

Name: \_\_\_\_\_ NetID: \_\_\_\_\_

S&DS 365 / 665

## **Intermediate Machine Learning**

Final Exam (Sample Solution)

Tuesday, December 20, 2022

Complete all of the problems. You have 2.5 hours (150 minutes) to complete the exam.

The exam is closed book, computer, phone, etc. You are allowed one double-sided  $8\frac{1}{2} \times 11$  sheet of paper with hand-written notes.

Please use a black pen or dark pencil so that your exam scans clearly.

1. **Multinomial choice** (25 points)

For each of the following questions, circle the *single best* answer.

1.1. Suppose the lasso is carried out for a response vector  $Y$  on a design matrix  $\mathbb{X}$  that satisfies  $\mathbb{X}\mathbb{X}^T = I_n$ , where  $I_n$  is the  $n \times n$  identity matrix. Let  $\text{Soft}_\lambda(\cdot)$  denote the soft thresholding operator with threshold  $\lambda$ . If the  $\ell_1$  penalty of the lasso is  $\lambda$ , then the solution  $\hat{\beta}$  satisfies which of the following?

- (a)  $\hat{\beta} = \text{Soft}_\lambda(Y)$
- ☒ (b)  $\hat{\beta} = \text{Soft}_\lambda(\mathbb{X}^T Y)$
- (c)  $\hat{\beta} = \text{Soft}_\lambda(\mathbb{X}Y)$
- (d)  $\hat{\beta}$  does not have a closed form
- (e) None of the above

1.2. Consider the regression estimator computed as

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2}(Y - \beta)^2 + \frac{1}{2}\beta^2 + |\beta|$$

where  $Y$  is a random variable and  $\beta$  is a scalar. If  $Y = -2$  the solution is

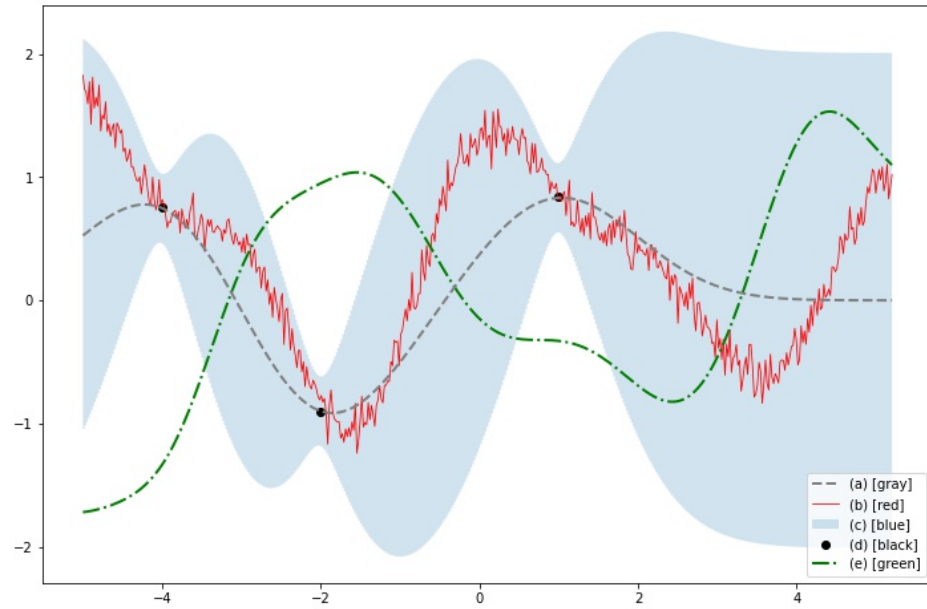
- (a)  $\hat{\beta} = 0$
- (b)  $\hat{\beta} = -1$
- ☒ (c)  $\hat{\beta} = -\frac{1}{2}$
- (d)  $\hat{\beta} = 1$
- (e) None of the above

1.3. Suppose that we have a kernel regression technique in one dimension with bandwidth parameter  $h$  for which the squared bias scales as  $O(\sqrt{h})$  and the variance scales as  $O\left(\frac{1}{\sqrt{nh}}\right)$  as  $h \rightarrow 0$  with  $nh \rightarrow \infty$ , for a sample of size  $n$ , under certain assumptions. What is the fastest rate at which the risk (expected squared error) will decrease with sample size for this technique?

