

Exam Answer Sheet

A.D. van der Merwe
Department of Computer Science
University of Stellenbosch
24923273
24923273@sun.ac.za

October 29, 2024

Section A: Introduction to Machine Learning

Question 1

Ockham's razor is a principle that states if there are multiple ways to explain a phenomenon, the simplest explanation that still effectively explains the phenomenon should be used.

Regarding predictive models, Ockham's razor is crucial in model selection and development. A simple model that is less computationally expensive or a model with fewer parameters and assumptions should be favoured if the prediction accuracy is sufficient enough for the task at hand.

For example, if a classification decision tree model achieves similar results to a more complex random forest classification model, the simple classification tree is preferred to perform predictions for the specific task. Simpler models are also easier to interpret, are less prone to overfitting on the training data, and are much less computationally expensive.

Question 2

Variance in a model refers to the extent to which a prediction varies when different test sets are used to construct the model [1]. It is almost impossible to find real-world data where there is no variability in the dataset. Therefore there will always be variability in each training dataset and the prediction of a model will almost always result in different values. Ideally should the prediction of an instance not vary too much between training sets. High variance indicates that the model may be overfitting, meaning it captures noise and specific patterns from the training data and generalises poorly to unseen data. This sensitivity to small changes in the training data can result in significant fluctuations in predictions. To calculate the variance of a regression model, the following equation is used.

$$Var(\hat{y}) = \mathbb{E} \left[(\hat{y} - \mathbb{E}[\hat{y}])^2 \right] \quad (1)$$

where \hat{y} is the predicted value of the regression model and $\mathbb{E}[\hat{y}]$ is the expected predicted value across different training sets.

Bias in a model refers to the error introduced by approximating a real-world problem, which may be extremely complicated, by a much simpler model. A model with high bias can lead to inaccurate predictions by underfitting the data as a simple model might not be able to capture the complex relationships between the descriptive features and the target feature. A high bias indicates that the model's predictions deviate significantly from the actual values, often due to an overly simplistic model that cannot capture the underlying complexities of the data. To calculate the bias of a regression model, the following equation is used.

$$Bias(\hat{y})^2 = (\mathbb{E}[\hat{y}] - y)^2 \quad (2)$$

where y is the actual value that the model aims to approximate.

The expected test mean squared error (MSE) can be decomposed into the sum of three fundamental quantities which is given in the following equation.

$$\mathbb{E} (y - \hat{y})^2 = Var(\hat{y}) + Bias(\hat{y})^2 + \sigma^2 \quad (3)$$

where σ^2 is the irreducible error, which is the noise in the observations that cannot be reduced by any model.

From equation (3) it is clear that both the variance and the bias of a model has to be small values to minimise the test MSE. The bias-variance dilemma arises from this. The more complex a model becomes to capture relationships between the descriptive features and the target feature, the more susceptible the model becomes to overfitting on the data, therefore the reduction in bias may lead to an increase in variance. Conversely, if a model is kept simple to prevent overfitting, the model may struggle to capture complex relationships between the descriptive features and the target feature, therefore the reduction in variance may lead to an increase in bias.

Question 3

(a)

Underfitting occurs when a predictive model is too simplistic to capture the underlying relationships between the descriptive features and the target feature resulting in poor predictive performance for both the training and test datasets.

(b)

Overfitting occurs when the predictive model is overly complex, leading it to become sensitive to noise and fluctuations in the training dataset. This results in the model fitting too closely to the training data, where the complex underlying relationships between the descriptive features and target features are captured, but the model learns the noise present in the training dataset. This then results in a good predictive performance for the training set, but the generalisation to the unseen test dataset will be poor.

(c)

When a model is underfitting, it fails to learn underlying relationships between the descriptive features and the target feature, thus resulting in a high bias and a low variance. Conversely, when a model is overfitting, it accurately learns the underlying relationships between the descriptive features and the target feature but also captures noise within the dataset. This leads to the model becoming sensitive to small fluctuations in the training dataset, thus resulting in a high variance and a low bias.

Question 4

When a dataset contains outliers, it is generally a bad idea to impute missing numerical values by use of the mean as the mean is influenced by the outliers that leads to an inaccurate representation of the central tendency of the data. It is generally better to use the median of the numerical values of a feature when outliers are present in the feature values. The median is less affected by outliers and provides a better measure of central tendency of the feature values if these values are not normally distributed.

Additionally considering more advanced imputation methods that accounts for the relationships between the descriptive features can yield better results. This is especially the case for when the data structure is complex. A k-nearest neighbours (KNN) classifier technique can be used to identify the k-nearest instances to the observation with a missing value based on the other features similarity. The missing value is then imputed using the mean or median of the feature values from the k-nearest instances, which often results in a more robust and representative imputed value.

