# Random Forest Classifier

Damilola O.Said
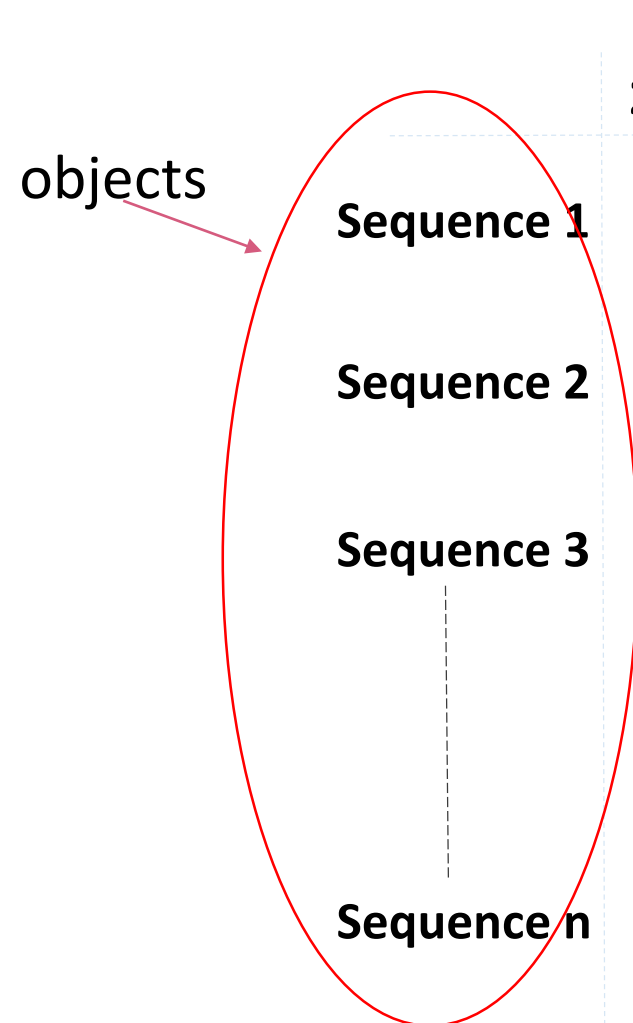
Advanced Bioinformatics

20th of April, 2017

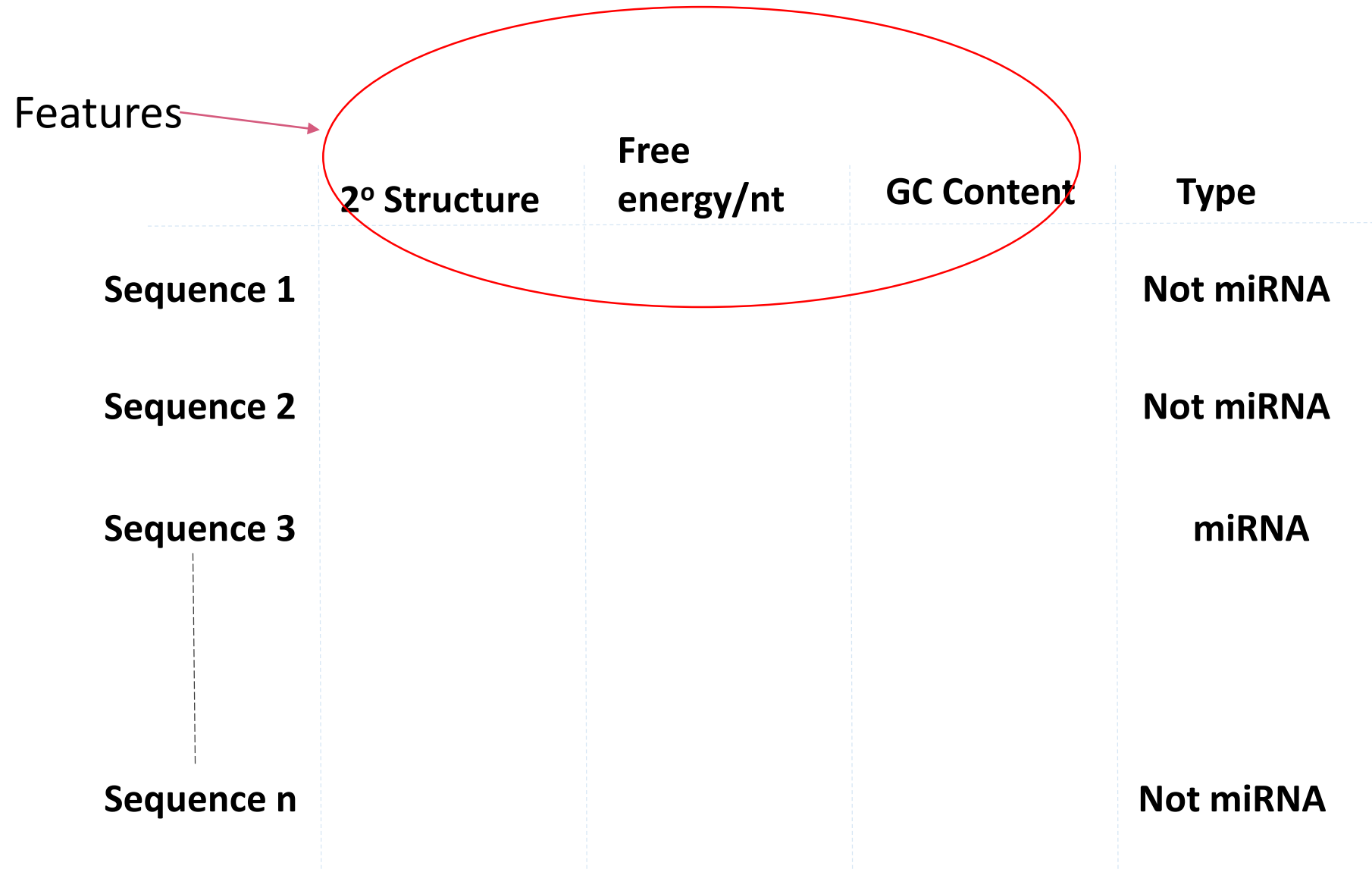# DATA STRUCTURE

objects →

| | 2° Structure | Free energy/nt | GC Content | Type |
|---|---|---|---|---|
| **Sequence 1** | | | | **Not miRNA** |
| **Sequence 2** | | | | **Not miRNA** |
| **Sequence 3** | | | | **miRNA** |
| **Sequence n** | | | | **Not miRNA** |

| | 2° Structure | Free energy/nt | GC Content | Type |
|---|---|---|---|---|
| Features | | | | |
| Sequence 1 | | | | Not miRNA |
| Sequence 2 | | | | Not miRNA |
| Sequence 3 | | | | miRNA |
| Sequence n | | | | Not miRNA |

# DATA STRUCTURE

Targets/Labels

| | 2° Structure | Free energy/nt | GC Content | Type |
|---|---|---|---|---|
| Sequence 1 | | | | Not miRNA |
| Sequence 2 | | | | Not miRNA |
| Sequence 3 | | | | miRNA |
| Sequence n | | | | Not miRNA |

# SPLITTING THE DATA

DATA (SAMPLES)

TRAINING SET (80%)

TRAINING SET FEATURES

TRAINING SET LABELS

VALIDATION SET (20%)

VALIDATION SET FEATURES

VALIDATION SET LABELS

# TRAINING THE CLASSIFIER

| TRAINING SET FEATURES | | VALIDATION SET FEATURES |

Fit model to training data → **RANDOM FOREST CLASSIFIER** → Validate → **ACCURACY SCORE**

| TRAINING SET LABELS | | VALIDATION SET LABELS |

OPTIMIZE

# THE CODE: IRIS DATASET

```python
7
8  import pandas as pd
9  from pandas.tools.plotting import scatter_matrix
10 import matplotlib.pyplot as plt
11 from sklearn import model_selection
12 from sklearn import metrics
13 from sklearn.cross_validation import train_test_split
14 from sklearn.ensemble import RandomForestClassifier
15 from sklearn.metrics import accuracy_score
16 import skfuzzy as fuzz
17 import numpy as np
18
19
20
21 #url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
22
23 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
24
25
26 dataset = pd.read_csv('C:\Users\doyewolesaid1\Downloads/iris.data.txt', names=names)
27
28 #80:20 slpit for test and validation
29 #train data excludes class
30 #train_features includes ONLY class
31 train_data= dataset.values
32 train_features = train_data[:, :4]
33 train_target = train_data[:, 4]
34
35 seed =10
36 train_x, test_x, train_y, test_y = train_test_split(train_features, train_target, test_size=0.20, random_state=seed)
37
38 classify = RandomForestClassifier(n_estimators=100)
39
40 classify = classify.fit(train_x, train_y)
41 predict_y = classify.predict(test_x)
42
43 print ("Accuracy = %.2f" % (accuracy_score(test_y, predict_y)))
```

**Fisher's *Iris* Data**

| Sepal length ⬦ | Sepal width ⬦ | Petal length ⬦ | Petal width ⬦ | Species ⬦ |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | I. setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | I. setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | I. setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | I. setosa |
| 5.0 | 3.6 | 1.4 | 0.2 | I. setosa |
| 5.4 | 3.9 | 1.7 | 0.4 | I. setosa |
| 4.6 | 3.4 | 1.4 | 0.3 | I. setosa |
| 5.0 | 3.4 | 1.5 | 0.2 | I. setosa |
| 4.4 | 2.9 | 1.4 | 0.2 | I. setosa |
| 4.9 | 3.1 | 1.5 | 0.1 | I. setosa |
| 5.4 | 3.7 | 1.5 | 0.2 | I. setosa |
| 4.8 | 3.4 | 1.6 | 0.2 | I. setosa |
| 4.8 | 3.0 | 1.4 | 0.1 | I. setosa |

(150 samples, 4 features + label)

# REFERENCES

1. Brownlee, J. http://machinelearningmastery.com/machine-learning-in-python-step-by-step/
2. Martin, D. http://nbviewer.jupyter.org/github/donnemartin/data-science-ipython-notebooks/blob/master/kaggle/titanic.ipynb