

Lab 2: Basic Classification Algorithms

First Name: _____ Last Name: _____ NetID: _____

Lab 2 will be graded. Upload your work to Gradescope by Feb. 18th at 11:59pm. Make sure to upload a single .pdf containing your answers and the R code you wrote to obtain them. Regrade requests should be made via Gradescope.

In this lab we will learn the basics of logistic regression and the k -nearest neighbor algorithm, and learn how to measure the accuracy of the models we build.

Part (a): Logistic regression

Go to Canvas and download the `Titanic` data set.

The `Titanic` data set contains information(age, gender, passenger class etc.) about 1309 passengers that travelled on the ship `Titanic` and whether or not they survived. We will fit a logistic regression model to predict what sorts of passengers were more likely to survive.

We will divide the `Titanic` data set into 2 parts: the training data set and the test data set. We will fit the model and learn the parameters with the training set, and then test the accuracy of our model *using the test set* by comparing the predicted outcomes of the data in the test set using our model and their true outcomes. The training and test sets must be created from the original data without bias. In this lab, we'll randomly sample 2/3 of the original data to be the training data and let the rest be the test data. Why must we use the test set to calculate the accuracy?

Load the `titanic` data set into R. There are 11 variables in the data set except the outcome variable. We suspect that `Passenger ID`, `Name`, `Ticket` and `Cabin` are not useful in predicting survival chance so we can remove them. In addition, the `Cabin` column has lots of missing data. Check if there is still any other missing data in the data set. You'll find the `Age` column also has many rows where the data is missing. Can you think of two approaches to deal with these rows?

After figuring out how to deal with missing data, we need to check if all categorical variables have been identified by R and coerce them if not. Finally, we can create the training and test sets. Note that you need to store the outcome value of the test set separately and remove the outcome column from the test set.

Now you can fit the logistic regression model to the training data. Report the summary of your model. Once you have fitted the model, use your model to predict the outcomes of the test data. Since 'Survived' is already binary valued, you do not need to coerce the 'Survive' column to a factor. Also DO NOT remove the "Survived" column from df.test.

(Note that the program may reports error when using predict function due to some factor levels only being present in the test samples. In this case, you can do one of the following:

- remove test samples containing the offending factors
- do not coerce the variable to factor
- retry with a new random seed until you get a split that works, and keep that seed

Depends on the split of your data set, you may or may not encounter this issue. see also: <https://www.r-bloggers.com/linear-logistic-regression-in-r-dealing-with-unknown-factor-levels-in-test-data/>)

The logistic regression model we built will output the probability of surviving of each passenger in the test data set. To get a binary prediction of whether a passenger survives or not, we need to pick a threshold. For now we simply predict a passenger would survive if the output survival probability is greater than or equal to 0.5. What is the (overall) accuracy of your model? In addition, report the true positive and true negative rates.

Now that you know how to measure accuracy, can you find a very simple classifier that would give you a decent overall accuracy (about 65% or more) by just looking at the data?

As you can see, the data is biased. As the data gets more biased, our formula for overall accuracy becomes less indicative of how good our model is. A way to solve this is to use the following formula:

$$\text{balanced_accuracy} = \frac{1}{2} \times \frac{t_p}{(t_p + f_n)} + \frac{1}{2} \times \frac{t_n}{(t_n + f_p)}$$

Using this formula, calculate the balanced accuracy of our logistic regression model on the test dataset.

Identify which predictors are most influential to the final prediction.

Part (b): K Nearest Neighbors

We will continue analyzing the `Titanic` data set, but with a different classification method. In particular, We will learn how to do KNN in R, and compare the performance of KNN with the logistic regression approach used in part (a).

Do the same procedures to the Titanic data set as we did in part(a), i.e.

1. Delete columns `Passenger ID`, `Name`, `Ticket` and `Cabin`.
2. Deal with missing data by removing all rows associated with a missing value (alternatively, you may replace missing values with mean/median).
3. Check for factor variables and convert them into dummy variables.

Note that we did not need the third step in part(a). Neither KNN nor linear/logistic regression deals with factor variables directly. The `lm()` and `glm()` functions in R automatically convert factors to dummies, but the built-in function `knn()` for KNN does not.

Can you think of the reason why KNN cannot deal with factor variables directly?

To convert factor variables into dummy variables, we use the function `model.matrix()`.

```
> titanic$Sex <- model.matrix( ~ Sex - 1, data=titanic )
> titanic$Embarked <- model.matrix( ~ Embarked - 1, data=titanic)
```

Why don't we just use the `as.numeric()` function to do the conversion?

