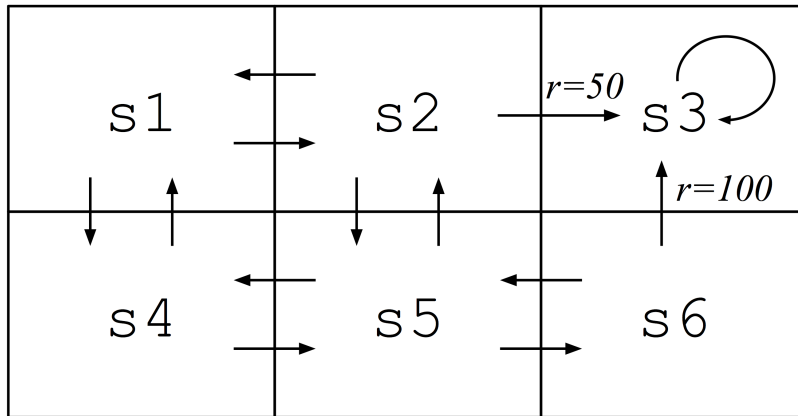
The SGD algorithm needs a *random* data point to be selected for the gradient computation. Thus, the `Xt` and `Yt` assignments should draw from a randomly chosen `j`, e.g.:

```
1     for k in range(iters):
2         j = np.random.randint(N)
3         Xt = X[:, j:j+1]
4         Yt = Y[:, j:j+1]
5         ...
```

Note that Conne's code may throw errors when `iters` $\geq N$.

Robots

5. (16 points) Consider the following deterministic Markov Decision Process (MDP), describing a simple robot grid world. Notice that the values of the immediate rewards r for **two** transitions are written next to them; the other transitions, with no value written next to them, have an immediate reward of $r = 0$. **Assume the discount factor γ is 0.8.**



- (a) For states $s \in \{s6, s5, s2\}$, write the value for $V_{\pi^*}(s)$, the discounted infinite horizon value of state s using an optimal policy π^* . It is fine to write a numerical expression—you don't have to evaluate it—but it shouldn't contain any variables.

Solution:

$$V_{\pi^*}(s6) = 100$$

Solution:

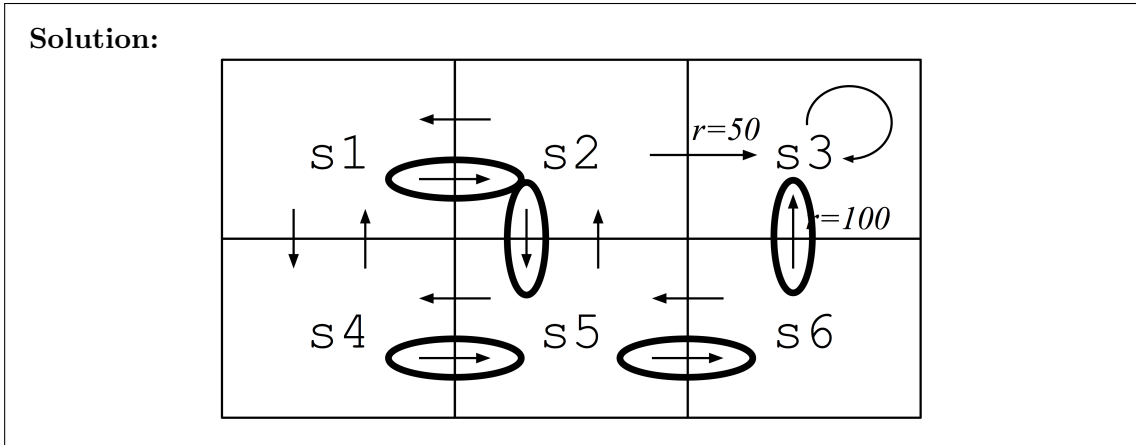
$$V_{\pi^*}(s5) = \gamma V_{\pi^*}(s6) = 80$$

Solution:

$$V_{\pi^*}(s2) = \gamma V_{\pi^*}(s5) = 64$$

Name: _____

- (b) For each state in the state diagram below, circle exactly one outgoing arrow, indicating an optimal action $\pi^*(s)$ to take from that state. If there is a tie, it is fine to select any action with optimal value.



- (c) Give a value for γ (constrained by $0 < \gamma < 1$) that results in a different optimal policy, and describe the resulting policy by indicating which $\pi^*(s)$ values (i.e., which policy actions) change.

