

---

# ***SUPERVISED LEARNING: CLASSIFICATION***

---

Main objective of the analysis that specifies whether your model will be focused on prediction or interpretation and the benefits that your analysis provides to the business or stakeholders of this data.

Main objective with this work is to model for prediction purposes.

Benefits:

- Main aim for this project from a business perspective is to reduce cost the company.
- Context here is that a company is hiring people for sales and training them so a huge upfront cost is involved.
- This model tries to predict during the time of application itself whether an applicant will turn out to be a good sales agent or not.

Brief description of the data set you chose, a summary of its attributes, and an outline of what you are trying to accomplish with this analysis.

The data I took is about hiring new associates/ salespeople by managers in sales division. Since the company is making upfront investment in hiring and training, it tries to predict based on past data if an applicant will produce business after getting hired in future.

So, the data is made up of 14572 rows (each of a distinct applicant) and 22 features for prediction and one binary target variable of which training data has 9527 entries, and test has rest. So we have 14572 people who applied in the past and got hired and whether they sold a product within three months is reflected in business sourced target column. The features/ variables are:

- Identification: 'ID'
- Location of Offices: 'Office\_PIN'
- Applicant Details: 'Application\_Receipt\_Date', 'Applicant\_City\_PIN', 'Applicant\_Gender', 'Applicant\_BirthDate', 'Applicant\_Marital\_Status', 'Applicant\_Occupation', 'Applicant\_Qualification'
- Manager details: 'Manager\_DOJ', 'Manager\_Joining\_Designation', 'Manager\_Current\_Designation', 'Manager\_Grade', 'Manager\_Status', 'Manager\_Gender', 'Manager\_DoB', 'Manager\_Num\_Application', 'Manager\_Num\_Coded', 'Manager\_Business', 'Manager\_Num\_Products', 'Manager\_Business2', 'Manager\_Num\_Products2'
- Target: 'Business\_Sourced'<sup>1</sup>

Here, variable types are:

- Date: 'Application\_Receipt\_Date', 'Applicant\_BirthDate', 'Manager\_DOJ' and 'Manager\_DoB'
- Numerical: 'Manager\_Num\_Application', 'Manager\_Business', 'Manager\_Num\_Products', 'Manager\_Business2', 'Manager\_Num\_Products2'

---

<sup>1</sup> For the rest of the document, successful means Business sourced is 1, otherwise 0.

Categorical: rest all

I am trying to model the data to predict whether an applicant now will become successful in future or not, ie, Business sourced is 1 or not.

## Brief summary of data exploration and actions taken for data cleaning and feature engineering.

1. Handling date variables
  - a. Since we cannot handle date variables directly, we can get the information out of it by creating new columns.
  - b. To see if month or year of application date has any effect in determining target, creating separate columns each for month and year instead of date directly
  - c. From applicant's DoB, creating a variable for applicant's age since I believe month and date wont have any meaningful impact on their success of work
  - d. Similarly using 'Manager\_DOJ' and 'Manager\_DoB' to calculate manager's age and experience in the company
2. Handling non-numerical variables:
  - a. PIN codes were used as geocodes to identify distance between applicant address and office location.
  - b. Gender and Manager Status (whether probation or confirmed) were binary encoded.
  - c. One Hot encoding (encoding with indicator variables) was performed for: 'Applicant\_Marital\_Status', 'Applicant Occupation'
  - d. Label encoding (Rank ordering) was performed for: 'Manager\_Joining\_Designation', 'Manager\_Current\_Designation', 'Manager\_Grade'. Grade was already ordinal; designation was of format "Level xx" was converted to ran based on xx.
  - e. Target encoding was performed for: 'Applicant\_Qualification'
3. Handling missing values
  - a. I noticed something in the data wrt missing values, ie, for 683 observations there was no information on any column describing manager details. So I figured these entries must mean, they had no manager, So created a new column for indicating the presence of manager.
  - b. Since every other variable with missing values except birthdate was a categorical variable and #missing is less than half a percent of overall data, observations with birthdate missing were dropped and others excluding Applicant\_Occupation were imputed with KNN iteratively with Applicant\_Gender having least missing imputed first. Due to large number of missing values, Applicant\_Occupation missing was treated as a new category.
4. Handling outliers
  - a. 'Manager\_Num\_Products' had unproportionally large count of 0 and very low counts of extremely high values. Hence, they were floored and capped at 10, 90 percentiles respectively and indicator for flooring and capping was added as a new feature
  - b. 'Manager\_Business' also had a couple of negative values and some extremely high values. Hence, they were capped and floored at 95, 5 percentiles respectively.
5. Variable Transformation
  - a. 'Manager\_Business' was heavily skewed towards left hence a log transformation was applied.
  - b. 'Manager\_Business2', 'Manager\_Num\_Products2' are essentially second line of business of a manager if exists. So, when it doesn't the data replicates the first and replicates them which leads to extremely high correlation between them. So instead of these, their difference from 1<sup>st</sup> line of business, ie, 'Manager\_Business', 'Manager\_Num\_Products' were calculated and used to create new columns and dropping these.

