$\begin{array}{cccc} 10\text{-}601 \text{ Machine Learning} & \text{Name:} \\ \text{Spring 2025} & \text{Andrew ID:} \\ \text{Practice Problems} & \text{Room:} \\ \text{Updated: April 22, 2025} & \text{Seat:} \\ \text{Time Limit: N/A} & \text{Exam Number:} \end{array}$

Instructions:

• Verify your name and Andrew ID above.

- This exam contains 29 pages (including this cover page). The total number of points is 0.
- Clearly mark your answers in the allocated space. If you have made a mistake, cross out the invalid parts of your solution, and circle the ones which should be graded.
- Look over the exam first to make sure that none of the 29 pages are missing.
- No electronic devices may be used during the exam.
- Please write all answers in pen or darkly in pencil.
- You have N/A to complete the exam. Good luck!

Question	Points
1. CNNs and RNNs	0
2. Transformers	0
3. Reinforcement Learning	0
4. Ensemble Methods	0
5. Recommender Systems	0
6. K-Means	0
7. Principal Component Analysis	0
Total:	0

Instructions for Specific Problem Types

For "Select One" questions, please fill in the appropriate bubble completely:

Select One: Who taught this course?

- Matt Gormley
- O Marie Curie
- Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

Select One: Who taught this course?

- Henry Chai
- O Marie Curie
- Noam Chomsky

For "Select all that apply" questions, please fill in all appropriate squares completely:

Select all that apply: Which are instructors for this course?

- Matt Gormley
- Henry Chai
- □ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

Select all that apply: Which are the instructors for this course?

- Matt Gormley
- Henry Chai
- I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

Fill in the blank: What is the course number?

10-601

10-8301

1 CNNs and RNNs (0 points)

1.1. Let's begin by considering some of the high-level components of a CNN kernel along with the basic motivation.

(a) What is a kernel?

(b) Why do we need stride, and what benefits/tradeoffs might different values of stride have on the output?

(c) What functionality does padding add to the kernel? Why might we want to use it?

1.2. Consider the following image, filter, and output shape, which you have seen in prior homework in this course.

$$F = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \qquad Y = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & j & k & l \\ m & n & o & p \end{bmatrix}$$

The shape of this particular Y is that of an kernel using no padding and a stride of 1.

(a) Suppose we decide that, instead of having our output shape be (4,4), we want

1.5.		I that apply: In which of the following settings is it more appropriate to NN over a CNN?
		Speech recognition
		Facial recognition
		Music composition
		Autocorrect system
		None of the above
1.6.		${\bf False} \colon$ RNN's are helpful in analyzing time series data. Explain your in 1-2 sentences.
	\bigcirc	True
	\bigcirc	False

Transformers (0 points) 2

2.1.	True or False: Masking in transformers is used to enhance the model's focus on specific parts of the input sequence during training.
	○ True
	○ False
2.2.	Short answer: Explain why padding and truncation are important in training transformer models. Provide a scenario where padding is necessary and another where truncation is required, explaining the implications of each in the context of a natural language processing task.
2.3.	In the context of sequence learning with classical RNN architectures, what is the primary challenge associated with forgetting, and why does it occur?
	□ Forgetting occurs due to the network's inability to store information permanently, making it difficult to recall specific details from long sequences.
	☐ The vanishing gradient problem leads to forgetting, as it causes the gradients to become very small, reducing the network's ability to learn dependencies between distant sequence elements.
	□ Forgetting is caused by the network's overemphasis on recent inputs, which overshadows earlier inputs in long sequences.
	$\hfill\Box$ The issue arises from the network's fixed-size memory, which is insufficient for storing all the information from long sequences.
	□ None of the above
2.4.	Short answer: What is the major problem that is addressed by layer normaliza-

3 Reinforcement Learning (0 points)

3.1 Markov Decision Process

Environment Setup (may contain spoilers for Shrek 1)

Lord Farquaad is hoping to evict all fairytale creatures from his kingdom of Duloc, and has one final ogre to evict: Shrek. Unfortunately all his previous attempts to catch the crafty ogre have fallen short, and he turns to you, with your knowledge of Markov Decision Processes (MDP's) to help him catch Shrek once and for all.