Solution: A small $\gamma = 0.001$ will make it not worthwhile to defer gains for very long. In this problem, if $\gamma^2 100 < 50$, then it will be better to directly take the 50 reward. So valid answers here are $0 < \gamma < \frac{\sqrt{2}}{2}$.

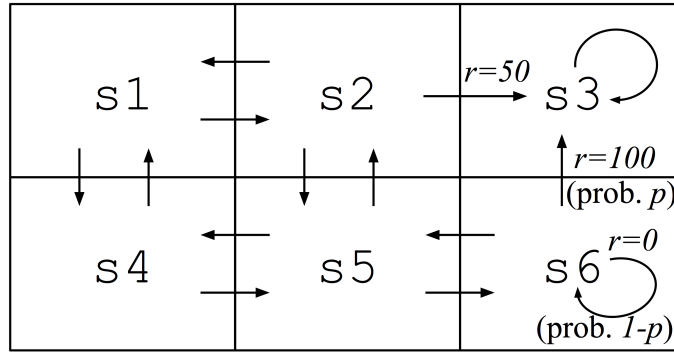
Solution: Now $\pi^*(s2)$ is to go right (east).

- (d) Is it possible to change the immediate reward for each state in such a way that V_{π^*} changes but the optimal policy π^* remains unchanged? If yes, provide a new reward function, and explain how the resulting V_{π^*} changes but π^* does not. Otherwise, explain in at most two sentences why this is impossible.

Solution: Yes. We can establish small immediate rewards, say $r = 1$, for all of the transitions currently with $r = 0$. These are not enough to change the π^* decisions, but do change V_{π^*} for all of these states.

Name: _____

When winter comes, snow also appears on one path in the grid world, making exactly one of the actions non-deterministic. The resulting MDP is shown below. Specifically, the change is that now the result of the action “go north” from state **s6** results in one of two outcomes. With probability p , the robot succeeds in transitioning to state **s3** and receives immediate reward 100. However, with probability $(1 - p)$ it slips on the ice, and remains in state **s6** with 0 immediate reward. **Assume again that the discount factor $\gamma = 0.8$.**



- (e) Assume $p = 0.75$. For each of the states $s \in \{s2, s5, s6\}$, write the value for $V_{\pi^*}(s)$. It is fine to write a numerical expression, but it shouldn't contain any variables.

Solution:

$$V_{\pi^*}(s6) = 100p + (1 - p)\gamma V_{\pi^*}(s6)$$

$$V_{\pi^*}(s6)(1 - (1 - p)\gamma) = 100p$$

$$V_{\pi^*}(s6) = \frac{100p}{1 - (1 - p)\gamma} = 93.75$$

Solution:

$$V_{\pi^*}(s5) = \gamma V_{\pi^*}(s6) = 75$$

Solution:

$$V_{\pi^*}(s2) = \gamma V_{\pi^*}(s5) = 60$$

- (f) How bad does the ice have to get before the robot will prefer to completely avoid the ice? Let us answer the question by giving a value for p for which the optimal policy chooses actions that completely avoid the ice, i.e., choosing the action “go left” over “go up” when

Name: _____

the robot is in the state s_6 . Approach this in four parts. The answer to each of the first three parts can be a numerical expression; the answer to the last part can be an expression involving numbers and p .

- i. What is the value V of going right in state s_2 ?

Solution: 50

- ii. What is the value V of going up in state s_5 , if you're going to go right in state s_2 ?

Solution: $\gamma \cdot 50 = 40$

- iii. What is the value V of going left in state s_6 , if you're going to go up in state s_5 and right in state s_2 ?

Solution: $\gamma^2 \cdot 50 = 32$

- iv. Under what condition on p is it better to go left in state s_6 (then up in state s_5 and right in state s_2) than it is to go up in state s_6 ?

Solution:

$$\frac{p \cdot 100}{1 - (1 - p) \cdot 0.8} < 32$$
$$p < \frac{8}{93} \approx 0.086$$

Recommender System

6. (14 points) After taking 6.036, Bob decides to train a recommender system to predict what ratings different customers will give to different movies. Currently, he knows of three really popular movies, and he knows of two potential customers who have ranked some of these movies. The data matrix currently looks like: $Y = [[2, ?, 3], [4, 2, ?]]$ where, as in class, rows correspond to customers and columns correspond to movies, and ? indicates a missing or unknown ranking. He decides to find a low rank factorization of Y using the alternating least squares algorithm implemented in class. Assume for this question that offsets are set to 0.

