

STROKE PREDICTION

THOMAS CHAMPION
UIS CSC 532B
SPRING 2021

PROJECT DEFINITION AND GOALS

- STROKES ACCOUNT FOR APPROXIMATELY 11% OF DEATHS WORLDWIDE
- SURVIVORS MAY EXPERIENCE LASTING EFFECTS AFTER A STROKE
 - DEPRESSION
 - PHYSICAL IMPAIRMENTS
- ASIDE FROM PHYSICAL IMPACT, THERE ARE FINANCIAL IMPACTS ASSOCIATED WITH A STROKE

PROJECT DEFINITION AND GOALS

- GOAL IS TO FIND A MODEL THAT WILL PREDICT A STROKE BASED UPON DATASET CONTAINING INFORMATION OF INDIVIDUALS THAT HAVE AND HAVE NOT EXPERIENCED A STROKE

STROKE PREDICTION DATASET

- DATASET WAS OBTAINED FROM KAGGLE.COM
- ORIGINALLY POSTED IN JANUARY 2021
- ~177K VIEWS
- ~26.8K DOWNLOADS
- 346 UNIQUE CONTRIBUTORS
 - DATA VISUALIZATION
 - DATA CLEANING AND BALANCING
 - STROKE PREDICTION

DATA EXPLORATION

- 5110 OBSERVATION
- 12 VARIABLES
 - 1 BINARY RESPONSE VARIABLE
 - 11 PREDICTIVE VARIABLES

Feature Name	Data Type	Description
id	Numeric	unique identification number
gender	Character	gender of individual with three values:: "Male", "Female", "Other"
age	Numeric	age of individual
hypertension	Numeric	Indication if the individual was diagnosed with hypertension. 0=No, 1=Yes
heart_disease	Numeric	indication if an individual has heart disease. 0=No, 1=Yes
ever_married	Character	indication if the individual has ever been married. "Yes" or "No"
work_type	Character	Indication of the individual's category of work. Five values: "children", 'Govt_job", "Never_worked", "Private" and "Self-employed"
Residence_type	Character	Indication of general area an individual resides. Two values: "Rural" and "Urban"
avg_glucose_level	Numeric	Individual's average blood glucose level
bmi	Character	Individual's body mass index
smoking_status	Character	Indication if individual smokes or has ever smoked. Four values: "formerly smoked", 'never smoked", "smokes", "Unknown"
stroke	Numeric	Binary response variable indicating whether an individual experienced a stroke. 0=No, 1=Yes

DATA EXPLORATION AND PREPROCESSING

TYPE CONVERSIONS

- CONVERTED EVER_MARRIED TO NUMERIC
 - 0 = NO, 1 = YES
- CONVERTED RESIDENCE_TYPE TO NUMERIC
 - 0 = RURAL, 1 = URBAN
- REMOVED 1 OBSERVATION WITH GENDER = 'OTHER'
- CONVERTED GENDER TO NUMERIC
 - 0 = FEMALE, 1 = MALE
- CONVERTED WORK_TYPE TO FACTOR
- CONVERTED SMOKING_STATUS TO FACTOR

DATA EXPLORATION AND PREPROCESSING

STROKE RESPONSE VARIABLE

- HIGHLY IMBALANCED DATA
 - 4861 OBSERVATIONS = NO STROKE
 - 249 OBSERVATIONS = STROKE
- WILL BE NECESSARY TO INCLUDE BALANCING AS PART OF THE MODELING PROCESS

DATA EXPLORATION AND PREPROCESSING

BMI

- BMI IS CHARACTER FIELD STORING NUMERIC VALUES
 - CONTAINS 'N/A' AS VALUES
- CONVERT 'N/A' TO NA
- CONVERT FIELD TO NUMERIC
- 201 MISSING VALUES