Consider the following MDP environment where the agent is Lord Farquaad:

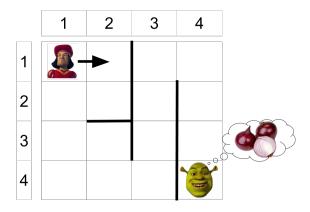


Figure 1: Kingdom of Duloc, circa 2001

Here's how we will define this MDP:

- S (state space): a set of states the agent can be in. In this case, the agent (Farquaad) can be in any location (row, col) and also in any orientation $\in \{N, E, S, W\}$. Therefore, state is represented by a three-tuple (row, col, dir), and S = all possible of such tuples. Farquaad's start state is (1, 1, E).
- A (action space): a set of actions that the agent can take. Here, we will have just three actions: turn right, turn left, and move forward (turning does not change row or col, just dir). So our action space is $\{R, L, M\}$. Note that Farquaad is debilitatingly short, so he cannot travel through (or over) the walls. Moving forward when facing a wall results in no change in state (but counts as an action).
- R(s, a) (reward function): In this scenario, Farquand gets a reward of 5 by moving into the swamp (the cell containing Shrek), and a reward of 0 otherwise.
- p(s'|s,a) (transition probabilities): We'll use a deterministic environment, so this will bee 1 if s' is reachable from s and by taking a, and 0 if not.

3.2. Why is it called a "Markov" decision process? (Hint: what is the assumption made with p?)

3.3. What are the following transition probabilities?

$$p((1,1,N)|(1,1,N),M) =$$

$$p((1,1,N)|(1,1,E),L) =$$

$$p((2,1,S)|(1,1,S),M) =$$

$$p((2,1,E)|(1,1,S),M) =$$

3.4. Given a start position of (1,1,E) and a discount factor of $\gamma=0.5$, what is the expected discounted future reward from a=R? For a=L? (Fix $\gamma=0.5$ for following problems).

3.5. What is the optimal action from each state, given that orientation is fixed at E? (if there are multiple options, choose any)

3.6. Farquaad's chief strategist (Vector from Despicable Me) suggests that having $\gamma = 0.9$ will result in a different set of optimal policies. Is he right? Why or why not?

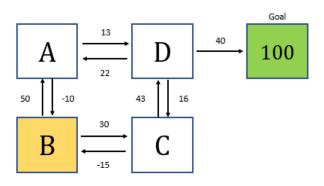
3.7. Vector then suggests the following setup: R(s, a) = 0 when moving into the swamp, and R(s, a) = -1 otherwise. Will this result in a different set of optimal policies? Why or why not?

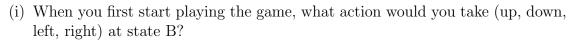
 \Box Unknown Reward Function

 $\hfill\square$ None of the Above

3.8.	Vector now suggests the following setup: $R(s,a)=5$ when moving into the swamp, and $R(s,a)=0$ otherwise, but with $\gamma=1$. Could this result in a different optimal policy? Why or why not?
3.9.	Surprise! Elsa from Frozen suddenly shows up. Vector hypnotizes her and forces her to use her powers to turn the ground into ice. The environment is now stochastic: since the ground is now slippery, when choosing the action M , with a 0.2 chance, Farquaad will slip and move two squares instead of one. What is the expected future-discounted rewards from $s = (2, 4, S)$?
0.0	
3.2	Value and Policy Iteration
3.1.	Select all that apply: Which of the following environment characteristics would increase the computational complexity per iteration for a value iteration algorithm? Choose all that apply:
	□ Large Action Space
	□ A Stochastic Transition Function
	□ Large State Space
	☐ Unknown Reward Function
	\square None of the Above
3.2.	Select all that apply: Which of the following environment characteristics would increase the computational complexity per iteration for a policy iteration algorithm? Choose all that apply:
	☐ Large Action Space
	□ A Stochastic Transition Function
	□ Large State Space

3.3. In the image below is a representation of the game that you are about to play. There are 5 states: A, B, C, D, and the goal state. The goal state, when reached, gives 100 points as reward (that is, you can assume $R(\mathtt{D},\mathtt{right})=140$). In addition to the goal's points, you also get points by moving to different states. The amount of points you get are shown next to the arrows. You start at state B. To figure out the best policy, you use asynchronous value iteration with a decay (γ) of 0.9. You should initialize the value of each state to 0.







(ii) What is the total reward at state B at this time?



(iii) Let's say you keep playing until your total values for each state has converged. What action would you take at state B?