- (a) Bob starts out by trying a rank 1 factorization of Y as UV^T . He initializes $U = [1, 2]^T$. Assume there is no regularization. In the first iteration of alternating least squares, we will find the best V given the current U . What is the objective function $J(V)$ in terms of V ? Write it in terms of V_1, V_2, V_3 and specific numerical values from Y .

Solution:

$$J(V) = (1 \cdot V_1 - 2)^2 + (1 \cdot V_3 - 3)^2 + (2 \cdot V_1 - 4)^2 + (2 \cdot V_2 - 2)^2$$

- (b) What is the optimal value of V ?

Solution: The optimal value is $V = [2, 1, 3]^T$. We are fortunate in being able to exactly match all of the non-empty Y elements.

- (c) What is the associated overall training error?

Solution: The training error is 0.

- (d) Bob is happy about what he has accomplished, until he realizes that there are a bunch of movies and users that he still needs to add to his database! He sees that his database will slowly grow over time, and that it will be time-consuming to train a completely new model every single time he updates his database. If Bob has an $m \times n$ data matrix which he wants to find a rank k factorization for, his analysis indicates that the worst-case run-time (in terms of number of expensive multiplications) of performing alternating least squares for t iterations (where each iteration updates both U and V) will be $O(k^2mnt)$.

Instead, Bob comes upon the following idea: whenever he gets information about a new movie, he adds an extra row to V but does not alter the existing entries of U or V . He then finds the values of the entries in that extra row that minimize the objective function (with no regularization). He performs a similar procedure when he gets a new user, but instead adds an extra row to U . Working from the dataset in the first part, say Bob receives a new movie to which his first user has given the rating 4. What is the updated value of V ?

Name: _____

Solution: Adding an extra row to V makes $V^T = [2, 1, 3, a]$. Then, multiplying out UV^T , the rating the first user gives to the new movie is a . To minimize the objective function, we can thus set $a = 4$.

Thus, $V = [2, 1, 3, 4]^T$

- (e) With this updated V , what rating does Bob predict that the second user will give this movie?

Solution: $(UV^T)_{14} = 2 \cdot a = 2 \cdot 4 = 8$ from above.

- (f) Bob continues using this update scheme whenever he adds new movies and users. Does the order in which Bob receives new information affect the final values of U and V that he learns? If your answer is yes, explain in detail why they are different. If your answer is no, explain in detail why they are the same.

Solution: Yes. Let us say Bob gets information about movie k and person a in that order. Based on this new update scheme, the row in V corresponding to movie k will be frozen after the information is received, and will not be updated when the information about person a is received. On the other hand, the learned row in U corresponding to person a will depend in part on the previously updated row in V corresponding to movie k .

If the information was received in the opposite order, we would have the opposite result. The row in U corresponding to person a would be frozen after the first piece of information was received, and not be influenced by the information about movie k . Meanwhile, the row in V corresponding to movie k would be learned in part based on the information gained about person a previously.

Thus, the order of new information matters a lot in this new scheme, because U and V aren't jointly optimized completely every time new information is received.

- (g) Let us say that Bob modifies this procedure, so that he still adds new movies and users in this way, but after every 100 new additions, he retrains U and V from scratch using alternating least squares. Would you expect that this method would make better predictions than if we just used Bob's original procedure? Explain.

Solution: Yes.

Whenever we retrain U and V from scratch, we are minimizing the objective function over all variables in the problem (all entries of U and V) so the minimum of the objective will be lower than we could obtain by just retraining a subset of variables, as we were doing in the previous part to lower computational costs.

Name: _____

Thus, this method will make better predictions.

- (h) After having added a few thousand users and movies to his database, Bob wants to try analyzing the user and movie vectors that he has learned, in order to see whether he can interpret what is causing customers to like certain movies over others. However, some of the numbers in U and V have a very high magnitude, which may lead to problems with numerical precision. How might Bob adjust his training process to fix the problem of high magnitude numbers in U and V ?

Solution: In order to have fewer numbers of large magnitude, Bob can employ regularization of both U, V .