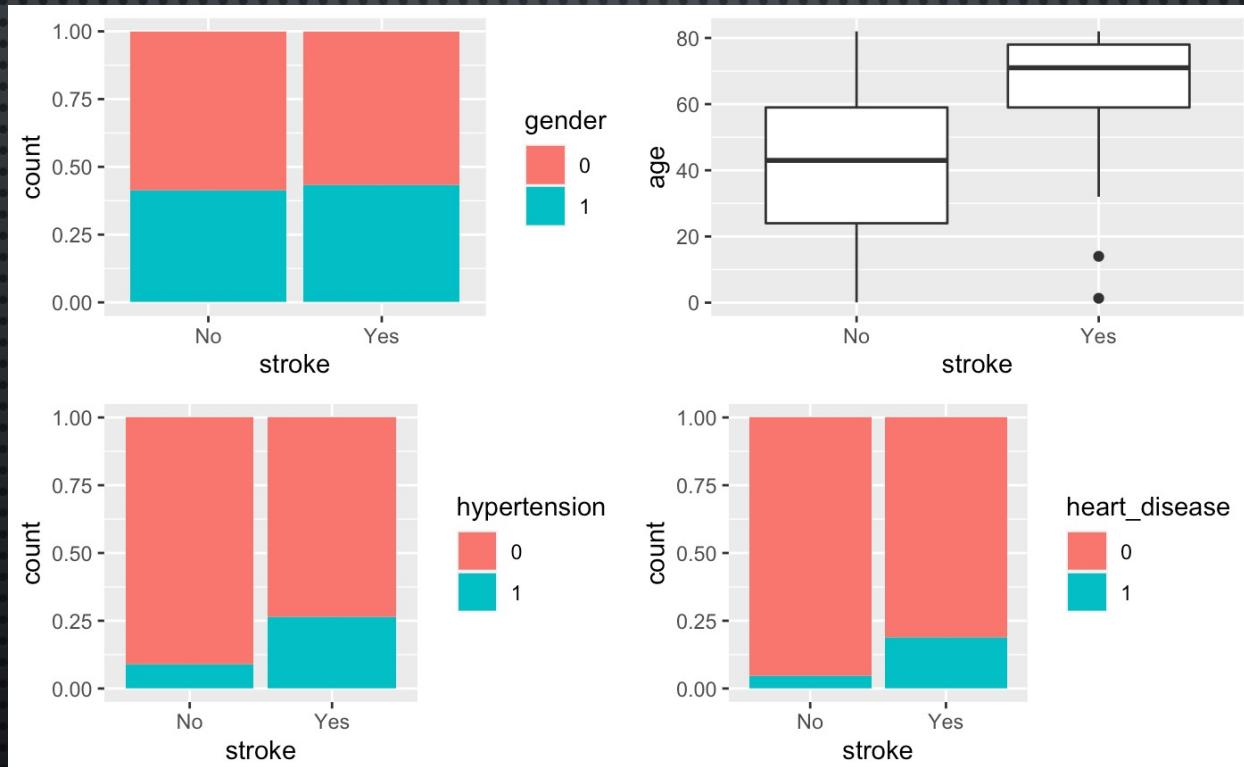
DATA EXPLORATION AND PREPROCESSING

BMI

- IMPUTE MISSING VALUES
 - CREATE 10 BINS USING CLUSTERING METHOD ON AGE VARIABLE
 - SET MISSING BMI VALUES EQUAL TO THE MEDIAN AGE VALUE OF THE BIN THEY BELONG TO

