

Final

● Graded

Student

TIYA CHOKHANI

Total Points

87 / 100 pts

Question 1

1.1

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 2

1.2

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 3

1.3

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 4

1.4

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 5

1.5

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 6

1.6

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 7

1.7

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 8

1.8

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 9

1.9

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 10

1.10

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 11

1.11

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 12

2.0

0 / 0 pts

- 0 pts Correct

- 2 pts Incorrect

Question 13

2.1

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 14

2.2

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 15

2.3

2 / 2 pts

- 0 pts Correct

- 1 pt Partially Correct

- 2 pts Incorrect

Question 16

2.4

3 / 3 pts

- 0 pts Correct

- 3 pts Incorrect

Question 17

2.5

3 / 3 pts

- 0 pts Correct

- 3 pts Incorrect

Question 18

2.6

3 / 3 pts

- 0 pts Correct

- 3 pts Incorrect

Question 19

2.7

3 / 3 pts

- 0 pts Correct

- 3 pts Incorrect

Question 20

3.1

0 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 21

3.2

0 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 22

3.3

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 23

3.4

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 24

3.5

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 25

4.1

0 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 26

4.2

0 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 27

4.3

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 28

4.4

2 / 2 pts

- 0 pts Correct

- 2 pts Incorrect

Question 29

5.1

3 / 3 pts

- 0 pts Correct

- 3 pts Incorrect

Question 30**5.2****5 / 5 pts****✓ - 0 pts** Correct**- 5 pts** Incorrect**Question 31****5.3****0 / 3 pts****- 0 pts** Correct**✓ - 3 pts** Incorrect**Question 32****5.4****2 / 2 pts****✓ - 0 pts** Correct**- 2 pts** Click here to replace this description.**Question 33****5.5****1 / 1 pt****✓ - 0 pts** Correct**- 1 pt** Click here to replace this description.**Question 34****5.6****1 / 1 pt****✓ - 0 pts** Correct**- 1 pt** Click here to replace this description.**Question 35****5.7****1 / 1 pt****✓ - 0 pts** Correct**- 1 pt** Click here to replace this description.**Question 36****5.8****1 / 1 pt****✓ - 0 pts** Correct**- 1 pt** Click here to replace this description.**Question 37****6.1****3 / 3 pts****✓ - 0 pts** Correct**- 3 pts** Click here to replace this description.

Question 38**6.2****3 / 3 pts****✓ - 0 pts** Correct**- 3 pts** Click here to replace this description.**Question 39****6.3****3 / 3 pts****✓ - 0 pts** Correct**- 3 pts** Click here to replace this description.**Question 40****6.4****2 / 4 pts****- 0 pts** Correct**✓ - 2 pts** 1 Correct + 1 Incorrect**- 4 pts** Not correct answer/more than 2 incorrect answers**Question 41****6.5****2 / 2 pts****✓ - 0 pts** Correct**- 2 pts** Click here to replace this description.**Question 42****7.1****2 / 2 pts****✓ - 0 pts** Correct**- 2 pts** Click here to replace this description.**Question 43****7.2****2 / 2 pts****✓ - 0 pts** Correct**- 2 pts** Click here to replace this description.**Question 44****7.3****2 / 2 pts****✓ - 0 pts** Correct**- 2 pts** Click here to replace this description.**- 0 pts** Click here to replace this description.

Question 45

7.4

2 / 2 pts

 **- 0 pts** Correct

- 2 pts Click here to replace this description.

Question 46

7.5

2 / 2 pts

 **- 0 pts** Correct

- 2 pts Click here to replace this description.

Write your name and UID:

Tiya Chokhani

305933966

Note 1: If you find a question difficult, move on with the rest of the questions and come back to it in the end!

Note 2: Your final grade will be curved, if necessary.

Note 3: There are 7 questions. Only answers written in the boxes will be graded. If you need extra space, please ask for an extra sheet. Good Luck! :-)

1 Linear Regression & Bias: True or False (22 Points)

Write 'T' or 'F' in the box corresponding to each of the statements below.