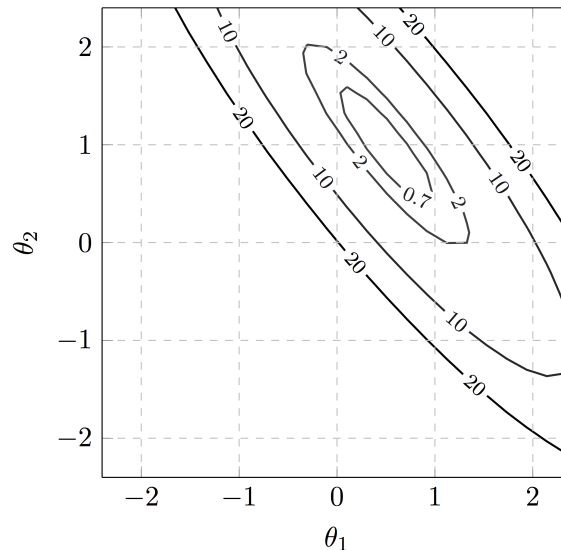
Linear Regression with Regularization

7. (14 points) In this problem, we consider using linear regression with a regularization term. Assume a dataset of n samples $\{(x^{(i)}, y^{(i)})\}$ with $x^{(i)} \in \mathbb{R}^2$ and output values $y^{(i)} \in \mathbb{R}$. Recall that the ridge regression objective is defined as follows:

$$J_{\text{ridge}}(\theta) = J_{\text{data}}(\theta) + J_{\text{reg}}(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 + \lambda \|\theta\|^2$$

where $\theta = [\theta_1, \theta_2]$ and λ is the regularization trade-off parameter.

Chris would like to solve the problem of computing θ that minimizes the ridge regression objective. He will employ graphical methods to obtain the solution. When plotting *just* the data error term, $J_{\text{data}}(\theta)$, as a function of θ_1 and θ_2 , the following set of isocontour lines (curves connecting sets of θ_1, θ_2 for which the objective value is constant) is obtained, for his dataset:



- (a) What is the optimum solution θ^* when you minimize only the data error term, $J_{\text{data}}(\theta)$, i.e., for $\lambda = 0$? Give an approximate value, for Chris's data.

Solution: [0.5, 1.0]

- (b) In general, is the data error term $J_{\text{data}}(\theta^*)$ guaranteed to be zero for the optimal value of θ , for the case when $\lambda = 0$? Explain.

Solution: No. Since we are not likely to perfectly fit all of the data, the data term error is likely to be larger than zero even for the optimal θ value.

Name: _____

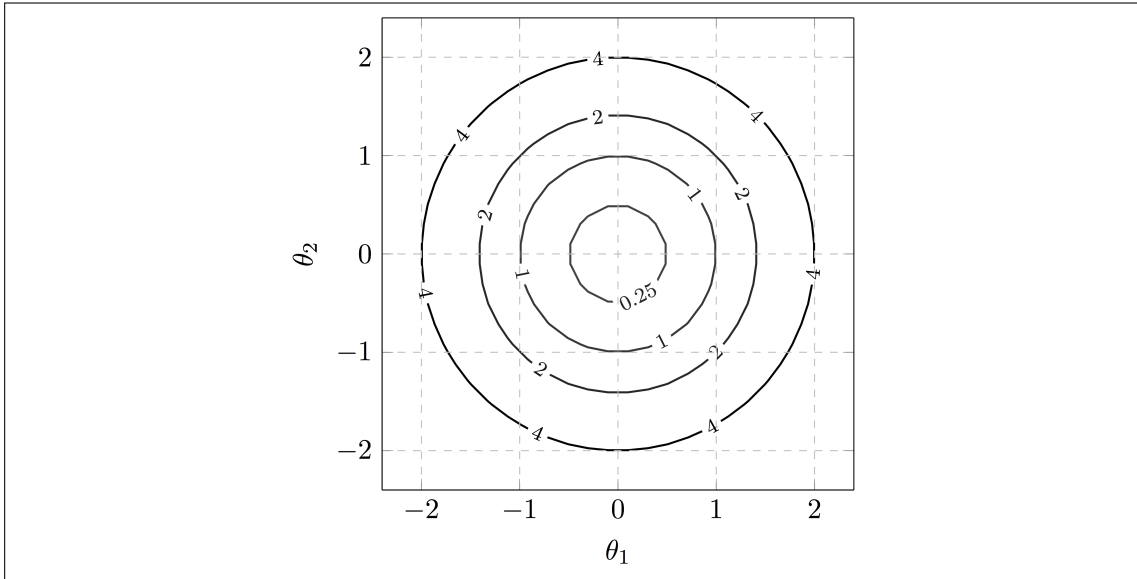
- (c) Recall that $\nabla J_{data}(\theta)$ is a vector in 2D. In general, at any parameter vector θ , describe the geometric relationship between $\nabla J_{data}(\theta)$ and the isocontour line of the data error term $J_{data}(\theta)$ that passes through θ .