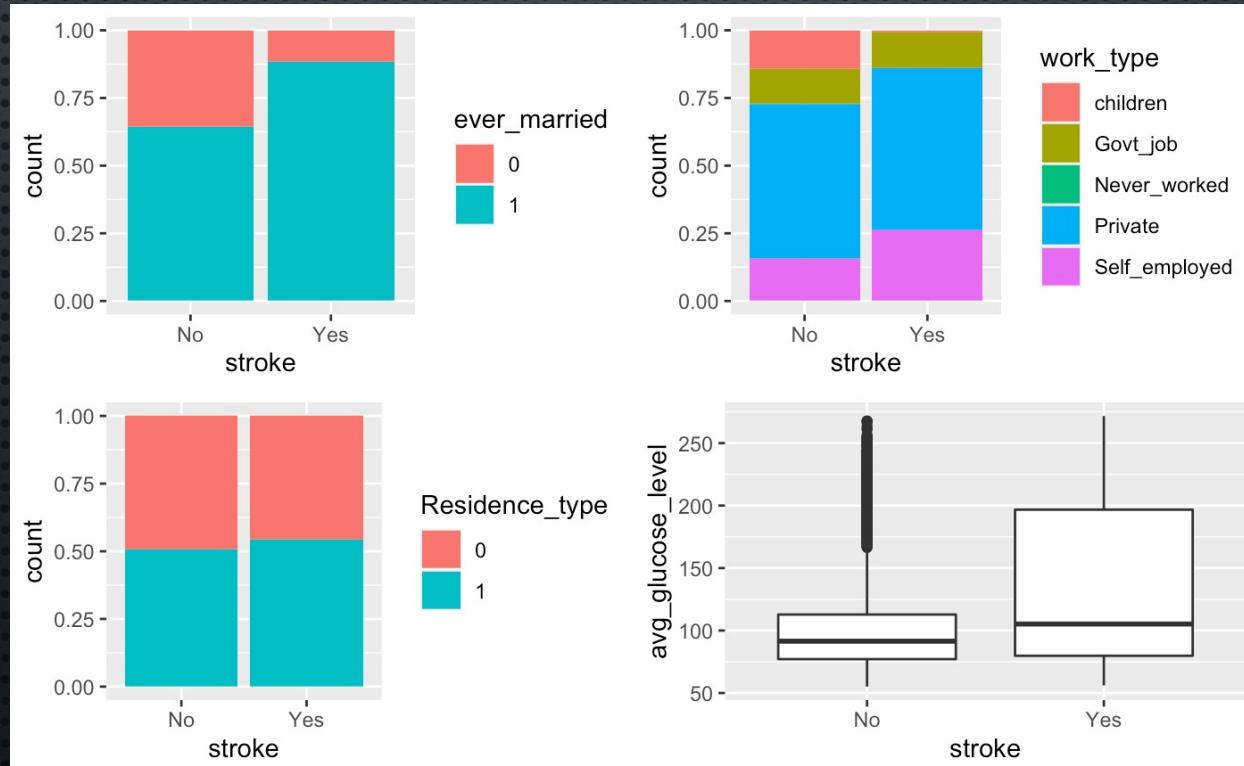
DATA EXPLORATION AND PREPROCESSING

ASSOCIATION ANALYSIS



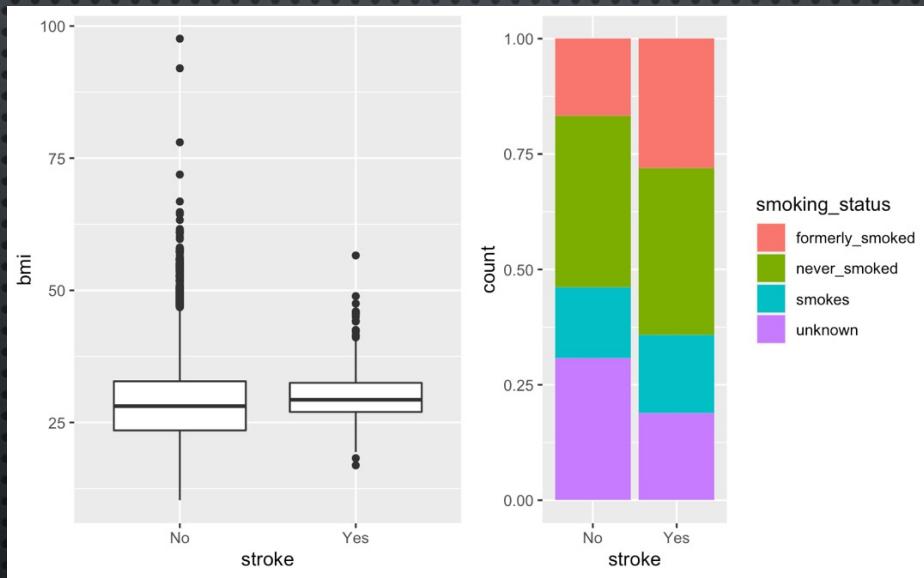
DATA EXPLORATION AND PREPROCESSING

ASSOCIATION ANALYSIS



DATA EXPLORATION AND PREPROCESSING

ASSOCIATION ANALYSIS



DATA EXPLORATION AND PREPROCESSING

ASSOCIATION ANALYSIS

Chi-Square Tests

chisq_labels	chisq_values
smoking_status	2.008e-06
gender	0.5598
hypertension	1.689e-19
heart_disease	2.121e-21
ever_married	1.686e-14
work_type	5.409e-10
residence_type	0.2998