(iv) What is the total reward at state B at this time?



3.4. **Select one:** Let $V_k(s)$ indicate the value of state s at iteration k in (synchronous) value iteration. What is the relationship between $V_{k+1}(s)$ and $\sum_{s' \in S} P(s'|s,a)[R(s,a,s') + \gamma V_k(s')]$, for any $a \in A$? Indicate the most restrictive relationship that applies. For example, if x < y always holds, use < instead of \le . Selecting ? means it's not possible to assign any true relationship. Assume $R(s,a,s') \ge 0 \ \forall s,s' \in S, \ a \in A$.

 $V_{k+1}(s) \ \Box \ \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V_k(s')]$

- \bigcirc =
- \bigcirc <
- O >
- $\bigcirc \leq$
- $\bigcirc \ge$
- \bigcirc ?

3.3 Q-Learning

3.1. For the following true/false, circle one answer and provide a one-sentence explanation:

(i) One advantage that Q-learning has over Value and Policy iteration is that it can account for non-deterministic policies.

Circle one: True False

(ii) You can apply Value or Policy iteration to any problem that Q-learning can be applied to.

Circle one: True False

(iii) Q-learning is guaranteed to converge to the true value Q* for a greedy policy.

Circle one: True False

3.2. For the following parts of this problem, recall that the update rule for Q-learning is:

 $\mathbf{w} \leftarrow \mathbf{w} - \alpha \left(q(\mathbf{s}, a; \mathbf{w}) - (r + \gamma \max_{a'} q(\mathbf{s'}, a'; \mathbf{w})) \nabla_{\mathbf{w}} q(\mathbf{s}, a; \mathbf{w}) \right)$

(i) From the update rule, let's look at the specific term $X = (r + \gamma \max_{a'} q(\mathbf{s}', a'; \mathbf{w}))$ Describe in English what is the role of X in the weight update.

(ii) Is this update rule synchronous or asynchronous?

(iii) A common adaptation to Q-learning is to incorporate rewards from more time steps into the term X. Thus, our normal term $r_t + \gamma * max_{a_{t+1}}q(s_{t+1}, a_{t+1}; w)$

would become $r_t + \gamma * r_{t+1} + \gamma^2 \max_{a_{t+2}} q(\mathbf{s}_{t+2}, a_{t+2} : \mathbf{w})$ What are the advantages of using more rewards in this estimation?

3.3. Select one: Let Q(s,a) indicate the estimated Q-value of state-action pair $(s,a) \in |S| \times |A|$ at some point during Q-learning. Suppose you receive reward r after taking action a at state s and arrive at state s'. Before updating the Q values based on this experience, what is the relationship between Q(s,a) and $r+\gamma \max_{a'\in A} Q(s',a')$? Indicate the most restrictive relationship that applies. For example, if x < y always holds, use < instead of \le . Selecting ? means it's not possible to assign any true relationship.

Q(s,a)	$r + \gamma \max_{a'} Q(s', a')$
\bigcirc	=
\bigcirc	<
\bigcirc	>
\bigcirc	\leq
\bigcirc	\geq
\bigcirc	?

3.4. During standard (not deep) Q-learning, you get reward r after taking action North from state A and arriving at state B. You compute the sample $r + \gamma Q(B, South)$, where $South = \arg \max_a Q(B, a)$.

Which of the following Q-values are updated during this step? (Select all that apply)

Q(A, North)
Q(A, South)
Q(B, North)
Q(B, South)
None of the above

3.5. In general, for Q-Learning (standard/tabular Q-learning, not approximate Q-learning) to converge to the optimal Q-values, which of the following are true?

True or False: It is necessary that every state-action pair is visited infinitely often.

TrueFalse

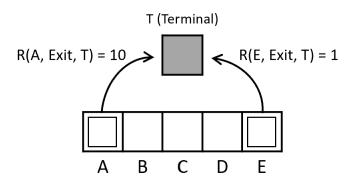
True or False: It is necessary that the discount γ is less than 0.5.

O True

O False

True or False: It is necessary that actions get chosen according to $\arg \max_a Q(s, a)$.

- False
- 3.6. Consider training a robot to navigate the following grid-based MDP environment.