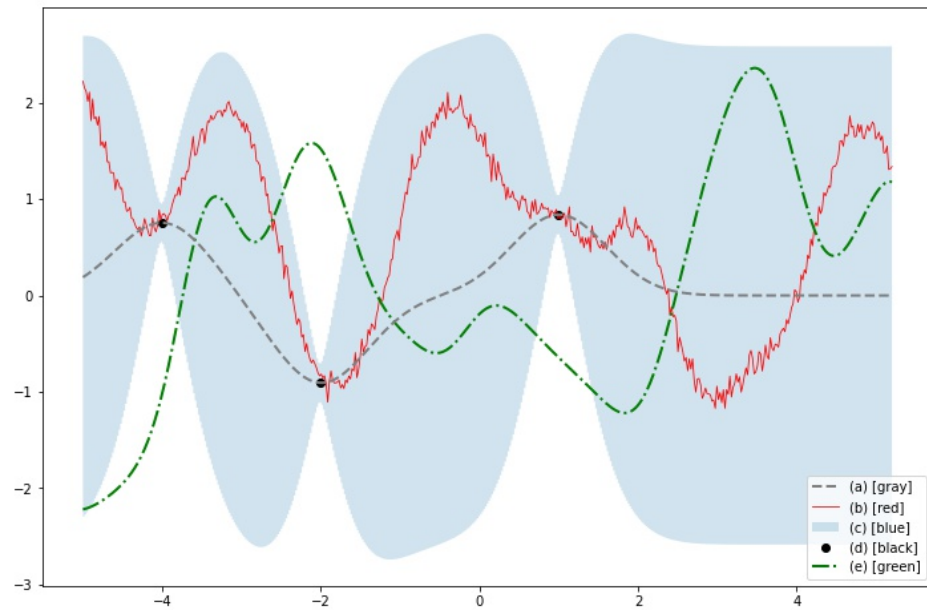
- (a)  $O(n^{-1/3})$
- ☒ (b)  $O(n^{-1/4})$
- (c)  $O(n^{-1/2})$
- (d)  $O(n^{-1})$
- (e) None of the above

1.4. Consider the following two plots associated with a Gaussian process. The data are the same in the two plots; only the hyperparameters are different.

Plot A



Plot B



- (1) What quantities are shown in the plots? Write the appropriate label (a), (b), (c), (d), or (e) next to each of the following.

(e) Sample from prior

(b) Sample from posterior

(d) Data

(a) Posterior mean

(c) Posterior confidence interval

- (2) The two plots are generated with different hyperparameters, kernel bandwidths  $h_A$  and  $h_B$  and measurement noise levels  $\sigma_A^2$  and  $\sigma_B^2$ . Circle the appropriate relations:

Kernel bandwidth:

$$h_A < h_B$$

$$h_A > h_B$$

Measurement noise:

$$\sigma_A^2 < \sigma_B^2$$

$$\sigma_A^2 > \sigma_B^2$$

- (3) Briefly explain your answers to (2) above.

*Bandwidths: The sample from the prior for B is slightly wigglier, and the confidence intervals are larger, both indicating that there is greater variance in the prior. This corresponds to smaller  $h$ .*

*Noise level: The sample from the posterior have larger fluctuations in A, and the confidence intervals at the data points are wider, bot indicating that there is a greater variance in the likelihood term. This corresponds to larger  $\sigma^2$ .*

- 1.5. A convolutional neural network is constructed to classify color images, each with 3 color channels (red, blue, green). The network has two convolutional layers followed by two dense layers.