T-Tests

ttest_labels	ttest_values
age	2.176e-95
glucose	2.373e-11
bmi	0.0003064

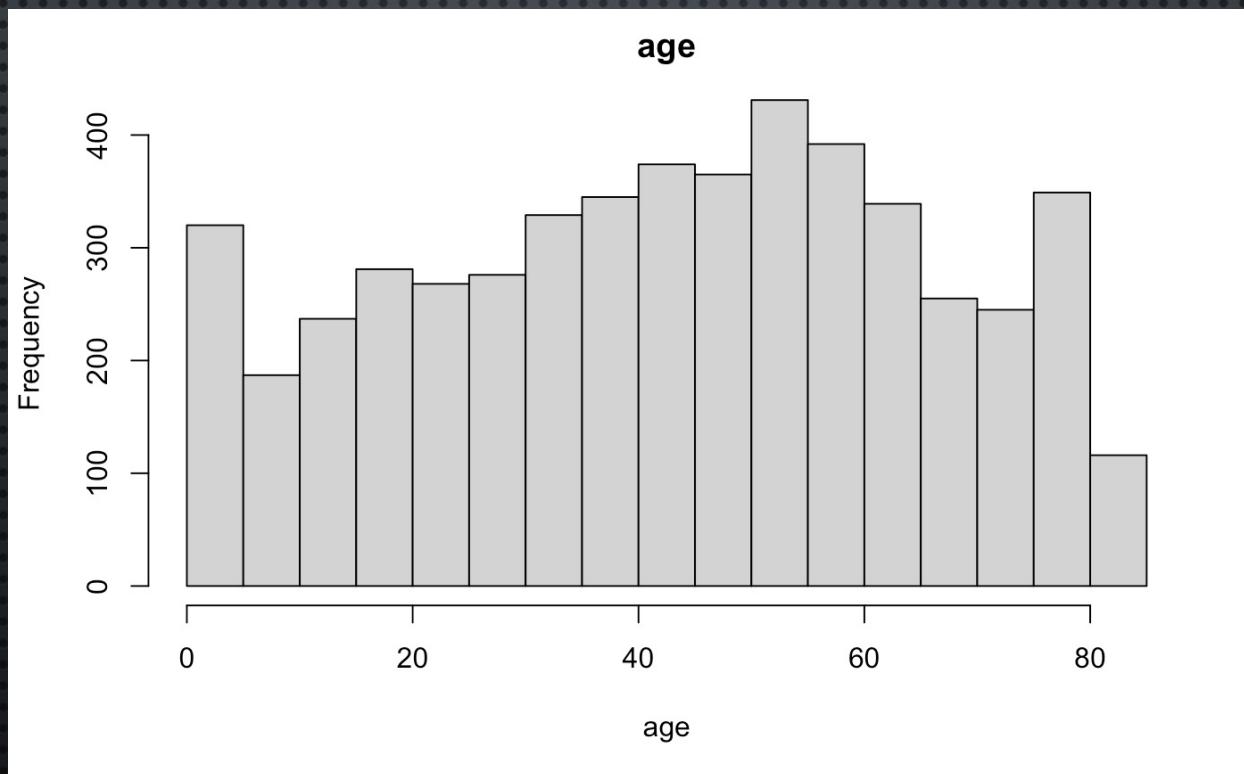
DATA EXPLORATION AND PREPROCESSING

ASSOCIATION ANALYSIS

- ALL VARIABLES SHOWED A STATISTICALLY SIGNIFICANT RELATION WITH STROKE EXCEPT FOR RESIDENCE_TYPE AND GENDER
 - BOTH VARIABLES REMOVED