- There are six states, A, B, C, D, E, and a terminal state T.
- Actions from states B, C, and D are Left and Right.
- The only action from states A and E is Exit, which leads deterministically to the terminal state

The reward function is as follows:

- R(A, Exit, T) = 10
- R(E, Exit, T) = 1
- The reward for any other tuple (s, a, s') equals -1

Assume the discount factor is 1. When taking action Left, with probability 0.8, the robot will successfully move one space to the left, and with probability 0.2, the robot will move one space in the opposite direction. When taking action Right, with probability 0.8, the robot will successfully move one space to the right, and with probability 0.2, the robot will move one space in the opposite direction. Run synchronous value iteration on this environment for two iterations. Begin by initializing the value of all states to zero.

Write the value of each state after the first (k = 1) and the second (k = 2) iterations. Write your values as a comma-separated list of 6 numerical expressions in the alphabetical order of the states, specifically V(A), V(B), V(C), V(D), V(E), V(T). Each of the six entries may be a number or an expression that evaluates to a number. Do not include any max operations in your response.

 $V_1(A), V_1(B), V_1(C), V_1(D), V_1(E), V_1(T)$ (Values for 6 states):



$V_2(A), V_2(B), V_2(C), V_2(C)$	$_{2}(D),V_{2}(E),V_{2}(T)$	$({\rm values}$	for 6	states):

$\pi(B), \pi(C), \pi(D)$	based	on	V_2 :

4 Ensemble Methods (0 points)

4.1 AdaBoost

1. In the AdaBoost algorithm, if the final hypothesis makes no mistakes on the training data, which of the following is correct?

Select all that apply:

- $\hfill \Box$ Additional rounds of training can help reduce the errors made on unseen data.
- ☐ Additional rounds of training have no impact on unseen data.
- ☐ The individual weak learners also make zero error on the training data.
- □ Additional rounds of training always leads to worse performance on unseen data.
- 2. **True or False:** In AdaBoost weights of the misclassified examples go up by the same multiplicative factor.
 - \bigcirc True

Round	$D_t(A)$	$D_t(B)$	$D_t(C)$	$D_t(D)$	$D_t(E)$	$D_t(F)$
1	?	?	$\frac{1}{6}$?	?	?
2	?	?	?	?	?	?
219	?	?	?	?	?	?
220	$\frac{1}{14}$	$\frac{1}{14}$	$\frac{7}{14}$	$\frac{1}{14}$	$\frac{2}{14}$	$\frac{2}{14}$
221	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{7}{20}$	$\frac{1}{20}$	$\frac{1}{4}$	$\frac{1}{10}$
			•••			
3017	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{16}$	0
			•••			
8888	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{1}{8}$	$\frac{2}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

3. In the last semester, someone used AdaBoost to train some data and recorded all the weights throughout iterations but some entries in the table are not recognizable. Clever as you are, you decide to employ your knowledge of Adaboost to determine some of the missing information.

Below, you can see part of table that was used in the problem set. There are columns for the Round # and for the weights of the six training points (A, B, C,

D, E, and F) at the start of each round. Some of the entries, marked with "?", are impossible for you to read.

In the following problems, you may assume that non-consecutive rows are independent of each other, and that a classifier with error less than $\frac{1}{2}$ was chosen at each step.

(a) The weak classifier chosen in Round 1 correctly classified training points A, B, C, and E but misclassified training points D and F. What should the updated weights have been in the following round, Round 2? Please complete the form below.

Round	$D_2(A)$	$D_2(B)$	$D_2(C)$	$D_2(D)$	$D_2(E)$	$D_2(F)$
2						

(b)	During Round 219, which of the training points (A, B, C, D, E, F) must
	have been misclassified, in order to produce the updated weights shown at the
	start of Round 220? List all the points that were misclassified. If none were
	misclassified, write 'None'. If it can't be decided, write 'Not Sure' instead.

		1
		ı
		ı
		ı
		ı
		ı
		,

- (c) You observes that the weights in round 3017 or 8888 (or both) cannot possibly be right. Which one is incorrect? Why? Please explain in one or two short sentences.
 - O Round 3017 is incorrect.
 - O Round 8888 is incorrect.
 - O Both rounds 3017 and 8888 are incorrect.