- c. A new column for growth of manager was calculated as a difference of 'Manager\_Joining\_Designation' & 'Manager\_Current\_Designation'
- d. Applicant\_Qualification was re-grouped into Class X, Class XII, Graduate, Certified Associate, Masters & Others and target encoded.

## Summary of training at least three different classifier models, preferably of different nature in explainability and predictability.

The modelling methods which I too into consideration were: 'Logistics regression', 'Decision tree', 'Support Vector Machine', 'Random Forest', 'K Nearest Neighbours'. Note that all of these we modelled and predicted on same train & test data.

The exact hyper-parameters used for each method were:

1. **LogisticRegression**(solver='newton-cg',max\_iter=500,n\_jobs=-1)
2. **DecisionTreeClassifier**(max\_depth=7, min\_samples\_split=50, min\_samples\_leaf=10,min\_impurity\_decrease=0.0001)
3. **SVC**(probability = True,max\_iter=500,break\_ties=True)
4. **RandomForestClassifier**(oob\_score=True,max\_depth=7, min\_samples\_split=50,min\_samples\_leaf=10,min\_impurity\_decrease=0.0001,n\_jobs=-1)
5. **KNeighborsClassifier**(n\_neighbors=10,n\_jobs=-1)

The results I got were:

### -----Training Data-----

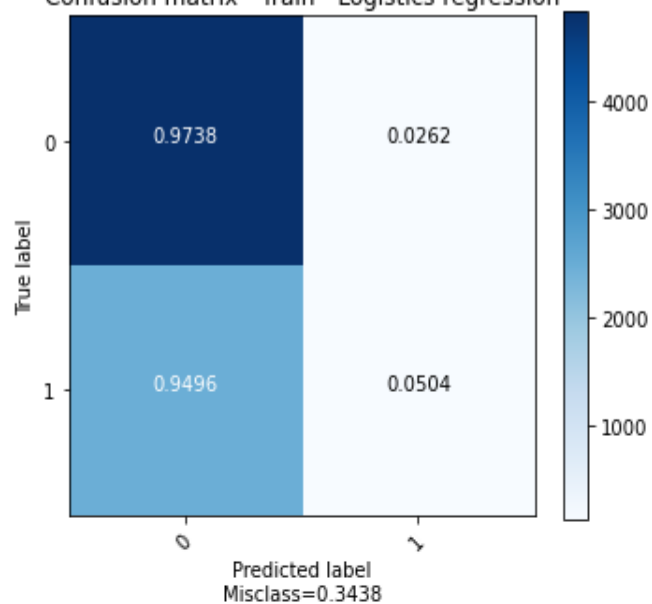
	accuracy	cohen_kappa_score	f1_score	precision_score	recall_score	roc_auc_score
<b>Logistics regression</b>	0.656221	0.030747	0.091544	0.501916	0.050365	0.512083
<b>Decision tree</b>	0.680550	0.157208	0.300116	0.608696	0.199154	0.566022
<b>SVM</b>	0.346291	-0.005482	0.506587	0.342095	0.975779	0.496051
<b>Random Forest</b>	0.666006	0.039328	0.061664	0.912088	0.031911	0.515149
<b>KNN</b>	0.683327	0.158393	0.292467	0.631378	0.190311	0.566034