1. (2 points) F In a linear regression model, a large value for β and p -value < 0.05 shows that the corresponding predictor causes the response. F
2. (2 points) T A confounding factor causes a high correlation between multiple other predictors and the response.
3. (2 points) T If R^2 is large and the residual plot has a normal distribution centered at zero, the model complexity is enough for the given data.
4. (2 points) F A relatively small value for R^2 means that increasing the model complexity improves the performance. F
5. (2 points) F If the 95% confidence interval for a predictor is $[-0.2, 100]$, there is a significant relationship between the predictor and the response.
6. (2 points) T Multicollinearity doesn't affect the performance of the model, but makes interpretation of coefficients unreliable.
7. (2 points) T Stratified sampling can reduce bias of data collection, when the population is imbalanced.
8. (2 points) T Any type of regularization with appropriate regularization coefficient yields a higher training error but a lower validation error.
9. (2 points) T Lasso (L1 regularization) makes some coefficients exactly equal to zero and can make interpretation easier.
10. (2 points) T Larger dataset yields a smaller confidence interval for the predictors, compared to smaller dataset coming from the same distribution.
11. (2 points) T Mini-batch stochastic gradient descent with appropriate batch size can train a neural network to a better performance compared to gradient descent.

2 Logistic Regression (18 Points)

Fig. 1 shows the log odds of success, based on X_1 (real-valued) and X_2 (categorical with 3 categories $\{C_1, C_2, C_3\}$). We use binary variable(s) to model X_2 and model C_3 by making all binary variable(s) equal to 0.

$$w = \beta_0 + \beta_1 x_1 + \beta_2 x_{21} + \beta_3 x_{22}$$

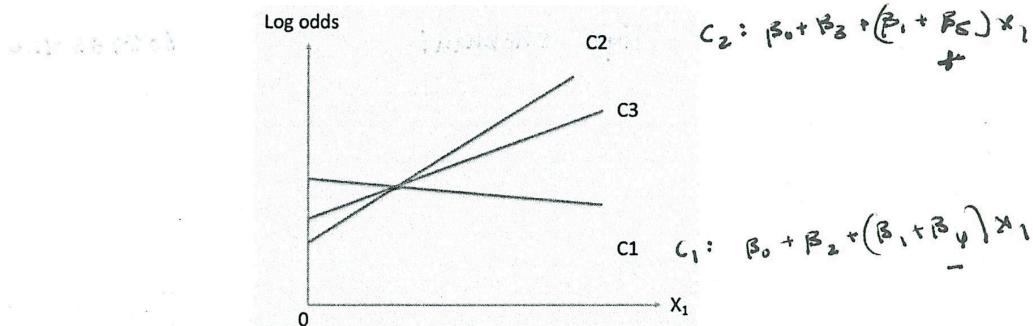


Figure 1: Log odds of success

Write the logistic model corresponding to Fig. 1. This part has 0 credit, but we won't grade the subsequent parts of this question if your logistic model is missing.

$$\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2^{c_1} + \beta_3 X_2^{c_2} + \beta_4 X_1 X_2^{c_1} + \beta_5 X_1 X_2^{c_2}$$

when $c_3 X_2^{c_1} = 0 = X_2^{c_2}$ so $\ln\left(\frac{P(Y=1)}{P(Y=0)}\right) = \beta_0 + \beta_1$

Answer the following questions based on Fig. 1 and your logistic model. For parts (d), (e), (f), choose one of the options (from 1 to 10) at the end of the question.

(a) (2 points) How many binary variables do you use for X_2 to have an interpretable logistic regression model? 2

(b) (2 points) How many interaction terms are used in the logistic regression model corresponding to Fig. 1? 2

(c) (2 points) Write the sign of the coefficient β_i of the interaction term(s). If you have more than 1 interaction term, write the sign for all (ordering doesn't matter). $\beta_4: - \beta_5: +$

(d) (3 points) What's the interpretation of β_i , where β_i is the intercept of the logistic model? 3

(e) (3 points) What's the interpretation of β_i , where β_i is the coefficient of X_1 ? 1

(f) (3 points) What's the interpretation of the β_i , where β_i is coefficient of a binary variable used for modeling X_2 ? 7

(g) (3 points) What's the interpretation of the β_i , where β_i is coefficient of an interaction term between X_1 and a binary variable used for modeling X_2 ? 6

✓ 1. One unit increase in X_1 changes the odds of success for C3 by $(e^{\beta_i} - 1)\%$.