4. What condition must a weak learner satisfy in order for boosting to work? Short answer:

	Short answer:
	train the next weak learner? (Provide an intuitive answer with no math symbols.)
5.	After an iteration of training, AdaBoost more heavily weights which data points to

6. Extra credit Do you think that a deep neural network is nothing but a case of boosting? Why or why not? Impress us.
Answer:

4.2 Random Forests

1. Consider a random forest ensemble consisting of 5 decision trees DT1, DT2 ... DT5 that has been trained on a dataset consisting of 7 samples. Each tree has been trained on a random subset of the dataset. The following table represents the predictions of each tree on its out-of-bag samples.

Tree	Sample Number	Prediction	Actual
DT1	6	No	Yes
DT1	7	No	Yes
DT2	2	No	No
DT3	1	No	No
DT3	2	Yes	No
DT3	4	Yes	Yes
DT4	2	Yes	No
DT4	7	No	Yes
DT5	3	Yes	Yes
DT5	5	No	No

(a)	What	is th	he O()B er	or o	f the	above	e rando	om f	forest	classi	fier?

(b) In the above random forest classifier, which Decision tree(s) will be given the highest weight in inference? If there are multiple trees, mention them all

(c)	To reduce the error of each individual decision tree, Neural uses all the features to train each tree. How would this impact the generalisation error of the random forest?
	\bigcirc The generalisation error would decrease as each tree has lower generalisation error
	O The generalisation error would increase as each tree has insufficient training data
	The generalisation error would increase as the trees are highly correlated

Recommender Systems (0 points) **5**

1.	Applied to the Netflix Prize problem, which of the following methods does NOT always require side information about the users and the movies?
	Select all that apply:
	□ Neighborhood methods
	□ Content filtering
	☐ Latent factor methods
	□ Collaborative filtering
	\square None of the above
2.	Select all that apply:
	\Box Using matrix factorization, we can embed both users and items in the same space
	\square Using matrix factorization, we can embed either solely users or solely items in the same space, as we cannot combine different types of data
	\square In a rating matrix of users by books that we are trying to fill up, the best-known solution is to fill the empty values with 0s and apply PCA, allowing the dimensionality reduction to make up for this lack of data
	\Box Alternating minimization allows us to minimize over two variables
	\Box Alternating minimization avoids the issue of getting stuck in local minima
	\Box If the data is multidimensional, then overfitting is extremely rare
	\Box Nearest neighbor methods in recommender systems are restricted to using euclidian distance for their distance metric
	\square None of the above
3.	Your friend Duncan wants to build a recommender system for his new website DuncTube, where users can like and dislike videos that are posted there. In order to build his system using collaborative filtering, he decides to use Non-Negative Matrix Factorization. What is an issue with Duncan's approach, and what could he change about the website <i>or</i> the algorithm in order to fix it?

6 K-Means (0 points)

		True or False questions, circle your answer and justify it; for $\mathbf{Q}\mathbf{A}$ questions, se down your answer.
	(i)	For a particular dataset and a particular k, k-means always produce the same result, if the initialized centers are the same. Assume there is no tie when assigning the clusters.
		○ True
		○ False
		Justify your answer:
	(ii)	k-means can always converge to the global optimum.
		○ True
		○ False
		Justify your answer:
((iii)	k-means is not sensitive to outliers.
		○ True
		○ False
		Justify your answer:
((iv)	k in k-nearest neighbors and k-means have the same meaning.
		○ True
		○ False
		Justify your answer:
	(v)	What's the biggest difference between k-nearest neighbors and k-means?
		Write your answer in one sentence:

(vi) In k-means, the cost always drops after one update step.

- \bigcirc True
- False
- (ix) When α in k-means++ becomes 0, it means random sampling.
 - True
- 6.2. In k-means, random initialization could possibly lead to a local optimum with very bad performance. To alleviate this issue, instead of initializing all of the centers completely randomly, we decide to use a smarter initialization method. This leads us to k-means++.

The only difference between k-means and k-means++ is the initialization strategy, and all of the other parts are the same. The basic idea of k-means++ is that instead of simply choosing the centers to be random points, we sample the initial centers iteratively, each time putting higher probability on points that are far from any existing center. Formally, the algorithm proceeds as follows.