### -----Test Data-----

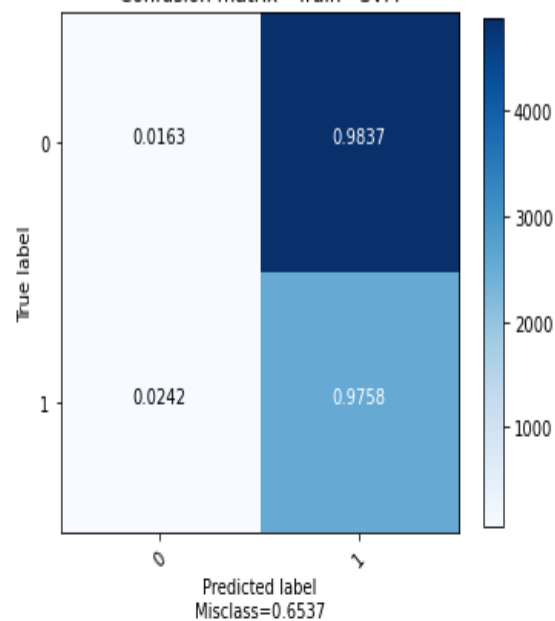
	accuracy	cohen_kappa_score	f1_score	precision_score	recall_score	roc_auc_score
<b>Logistics regression</b>	0.643575	0.017638	0.081744	0.441176	0.045045	0.507012
<b>Decision tree</b>	0.632470	0.043491	0.202067	0.429268	0.132132	0.518311
<b>SVM</b>	0.352195	-0.008273	0.514467	0.349488	0.974474	0.494176
<b>Random Forest</b>	0.645690	0.000207	0.014706	0.357143	0.007508	0.500080
<b>KNN</b>	0.631941	0.029079	0.169451	0.412791	0.106607	0.512079

And the confusion matrix for each data sample for each method is shown below:

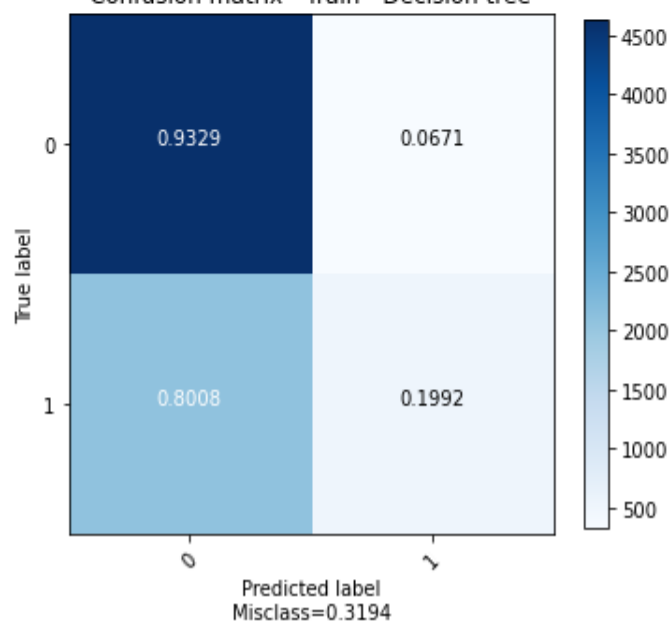
Confusion matrix - Train - Logistics regression



