

# 17. Random Forest:-

M T W T F S

DATE: \_\_\_\_\_

- Random forest have the ability to greatly increase the performance based on expanding ideas from the Decision Tree.
- Random Forests are known as ensemble learners, since they rely on an ensemble of models (multiple decision tree)

## Section Overview

- Motivation and History
- Hyperparameter
- Random Forest Classification
- Random Forest Regression
  - Comparing many regression models (SVR, DTR, etc)

## 02 Random Forest

### Theory & Intuition

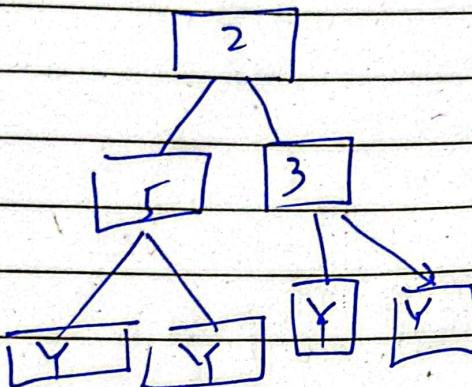
Question we can ask

- Why not just use Decision Tree?
- A motivation behind ~~Decision Tree~~ <sup>Random Forest</sup> and how they improve decision tree -

.

- Imagine a dataset

|   |   |   |   |   |     |
|---|---|---|---|---|-----|
| 1 | 2 | 3 | 4 | 5 | $y$ |
|   |   |   |   |   |     |
|   |   |   |   |   |     |
|   |   |   |   |   |     |
|   |   |   |   |   |     |



Decision Tree

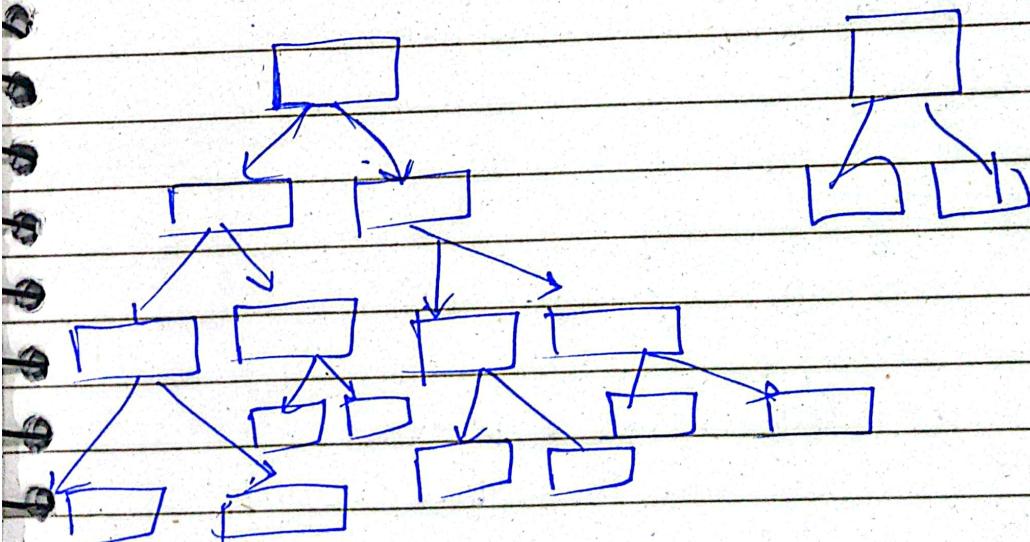
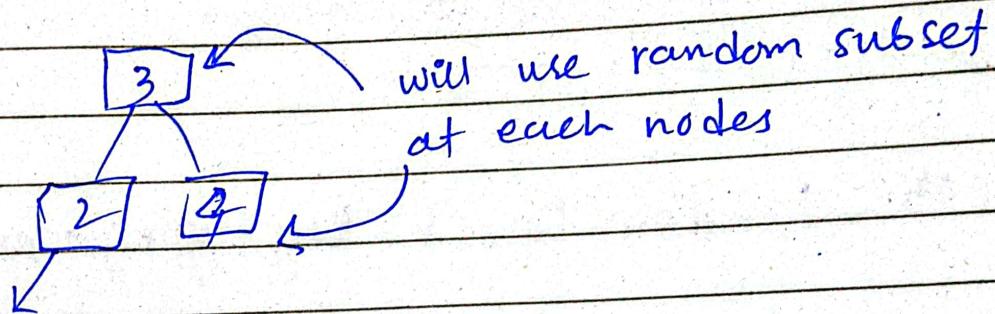
we will always use certain feature every time we try to split mean we will have same node in decision tree unless we