Given: Data set $x^{(i)}, i = 1, \dots, N$ Initialize:

$$\mu^{(1)} \sim \text{Uniform}(\{x^{(i)}\}_{i=1}^N)$$

For $j = 2, \dots, k$

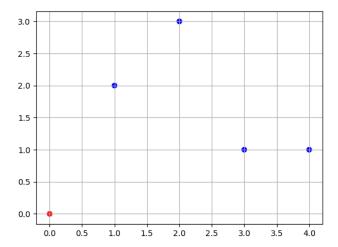
Computing probabilities of selecting each point

$$p_i = \frac{\min_{j' < j} \|\mu^{(j')} - x^{(i)}\|_2^2}{\sum_{i'=1}^N \min_{j' < j} \|\mu^{(j')} - x^{(i')}\|_2^2}$$

Select next center given the appropriate probabilities $\mu^{(j)} \sim \text{Categorical}(\{x^{(i)}\}_{i=1}^{N}, \mathbf{p}_{1:N})$

Note: n is the number of data points, k is the number of clusters. For cluster 1's center, you just randomly choose one data point. For the following centers, every time you initialize a new center, you will first compute the distance between a data point and the center closest to this data point. After computing the distances for all data points, perform a normalization and you will get the probability. Use this probability to sample for a new center.

Now assume we have 5 data points (n=5): (0, 0), (1, 2), (2, 3), (3, 1), (4, 1). The number of clusters is 3 (k=3). The center of cluster 1 is randomly choosen as (0, 0). These data points are shown in the figure below.



(i) What is the probability of every data point being chosen as the center for cluster 2? (The answer should contain 5 probabilities, each for every data point)



(ii) Which data point is mostly liken chosen as the center for cluster 2?

(vii) According to the result of (ii) and (iv), explain how does k-means++ alleviate the local optimum issue due to initialization?

6.3. Consider a dataset with seven points $\{x_1, \ldots, x_7\}$. Given below are the distances between all pairs of points.

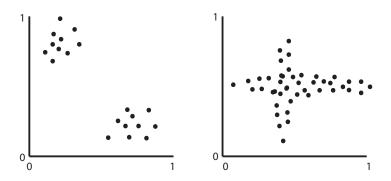
	x_1	x_2	x_3	x_4	x_5	x_6	x_7
x_1	0	5	3	1	6	2	3
x_2	5	0	4	6	1	7	8
x_3	3	4	0	4	3	5	6
x_4	1	6	4	0	7	1	2
x_5	6	1	3	7	0	8	9
x_6	2	7	5	1	8	0	1
x_7	3	8	6	2	9	1	0

Assume that k = 2, and the cluster centers are initialized to x_3 and x_6 . Which of the following shows the two clusters formed at the end of the first iteration of k-means? Circle the correct option.

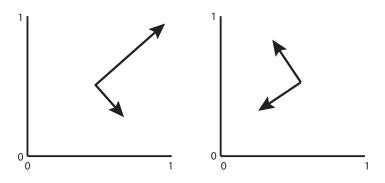
- $\bigcirc \{x_1, x_2, x_3, x_4\}, \{x_5, x_6, x_7\}$
- $\bigcirc \{x_2, x_3, x_5\}, \{x_1, x_4, x_6, x_7\}$
- $\bigcirc \{x_1, x_2, x_3, x_5\}, \{x_4, x_6, x_7\}$
- $\bigcirc \{x_2, x_3, x_4, x_7\}, \{x_1, x_5, x_6\}$

7 Principal Component Analysis (0 points)

7.1. (i) Consider the following two plots of data. Draw arrows from the mean of the data to denote the direction and relative magnitudes of the principal components.



(ii) Now consider the following two plots, where we have drawn only the principal components. Draw the data ellipse or place data points that could yield the given principal components for each plot. Note that for the right hand plot, the principal components are of equal magnitude.



7.2. Circle one answer and explain.

In the following two questions, assume that using PCA we factorize $X \in \mathbb{R}^{n \times m}$ as $Z^T U \approx X$, for $Z \in \mathbb{R}^{m \times n}$ and $U \in \mathbb{R}^{m \times m}$, where the rows of X contain the data points, the rows of U are the prototypes/principal components, and $Z^T U = \hat{X}$.

(i) Removing the last row of U and Z will still result in an approximation of X, but this will never be a better approximation than \hat{X} .

Circle one: True False

(ii) $\hat{X}\hat{X}^T = Z^TZ$.

Circle one: True False

(iii) The goal of PCA is to interpret the underlying structure of the data in terms of the principal components that are best at predicting the output variable.

Circle one: True False

(iv) The output of PCA is a new representation of the data that is always of lower dimensionality than the original feature representation.

Circle one: True False