2. One unit increase in X_1 changes the odds of success for C3 by β_i .

3. e^{β_i} is the odds of success at $X_1 = 0$ for C3. ✓

4. e^{β_i} is the odds of success for C3.
5. β_i is the odds of success for C3.
6. For one unit increase in X_1 , e^{β_i} is the ratio of the multiplicative change in odds of success for C1 or C2 over that of C3.
7. At $X_1 = 0$, e^{β_i} is the ratio of the multiplicative change in odds of success for C1 or C2 over that of C3.
8. At $X_1 = 0$, β_i is the additional change in odds of success for C1 or C2 over that of C3.
9. Odds ratio of success for C1 or C2 over that of C3 at $X_1 = 0$.
10. Odds ratio of success for C1 or C2 over that of C3.

3 Classification Metrics (10 Points)

For a trained binary classifier, we sort examples based on their probability of having label 1. Write 'T' or 'F' in the box corresponding to each of the statements below.

Hint: AUC stands for Area Under the RoC Curve.

$\frac{T \quad F}{V \rightarrow N}$

1. (2 points) T If the classifier sorts the examples correctly, the model has a good accuracy.
2. (2 points) F If the classifier sorts the examples correctly, the model has a good AUC.
3. (2 points) F If the classifier sorts the examples correctly, the model has a good accuracy and a good AUC.
4. (2 points) T If the model has a good AUC but a poor accuracy for class 1, changing the threshold for predicting class 1 improves the accuracy.
5. (2 points) F If the model has a good accuracy for class 1 but a poor AUC, changing the threshold for predicting class 1 improves the accuracy.

4 Multi-class Logistic Regression (8 Points)

- (a) (4 points) Fig. 2 shows the binary classifiers used for multinomial and One-vs-Rest (OvR) Logistic Regressions. Remember that for one of the logistic models, one of the classifiers is derived from the other two. Which figure corresponds to multinomial and OvR models? Write (a) or (b) in the boxes below:

1. Multinomial b

2. OVR a

- (b) (4 points) We want to do greedy (forward) model selection (greedy step-wise variable selection) for *each of the linear classifiers* in a multinomial or OvR logistic regression. If we have c classes, and p predictors, how many logistic regression models should we fit in total to choose $q \leq p$ predictors for each model? Assume, we only consider the original predictors in the model selection without including

$$c > p(p-1)(p-2) \cdots (q)$$

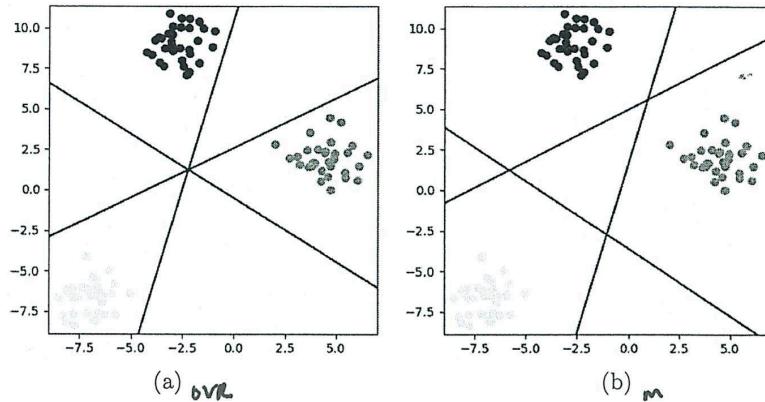


Figure 2: Multiclass Logistic Regression

any higher order terms. In each of the boxes below, choose one option from (a) to (f) at the end of the question.

Multinomial: F

$$p = 10$$

$$q = 6$$

OvR: E

(A) $c \times [p! + 1]$

(B) $(c-1) \times [p! + 1]$

(C) $c \times [\frac{p!}{q! \times (p-q)!} + 1]$

(D) $(c-1) \times [\frac{p!}{q! \times (p-q)!} + 1]$

(E) $c \times [(p+p-1+\dots+p-q+1)+1]$

(F) $(c-1) \times [(p+p-1+\dots+p-q+1)+1]$

$$P \times P-1 \dots \frac{P-P}{P-P} \times P-2 \times P-3 = 2!$$

$$cP \rightarrow$$

$$cP + c(p-1) \dots + c(p-q) + \dots + c(c)$$

$$c(10) + c(9) + c(8) + \dots + c(3) + c(2) + c(1)$$

$$c \times 10 \times 9 \times 8 \times 7 \times 6 \Rightarrow$$

$$\frac{10!}{6! 4!} + 1 \Rightarrow \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2 \times 1}$$

5 Neural Networks (17 Points)

Consider the following neural network with sigmoid activation functions in hidden nodes h_1, h_2 , and sigmoid output p .