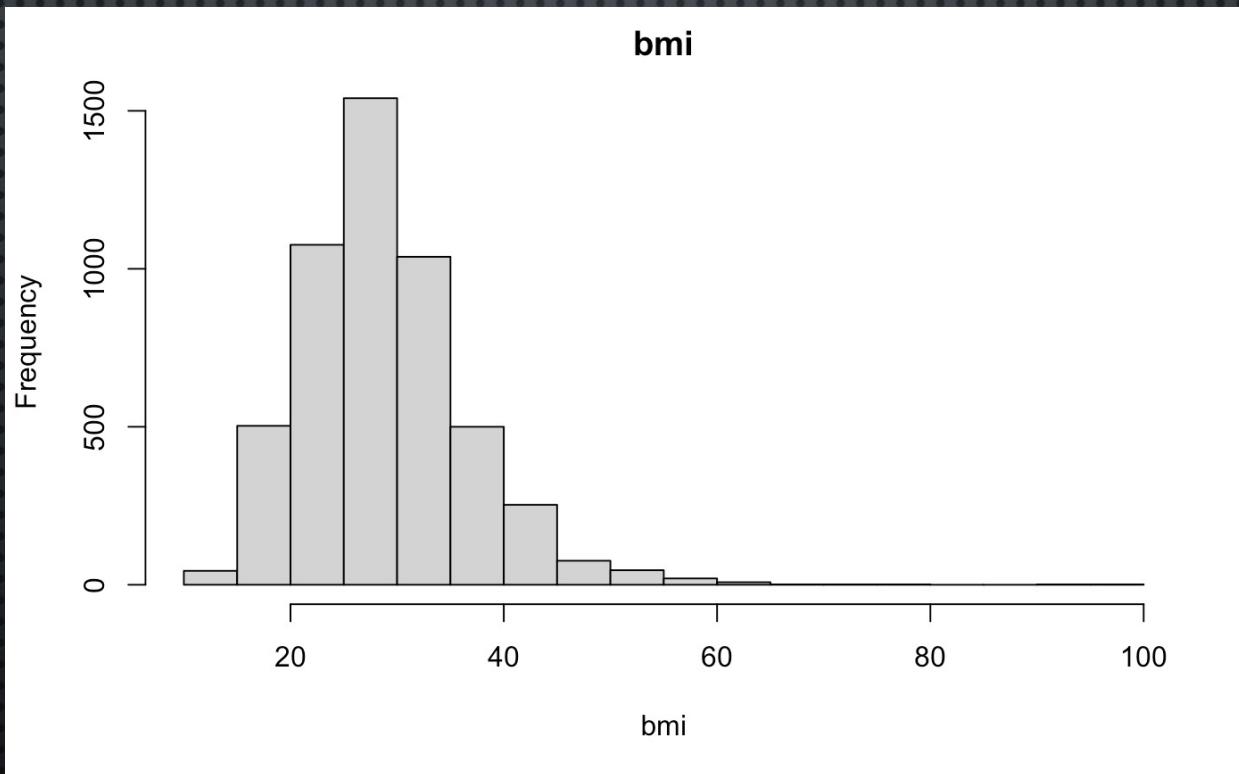
DATA EXPLORATION AND PREPROCESSING

DISTRIBUTION OF NUMERIC VARIABLES



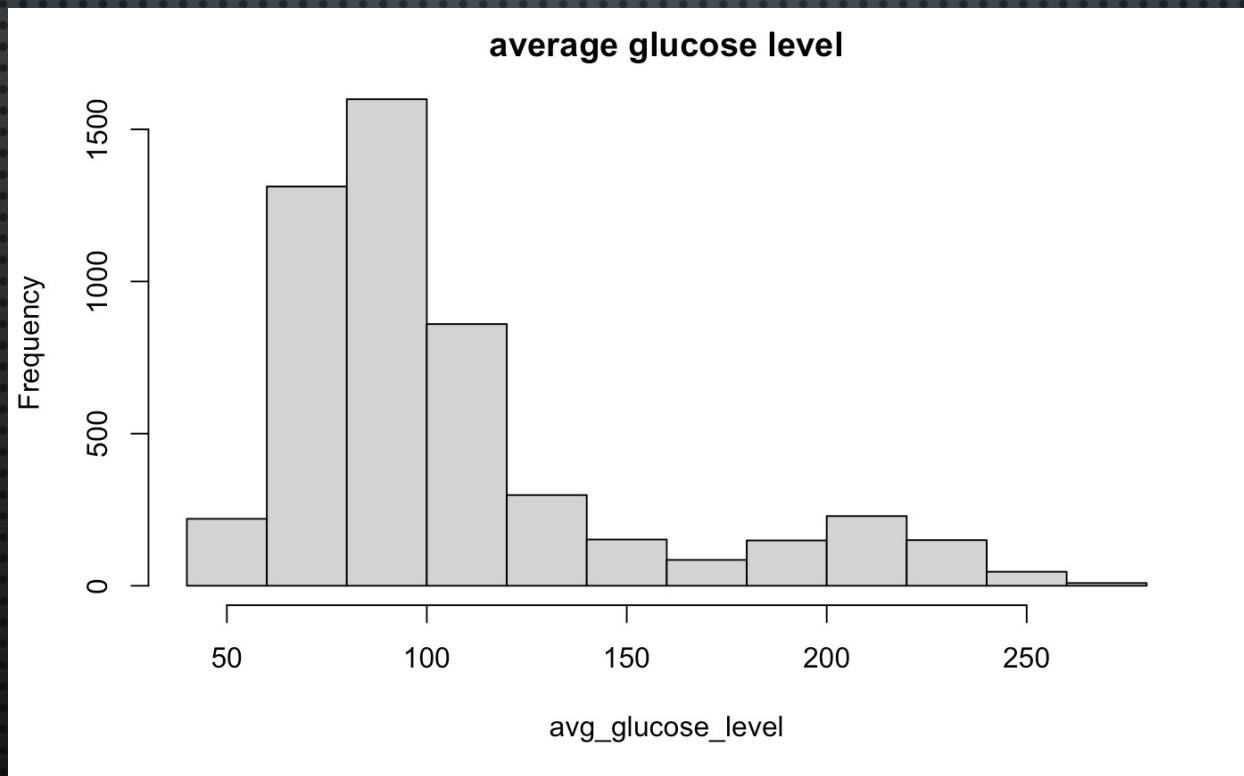
DATA EXPLORATION AND PREPROCESSING

DISTRIBUTION OF NUMERIC VARIABLES



DATA EXPLORATION AND PREPROCESSING

DISTRIBUTION OF NUMERIC VARIABLES



DATA EXPLORATION AND PREPROCESSING

DISTRIBUTION OF NUMERIC VARIABLES

- AVERAGE_GLUCOSE_LEVEL AND BMI ARE RIGHT-SKEWED
- BOTH VARIABLES WERE LOG TRANSFORMED

DATA EXPLORATION AND PREPROCESSING

DATASET PREPARATION

- INITIAL SPLIT OF DATASET INTO FULLTRAIN AND TEST DATASETS
 - 80% FULLTRAIN, 20% TEST
- FULLTRAIN SPLIT INTO TEST AND TRAIN DATASETS
 - 80% TRAING, 20% TEST
- SECOND COPY OF EACH NEW DATASET MADE
 - ONE SET FOR TRAINING WITH CARET
 - ONE SET FOR TRAINING WITH KERAS/TFRUNS

DATA EXPLORATION AND PREPROCESSING

DATASET PREPARATION: CARET

- FULLTRAIN DATASET BALANCED USING ROSE PACKAGE
- FULLTRAIN DATASET CENTERED AND SCALED
 - SAME INSTANCE USED TO CENTER AND SCALE TEST DATASET
- NO CENTERING, SCALING OR BALANCING FOR TRAIN OR VALIDATION DATASETS
- COPIES OF ALL FOUR DATASETS MADE AND ONE-HOT ENCODED USING DUMMYVARS FUNCTION

DATA EXPLORATION AND PREPROCESSING

DATASET PREPARATION: KERAS/TFRUNS

- NNTRAIN DATASET WAS SCALED USING SCALE FUNCTION
 - NNVALIDATION DATASET SCALED USING SAME COLUMNS MEANS AND STANDARD DEVIATIONS
- NNFULLTRAIN DATASET WAS SCALED USING SCALE FUNCTION
 - NNTTEST DATASET SCALED USING SAME COLUMNS MEANS AND STANDARD DEVIATIONS
- COPIES OF ALL FOUR DATASETS MADE AND ONE-HOT ENCODED USING DUMMYVARS FUNCTION

