a. Please compute the Residual Sum of Squares (RSS) for the given data set where $RSS = \sum_{i=1}^{N} (\hat{f}(x_i) - y_i)^2. \text{ Please show the steps involved. [Note that MSE = RSS/N where N is the number of data points.]}$

$$RSS = \sum_{k=1}^{6} (\hat{f}(x_1) - y_1)^2 = (1.01 - 0.9)^2 + (2.04 - 2.2)^2 + (3.04 - 3)^2 + (4.04)^2 + (5.25 - 4.8)^2 + (6.36 - 6)^2 + (6.36 - 6)^2 = 0.004 + 0.026 + 0.01 + 0.004 + 0.202 + 0.129$$

$$= 0.383$$

b. Your goal is to minimize RSS. If you know the ground truth that values in decimal are due to noise component - please design a regularization function to reduce the effect of \(\chi^2\) i.e., present the new RSS function you plan to minimize. Please explain in detail how this would achieve the desired effect.

(BAND) > (INT)

Q1. Tick all correct answers, Cross all wrong answers. Negative marks of 25% will be awarded for each wrongly marks of 25% will be awarded for each wrongly marked answer.

- If Ram took the bus (B) or drove in his own car (D), then he arrived late (L) and missed the first session (M).
 - (a) ¬Bv¬DvLvM
 - (b) ¬BA¬DVLVM
 - $\mathcal{L}(e) (\mathsf{B} \mathsf{V} \mathsf{D}) \to (\mathsf{L} \mathsf{A} \mathsf{M})$

(d) BvDvLvM

(K=King castles) and the pawn advances (P) then either the bishop is blocked or the rook is pinned (R).

$$(B)'(\neg K \land P) \rightarrow (B \lor R)$$

- (b) KVPVBVR
- (c) ¬Kv¬PvBvR
- (d) Kv-PvBvR

Q2. Prove the resolution rule $(P \lor Q, \neg Q \lor R \Rightarrow P \lor R)$ using truth-tables.

[4]

[10]

P	0	R	BNB	-BUR	PURQ	PUBNGAUR) -> (KAK)
7	T	T	T	下	17	T	
+	Ť	F	T	F	T	# 1	
T	F	T	T	T		1	
T	E	F	T	T	-	<u>'</u>	
	T	T	1 7	F	E	T6	
-	十	F		7	1	1	
r	C	1-	F		1 6	T	
-	1	P	F				
F	1	1					

Q3. Prove using resolution: $P \rightarrow \neg Q$, $\neg Q \rightarrow R \implies P \rightarrow R$. (Represent a clause with its number.) [10]

Clauses being resolved	Substitution required	New resulting clause	(of new clause)
		Physical Company	

Two popular general regularization functions in literature are Ridge and Lasso regularization. Ridge regularization also called an L2 penalty, is going to square your coefficients. Lasso regularization or an L1 penalty, is going to take the absolute value of your coefficients.

RSS with Ridge regularization: RSS = $\sum_{i=1}^{n} (\hat{f}(x_i) - y_i)^2 + \alpha \sum_{j=1}^{n} (\theta_j)^2$

RSS with Lasso regularization: RSS = $\sum_{i=1}^{n} (\hat{f}(x_i) - y_i)^2 + \alpha \sum_{j=1}^{n} (\theta_j)$

Assuming $\alpha = 1$, which of the two regularization functions would be better for your $\hat{f}(x)$? Please explain your choice with mathematical reasoning.

Q4. Suppose a training set consists of points X_1, X_2, \dots, X_n and real values Y_i associated with each point X_i . We assume there is a function with noise $y = f(x) + \varepsilon$, where the noise ε has zero mean and variance σ^2 . Please provide all steps of derivation for

 $E\left[\left(y-\hat{f}(x)\right)^{2}\right]=\left(Bias\left[\hat{f}(x)\right]\right)^{2}+Var\left[\hat{f}(x)\right]+\sigma^{2}$

where $\hat{f}(x)$ is the best approximation for f(x) identified by the machine learning algorithm.

[8]

$$\begin{split} E \left[(g - f(\alpha))^2 \right] &= E \left[(f + \varepsilon - \hat{f})^2 \right] \\ &= E \left[(f + \varepsilon - \hat{f} + E(\hat{f}) - E(\hat{f}))^2 \right] \\ &= E \left[(f - E(\hat{f})) + \mathbf{G} + (E(\hat{f}) - \hat{f})^2 \right] \\ &= E \left[(f - E(\hat{f})) + E(e^2) \right] + E \left(E(\hat{f}) - \hat{f} \right) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] + 2 E \left((f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] - E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right] \\ &+ 2 E \left[(f - E(\hat{f})) \right] - 2 E \left[(f - E(\hat{f})) \right]$$

Q6. Please answer the following questions:

What is the AlphaGo system? Please describe idea behind the system in 2-3 sentences? You should mention atleast three specific ideas or techniques used as part of this system.

Alpha Goo is a Madrine Learning system trained on the Monte Carolo tree method to play chess. Within 4 hours of taining It bear the reigning chess robot champion.

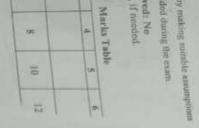
. What does dimensionality reduction mean ? Please explain how it can be used to reduce overfitting?

Dimensionality seduction deals with reducing the features used to train the model. This is done because disregarding features not relevant to the seasons for the aning the model can better help reseate a generalised model which doesn't memosize the data training set.

What is feature engineering? Please explain how binning can be used to perform feature

Feature engineering is the process of creating new features to train your model based or poe-existing ones.

What is the goal of Power transformation? Please present a diagram with some description that [1.5+1.5] conveys the essence of your explanation?



Transformations are a way of making your features more easily possable as an import to your model.

Power transformations involve modifying the power of your features to addieve the same.

Q7. Given the following frequent itemsets what candidates will Apriori compute for the next database scan? Show your steps.

(i) AB, AC, AD, BC, BD, CD, AE

[6]

i) We first create a set of sets of three by matching the given sets so that they have at most one substem defferent. Then, we have:

[ABC, ABD, ABE, ACD, ACE, ADG, BCD]

11.) From this set, we eliminate itemsets who do not have sobsets in the given set. So, we have,

[ABC, ABD, ACD, BCD]
This is own final answers.

(ii) ABC, ABD, ACD, BCD, BCE, CDE

[6]

9) Similar to above, the first step produces the following:

[ABCD, ABCE, ACDE, BCDE]

11) Now, eliminating results in our final answer:

[ABCD]