Remember: $\frac{\partial \log z}{\partial z} = \frac{1}{z}$, and for a sigmoid function p , we have $\frac{\partial p(z)}{\partial z} = p(1-p)$.

$$a_1 = x_1 w_1 + x_2 w_3$$

$$a_2 = x_1 w_2 + x_2 w_4$$

$$h_1 = \sigma(a_1)$$

$$h_2 = \sigma(a_2)$$

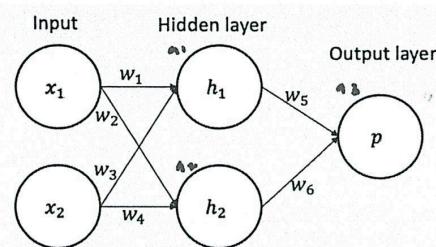


Figure 3: Neural Network

- (a) (3 points) Consider a mean squared loss, i.e. $\mathcal{L} = \frac{1}{2} \sum_i (p_i - y_i)^2$. Use the chain rule to write the derivative of the loss w.r.t. w_6 .

- (A) $\frac{\partial \mathcal{L}}{\partial w_6} = (p - y)h_2$
- (B) $\frac{\partial \mathcal{L}}{\partial w_6} = (p - y)p(1 - p)h_2$
- (C) $\frac{\partial \mathcal{L}}{\partial w_6} = 2(p - y)ph_2$
- (D) $\frac{\partial \mathcal{L}}{\partial w_6} = 2(p - y)(1 - p)h_2$
- (E) $\frac{\partial \mathcal{L}}{\partial w_6} = (p - y)p(1 - p)$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial p} &= (p - y) \\ \frac{\partial p}{\partial a_3} &= p(1 - p) \\ \frac{\partial a_3}{\partial w_6} &= h_2 \\ \frac{\partial \mathcal{L}}{\partial w_6} &= (p - y)(p(1 - p))h_2\end{aligned}$$

Answer: B

- (b) (5 points) Consider a binary cross entropy loss: $\mathcal{L} = \sum_i y_i \log p_i + (1 - y_i) \log(1 - p_i)$. Use the chain rule to write the derivative of the loss w.r.t. w_2 .

- (A) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} + \frac{1-y}{1-p})p(1-p)w_6h_2(1-h_2)x_1$
- (B) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)h_2(1-h_2)x_1$
- (C) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)w_6x_1$
- (D) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)w_6h_2(1-h_2)x_1$
- (E) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)w_6h_2x_1$
- (F) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)h_2(1-h_2)x_1$
- (G) $\frac{\partial \mathcal{L}}{\partial w_2} = (\frac{y}{p} - \frac{1-y}{1-p})p(1-p)h_2(1-h_2)$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial p} &= \frac{y}{p} - \frac{1-y}{1-p} & \frac{\partial \mathcal{L}}{\partial w_2} &= \left(\frac{y}{p} - \frac{1-y}{1-p}\right)p(1-p)w_6 \\ \frac{\partial p}{\partial a_3} &= p(1-p) & \frac{\partial a_3}{\partial w_2} &= w_6 \\ \frac{\partial h_2}{\partial a_2} &= & \frac{\partial h_2}{\partial a_2} &= h_2(1-h_2) \\ \frac{\partial a_2}{\partial w_2} &= x_1 & \frac{\partial a_2}{\partial w_2} &= x_1\end{aligned}$$

Answer: D

- (c) (3 points) Write one step of *mini-batch stochastic gradient descent* update for w_2 , using η as learning rate. Assume there are k examples in each mini-batch.