MODEL TRAINING AND EVALUATION

- MULTIPLE MODELS TRAINED AND EVALUATED
 - K-NEAREST NEIGHBORS
 - LOGISTIC REGRESSION WITH ELASTIC NET REGULARIZATION
 - DECISION TREE
 - RANDOM FOREST
 - LINEAR SUPPORT VECTOR MACHINE
 - STOCHASTIC GRADIENT BOOST
 - NEURAL NETWORK
- MODELS WERE TRAINED USING THE CARET PACKAGE FOR ALL MODELS ASIDE FROM THE NEURAL NETWORK

MODEL TRAINING AND EVALUATION

CARET PACKAGE

ALL MODELS TRAINED USING THE CARET PACKAGE WERE TRAINED AS FOLLOWS :

- 10-FOLD CROSS-VALIDATION REPEATED 5 TIMES
- CLASSPROBS WERE SET TO BE CALCULATED
- TWOCLASSSUMMARY AS THE SUMMARYFUNCTION
- ONESE AS THE SELECTIONFUNCTION
- SAMPLING WAS DONE USING THE ROSE METHOD
- DATA WAS CENTERED AND SCALED USING THE PREPROC OPTION
- METRIC FOR EVALUATION WAS ROC

MODEL TRAINING AND EVALUATION

CARET PACKAGE

MODEL TUNING PARAMETERS

Model	Method	Tuning Parameters
K-Nearest Neighbors	knn	<code>k = seq(1,100,by=2)</code>
Logistic Regression with Elastic Net Regularization	glmnet	<code>alpha: seq(0,1,length=10)</code> <code>lambda: 10*seq(-3,3,length=100)</code>
Decision Tree	rpart	<code>cp = seq(0,1,by=0.001)</code>
Random Forest	rf	<code>ntry = c(1,2,4,6,8)</code>
Linear Support Vector Machine	svmLinear	<code>C = 3**(-7:7)</code>
Stochastic Gradient Boost	gbm	Auto-tuned with <code>tuneLength = 10</code>

MODEL TRAINING AND EVALUATION KERAS/TFRUNS

- NEURAL NETWORK WAS TUNED USING THE TFRUNS PACKAGE
- MODEL WAS DESIGNED WITH TWO HIDDEN LAYERS, EACH FOLLOWED BY A DROPOUT LAYER
- LOSS FUNCTION WAS BINARY CROSSENTROPY
- METRIC WAS ACCURACY
- FINAL LAYER USED A SIGMOID ACTIVATION

MODEL TRAINING AND EVALUATION KERAS/TFRUNS

```
MODEL <- KERAS_MODEL_SEQUENTIAL()
MODEL %>%
  LAYER_DENSE(UNITS = FLAGS$NODES1, ACTIVATION = FLAGS$ACTIVATION, INPUT_SHAPE = 15) %>%
  LAYER_DROPOUT(FLAGS$DROP1) %>%
  LAYER_DENSE(UNITS = FLAGS$NODES2, ACTIVATION = FLAGS$ACTIVATION) %>%
  LAYER_DROPOUT(FLAGS$DROP2) %>%
  LAYER_DENSE(UNITS = 1, ACTIVATION = 'SIGMOID')

MODEL %>% COMPILE(
  OPTIMIZER = OPTIMIZER_Adam(LR=FLAGS$LEARNING_RATE),
  LOSS = "BINARY_CROSSENTROPY",
  METRICS=C('ACCURACY'))

MODEL %>% FIT(AS.MATRIX(NNTRAINFEATURES), AS.MATRIX(NNTRAINLABELS), EPOCHS = FLAGS$EPOCHS, BATCH_SIZE =
  FLAGS$BATCH_SIZE, VALIDATION_DATA=LIST(AS.MATRIX(NNVALFEATURES), AS.MATRIX(NNVALLABELS)))
```

MODEL TRAINING AND EVALUATION

KERAS/TFRUNS

TUNING PARAMETERS

Parameter	Value
nodes1	c(20,50,75)
nodes2	c(20,50,75)
drop1	c(0.2,0.4,0.6)
drop2	c(0.2,0.4,0.6)
learning_rate	c(0.01,0.05,0.001,0.0001)
epochs	c(25,50,100)
batch_size	c(10,25,50)
activation	c('relu', 'sigmoid','tanh')

MODEL TRAINING AND EVALUATION

- AFTER EACH MODEL WAS TRAINED AND TUNED IT WAS FIT TO THE APPROPRIATE VALIDATION DATASET
- MODELS WERE EVALUATED ON AUC, SPECIFICITY AND SENSITIVITY