We will divide the Titanic data set into training and test sets, again by randomly sampling 2/3 of the data to be the training set and using the rest as the test set. We will set a random seed so the the sampling is reproducible.

```
> set.seed(1)
> train_ind <- sample(1:nrow(titanic), 2/3*nrow(titanic))
```

We will now perform KNN using the `knn()` function, which is part of the `class` library. This function works rather differently from the other model fitting functions that we have encountered thus far. Rather than a two-step procedure in which we first fit the model and then we use the model to make predictions, the `knn()` function computes the predictions directly in a single command. Running `?knn` to display the help page, we see that the function requires four inputs:

1. A matrix containing the predictors associated with the training data, labeled `titanic.train` below.
2. A matrix containing the predictors associated with the data for which we wish to make predictions, labeled `titanic.test` below.
3. A vector containing the outcomes (class labels) for the training observations, labeled `train.survive` below.
4. A value for k , the number of nearest neighbors to be used by the classifier.

Before applying the `knn()` function, we normalize each predictors to have unit variance. Here we use function `scale()` to make each variables to have 0 mean and variance 1.

```
> titanic_normal <- scale(titanic[,!names(titanic) %in% 'Survived'])
> titanic.train <- titanic_normal[train_ind, ]
> titanic.test <- titanic_normal[-train_ind, ]
```

Why do we need this step?

We next create two variables that store the outcome values of the training/test sets.

```
train.survive = titanic$Survived[train_ind]
test.survive  = titanic$Survived[-train_ind]
```

We can now apply the `knn()` function to predict the survival of the passengers. Note that in KNN, if several observations are tied as nearest neighbors, R will randomly break the tie.

```
> knn.pred = knn(titanic.train, titanic.test, train.survive, k=1)
> table(knn.pred, test.survive)
> mean(knn.pred == test.survive)
```

Here we choose $K = 1$, which gives an error rate of about 23% (according to my experiment).

We explore the performance of KNN with other values of K . Run KNN for $K = 1, 3, 5, 10, 20, 50$. (You may use a for loop for convenience). Report the corresponding error rates, and find the value of K that gives the lowest error rate.

Finally, compare the prediction performance of KNN (using your optimal K) with the logistic regression model from part(a). Which method gives a better accuracy rate in this case?

Take Home Questions

- (1) We introduced *balanced accuracy* in Part (a). Is balanced accuracy the same as overall accuracy? Prove or give a counterexample.

- (2) Give an example of a scenario when the balanced accuracy provides a better measure than the overall accuracy.

- (3) Recall in Lab 1 you tried to get the units of the ω coefficients of the linear regression model. What are the units of the coefficients of the logistic regression model. (Hint: Recall the logistic regression model assumes:

$$\mathbb{P}(y = 1|\mathbf{x}, \beta) = \frac{1}{1 + \exp(-\beta^T \mathbf{x})}$$

Try to find an expression $f(p)$ of $p = \mathbb{P}(y = 1|\mathbf{x}, \beta)$ which satisfies $f(p) = \beta^T \mathbf{x}$. We call $f(p)$ “logit” or “log-odds”. Recall the interpretation of β coefficients in the linear regression model is the change of outcome variable when a predictor changes by 1 unit. Interpret β coefficients of the logistic regression model in a similar way.)

For every one year of age, a passenger’s logit of surviving on Titanic will decrease (or increase) by a factor of what?

- (4) Go to Canvas and download the Cancer data set: `pros.xls` and `pros_meta.txt`. We want to build a logistic regression model to predict whether the tumor has penetrated the prostatic capsule.

Load the data into R. Note the ‘VOL’ column has lots of missing data and you need to remove these instances first. We will fit a logistic regression model with the data and test the accuracy of our model.

Divide the data set into training and test data. Fit a logistic regression model with the training data. Report the summary of your model.

- (5) Now let’s check the accuracy of your model in Q2. Remember the logistic regression model outputs a probability of penetration p for each patient. We could pick a threshold

p_0 so that we predict the outcome to be 1 if $p \geq p_0$. Hence different p_0 gives different classifiers.

Define $\text{FPR} = \frac{f_p}{f_p + t_n}$ to be the false positive rate: the proportion of negative instances that are falsely predicted to be positive instances. Moreover, define $\text{TPR} = \frac{t_p}{t_p + f_n}$ to be the true positive rate: the proportion of positive instances that are also predicted correctly. Each p_0 gives a new pair of (FPR, TPR).