change hyperparameters so some features might be neglected

- No guarantee of using all features
- Root node will always be same - (using same) rule
- Root node will have a huge influence on tree.
- We could try adjusting rules such as :
  - splitting criterion (Information Gain)
  - Minimum Gini Impurity Decrease
  - Setting Depth limits
  - Limit on number of terminal leaf nodes.
- However even with all these added hyperparameter adjustments the single decision tree is still limited!

- Single feature for root node
- splitting criteria can lead to some features not being used.
- potential for overfitting the data
- 1995: Tim Kam Ho present Random Decision Forest at 3rd International Conference on Document analysis and Recognition in Montreal.

↳ create subsets of randomly picked features at each potential split :



$$\text{Tree 1: } Y = 0 \quad | \quad \text{Tree 2: } Y = 0 \quad | \quad \text{Tree 3: } Y = 1$$

$$(0, 0, 1) \quad | \quad Y = 0$$

M T W T F S

We can also see probability of  
the classification.

DATE:

- Ho and other researcher shows that adding the random features subspace constraint could allow trees to grow much deeper without causing overfitting.
- Mathematical explanation of this was shown by Eugene Kleinberg's theory of Stochastic discrimination.
- 2001 : Leo Breiman published Random Forests in the journal Machine Learning.
- Introduces bootstrapping samples and out of bag error.

### 003 Random Forest - Key Hyperparameters

DecisionTreeClassifier(\*, criterion='gini',

)

RandomForestClassifier(n\_estimators,  
max\_features='auto',  
bootstrap=True,  
oob\_score=False)

Random Forest Hyperparameter different from:

- Number of Estimators → How many decision tree used
- Number of features → How many features to include in each subset

• Bootstrap sample

↳ Out-of-Bag Error

↳ Allow for bootstrap sampling of each training subset of features?

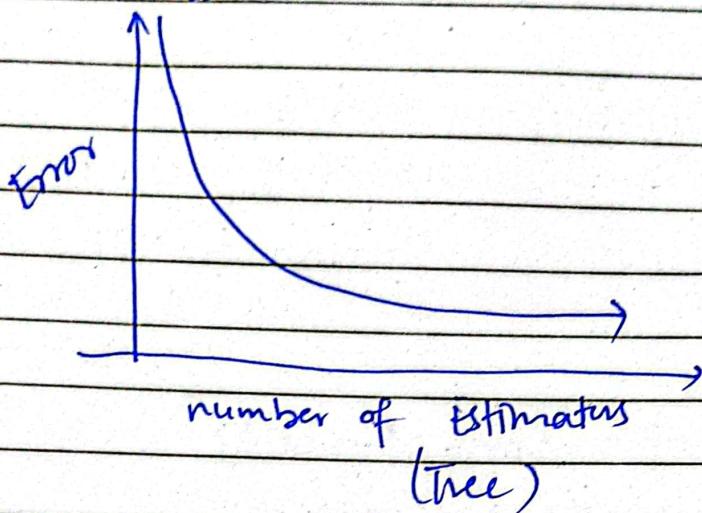
↳ Calculate OOB error during training?