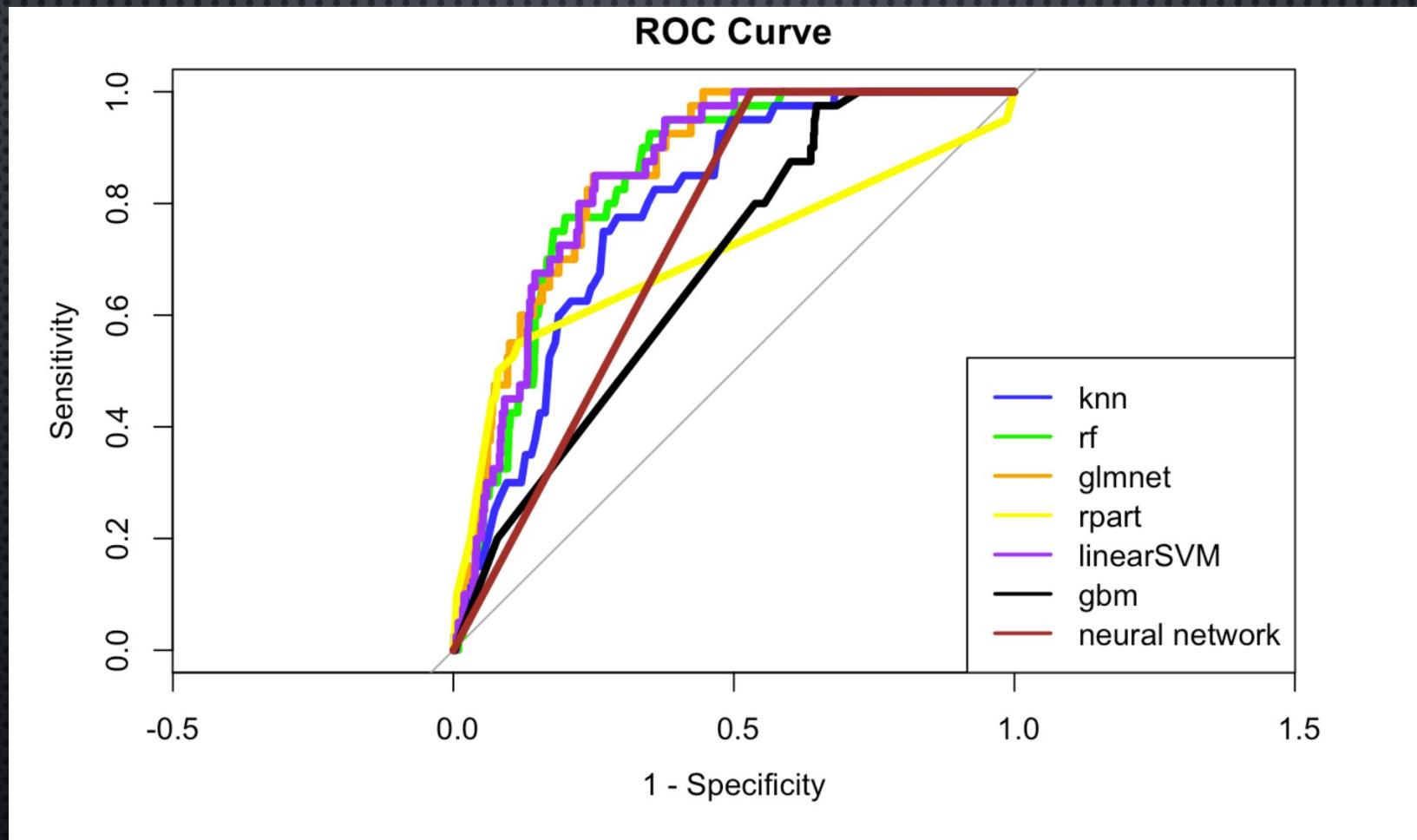
MODEL TRAINING AND EVALUATION

EVALUATION METRICS

	ROC	Spec	Sens
knn	0.7889	0.7773	0.625
glmnet	0.8555	0.7091	0.85
rpart	0.7064	0.9936	0.1
rf	0.8387	0.5354	0.95
linearSVM	0.8488	0.7516	0.825
gbm	0.6796	0.003861	1
neural network	0.7342	0.4646	1

MODEL TRAINING AND EVALUATION

EVALUATION METRICS



MODEL SELECTION

BASED UPON THE EVALUATION METRICS SELECTED, THE LOGISTIC REGRESSION WITH ELASTIC NET
REGULARIZATION WAS DETERMINED TO BE THE BETTER PERFORMING MODEL

FINAL MODEL EVALUATION

- FINAL MODEL WAS BUILT USING THE GLMNET PACKAGE
 - ALPHA AND LAMBDA WERE SET BASED UPON THE RESULTS OF THE CROSS-VALIDATION AND TUNING PROCESS WITH THE CARET PACKAGE
 - MODEL WAS BUILT USING THE FULLTRAIN DATASET
 - MODEL WAS EVALUATED AGAINST THE PREVIOUSLY UNUSED TEST DATASET