When a classification task is considered, it is generally better to impute the missing value of a feature with only the mean or median features related to the class which contains the missing value. This approach leverages the conditional relationships between the descriptive features and the target feature, which could lead to more accurate imputations.

Question 5

Feature selection is a crucial step when preprocessing data to construct a predictive machine learning. Reducing the number of features in the dataset will improve the performance of the model, reduce overfitting, result in faster training times, create a simpler model which is more interpretable and improve the data quality.

The model's performance improves, because only the most relevant descriptive features are chosen to perform predictions of the target feature, which often leads to improved accuracy and predictive power. Additionally, overfitting is reduced as the reduction in dimensionality of the original dataset leads to models that are less complex and reduces the chances of capturing noise in the dataset and increases the chance of the model learning the underlying relationships between the remaining descriptive features and the target feature. The dimensionality reduction also leads to faster training times as the model is constructed on less data, which is very beneficial when large datasets are used and the models are easier to interpret, which is beneficial when a model is used where the decision-making process is important to understand. Lastly, the quality of the data is improved as features with noise or irrelevant features are removed.

Question 6

(a)

Stationary data is data that has statistical properties which stays consistent over time, such as the mean and variance. This means that the behavior of the dataset does not significantly change over time and stays more or less constant. Conversely non-stationary data has statistical properties which changes over time. Non-stationary data may show increasing or decreasing means and changing variances, due to trends or even seasonal effects, making it less predictable.

(b)

A real world problem that exhibits non-stationary data is stock market data. This is because stock market data exhibit upward and downward trends over time, which is influenced by different factors such as economic growth and some stocks may exhibit seasonal patterns, such as retail stocks that perform better during the holiday.

(c)

The main implications are that the models can become unstable over time, the model needs to constantly learn and the models might overfit past historical patterns.

The first implications is that the model may become unstable overtime and lead to invalid predictions, as the relationships between the descriptive features and target feature may change over time. This makes previously learned patterns unreliable and justifies the need for the second implication, which is that the model needs to constantly be updated. The model has to incorporate some form of continuous learning or regular retraining to ensure that the models remain relevant and can adapt as soon as the underlying patterns between the descriptive features and the target feature changes over time. The models may also overfit to past historical, which leads to poor generalisation. This may lead to high error rates when the model is applied to current or future data.

Question 7

When binning is used to convert a numeric feature into a categorical feature, what are the consequences of too few bins versus too many bins?

The consequences of using too few bins is that important information may be lost with respect to the distribution of values in the original continuous feature. This loss of information may lead to reduced predictive power of the model, as the small number of bins reduces the feature's ability to differentiate between classes or predict outcomes accurately. A small number of bins can also potentially hide variations or patterns that could be useful for modeling.

Conversely, the consequences of using too many bins, is that each bin will have a small number of instances contained within them, where some bins may end up having no instances. Having too many bins may lead to a more computationally expensive model. Additionally, too many bins may increase the complexity of the model, such as a classification tree that has to evaluate a greater number of categories.

Question 8

(a)

TODO

(b)

The distribution of ID (d) indicates a feature that should be removed, as the feature contains unique values. This feature has no predictive information relevant to the target feature and therefore it should be removed.

(c)

TODO

Section B: Information-based Learning

Question 1

e

Question 2

(a)

TODO

(b)

TODO

Question 3

TODO

Question 4

The first inductive bias is that the decision boundaries can only be parallel to the axis. The consequence of this is that more complex trees are needed if the true decision boundaries are not parallel to the axes. This inductive bias can be addressed by using oblique trees.

The second inductive bias is that the tree is induced to overfit the training data. Meaning that the tree creates decision boundaries until all of the classes are separated and achieves a perfect score on the training data. The consequence of this inductive bias is that the tree generalises poorly to unseen data. This inductive bias can be addressed by either pruning the tree while it is being induced or by pruning the induced tree.

The third inductive bias is that if the information gain is maximised, descriptive features with many outcomes are favoured higher up in the tree. This inductive bias results into an induced tree with many decision rules and trees with high branching factors high up in the tree. This inductive bias can be addressed by using the gain ratio criterion.

TODO

Question 5

An oblique tree algorithm would result in the smallest tree structure for this classification problem, as the classification problem represents the classification trees, such as C4.5 and ID3's, inductive bias of only creating axis-parallel decision boundaries.

Oblique trees, by contrast, allow splits that are not parallel to any axis, enabling the model to create boundaries at arbitrary angles. Fewer nodes are then required to split fully induce the tree to overfit on this particular classification task and therefore the tree resulting from the oblique tree algorithm would result in the smallest tree.

Question 6

(a)

TODO

(b)

TODO

Question 7

TODO

Question 8

TODO

Section C: Similarity-based Learning

Question 1

b

Question 2

b, c

Question 3

b, c, g TODO a???