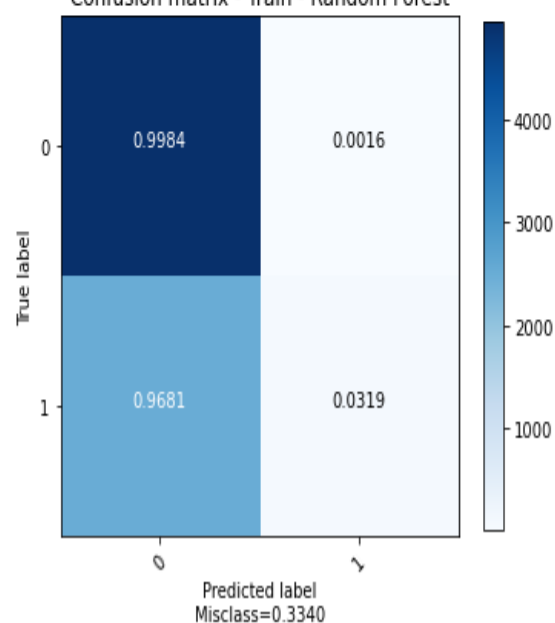
Confusion matrix - Train - SVM



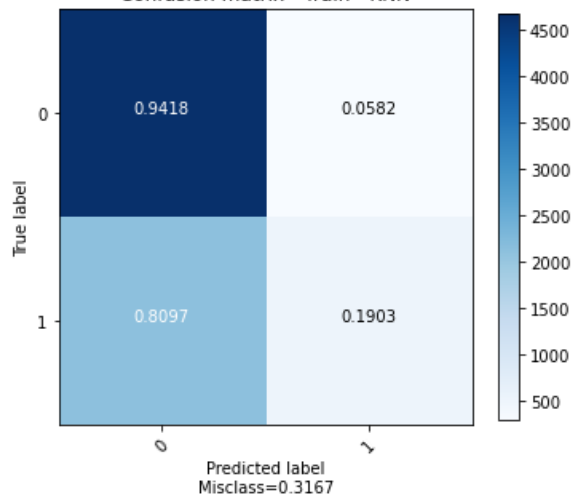
Confusion matrix - Train - Decision tree