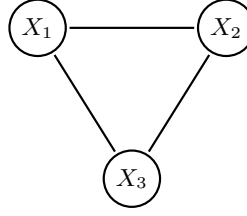
The first convolutional layer uses 10 filters, each  $3 \times 3$ , constructed with a command such as `model.add(layers.Conv2D(10, (3, 3)))`. How many trainable parameters does this layer have?

- (a) 90
- (b) 100
- (c) 270
- ☒ (d) 280
- (e) 900
- (f) 910
- (g) It depends on the number of pixels in the images
- (h) It depends on the pooling operation after the first layer

The second convolutional layer also uses 10 filters, each  $3 \times 3$ , constructed with a command such as `model.add(layers.Conv2D(10, (3, 3)))`. How many trainable parameters does this layer have?

- (a) 90
- (b) 100
- (c) 270
- (d) 280
- (e) 900
- ☒ (f) 910
- (g) It depends on the number of pixels in the images
- (h) It depends on the pooling operation after the first layer

1.6 Consider a graph neural network for binary classification of inputs  $x = (x_1, x_2, x_3)^T$  on the graph below:

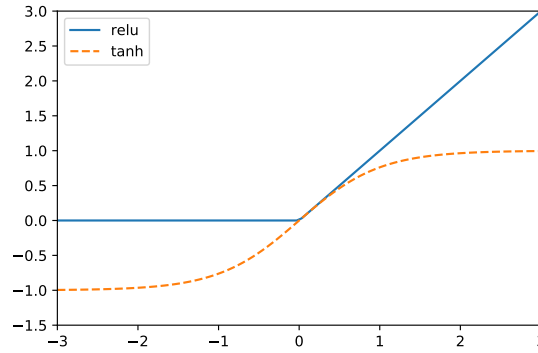


The graph neural network defines a classifier with discriminant function

$$\log \left( \frac{p(Y = 1 | x)}{p(Y = 0 | x)} \right) = \beta^T h(x)$$

with  $h(x) = \varphi(Lx)$  where  $\varphi$  is an activation function and  $L$  is the graph Laplacian (with edge weights 1), and  $\beta = (1, -1, 1)^T$ .

Recall the relu and tanh activation functions look like this, with  $\tanh(-x) = -\tanh(x)$ :



TRUE ☐ FALSE ☒

(1) If  $\varphi$  is relu and  $x = (1, 2, 1)$  then  $p(Y = 1 | x) > \frac{1}{2}$

