

Name/Student Registration Number: _____

1. (10 points) Please select which statement is correct or wrong:

- | | | |
|--|-------------------------------|-----------------------------|
| MLR requires collinear variables | <input type="radio"/> correct | <input type="radio"/> wrong |
| ANOVA tests for equal variances | <input type="radio"/> correct | <input type="radio"/> wrong |
| Please assign the dendrograms A, B, C, and D to the datasets 1, 2, 3, and 4: | <input type="radio"/> correct | <input type="radio"/> wrong |
| Forward selection can be used to find outliers | <input type="radio"/> correct | <input type="radio"/> wrong |
| The VIF of a particular variable has nothing to do with PRESS | <input type="radio"/> correct | <input type="radio"/> wrong |
| The confusion matrix can be used to calculate the probabilities of type 1 and type two error | <input type="radio"/> correct | <input type="radio"/> wrong |
| Confounding cannot be reduced by randomisation of the observations | <input type="radio"/> correct | <input type="radio"/> wrong |
| The order of a model is equal to the number of non-zero eigenvalues | <input type="radio"/> correct | <input type="radio"/> wrong |
| The kMeans algorithm can be used to estimate the means of k neighbors | <input type="radio"/> correct | <input type="radio"/> wrong |
| PLS/DA stands for "PLS-based Data Analysis" | <input type="radio"/> correct | <input type="radio"/> wrong |
| The F value of an MLR model can be calculated from the coefficient of determination if the number of variables and objects are known | <input type="radio"/> correct | <input type="radio"/> wrong |
| MLR can be used to determine parameter b in the model $y = a \cdot \sin(b \cdot x)$ | <input type="radio"/> correct | <input type="radio"/> wrong |
| The model error can be obtained by stepwise regression | <input type="radio"/> correct | <input type="radio"/> wrong |
| The reliability of an MLR model can be checked by cross validation | <input type="radio"/> correct | <input type="radio"/> wrong |
| Blocking can be used to control confounding variables | <input type="radio"/> correct | <input type="radio"/> wrong |
| The ratio of the number of variables to the number of object influences overfitting in MLR models | <input type="radio"/> correct | <input type="radio"/> wrong |
| Majority voting is used in connection with dendrograms | <input type="radio"/> correct | <input type="radio"/> wrong |
| ANOVA can be used to detect effects of factors on the independent variable | <input type="radio"/> correct | <input type="radio"/> wrong |
| Random generators have to be used for randomizations | <input type="radio"/> correct | <input type="radio"/> wrong |
| LDA is based on multilinear regression | <input type="radio"/> correct | <input type="radio"/> wrong |
| Variables which exert an undesired influence on the dependent variable are confounding variables | <input type="radio"/> correct | <input type="radio"/> wrong |
| Treatments are conditions for a given experiment | <input type="radio"/> correct | <input type="radio"/> wrong |
| The order of a model can be recognized in a scree plot | <input type="radio"/> correct | <input type="radio"/> wrong |
| PLS can be used for problems with more variables than observations | <input type="radio"/> correct | <input type="radio"/> wrong |
| Fractional factorial designs reduce the number of required experiments | <input type="radio"/> correct | <input type="radio"/> wrong |
| The significance of an MLR coefficient can be calculated from the ratio of its standard error to its value | <input type="radio"/> correct | <input type="radio"/> wrong |
| Linear discriminant analysis maximizes the t-value between two groups | <input type="radio"/> correct | <input type="radio"/> wrong |
| ANOVA is used to check the overall reliability of an MLR model | <input type="radio"/> correct | <input type="radio"/> wrong |
| PRESS is used to detect multi-collinearity | <input type="radio"/> correct | <input type="radio"/> wrong |
| Balanced experiments have the same number of cases per treatment | <input type="radio"/> correct | <input type="radio"/> wrong |
| Multi-collinearity means that the multi-collinear variables can be related to each other by linear equations | <input type="radio"/> correct | <input type="radio"/> wrong |
| Overfitting occurs if there are more objects than variables | <input type="radio"/> correct | <input type="radio"/> wrong |
| The inverse ratio of a regression parameter to its standard error is t-distributed | <input type="radio"/> correct | <input type="radio"/> wrong |
| The general Minkowski distance includes the Mahalanobis distance | <input type="radio"/> correct | <input type="radio"/> wrong |
| PLS1 calculates a model having only one input variable | <input type="radio"/> correct | <input type="radio"/> wrong |
| Heteroscedasticity implies that the variance of repeat measurements depends on y-hat | <input type="radio"/> correct | <input type="radio"/> wrong |
| Homoscedasticity means equal means of repeat measurements | <input type="radio"/> correct | <input type="radio"/> wrong |
| PCR stands for "Principal Component Reduction" | <input type="radio"/> correct | <input type="radio"/> wrong |
| MLR applied to principal component scores results in PCR | <input type="radio"/> correct | <input type="radio"/> wrong |
| Majority voting requires an odd number of nearest neighbors to be unambiguous | <input type="radio"/> correct | <input type="radio"/> wrong |
| For the analysis of variances the variances of the compared factor levels have to be equal | <input type="radio"/> correct | <input type="radio"/> wrong |
| Variable selection can be used to fight the curse of dimensionality | <input type="radio"/> correct | <input type="radio"/> wrong |
| The Mahalanobis distance takes correlations into account | <input type="radio"/> correct | <input type="radio"/> wrong |
| PRESS is the square root of RMSEP | <input type="radio"/> correct | <input type="radio"/> wrong |
| The Lance-Williams equation controls the type of hierarchical clustering | <input type="radio"/> correct | <input type="radio"/> wrong |
| Confounding variables are unimportant in experimental designs | <input type="radio"/> correct | <input type="radio"/> wrong |
| A Latin square can be used to set up an experimental design | <input type="radio"/> correct | <input type="radio"/> wrong |
| Cluster analysis is based on distances in the p-dimensional space | <input type="radio"/> correct | <input type="radio"/> wrong |
| Stepwise regression always results in the best set of variables | <input type="radio"/> correct | <input type="radio"/> wrong |
| Fractional factorial designs require less variables than full experimental designs | <input type="radio"/> correct | <input type="radio"/> wrong |

2. (3 P) Explain the advantages and drawbacks of kNN-based models

3. (3 P) What are the main assumptions of multiple linear regression? How can you check each of these assumptions?

4. (2P) What is a confusion matrix? Draw an example and explain the particular cells of the matrix.

5. (2 P) Draw the dendrogram of the following data:

