# Group 5

Code for this assignment at
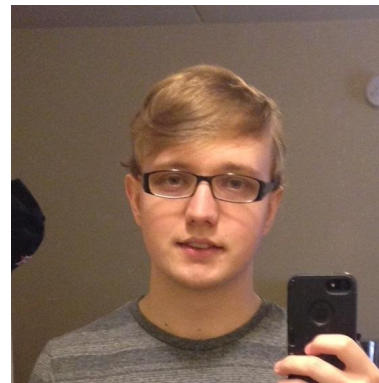https://github.com/TetroVolt/COSC3337Assign2

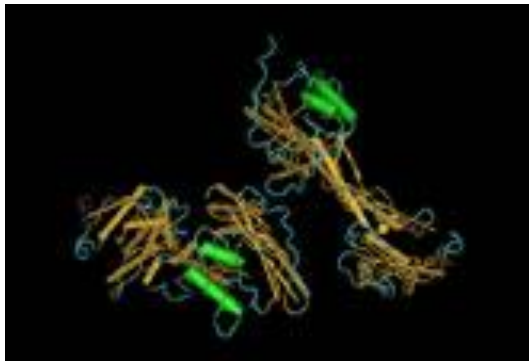# Who we are

Raymond Sutrisno

Niels Moeller

Gal Egozi

Colby Kuhnel

# Data Set I (Molecular Biology Data Set)

| Data Set Characteristics: | Sequential, Domain-Theory | Number of Instances: | 3190 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Categorical | Number of Attributes: | 61 | Date Donated | 1992-01-01 |
| Associated Tasks: | Classification | Missing Values? | No | Number of Web Hits: | 95406 |



- Introns do not code for proteins
- Exons code for proteins

# Data Set II: Preliminary analysis

1. Each Sequence is 60 base pairs long

2. 8 possible base pairs: ('A', 'C', 'D', 'G', 'N', 'R', 'S', 'T')

3. Three Class Labels

    a. Exon-Intron Boundaries (EI) (767 examples)

    b. Intron-Exon Boundaries (IE) (768 examples)

    c. Neither (N) (1655 examples)

# Data Preprocessing I (RAW data)

# Algorithms considered

Required:

1. Neural Networks (Gal Egozi)
2. Support Vector Machines (Colby Kuhnel)

Chosen freely:

1. Decision Tree (Niels Moeller)
2. Random Forest (Raymond Sutrisno)

# Data Processing II (considerations)

1. NN and SVM are "One vs. All" algorithms
2. Features are categorical in nature
3. Need One Hot Encoding for both class labels and features

# Data Processing III (procedure)

1.  One hot encode features
    a.  Each Sequence is 60 base pairs long and categorical
    b.  60 features
    c.  8 base pair categories ('A', 'C', 'D', 'G', 'N', 'R', 'S', 'T')
    d.  60 * 8 = 480 columns total after one hot encoding
2.  One hot encode classes for algorithms that need it
    a.  Three classes: 'EI' , 'IE', 'N'

# Training Procedure

1. Gridsearch model parameters to tune for the best model using stratified 10 fold cross validation
2. Report accuracies (CV mainly, training and test)
3. Optional metrics considered
   a. Confusion Matrices
   b. Learning Curves

# Implementation Info

1. Language used: Python
2. Utility Libraries:
   a. Pandas (data processing)
   b. Numpy (data processing)
   c. Scikit Learn (data processing)
3. Machine Learning Libraries
   a. Scikit Learn    (SVM, Random Forest, Decision Tree)
   b. Keras (Neural networks)

# Results

# Random Forest (trees based on C4.5)

# Random Forest I (model parameter space available)

1. Criterion ['gini', 'entropy']
2. Max_depth
3. Min sample splits
4. Min impurity decrease
5. ...

# Random Forest I (model parameter space used)

1. 'N_estimators':[25, 50, 100, 150]
2. 'max_depth': [2, 4, 8, 16, 32, 64]

# Random Forest II (Raw Output) (20 % Test)

```
Dataset characteristics:
  Number of examples in the dataset = 3190
  Number of examples reserved for test set = 638
  Number of examples reserved for training via 10 fold CV = 2552
  Class Distribution Ratio (N : EI : IE) = 2 : 1 : 1
  N features : 60, all categorical (DNA base pairs in 60 base pair long sequence)


  Grid Search parameter space for Random Forest =
    {'n_estimators': [25, 50, 100, 150], 'max_depth': [2, 4, 8, 16, 32, 64]}
  best estimator parameters found = {'max_depth': 32, 'n_estimators': 150}
  best estimator mean training score   = 0.9995646116203684
  best estimator mean validation score = 0.9114420062695925
  best estimator test score            = 0.9247648902821317
```

# Random Forest III (Raw Output) (50 % Test)

```
Dataset characteristics:
  Number of examples in the dataset = 3190
  Number of examples reserved for test set = 1595
  Number of examples reserved for training via 10 fold CV = 1595
  Class Distribution Ratio (N : EI : IE) = 2 : 1 : 1
  N features : 60, all categorical (DNA base pairs in 60 base pair long sequence)


  Grid Search parameter space for Random Forest =
     {'n_estimators': [25, 50, 100, 150], 'max_depth': [2, 4, 8, 16, 32, 64]}
  best estimator parameters found = {'max_depth': 16, 'n_estimators': 100}
  best estimator mean training score   = 0.9992337891743421
  best estimator mean validation score = 0.8915360501567398
  best estimator test score            = 0.8984326018808777
```

# Random Forest IV (Raw Output) (70 % Test)

```
Dataset characteristics:
  Number of examples in the dataset = 3190
  Number of examples reserved for test set = 2233
  Number of examples reserved for training via 10 fold CV = 957
  Class Distribution Ratio (N : EI : IE) = 2 : 1 : 1
  N features : 60, all categorical (DNA base pairs in 60 base pair long sequence)


  Grid Search parameter space for Random Forest =
    {'n_estimators': [25, 50, 100, 150], 'max_depth': [2, 4, 8, 16, 32, 64]}
  best estimator parameters found = {'max_depth': 16, 'n_estimators': 100}
  best estimator mean training score   = 1.0
  best estimator mean validation score = 0.8610240334378265
  best estimator test score            = 0.8669950738916257
```

# Random Forest Best Parameters

| Train : Test Ratio | Max Depth | N Estimators |
|---|---|---|
| 80:20 | 32 | 150 (MAX) |
| 50:50 | 16 | 100 |
| 30:70 | 16 | 100 |

# Random Forest Mean Scores

| Train : Test Ratio | Mean Train | Mean CV | Test Score |
|---|---|---|---|
| 80:20 | 0.9995 | 0.9114 | 0.9247 |
| 50:50 | 0.9992 | 0.8915 | 0.8984 |
| 30:70 | 1.0000 | 0.8610 | 0.8669 |

# Random Forest Mean Scores

| Train : Test Ratio | Mean Train | Mean CV | Test Score | Max Depth | N Estimators |
|---|---|---|---|---|---|
| 80:20 | 0.9995 | 0.9114 | 0.9247 | 32 | 150 |
| 50:50 | 0.9992 | 0.8915 | 0.8984 | 16 | 100 |
| 30:70 | 1.0000 | 0.8610 | 0.8669 | 16 | 100 |

# Decision Tree (SKLearn based on C4.5)

# Decision Trees 80/20

Dataset characteristics:

Number of examples in the dataset = 3190

Number of examples reserved for test set = 638
Number of examples reserved for training via 10 fold CV = 2552
Class Distribution Ratio (N : EI : IE) = 2 : 1 : 1
N features : 60, all categorical (DNA base pairs in 60 base pair long sequence)

Grid Search parameter space for Random Forest = {'max_depth': range(3, 20)}
best estimator parameters found = {'max_depth': 6}
best estimator mean training score   = 0.9573321284086651
best estimator mean validation score = 0.9455329153605015
best estimator test score            = 0.9404388714733543

# Decision Tree Scores

| Train : Test Ratio | Mean Train | Mean CV | Test Score |
|---|---|---|---|
| 80:20 | 0.9570 | 0.9447 | 0.9341 |
| 50:50 | 0.9484 | 0.9347 | 0.9373 |
| 30:70 | 0.9564 | 0.9362 | 0.9305 |

# Decision Tree Best Parameters

| Train : Test Ratio | Max Depth (3-20) | Min samples split [3,5,25] |
|---|---|---|
| 80:20 | 5 | 25 |
| 50:50 | 5 | 3 |
| 30:70 | 6 | 3 |

# Resulting Decision Tree

# SVM

# Support Vector Machines (20% Test)

```
Grid Search parameter space for SVM = [{'C': [100, 1000, 10000], 'gamma': [0.01, 0.001], 'kernel':
['poly', 'rbf'], 'degree': [2, 3, 4]}]
    best estimator parameters found = {'C': 100, 'degree': 3, 'gamma': 0.01, 'kernel': 'poly'}
    best estimator mean training score      = 0.9996081466330754
    best estimator mean validation score    = 0.9678683385579937
    best estimator test score               = 0.9702194357366771
    Process Time                            = 8.436292 Minutes
```

# Support Vector Machines (50% Test)

```
Grid Search parameter space for SVM = [{'C': [100, 1000, 10000], 'gamma': [0.01, 0.001], 'kernel':
['poly', 'rbf'], 'degree': [2, 3, 4]}]
    best estimator parameters found = {'C': 100, 'degree': 3, 'gamma': 0.01, 'kernel': 'poly'}
    best estimator mean training score      = 0.9994427513515088
    best estimator mean validation score    = 0.9605015673981191
    best estimator test score               = 0.9661442006269593
    Process Time                            = 3.845329 Minutes
```

# Support Vector Machines (70% Test)

```
Grid Search parameter space for SVM = [{'C': [100, 1000, 10000], 'gamma': [0.01, 0.001], 'kernel':
['poly', 'rbf'], 'degree': [2, 3, 4]}]
    best estimator parameters found = {'C': 100, 'degree': 3, 'gamma': 0.01, 'kernel': 'poly'}
    best estimator mean training score      = 1.0
    best estimator mean validation score    = 0.9540229885057471
    best estimator test score               = 0.961486789072996
    Process Time                            = 1.603256 Minutes
```

# Support Vector Machines Mean Scores

| Train : Test Ratio | Mean Train | Mean CV | Test Score |
|---|---|---|---|
| 80:20 | 0.9996 | 0.9679 | 0.9702 |
| 50:50 | 0.9994 | 0.9605 | 0.9661 |
| 30:70 | 1.0 | 0.9540 | 0.9615 |

# Neural Network

# Neural Network

- One layer, 50 nodes
- Scores 95.52±0.74% on stratified cross validation
- 20% test
  - 99.61% Train
  - 94.36% Test
- 50% test
  - 99.44% Train
  - 93.67 Test
- 70% test
  - 100% train
  - 91.49% test
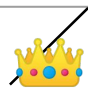
# Performance Comparison

# Top Dog for 20% Test

| Algorithm | Best Test Score | Accompanying CV Score | Accompanying Train Score |
| --- | --- | --- | --- |
| Random Forest | 0.9247 | 0.9114 | 0.9995 |
| Decision Tree | 0.9373 | 0.9347 | 0.9484 |
| SVM 🐕🔥💯 | 0.9702 👑 | 0.9679 👑 | 0.9996 👑 |
| Neural Network | 0.9436 | 0.9552 | 0.9961 |

# Top Dog for 50% Test

| Algorithm | Best Test Score | Accompanying CV Score | Accompanying Train Score |
|-----------|-----------------|------------------------|---------------------------|
| Random Forest | 0.8984 | 0.8915 | 0.9992 |
| Decision Tree | 0.9373 | 0.9347 | 0.9484 |
| SVM 🐕🔥💯 | 0.9661 👑 | 0.9605 👑 | 0.9994 👑 |
| Neural Network | 0.9436 | 0.9552 | 0.9961 |

# Top Dog for 70% Test

| Algorithm | Best Test Score | Accompanying CV Score | Accompanying Train Score |
|---|---|---|---|
| Random Forest | 0.8669 | 0.8610 | 1.0000 👔 |
| Decision Tree | 0.9373 | 0.9347 | 0.9484 |
| SVM 🐕🔥💯 | 0.9615 👑 | 0.9540 | 1.0000 👔 |
| Neural Network | 0.9436 | 0.9552 👑 | 0.9961 |

# What could have been done differently?

- Investigated Gini importances of features
- Investigated other pattern recognition neural network architectures
- Investigated learning curves of models to better understand generalization behaviors
- Investigated Confusion matrices to see if class distribution was an issue that skewed different types of errors
- Used stratified K fold rather than random sample K fold
- Played more with parameters in general

# Thank you!

End of presentation