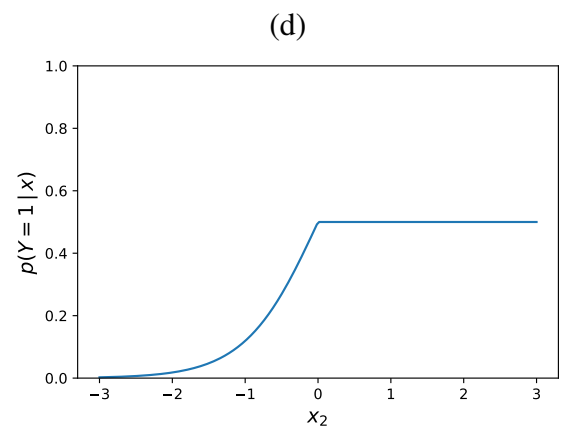
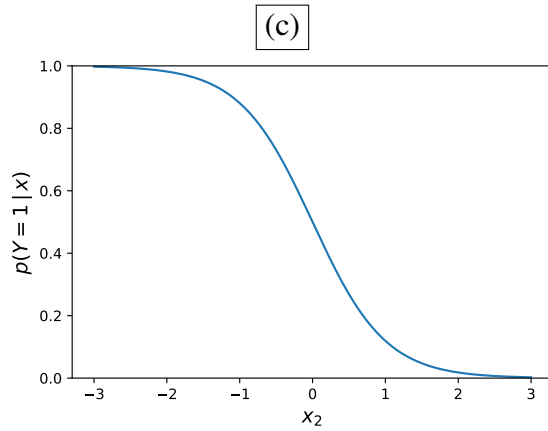
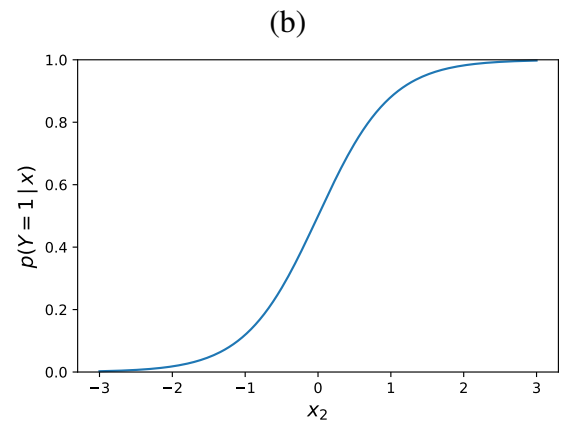
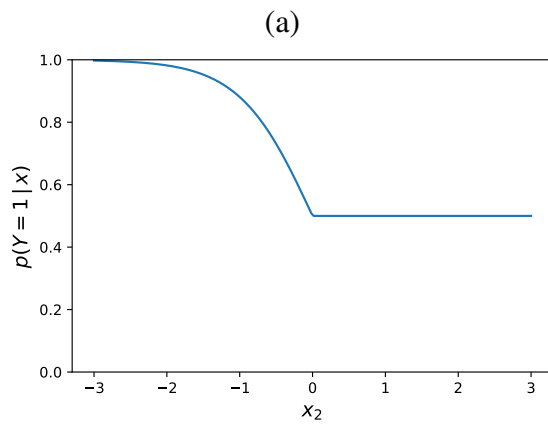
☒ TRUE ☐ FALSE

(2) If  $\varphi$  is tanh and  $x = (1, -2, 1)$  then  $p(Y = 1 | x) > \frac{1}{2}$

TRUE ☐ FALSE ☒

(3) If  $\varphi$  is tanh and  $x = (1, 1, 1)$  then  $p(Y = 1 | x) > \frac{1}{2}$

(4) If  $\varphi$  is relu and  $x = (0, x_2, 0)^T$ , which of the following is a plot of the probability  $p(Y = 1 | x)$  as a function of  $x_2$ ?



## 2. *Short answer* (18 points)

The following three subproblems ask you to explain important concepts associated with topics in the course.

- (a) ***Policy-gradient algorithm for reinforcement learning.*** (1) What problem is the policy gradient algorithm designed to solve, compared with the policy iteration algorithm? (2) How is it different from actor-critic and  $Q$ -learning? (3) Write down mathematical expressions for the parameter update equations used by the policy-gradient algorithm.

*(1) If the state space is large, the policy cannot be computed explicitly and stored in a table, as the policy iteration algorithm does. In a policy-gradient approach, the policy is parameterized in terms of a function of the state, for example using a neural network.*

*(2) Policy gradient only estimates a policy, and doesn't explicitly estimate the value function. In actor-critic, both a policy and a value function are estimated. In  $Q$ -learning, only the  $Q$ -function, and therefore the value function, is estimated.*

*(3) The algorithm is gradient ascent in the expected reward, where the parameters  $\theta$  are updated as*

$$\theta \leftarrow \theta + \eta \nabla_{\theta} J(\theta)$$

*where  $\eta$  is a small step size, and  $J(\theta) = \mathbb{E}_{\theta}(R)$  is the expected reward. To estimate the gradient  $\nabla_{\theta} J(\theta)$  we sample sequences  $\tau_i$  of states, actions, and rewards for a collection of  $N$  episodes using the current policy  $\pi(a | s)$ , and then compute*

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N R(\tau_i) \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

*or*

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T R_t(\tau_i) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

*where the discounted reward  $R_t(\tau_i)$  is computed using a backward recurrence*

$$R_t(\tau_i) = r_{t+1} + R_{t+1}(\tau_i).$$



- (b) **Gated recurrent units (GRUs)**. (1) Describe what problem GRUs are designed to solve.  
 (2) Explain the two types of gates used in a GRU, and how they are implemented.  
 (3) Write down mathematical expressions for the update equations for computing the hidden state vector, assuming one hidden layer.