## 004 Random Forests Hyperparameters Estimator and Features

- Number of Estimators:
  - Intuitively, we know the more decision tree, the more opportunities to learn from variety of features subset combinations.
  - Is there a limit to adding more trees?
  - Is there a danger of overfitting?
- How to choose number of trees?
  - Reasonable Default values: 100
  - publication suggest 64-128 trees.
  - cross validate a grid search of trees.

- plot Error versus number of trees  
(similar to elbow method of KNN)
- Should notice diminishing error reduction after some N trees.

Error vs Trees :



- After a certain number of trees, two things that can occur:
  - Different random selection don't reveal any more information
    - Trees become highly correctly
  - ~~Different~~ Different random selections are simply duplicating trees that have already been created.
- This allows us to be quite lenient in setting the number of estimators hyperparameters, as overfitting is of minimal concern.

\* Now let's us discuss how to choose number of features to randomly select at each split.

- Number of Feature
  - how many features to include in each subset when splitting at a node?
  - original publication suggested subset of  $\log_2(N+1)$  random features in subset given a set of N total features.
- From Leo Breiman's official page on Random Forests:
  - An interesting difference between regression and classification is that the correlation increase quite slowly as the number of features used increase,".
- current suggested convention is  $\sqrt{N}$  in the subset given N features
- later suggested by Breiman indicated  $N/3$  may be more suitable for regression task, typically larger

M T W T F S

DATE: \_\_\_\_\_

than  $\sqrt{N}$ ).

- ISLR indicates these to be as tuning parameter with  $\sqrt{N}$  as a good starting point
  - o All priors were chosen using empirical method -
- Start with 100 as default estimated and then grid search
- Start with  $\sqrt{N}$ , grid search for other possible values ( $N/3$ )

## OOS Random Forest - Bootstrapping and out-of-Bag Error

(N DIVIDEDS)

DATE: \_\_\_\_\_

- Random Forest Hyperparameters:
  - Bootstrap Samples
    - Allow for bootstrap sampling of each training subset of features?
  - Out-of-Bag Error
    - calculate OOB error during training ?

What is Bootstrapping?

- A term used to describe "random sampling with replacement".
- Let see a quick example of given a set of letters.

|   |            |
|---|------------|
| A | 10 samples |
| B |            |
| C |            |
| D |            |
| E |            |

→ [ C | A | D | C | E | D | A | B | B | C ]

Bootstrap allow for use of 1 value/feature in multiple places

default = True

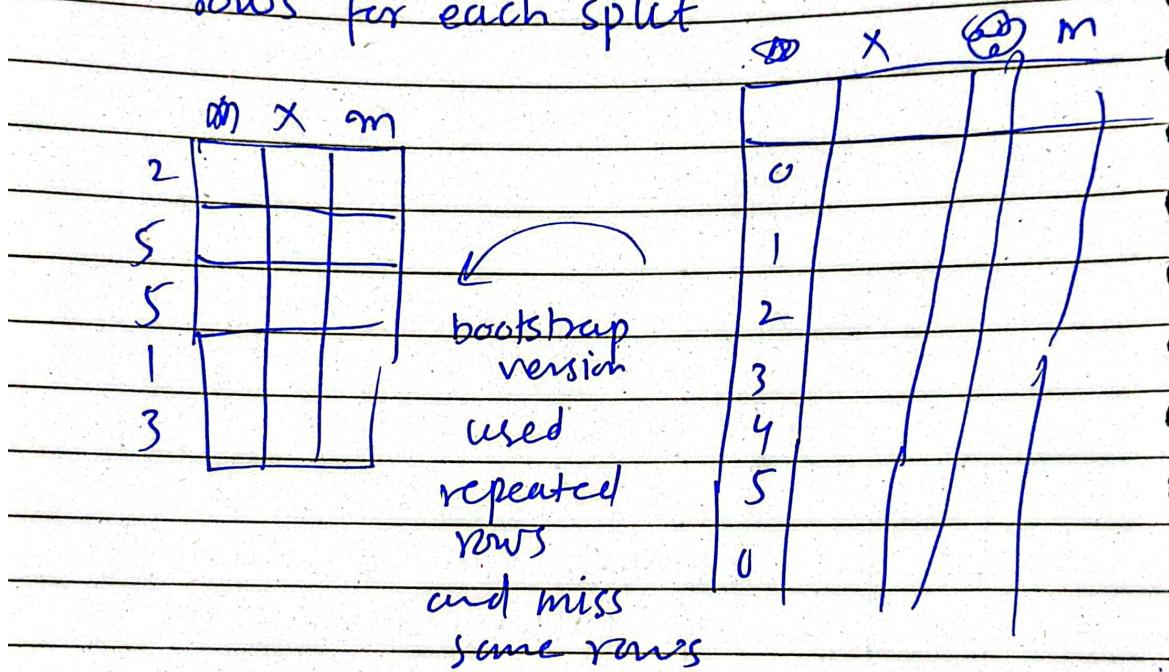
Bootstrapping is using same rows multiple time and not using same row in Decision tree to have more non correlation between the decision trees,