Solution: The vector $\nabla J_{data}(\theta)$ is locally perpendicular to the isocontour line of the data error term $J_{data}(\theta)$ at θ . The gradient points in the “uphill” direction.

- (d) What is $\nabla J_{data}(\theta^*)$ at the optimum θ^* , when $\lambda = 0$?

Solution: $\nabla J_{data}(\theta^*) = 0$.

- (e) Now we consider regularization. Sketch the isocontour lines for just the regularization term, $J_{reg}(\theta)$. Clearly label the contour line corresponding to the values of θ for which this term has value 1, when $\lambda = 1$.



Name: _____

- (f) What is the effect of the regularization trade-off parameter λ on the **shape** and **value** of the isocontour lines of the regularization term $J_{reg}(\theta)$?

Solution: The shape remains concentric circles centered at the origin. λ scales the isocontour value for each radius of these concentric circles. (Note that for a constant isocontour value, the radius then decreases.) Visualizing the shape as a bowl, larger λ makes the bowl steeper.

- (g) Now consider the gradient of the regularization term $\nabla J_{reg}(\theta)$. Towards what specific point does the $-\nabla J_{reg}(\theta)$ vector point to?

Solution: The origin, (0,0).

- (h) If λ is *very* large, what is the θ^* that minimizes $J_{ridge}(\theta^*)$? What approximate numerical value does $J_{data}(\theta^*)$ have for Chris's data?

Solution: λ being very large forces θ^* to be very nearly $[0, 0]$. Looking at the plot at the start of the problem, we see that the J_{data} at the origin is approximately 20.

- (i) Given a general optimal solution θ^* for $J_{ridge}(\theta)$ for a given (finite) λ , what is the algebraic relationship between $\nabla J_{data}(\theta^*)$ and $\nabla J_{reg}(\theta^*)$?

Solution: We know that $\nabla J_{ridge}(\theta^*) = 0$ at the optimal point. This forces $\nabla J_{data}(\theta^*) = -\nabla J_{reg}(\theta^*)$.

Training Neural Networks with Regularization

8. (8 points) In this problem we will investigate regularization for neural networks.

Kim constructs a fully connected neural network with $L=2$ layers using mean squared error (MSE) loss and ReLU activation functions for the hidden layer, and a linear activation for the output layer. The network is trained with a gradient descent algorithm on a data set of n points $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$.

Recall that the update rule for weights W^1 can be specified in terms of step size η and the gradient of the loss function with respect to weights W^1 . This gradient can be expressed in terms of the activations A^l , weights W^l , pre-activations Z^l , and partials $\frac{\partial L}{\partial A^2}$, $\frac{\partial A^l}{\partial Z^l}$, for $l = 1, 2$:

$$W^1 := W^1 - \eta \sum_{i=1}^n \frac{\partial L(h(x^{(i)}; W), y^{(i)})}{\partial W^1},$$

where $h(\cdot)$ is the input-output mapping implemented by the entire neural network, and

$$\frac{\partial L}{\partial W^1} = \frac{\partial Z^1}{\partial W^1} \cdot \frac{\partial A^1}{\partial Z^1} \cdot W^2 \cdot \frac{\partial A^2}{\partial Z^2} \cdot \frac{\partial L}{\partial A^2}.$$

- (a) Derive a new update rule for weights W^1 which also penalizes the sum of squared values of all individual weights in the network:

$$L^{new} = L(h(x^{(i)}; W), y^{(i)}) + \lambda \|W\|^2$$

where λ denotes the regularization trade-off parameter. You can express the new update rule as follows:

$$W^1 := \alpha W^1 - \eta \sum_{i=1}^n \frac{\partial L(h(x^{(i)}; W), y^{(i)})}{\partial W^1}$$

where $L(\cdot)$ represents the previous prediction error loss.

What is the value of α in terms of λ and η ?

Solution:

$$L^{new} = L + \lambda \sum_{i,j,l} (W_{i,j}^l)^2$$

$$\frac{\partial L^{new}}{\partial W^1} = \frac{\partial L}{\partial W^1} + 2\lambda W^1$$

$$W^1 := W^1 - \eta \sum \frac{\partial L^{new}}{\partial W^1}$$

$$W^1 := (1 - 2\lambda\eta)W^1 - \eta \sum \frac{\partial L}{\partial W^1}$$

Thus $\alpha = 1 - 2\lambda\eta$.