Pick a list of p_0 and plot the corresponding (FPR, TPR) pairs on a graph. We obtain the FPR vs. TPR plot (i.e., the ROC curve) for our logistic regression model as shown in Figure 1. (If you generate this plot yourself, it will look slightly different due to the randomness in splitting in to training and test datasets. But the general trend of the curve will be the same.)

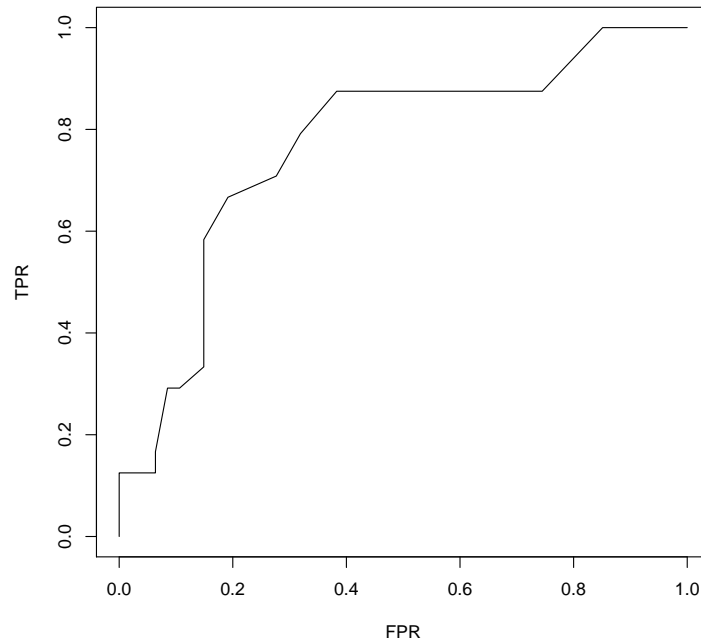


Figure 1: ROC curve: False Positive Rate vs. True Positive Rate

Next, assume we have a balanced data set, i.e. with 50% positive instances and 50% negative ones. What would the FPR vs. TPR plot look like for the following two classifiers:

- You were guessing the prediction randomly.
- You were predicting only one of the classes for everything.

What would the plots in 5a and 5b look like if we have an unbalanced data set, say $p\%$ positive class?

- (6) In predicting cancer diagnosis, we would like to be more conservative and prefer to error on the side of a false diagnosis rather than missing cancer. What would you do in this case when making predictions according to the probabilities output by your logistic regression model?

More on Logistic Regression

1. Will the fitted coefficients and the classification change if we duplicate the whole training set while using the logistic regression? What happens if we only duplicate a subset of the training set? Give intuitive explanations and verify your result using the **Titanic** data set.

2. Duplicating a subset of the training set may be used to fix an imbalanced data set. A common approach is called the oversampling, where we resample or duplicate the minority class to balance the data set. State one potential limitation of the oversampling method.

3. Recall from the lecture that the logistic regression can be extended to the multi-class classification tasks. Consider the following version of the **One-Versus-All** method. Suppose we have K classes, for each class i where $0 \leq i \leq K-1$, we perform a logistic regression and determine the weights $\beta_0^{(i)}$ and $\beta^{(i)}$. We define

$$\mathbb{P}(y = i \mid \mathbf{x}, \beta) = \frac{e^{\beta_0^{(i)} + \beta^{(i)} \cdot \mathbf{x}}}{\sum_{j=0}^{K-1} e^{\beta_0^{(j)} + \beta^{(j)} \cdot \mathbf{x}}} \quad (1)$$

In the case of binary classification, i.e. $K = 2$, how does the predicted probability using **One-Versus-All** differ from the original logistic regression?

Recall that logistic regression for 2 classes has

$$\mathbb{P}(y = 1 \mid \mathbf{x}, \beta) = \frac{1}{1 + e^{-(\beta_0 + \beta \cdot \mathbf{x})}}$$

What are the relationships between β_0 , β , $\beta_0^{(i)}$ and $\beta^{(i)}$ in the case of binary classification?

In practice, when we perform the classification using (1), we will fix a class i as the base class and set $\beta^{(1)} = 0$ and $\beta_0^{(i)} = 0$. Explain briefly why this step is performed. Is the choice of i important for classification?

4. For simplicity, we consider the logistic regression with threshold 0.5. What is the formula for the classification boundary?