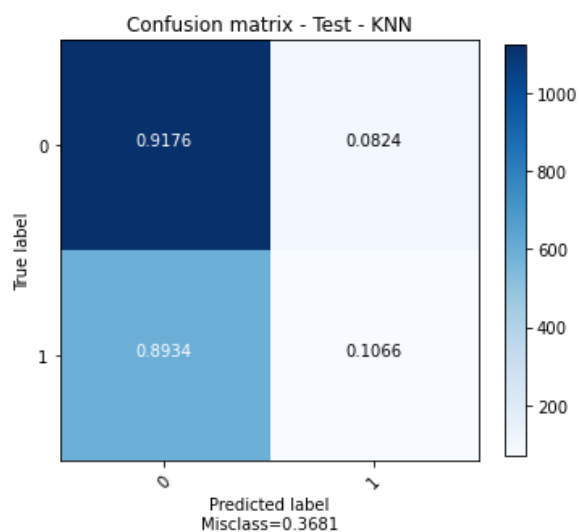
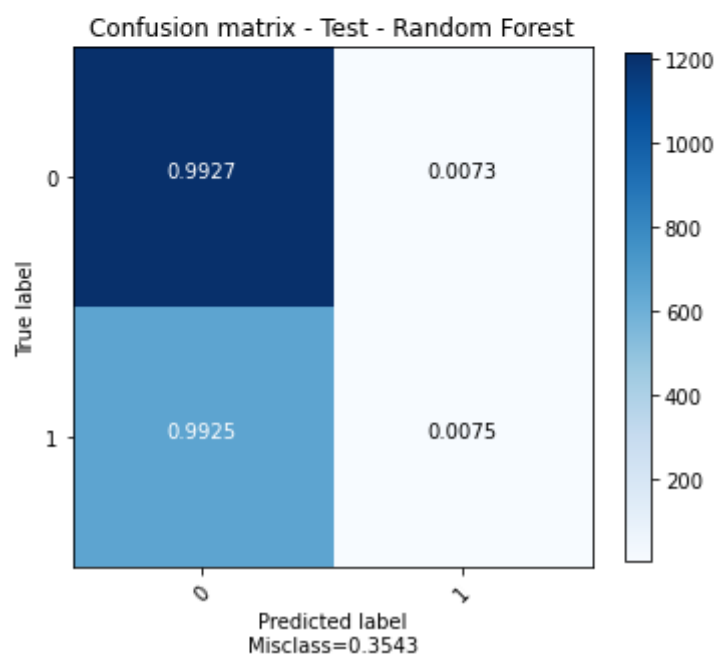
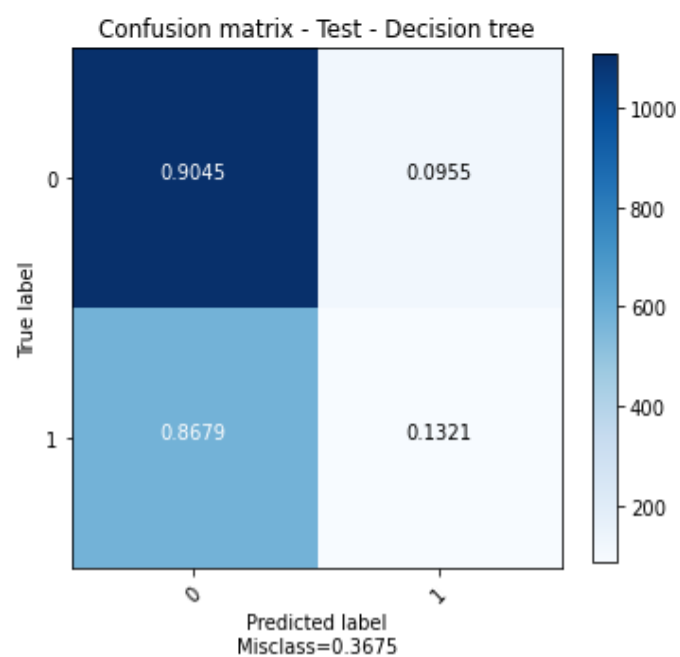
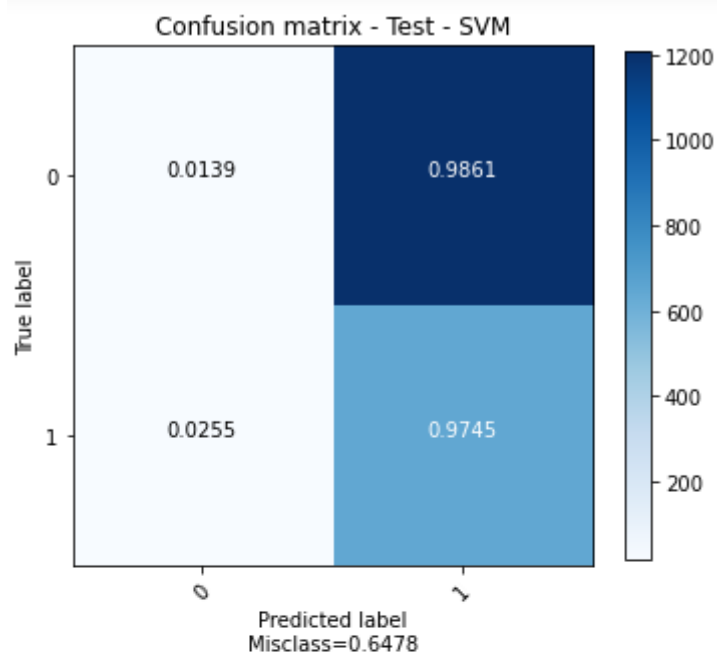
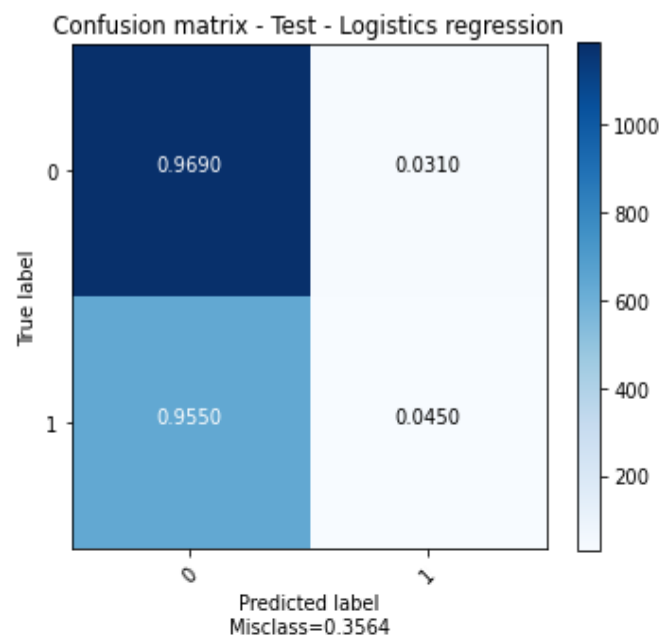


Confusion matrix - Train - Random Forest



Confusion matrix - Train - KNN





## A paragraph explaining which of your classifier models you recommend as a final model that best fits your needs in terms of accuracy and explainability.