Name: _____

- (b) Explain how this new update rule helps the neural network reduce overfitting to the data.

Solution: For reasonable λ and η , the weights are scaled by a factor less than 1 at each iteration. (If $1 - 2\lambda\eta > 1$, the weights will rapidly grow and diverge.) A value of $|\alpha| < 1$ pushes the weights toward zero in general, except those weights that are needed to fit substantial subsets of the data (i.e., those weights that are needed to keep the data loss term L low).

- (c) Given that we are training a neural network with gradient descent, what happens when we increase the regularization trade-off parameter λ too much, while holding the step size η fixed?

Solution: With too large a λ , α may approach zero and the weights would be too strongly penalized and thus tend to zero, preventing the neural network from fitting the available training data. That is to say, the network is pushed towards an overly “generalized” constant output based on zero or near-zero weights. With even larger values of λ , α may become negative causing oscillations in weights. With $|\alpha|$ larger than 1, the weights will grow in magnitude without bound.

Lost in Translation

9. (10 points) We want to make an RNN to translate English to Martian. We have a training set of pairs $(e^{(i)}, m^{(i)})$, where $e^{(i)}$ is a sequence of length $J^{(i)}$ of English words and $m^{(i)}$ is a sequence of length $K^{(i)}$ of Martian words. The sequences, even within a pair, do not need to be of the same length, i.e., $J^{(i)}$ need not equal $K^{(i)}$. We are considering two different strategies for turning this into a transduction or sequence-to-sequence learning problem for an RNN.

Method 1: Construct a training-sequence pair (x, y) from an example (e, m) by letting

$$\begin{aligned}x &= (e_1, e_2, \dots, e_L, \text{stop}) \\y &= (m_1, m_2, \dots, m_L, \text{stop})\end{aligned}$$

In Method 1, we assume that if the original e and m had different numbers of words, then the shorter sentence is padded with enough time-wasting words (“ummm” for English, “grlork” for Martian) so that they now have equal length, L . Any needed padding words are inserted at the end of $e^{(i)}$, and at the start of $m^{(i)}$.

Method 2: Construct a training-sequence pair (x, y) from an example (e, m) by letting

$$\begin{aligned}x &= (e_1, e_2, \dots, e_J, \text{stop}, \text{blank}, \dots, \text{blank}) \\y &= (\text{blank}, \dots, \text{blank}, m_1, m_2, \dots, m_K, \text{stop})\end{aligned}$$

In Method 2, blanks are inserted at the end of e and start of m such that the length of x and y are now both $J + K + 1$.

- (a) Assume an element-wise loss function $L_{elt}(p, y)$ on predicted versus true Martian words. What is an appropriate sequence loss function for **Method 1**? Assume that the predicted sequence p has the same length as the target sequence y .

Solution:

$$L_{seq} = \sum_{i=1}^{L+1} L_{elt}(p_i, y_i)$$

The RNN should seek to output the correct Martian words, as well as the *stop* indicator.

- (b) Assume an element-wise loss function $L_{elt}(p, y)$ on predicted versus true Martian words. What is an appropriate sequence loss function for **Method 2**? Assume the predicted sequence p has the same length as the target sequence y .

Solution:

$$L_{seq} = \sum_{i=J+1}^{J+K+1} L_{elt}(p_i, y_i)$$

It's really only necessary that the RNN correctly outputs the whole Martian sequence and the final *stop* indicator. But, it's okay if you sum starting from the first token, $i = 1$.

Name: _____

- (c) Which method is likely to need a higher dimensional state? Explain why.

Solution: Method 2 likely needs to have a larger state to hold a representation of the full input sentence e , while Method 1 might have a shorter state that enables mapping of individual words or shorter sub-sequences of words to corresponding output words or sub-sequences.

- (d) Which method is better if English and Martian have very different word order? Explain why.

Solution: Method 2 since it can first parse the entire input sentence, and then output in a different word order.

- (e) Martian linguist Grlymp thinks it is also important to pad the original English and Martian sentences with time-wasting word to be of the same length for Method 2 (i.e., so that $J = K$), but English linguist Chome Nimsky disagrees. Who is correct, and why?

Solution: Chome Nimsky is right: Method 2 already has full flexibility in processing the entire sentence e before outputting m , so additional time-wasting words would not help (and may hurt) in expressiveness and/or training.