

Lab1

| Category | Points | Description |
|------------------------------|--------|---|
| Business Understanding | 10 | Describe the purpose of the data set you selected (i.e., why was this data collected in the first place?). Describe how you would define and measure the outcomes from the dataset. That is, why is this data important and how do you know if you have mined useful knowledge from the dataset? How would you measure the effectiveness of a good prediction algorithm? Be specific. |
| Data Meaning Type | 10 | Describe the meaning and type of data (scale, values, etc.) for each attribute in the data file. |
| Data Quality | 15 | Verify data quality: Explain any missing values, duplicate data, and outliers. Are those mistakes? How do you deal with these problems? Give justifications for your methods. |
| Simple Statistics | 10 | Visualize appropriate statistics (e.g., range, mode, mean, median, variance, counts) for a subset of attributes. Describe anything meaningful you found from this or if you found something potentially interesting. Note: You can also use data from other sources for comparison. Explain why the statistics run are meaningful. |
| Visualize Attributes | 15 | Visualize the most interesting attributes (at least 5 attributes, your opinion on what is interesting). Important: Interpret the implications for each visualization. Explain for each attribute why the chosen visualization is appropriate. |
| Explore Joint Attributes | 15 | Visualize relationships between attributes: Look at the attributes via scatter plots, correlation, cross-tabulation, group-wise averages, etc. as appropriate. Explain any interesting relationships. |
| Explore Attributes and Class | 10 | Identify and explain interesting relationships between features and the class you are trying to predict (i.e., relationships with variables and the target classification). |
| New Features | 5 | Are there other features that could be added to the data or created from existing features? Which ones? |
| Exceptional Work | 10 | You have free reign to provide additional analyses. One idea: implement dimensionality reduction, then visualize and interpret the results. |

ML1- Ruberics

Mini Lab

| Category | Available | Requirements |
|------------------------------|-----------|---|
| Total Points | 100 | |
| Create Models | 50 | Create a logistic regression model and a support vector machine model for the classification task involved with your dataset. Assess how well each model performs (use 80/20 training/testing split for your data). Adjust parameters of the models to make them more accurate. If your dataset size requires the use of stochastic gradient descent, then linear kernel only is fine to use. That is, the SGDClassifier is fine to use for optimizing logistic regression and linear support vector machines. For many problems, SGD will be required in order to train the SVM model in a reasonable timeframe. |
| Model Advantages | 10 | Discuss the advantages of each model for each classification task. Does one type of model offer superior performance over another in terms of prediction accuracy? In terms of training time or efficiency? Explain in detail. |
| Interpret Feature Importance | 30 | Use the weights from logistic regression to interpret the importance of different features for the classification task. Explain your interpretation in detail. Why do you think some variables are more important? |
| Interpret Support Vectors | 10 | Look at the chosen support vectors for the classification task. Do these provide any insight into the data? Explain. If you used stochastic gradient descent (and therefore did not explicitly solve for support vectors), try subsampling your data to train the SVC model— then analyze the support vectors from the subsampled dataset. |

Lab2

| Category | Available | Requirements |
|---------------------------|------------|---|
| Total Points | 100 | |
| | | |
| Data Preparation Part 1 | 10 | Define and prepare your class variables. Use proper variable representations (int, float, one-hot, etc.). Use pre-processing methods (as needed) for dimensionality reduction, scaling, etc. Remove variables that are not needed/useful for the analysis. |
| Data Preparation Part 2 | 5 | Describe the final dataset that is used for classification/regression (include a description of any newly formed variables you created). |
| Modeling and Evaluation 1 | 10 | Choose and explain your evaluation metrics that you will use (i.e., accuracy, precision, recall, F-measure, or any metric we have discussed). Why are the measure(s) appropriate for analyzing the results of your modeling? Give a detailed explanation backing up any assertions. |
| Modeling and Evaluation 2 | 10 | Choose the method you will use for dividing your data into training and testing splits (i.e., are you using Stratified 10-fold cross validation? Why?). Explain why your chosen method is appropriate or use more than one method as appropriate. For example, if you are using time series data then you should be using continuous training and testing sets across time. |
| Modeling and Evaluation 3 | 20 | Create three different classification/regression models for each task (e.g., random forest, KNN, and SVM for task one and the same or different algorithms for task two). Two modeling techniques must be new (but the third could be SVM or logistic regression). Adjust parameters as appropriate to increase generalization performance using your chosen metric. You must investigate different parameters of the algorithms! |
| Modeling and Evaluation 4 | 10 | Analyze the results using your chosen method of evaluation. Use visualizations of the results to bolster the analysis. Explain any visuals and analyze why they are interesting to someone that might use this model. |

ML1- Ruberics

| | | |
|---------------------------|----|---|
| Modeling and Evaluation 5 | 10 | Discuss the advantages of each model for each classification task, if any. If there are not advantages, explain why. Is any model better than another? Is the difference significant with 95% confidence? Use proper statistical comparison methods. You must use statistical comparison techniques—be sure they are appropriate for your chosen method of validation as discussed in unit 7 of the course. |
| Modeling and Evaluation 6 | 10 | Which attributes from your analysis are most important? Use proper methods discussed in class to evaluate the importance of different attributes. Discuss the results and hypothesize about why certain attributes are more important than others for a given classification task. |
| Deployment | 5 | How useful is your model for interested parties (i.e., the companies or organizations that might want to use it for prediction)? How would you measure the model's value if it was used by these parties? How would you deploy your model for interested parties? What other data should be collected? How often would the model need to be updated, etc.? |
| Exceptional Work | 10 | You have free reign to provide additional analyses. One idea: grid search parameters in a parallelized fashion and visualize the performances across attributes. Which parameters are most significant for making a good model for each classification algorithm? |