Question 4

The inductive bias of the k -nearest neighbour algorithm is instances that have similar descriptive feature values also have the same target feature values. Meaning that if two instances are close in terms of their descriptive feature values, the k -nearest neighbours algorithm assumes that these instances belongs to the same class for classification tasks or similar outputs for regression tasks.

Question 5

For algorithms like k -nearest neighbours that makes use of similarity based learning it is important to normalise the input descriptive features so that each descriptive feature has exactly the same range. The features should be normalised to the same range as as input features with larger differences in the range has a stronger impact on the distance calculations used in the k -nearest neighbours. A descriptive feature with

a range of [1000,100000] will have a much stronger impact on the distance calculation than a feature with a range of [10,100]. Therefore these features should be normalised to a specific range such as [0,1].

Question 6

The k -nearest neighbours algorithm can be used to impute a missing value of a feature by calculating the similarity between instances. For an instance that contains a missing value in one of the descriptive features, the k -nearest neighbours algorithm can be used to examine the k most similar instances in the dataset, where the similarity is calculated based on the features that does not contain missing values. After the k most similar neighbours of this instance has been found, the missing value is imputed by use of the mean or median of the descriptive feature that contains the missing value of these k neighbours for a numerical descriptive feature. Conversely for a categorical descriptive feature, the mode of the descriptive feature that contains the missing value of the k neighbours are used to impute the missing value.

Question 7

The k -nearest neighbour algorithm can be applied to problems with categorical descriptive features. The categorical descriptive features should either be one-hot encoded or ordinally encoded to ensure the categories are represented in numerical form if they are not already. If the categorical features have been one-hot encoded, a variety of distance calculations can be used to calculate the similarity between observations such as the Jaccard, Hamming, Manhattan, Euclidean and cosine distance metrics. For ordinally encoded features, the similarity between observations can be calculated by use of the Jaccard similarity, which calculates the number of features with the same value.

Question 8

Discuss the consequences of different values for k when k -nearest neighbours is applied to regression problems.

If the number of neighbours, k , in the algorithm is small, the algorithm will become sensitive to noise in the target feature. The predicted value of an instance can be negatively influenced if k is small and there are a few neighbours around a particular instance with an abnormal target value due to noise or outliers. Therefore, the outcome may be skewed as there are fewer neighbours to average out the noise.

Conversely, if the number of neighbours, k , in the algorithm is large, the predictions of the algorithm will be bad as the average of the target features of the k -nearest neighbours is calculated and used as the predicted value of the instance. Therefore, a larger value of k could lead to an algorithm that underfits the data, failing to capture local patterns effectively.

Section D: Error-based Learning

Question 1

a

Question 2

d

Question 3

c

Question 4

a, b, g

Question 5

(a)

Yes, normalising or scaling of the target feature is required. This is because the output of the sigmoid activation function returns a value in the range $(0,1)$. In the case of a regression problem, the target features should be scaled to the range of $(0,1)$ to match the output of the sigmoid activation function. In the case of binary classification, the two target classes should be ordinally encoded to 0 and 1, which represent the probability of the observation belonging to class 1. If the target features are not scaled, the neuron will always produce a large error signal, which leads to a continuous adjustment of the weights, meaning the model would never converge.

(b)

Yes, it is prudent to normalise or scale the input features in this scenario. It is not necessary to scale the input feature values, but the performance of the model can be improved if the inputs are scaled to the active domain of the activation function. For the sigmoid activation function, the active domain is $[-\sqrt{3}, \sqrt{3}]$. This corresponds to the parts of the sigmoid function for which weight changes in the input features has a relatively large change in output. The values beyond these points would have a very small influence on the weight updates.

(c)

Large weights and biases used in the gradient descent optimisation algorithm can lead to premature convergence. This occurs because the large weights and biases move to the asymptotic ends of the sigmoid activation function too quickly, which leads to extreme output values with associated derivatives being close to zero, meaning the weight updates are also close to zero. Absolute values of the weights and biases is also a poor strategy as the active domain of the sigmoid activation function is $[-\sqrt{3}, \sqrt{3}]$. If the weights and biases are initialised as only large positive values, the activations of the sigmoid activation function will be biased towards the positive end of the sigmoid's active domain.

Question 6

The momentum term is essential to the stochastic gradient descent (SGD) optimisation algorithm, as it smooths out the search trajectories and prevents the oscillation of the search trajectories. The idea of the momentum term is to average out the weight changes, thereby ensuring that the search path is in the average downhill direction, meaning the search direction does not prematurely change between epochs, ensuring that oscillation between search directions does not occur. The momentum term is then simply the previous weight change weighted by a scalar value α , which is defined as any value between the range of $\alpha \in (0,1)$. For larger values of α , the more strict the optimisation algorithm becomes of staying on the current search path, meaning that more significant evidence is needed to jump over to the other side.