- (A) For $i \in \{1, \dots, k\}$ do $w_2 = w_2 - \eta \frac{\partial \mathcal{L}_i}{\partial w_2}$
- (B) $w_2 = w_2 - \frac{\eta}{k} (\frac{\partial \mathcal{L}_1}{\partial w_2} + \frac{\partial \mathcal{L}_2}{\partial w_2} + \dots + \frac{\partial \mathcal{L}_k}{\partial w_2})$
- (C) For $i \in \{1, \dots, k\}$ do $w_2 = w_2 - \frac{\eta}{k} (\frac{\partial \mathcal{L}_1}{\partial w_2} + \frac{\partial \mathcal{L}_2}{\partial w_2} + \dots + \frac{\partial \mathcal{L}_k}{\partial w_2})$
- (D) For $i \in \{1, \dots, k\}$ do $w_2 = w_2 - \frac{\eta}{n} (\frac{\partial \mathcal{L}_1}{\partial w_2} + \frac{\partial \mathcal{L}_2}{\partial w_2} + \dots + \frac{\partial \mathcal{L}_n}{\partial w_2})$

Answer: A

- (d) (2 points) Assuming the data has a pattern that is learnable by the model, after one *gradient descent* iteration with appropriate learning rate, how do the model predictions change for majority of examples?

- (A) Gets close to their actual labels
- (B) Gets farther away from their actual labels
- (C) Both (a) or (b) may happen

Answer: A

| | | | | |
|-------|----|----|----|----|
| x_1 | 0 | 0 | 1 | 1 |
| x_2 | 0 | 1 | 0 | 1 |
| class | +1 | -1 | -1 | +1 |



- (e) (4 points) Consider the following dataset: For each of the following activations, indicate if the model can classify the dataset correctly (Y) or not (N). In each box, write Y or N.

- (a) ReLU
- (b) Sigmoid
- (c) Tanh
- (d) linear

6 Black-box Interpretability (15)

- (a) Shapley Value (9 points) You want to use Shapley values to explain why a particular image is classified as a 'label', by a neural network. To do so, we first divide the image to 3 super-pixels. The following table shows the probability for the image to be in class 'label' by the neural network, when we input different combinations of super-pixels (Super-pixels kept) to the network and replace the other ones with gray color.

Remember: multinomial coefficient can be calculated as $\binom{n}{m} = n! / m!(n-m)!$. Note that this is not the exact formula for computation of the Shapley value.

| Super-pixels kept | Probability of 'label' |
|-------------------|------------------------|
| {} | 0.0 |
| {1} | 0.4 |
| {2} | 0.5 |
| {3} | 0.4 |
| {1,2} | 0.8 |
| {1,3} | 0.5 |
| {2,3} | 0.8 |
| {1,2,3} | 0.9 |

- (a) (3 points) What's the contribution of super-pixel 2 to prediction of 'label' conditioned on (when added to) super-pixels {1, 3}?

- (A) 0.4/3
- (B) 0.4/6
- (C) 0.3/6
- (D) 0.4/3 + 0.4/6 + 0.3/6

Answer:

$$\frac{2! \cdot 0!}{3!} \times 0.4 \Rightarrow \frac{0.4}{3}$$

- (b) (3 points) What's the contribution of super-pixel 2 to prediction of 'label' conditioned on (when added to) super-pixel 1?

- (A) 0.4/3
- (B) 0.4/6
- (C) 0.3/6

$$\frac{1! \cdot 1!}{3!} \times (0.4) = \frac{0.4}{6}$$

- (D) $0.4/3 + 0.4/6 + 0.3/6$

Answer: B

- (c) (3 points) What's the contribution of super-pixel 2 to prediction of 'label'?

- (A) $0.4/3 + 0.5/6$
- (B) $0.4/6 + 0.5/3$
- (C) $0.4/6 + 0.4/6 + 0.4/6 + 0.5/6$
- (D) $0.4/3 + 0.4/3 + 0.4/3 + 0.5/3$
- (E) $0.4/6 + 0.4/3 + 0.4/3 + 0.5/6$
- ~~(F)~~ (F) $0.4/3 + 0.4/6 + 0.4/6 + 0.5/3$

Answer: F

$$\begin{aligned} \text{when 0 kept: } & \frac{6!}{2!} (0.5) = \frac{0.5}{3} \\ \frac{0.8}{6} & \quad \text{1 kept: } \frac{1!}{1!} (0.4+0.4) = \frac{0.8}{6} \frac{0.4}{3} \\ & \quad \text{2 kept: } \frac{2!}{0!} (0.4) = \frac{0.4}{3} \\ & = \frac{0.5 + 0.4 + 0.4}{3} \end{aligned}$$

- (b) LIME (6 points) Next, we consider using LIME for the above problem. After model selection, we find the following liner model, where X_i models super-pixel i :

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1 X_2,$$

with the following coefficients:

| | β_0 | β_1 | β_2 | β_3 | β_4 |
|---------|-----------|-----------|-----------|-----------|-----------|
| value | 0.1 | 5 | 0.1 | 0.1 | 5 |
| p-value | 0.05 | 0.1 | 0.1 | 0.1 | 0.01 |

$P < 0.05$

- (4 points) Which of the following statements are correct? You can choose multiple correct answers, but we will deduct 2 points for each incorrect choice.

- (A) super-pixel 1 contributes to prediction of 'label' only for this image
- ~~(B)~~ (B) super-pixel 1 contributes to prediction of 'label' for all images of this class
- (C) super-pixel 2 contributes to prediction of 'label' only for this image
- ~~(D)~~ (D) super-pixel 2 contributes to prediction of 'label' for all images of this class
- ~~(E)~~ (E) super-pixel 1 or 2 alone do not contributes to prediction of 'label' for this image, but the combination does
- ~~(F)~~ (F) super-pixel 1 and 2 together contributes to prediction of 'label' for all images of this class

Answer: F, E

- (c) (2 points) When do we use LIME and SHAP?

- (A) When the test performance is low and we want to know why
- (B) When the test performance is high and we want to make sure model behavior is reasonable
- ~~(C)~~ (C) Both (a) and (b) are correct

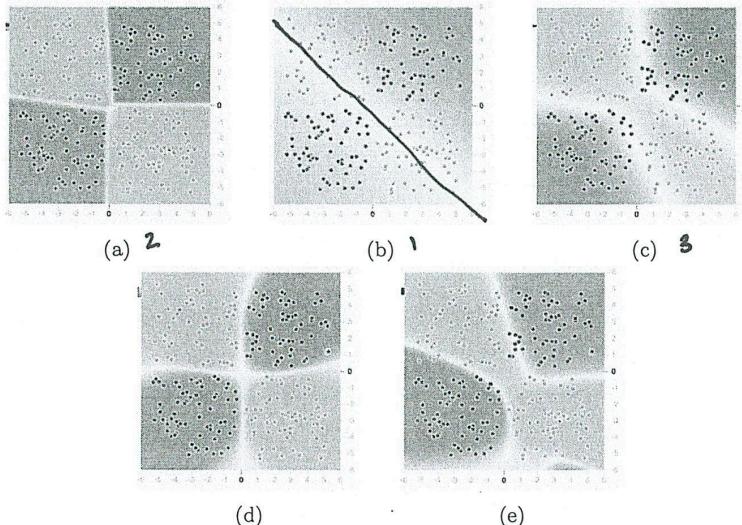
Answer: C

7 Decision Boundary (10 Points)

Consider the dataset with two classes in the figure below. In each of the plots (a)-(e), one of the following classification methods has been used, and the resulting decision boundary is shown:

- (1) Multi-layer Neural Network with linear activation functions **b**
- (2) Multi-layer ReLU network **a**
- (3) Regularized Multi-layer ReLU network **c**
- (4) Multi-layer tanh network **d**
- (5) Regularized Multi-layer tanh network **e**

Assign each of the previous methods to exactly one of the following plots (in a one to one correspondence), and explain briefly why you made each assignment. Write your choices (a, b, c ...) in the answer boxes below and your explanation in the corresponding explanation box. There is no partial credit for each part.



| | |
|--------------|---|
| (1) b | A multi layer network with linear activation functions no matter Explanation: how many hidden layers will produce a linear decision boundary as seen in (b) |
| (2) a | ReLU activation always results in a piecewise linear Explanation: decision boundaries & (a) shows this along with not having coeff tend to 0 due to regularization |
| (3) c | bc of RELU activation we see linear piecewise boundary but Explanation: also due to the regularization the boundary isn't as jagged & has fewer corners compared to (a) we also see the coeff tend more to 0 |

| | |
|-------|---|
| (4) d | tanh activation results in curvy lines in the decision boundary as seen in (d) but we can also tell that it hasn't been regularised as it's more squiggly than (e) & doesn't tend to 0. |
| (5) e | Due to the tanh activation we see curved lines in the boundary but also as compared to d it's smoother less wavy & tend coeff to 0 to give a more simpler boundary . |