Based on summary statistics of results, we can infer the following:

1. For all the methods, with respect to all comparison metrics, train and test data perform similarly.
2. SVM has uniformly poorer results across scores like accuracy, precision, cohen's kappa than other methods.
3. KNN and tree based methods gives better performance with training data over Logistic Regression.
4. Tree based methods gives better performance with test data over KNN and Logistic regression.

Therefore, now I am comparing between tree based methods.

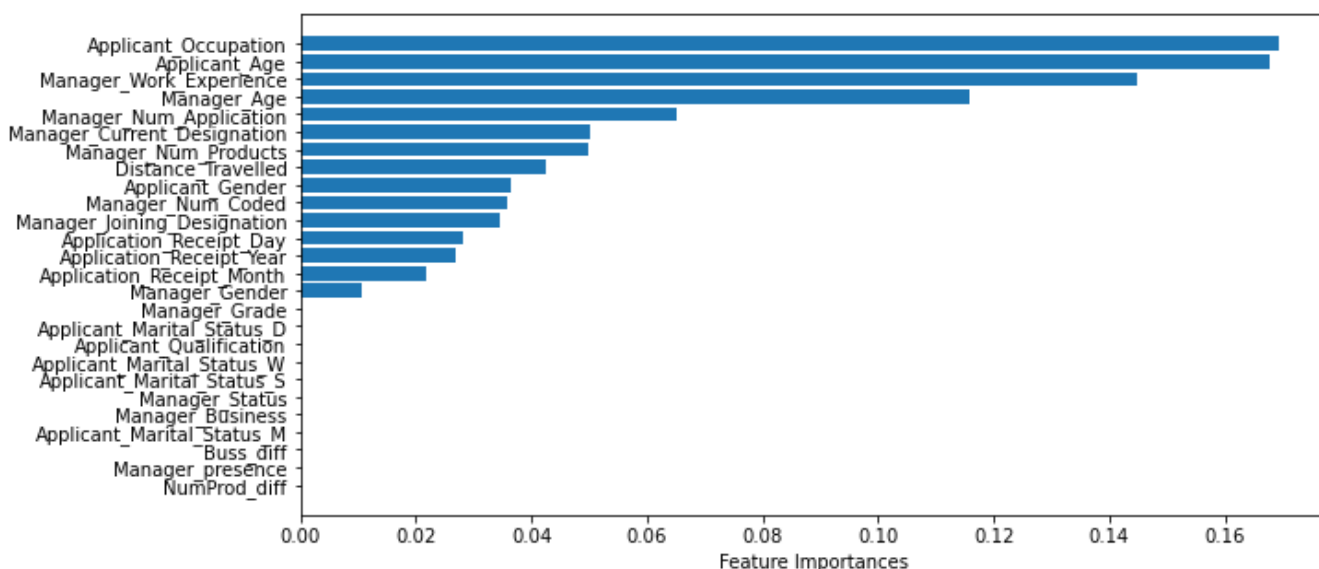
With respect to accuracy and roc-auc metrics, both have similar scores in both train and test whereas in precision, RF has better results than decision tree in train data and Decision tree has better results in both precision and recall in both train and test.

Now, looking at visual ease of understanding and from a perspective of explainability, Decision tree is better than random forest.

Hence, I will recommend to use Decision tree model with parameters as mentioned above.

## Summary Key Findings and Insights, which walks your reader through the main drivers of your model and insights from your data derived from your classifier model.

Feature importances based on the final tree is plotted here:



As we can infer from above plot, an applicants' current occupation and their age are the most important drivers in determining success of an applicant having very close importance values with Manager's work experience, age and number of applications they received coming next.

Suggestions for next steps in analysing this data, which may include suggesting revisiting this model after adding specific data features that may help you achieve a better explanation or a better prediction.

Next step would be to do Hyper parameter tuning and cost complexity pruning and also try a bagged or boosted version of decision tree to improve its overall accuracy and precision.

Additional data regarding performance of applicant in their occupation would give better insights into the overall capacity of the applicant. Also, data on whether that applicant has previous experience in sales would also be a good predictor I believe.