Question 7

This statement is false, because if the net input signal is calculated by use of product units, the single neuron can separate non-linearly separable classes. The mathematical equation used to calculate the product units is as follows.

$$net = \prod_{i=1}^I x_i^{w_i} \quad (4)$$

where net is the net input signal, I is the number of input signals, x_i is the i -th signal and w_i is the weight corresponding to the i -th input signal. Product units allow for higher-order combinations of inputs, having the advantage of increased information capacity. This increased information capacity and the product units allows the single neuron to separate non-linearly separable data.

Question 8

The main advantages of using the scaled conjugate gradient optimisation algorithm instead of the gradient descent optimisation algorithm is that the model converges faster due to the fast quadratic convergence of Newton's method. Additionally, the scaled conjugate gradient optimisation algorithm is less susceptible to the local minima as the algorithm restarts every n_w steps if there is no reduction in the test error, where n_w is the total number of weights and biases. Lastly, the scaled conjugate gradient optimisation algorithm performs automatic step size adjustment, eliminating the need for the learning rate used in the gradient descent optimisation algorithm, meaning there is less manual control parameters of which the optimal values should be found for.

Question 9

TODO

Question 10

TODO

Question 11

A linear logistic regression model can not separate the two classes. However a non-linear logistic regression model will be able to separate the two classes.

The non-linear logistic regression model can capture the underlying relationship between the two descriptive features x_1 and x_2 if they are transformed by use of sine and cosine as basis functions. Therefore the non-linear logistic regression model would become the following.

$$P(y = 1|\mathbf{x}_i) = \frac{1}{1 + e^{-(w_0 + w_1 \cdot \sin(x_{i,1}) + w_2 \cdot \cos(x_{i,1}) + w_3 \cdot \sin(x_{i,2}) + w_4 \cdot \cos(x_{i,2}))}} \quad (5)$$

where \mathbf{x}_i is a vector containing all of the descriptive features of the i -th observation, w_0 is the bias and $P(y = 1|\mathbf{x}_i)$ is the probability of the i -th observation belonging to class 1.

Therefore by utilising these basis functions to capture the relationship between the two descriptive features, will the non-linear logistic regression be able to separate these classes.

Question 12

Consider the classification problem depicted in the figure below. Explain in detail how logistic regression can be used to separate the classes.

Logistic regression can separate these classes, by making use of a multinomial logistic regression model, which is designed as an extension of the normal logistic regression model to handle multiclass classification problems.

The multinomial logistic regression model is constructed by making use of 5 one versus all logistic regression models for this scenario. These one-versus-all models distinguishes between one class of the target feature and all the other classes of the target feature. The 5 one-versus-all models would then be the following:

- Logistic regression model 1: classify class 0 against the other classes
- Logistic regression model 2: classify class 1 against the other classes
- Logistic regression model 3: classify class 2 against the other classes
- Logistic regression model 4: classify class 3 against the other classes
- Logistic regression model 5: classify class 4 against the other classes

This multinomial logistic regression model will then construct 4 decision boundaries combined from each models decision boundary to separated these classes into 5 segments, where each segment contains all the observations of one class.

Question 13

TODO

Question 14

TODO

Question 15

A neural network that deploys a sigmoid activation function as the activation function in the hidden layer, as opposed to using the linear activation function as the activation function in the hidden layer, will result in less hidden units for a highly non-linear mapping. This is because the linear activation function only perform linear transformations, which limits the model's capacity to capture non-linear patterns and more units in the hidden layer would have to be added for the model to capture these non-linear patterns. Conversely, the sigmoid activation function allows the hidden layer to capture non-linear patterns, which will result in less units that should be used in the hidden layer to achieve the same generalisation performance as the linear activation function deployed in the hidden layer with more hidden units.

Section E: Unsupervised Learning

Question 1

a, d

Question 2

TODO

Question 3

TODO

Question 4

TODO

Question 5

TODO

Question 6

TODO

Question 7

(a)

TODO

(b)

TODO

(c)

i TODO

ii TODO

iii TODO

Section F: Kernel-based Learning

Question 1

TODO

Question 2

TODO

Question 3

TODO

Question 4

TODO

Question 5

TODO

Question 6

TODO

Question 7

TODO

Question 8

TODO

Question 9

TODO

Question 10

TODO

Section G: Ensemble Learning

Question 1

TODO

Question 2

TODO

Question 3

TODO

Question 4

TODO

Question 5

TODO

Question 6

TODO

Question 7

TODO

Question 8

TODO

Question 9

TODO

Question 10

TODO

Section H: Reinforcement Learning

Question 1

TODO

Question 2

(a)

TODO

(b)

TODO

(c)

TODO

(d)

TODO

(e)

TODO

References

- [1] G. James, D. Witten, T. Hastie, R. Tibshirani *et al.*, *An introduction to statistical learning*. Springer, 2013, vol. 112.