Suppose the logistic regression gives β_0 and β as the fitted coefficients. How does the classification result change if we predict instead using the weights $c\beta_0$ and $c\beta$ for some constant $c > 0$? Will the classification change if $c < 0$.

(Optional) Use the observation above, explain heuristically why the logistic regression can perform badly when fitted by the maximum likelihood, when the data *can* be linearly separated?

(Optional) Test this phenomenon on some simulated data in R. Give a printout of the warning messages when you run the logistic regression.

KNN

1. Compare the performance of KNN with and without normalization of the `Titanic` data set.

Create a plot of the value of K vs. the associating error rate, where $K \in [1, 20]$, with normalization. Next, make the same plot but without normalizing the data set. Compare the two plots. (You may attach additional page for printing out the 2 plots.)

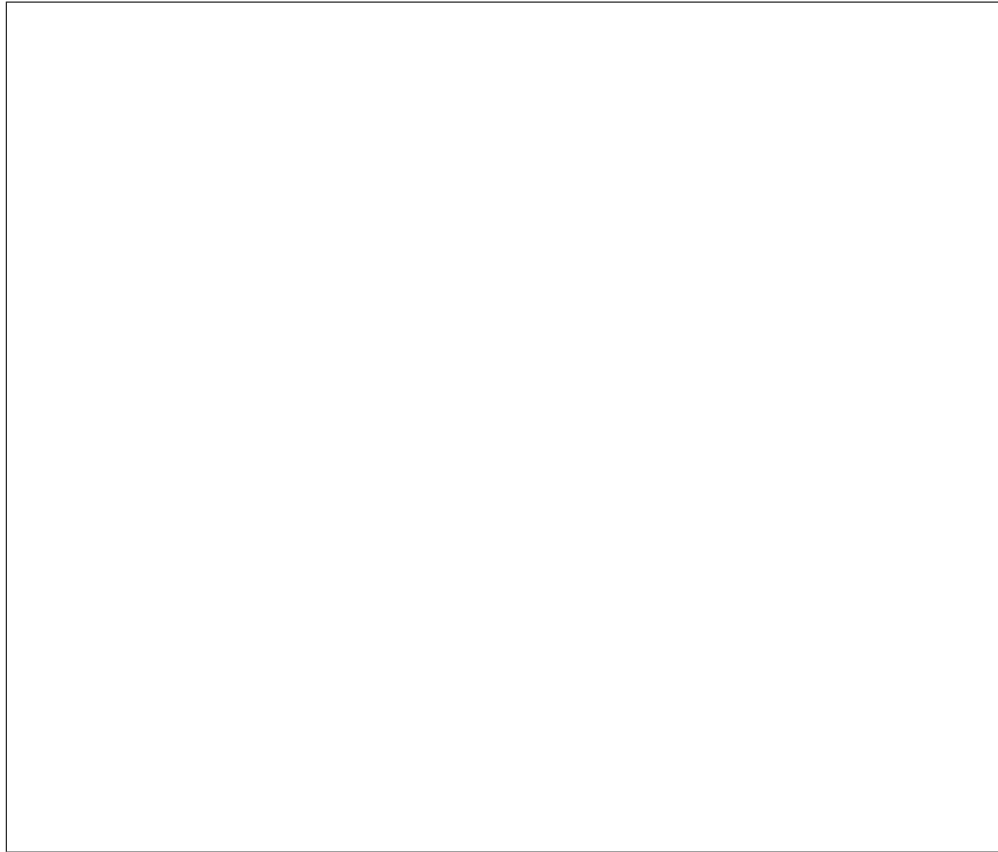
What can you say about the accuracy rates with and without normalization?



2. What will be the prediction result if we set K to be the total number of training data points?



3. Will the classification change if we duplicate the training set while keeping K constant? Someone claims that using $K_{\text{duplicated}} = 2k$ neighbors in the duplicated training set should always yield the same prediction with using $K_{\text{original}} = k$ neighbors in the original training set. Is this claim true? Illustrate it with the **Titanic** data set by plotting the error rates for $K_{\text{duplicated}} = 2k$ and $K_{\text{original}} = k$, $k \in [1, 30]$. (Set `use.all=FALSE` in `knn`).



4. Make a **concept map**¹ for the topics covered by this lab. While there is no single right answer, you are expected to demonstrate that you understand the relationships between various topics in the lab.

¹https://en.wikipedia.org/wiki/Concept_map