(1) *“Vanilla” RNNs have the vanishing gradient problem, where the derivatives become so tiny that the parameter updates do not allow information to be incorporated into the state over long sequences. This is a problem because predictions of symbols often depend on symbols far away—a basic example is balanced parentheses, where “)” is likely to follow “(”.*

(2) *GRUs use two types of gates. The update gate regulates information to be put into the state vector without going through the usual vanilla RNN state update; this is how vanishing gradients are avoided. The reset gate regulates when the state vector is started from fresh, by removing the dependence on the previous state. Both gates use linear transformations followed by a sigmoid activation function to keep the values between zero and one.*

(3) *The full update equations are*

$$\Gamma_t^u = \sigma(W_{ux}x_t + W_{uh}h_{t-1} + b_u)$$

$$\Gamma_t^r = \sigma(W_{rx}x_t + W_{rh}h_{t-1} + b_r)$$

$$c_t = \tanh(W_{hx}x_t + \Gamma_t^r \odot W_{hh}h_{t-1} + b_h)$$

$$h_t = (1 - \Gamma_t^u) \odot c_t + \Gamma_t^u \odot h_{t-1}$$

- (c) **Attention and transformers.** (1) Explain the idea of attention, and the problem it is designed to solve. (2) Write down how attention is implemented in terms of queries, keys, and values. (3) Write down a mathematical expression for how kernel regression can be expressed in terms of attention — what are the queries, keys, and values?

*(1) Attention, perhaps better called “alignment,” is meant to address the challenge that the relevant information needed to predict the next symbol can come from different previous locations. For example, the word meow might be aligned with cat in the sentence I heard a meow in the middle of the night and woke up realizing that I had locked the cat in the closet. This is the same problem that GRUs are designed to address, but the attention mechanism is different.*

*(2) Attention is computed using inner products of feature vectors. If  $\phi(w)$  is the vector of neurons corresponding to word  $w$  at a given stage of the network, then the attention between words  $w$  and  $w'$  will depend in the inner product  $\phi(w)^T \phi(w')$ . For example, we might expect that  $\phi(\text{meow})^T \phi(\text{cat})$  will be large but  $\phi(\text{meow})^T \phi(\text{closet})$  will be small. Transformers are neural networks that transform the input sequence in layers, using the usual feedforward layers, but also attention layers that weight neurons using attention. In general, attention layers in transformers are computed in terms of query, key, value triples as*

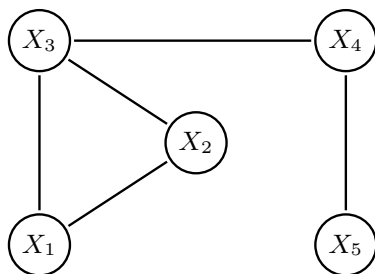
$$\text{Attn}(Q, K, V) = \text{Softmax}\left(QK^T/\sqrt{d}\right)V$$

*where  $Q$ ,  $K$  and  $V$  are linear transformations of neurons at the previous layer.*

*(3) Kernel regression can be viewed in terms of attention by taking the queries to be the test points  $x$ , the keys to be the training input points  $\{x_i\}$  and the values to be the training response values  $\{y_i\}$ .*

3. **Graphical models** (10 points)

Consider the following graph on five variables:

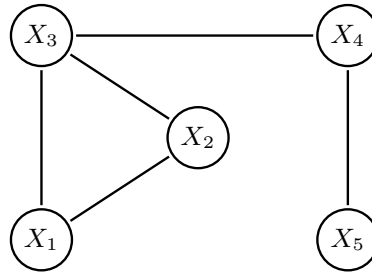


- (a) Write down the general form for a strictly positive probability density that factors (is Markov) with respect to this graph.

$$p(x_1, x_2, x_3, x_4, x_5) = \psi_a(x_1, x_2, x_3) \psi_b(x_3, x_4) \psi_c(x_4, x_5)$$

- (b) If a multivariate Gaussian has this graph, what is the form of the sparsity pattern of the precision matrix  $\Omega$ ? Enter 0 for a zero entry and \* for a non-zero entry.

$$\Omega = \begin{matrix} & \begin{matrix} X_1 & X_2 & X_3 & X_4 & X_5 \end{matrix} \\ \begin{matrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \end{matrix} & \begin{bmatrix} * & * & * & 0 & 0 \\ * & * & * & 0 & 0 \\ * & * & * & * & 0 \\ 0 & 0 & * & * & * \\ 0 & 0 & 0 & * & * \end{bmatrix} \end{matrix}$$



(c) Based on this graph (repeated above) which of the following independence relations hold?  
Circle all that apply.

(1)  $X_1 \perp\!\!\!\perp X_5$

☒ (2)  $X_3 \perp\!\!\!\perp X_5 \mid X_4$

(3)  $X_3 \perp\!\!\!\perp X_4 \mid X_5$

☒ (4)  $X_2 \perp\!\!\!\perp X_5 \mid X_3$

☒ (5)  $X_2 \perp\!\!\!\perp X_5 \mid X_3, X_4$

☒ (6)  $X_1 \perp\!\!\!\perp X_5 \mid X_2, X_4$

#### 4. *Coding: Deep Q-learning* (10 points)

In this problem you are asked to complete the implementation of some code for deep Q-learning, which only looks at the current reward, not future rewards. This corresponds to discount factor  $\gamma = 0$ .

The environment is specified by functions that begin `env_` such as `env_num_actions()` and `env_state_dim()` that specify the number of actions and dimension of the state. The function `env_step()` returns the reward and next state.

You need to provide the missing five lines in the code below.

```
def construct_q_network(state_dim, num_actions):
    inputs = layers.Input(shape=(state_dim,))
    hidden1 = layers.Dense(32, activation="relu")(inputs)
    hidden2 = layers.Dense(16, activation="relu")(hidden1)
    q_values = ... # line (a)
    deep_q_network = keras.Model(inputs=inputs, outputs=q_values)
    return deep_q_network

def loss_function(q_value, reward):
    loss = ... # line (b)
    return loss

exploration_rate = 0.1
num_steps = 2000

num_actions = env_num_actions()
q_network = construct_q_network(env_state_dim(), num_actions)
opt = tf.keras.optimizers.Adam(learning_rate=0.01)
state = env_init_state()

for i in range(num_steps):
    with tf.GradientTape() as tape:
        q_values = q_network(state)

        u = np.random.uniform()
        if u <= exploration_rate:
            action = np.random.choice(num_actions)
        else:
            action = ... # line (c)

        reward, next_state = env_step(action)
        q_value = q_values[0, action]
        loss = loss_function(q_value, reward)

        grads = ... # line (d)
        opt.apply_gradients(zip(grads, q_network.trainable_variables))

    state = ... # line (e)
```

Complete the implementation by giving *a single line of Python code* for each of the five missing lines above. For full credit your Python expressions should be syntactically correct.

```
line (a):  
    q_values = layers.Dense(num_actions)(hidden2)  
  
line (b):  
    loss = (q_value - reward) ** 2  
  
line (c):  
    action = np.argmax(q_values)  
  
line (d):  
    grads = tape.gradient(loss, q_network.trainable_variables)  
  
line (e):  
    state = next_state
```

## 5. Reinforcement learning (15 points)

The following three subproblems are on the topic of reinforcement learning using policy iteration.