FINAL MODEL EVALUATION

```
Confusion Matrix and Statistics

      Reference
Prediction   No  Yes
      No    610    5
      Yes   362   44

      Accuracy : 0.6405
      95% CI   : (0.6103, 0.67)
      No Information Rate : 0.952
      P-Value [Acc > NIR] : 1

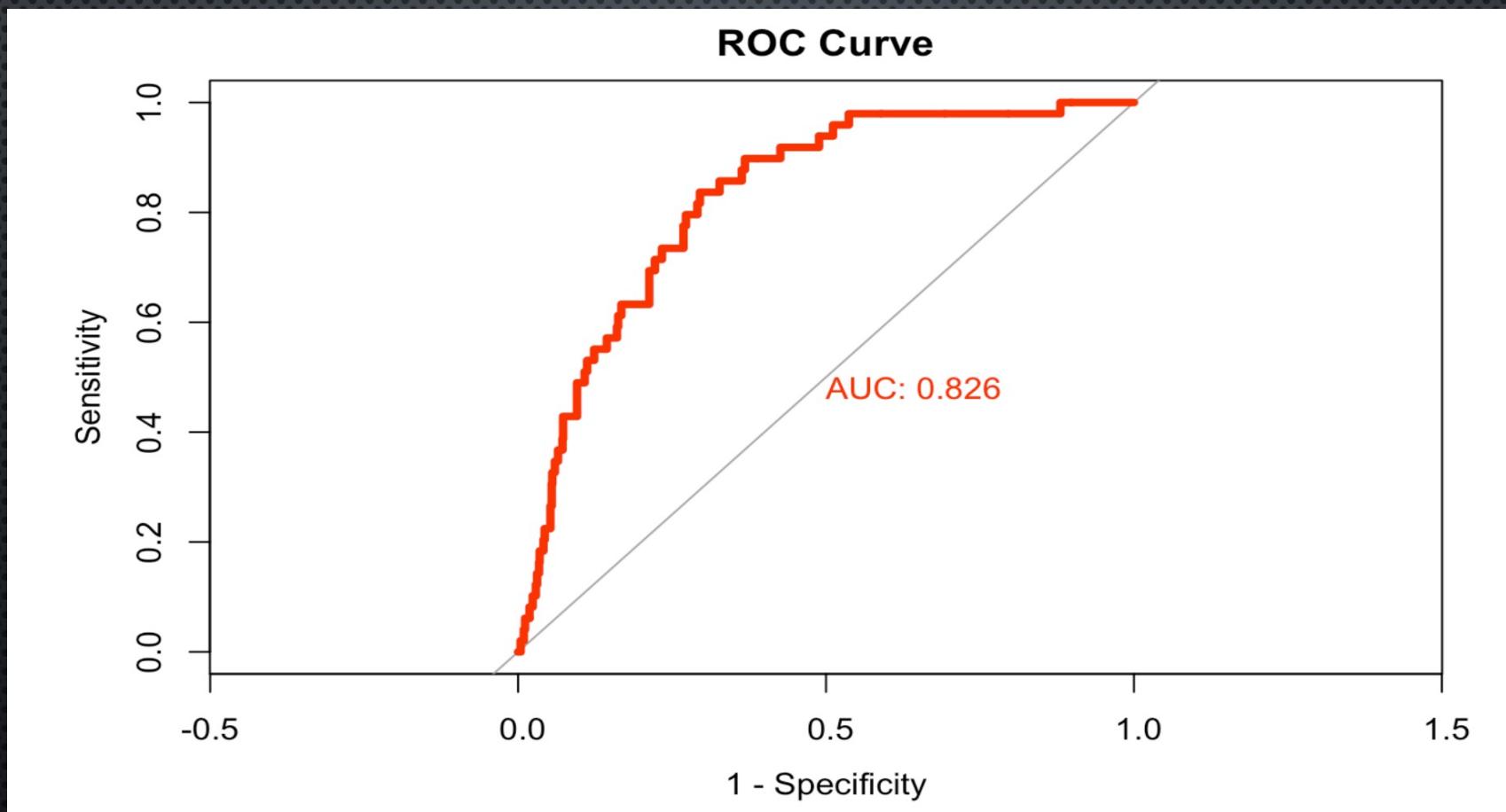
      Kappa   : 0.1179

      Mcnemar's Test P-Value : <2e-16

      Sensitivity   : 0.89796
      Specificity   : 0.62757
      Pos Pred Value : 0.10837
      Neg Pred Value : 0.99187
      Prevalence     : 0.04799
      Detection Rate : 0.04310
      Detection Prevalence : 0.39765
      Balanced Accuracy : 0.76277

      'Positive' Class : Yes
```

FINAL MODEL EVALUATION



FINAL MODEL EVALUATION

VARIABLE SELECTION SUMMARY

```
16 x 1 sparse Matrix of class "dgCMatrix"
                                         s0
(Intercept)          0.03373739
age                  0.28096432
hypertension         0.02208788
heart_disease        0.04157673
ever_married         0.04646855
work_type.children   -0.13217010
work_type.Govt_job   .
work_type.Never_worked .
work_type.Private    .
work_type.Self_employed .
avg_glucose_level   0.02937247
bmi                 .
smoking_status.formerly_smoked .
smoking_status.never_smoked .
smoking_status.smokes   .
smoking_status.unknown .
```

CONCLUSION

IN TRAINING AND EVALUATING A VARIETY OF MACHINE LEARNING MODELS FOR THE PREDICTION OF STROKE, THIS PROJECT FOUND THAT A LOGISTIC REGRESSION MODEL USING ELASTIC NET REGULARIZATION PROVIDED THE BEST BALANCE BETWEEN SENSITIVITY AND SPECIFICITY AS WELL AS PROVIDING A MORE EASILY INTERPRETABLE SOLUTION.