- (1) Give the algorithm for policy iteration, and outline a proof that it converges.

*The algorithm alternates between computing the value function for the current policy, and then updating the policy. At each policy evaluation step, we follow a policy  $\pi$  and update the value function  $V(s)$  for each state  $s$  as*

$$V(s) \leftarrow \sum_{s',r} p(s', r|s, \pi(s)) (r + \gamma V(s')).$$

*At each policy improvement step, we improve the current policy  $\pi(s)$  following the updated value function:*

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r|s, a) (r + \gamma V(s')).$$

*In Assignment 5 we outlined a proof that this converges using the following steps. First, supposing we have a policy  $\pi$  and a monotonically improved policy  $\pi'$ , such that for all states  $s \in S$*

$$Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s),$$

*it follows that the value function  $V_{\pi'}$  dominates  $V_{\pi}$ , i.e.  $V_{\pi'}(s) \geq V_{\pi}(s)$ . This is shown by expanding the value functions iteratively. Since the state space is finite and the value function is monotonically increasing, the algorithm converges. Upon convergence,  $Q_{\pi}(s, \pi(s)) = V_{\pi}(s)$  which is the Bellman equation optimality condition.*

- (2) Recall the Taxi problem discussed in class. A taxi navigates a  $5 \times 5$  grid, picking up passengers and delivering them to their desired destinations. There is a default per-step reward of  $-1$ , and a reward of  $+20$  for delivering the passenger. An illegal pickup or drop-off action has a reward of  $-10$ .

The state can be represented as a tuple

`(row, col, passenger_location, destination).`

The “ascii art” figure on the left below shows the grid environment the taxi navigates. Vertical lines `|` are barriers, which prevent the taxi from moving horizontally. The agent (taxi) can take the following six actions shown below.

```

+-----+
|R: | : G|
| : | : |
| : | : |
| | : | :
|Y| : B: |
+-----+

```

- 0: move south (down)
- 1: move north (up)
- 2: move east (right)
- 3: move west (left)
- 4: pickup passenger
- 5: drop off passenger

The passenger is either waiting to be picked up at one of four locations, or is in the taxi being driven to their desired destination. So, for a fixed passenger location and destination, an optimal policy  $\pi^*(row, col)$  assigns an action to each of the  $25 = 5 \times 5$  grid points.

For each of the two tables below, write an integer in each cell for the action to be taken by an optimal policy  $\pi^*$ . We have filled in two actions, which indicate whether the passenger is in the taxi or waiting to be picked up.

0	0	2	2	4
0	0	2	2	1
2	2	2	2	1
1	1	1	1	1
1	1	1	1	1

passenger is outside the taxi

0	0	0	0	0
0	0	0	0	0
0	3	3	3	3
0	1	1	1	1
5	1	1	1	1

passenger is inside the taxi



- (3) Suppose that an environment has three states,  $s = 0$ ,  $s = 1$ , and  $s = 2$ , and an agent can take two actions  $a = 0$  and  $a = 1$ . The environment is deterministic and has the following `reward` and `next` functions:

$$\text{next}(s, a) : \begin{matrix} & a=0 & a=1 \\ \begin{matrix} s=0 \\ s=1 \\ s=2 \end{matrix} & \begin{pmatrix} 1 & 1 \\ 2 & 2 \\ 0 & 0 \end{pmatrix} \end{matrix} \quad \text{reward}(s, a) : \begin{matrix} & a=0 & a=1 \\ \begin{matrix} s=0 \\ s=1 \\ s=2 \end{matrix} & \begin{pmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \end{matrix}$$

Give an optimal (deterministic) policy  $\pi^*(s)$  using discount  $\gamma = \frac{2}{3}$ . Write your answer in the following table, showing which action to take in each state. Show your work for partial credit.

$$\pi^*(s) : \begin{matrix} & s=0 & s=1 & s=2 \\ \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \end{matrix}$$