M T W T F S

DATE:

## Bootstrapping in Random Forest:

- Recall for each split we are randomly selecting a subset of features.
- This randomly subset of features help create more diverse trees that are not correlated -
- To further differentiate trees, we could bootstrap a selection of rows for each split



- Bootstrap can be set to False during training (it is True by default)
- Bootstrapping is yet another hyperparameter meant to reduce correlation between trees, since trees are then trained on different subsets of features columns and data rows.

## Calculate out-of-bag error referring to bagging

What is Bagging?

Recall to ~~construct~~ actually use a Random forest, we use bootstrapped data and then calculate a prediction based on the aggregated prediction of the trees: (average)

- Classification: Most Voted Y class
- Regression: Average predicted Ys

- If we performed bootstrapping when building out trees, this means that for certain trees, certain rows of data were not used in training. Those not used rows are called out-of-bag samples.

- Not used for constructing some trees.
- We could use these to get performance test metrics on trees that did not use these rows!

Then we can use that out-of-bag sample to find OOB error on each decision tree.

M T W T F S

DATE:

M.I

fn

fe

- Note the OOB score does not effect training. It's just calculate OOB error when bootstrapping is true mean. - OOB setting OOB to true doesn't change Decision Trees but Setting bootstrapping to True does change the structure of decision trees.
- Note that OOB score is also limited to not using all the trees in the random forest, it can only be calculated on trees that did not use the OOB data.
- Due to not using of the entire random forest, the default value of OOB score hyperparameters is set to False.

(Random Forest Classifier)

## OOB - Coding Classification - RFC

### part 1

```
df = pd.read_csv("pgm.csv")
```

```
df = df.dropna()
```

```
df.head()
```

```
X = pd.get_dummies(df.drop(["species"], axis=1), drop_first=True)
```

```
y = df["species"]
```

```
train-test-split(X, y, test_size=0.3, random_state=101)
```

```
from sklearn.ensemble import RandomForestClassifier
```

```
rfc = RandomForestClassifier(n_estimators=10)
```

max\_features = {  
 'auto',  
 'sqrt',  
 "log2"}

random\_state=101}

int = max

float = fraction

int(max\_features \* n\_features)

max\_features

auto = sqrt(N)

sqrt =  $\sqrt{N}$

$\log_2 N = \Theta(\log_2(N))$

```
rfc.fit(X-train, y-train)
```

```
preds = rfc.predict(X-test)
```

# predict

```
plot_confusion_matrix(model, X-test, y-test)
```

```
print(classification_report(y-test, y-preds))
```

# 007 Coding classification with Random Forest - part two

DATE: \_\_\_\_\_

NUDIBUL  
rf

```
df = pd.read_csv("ru data-banknote.csv")  
df.head()
```

```
sns.pairplot(df, hue = "Class")
```

```
X = df.drop('Class', axis=1)
```

```
y = df['Class']
```

```
from sklearn
```

```
train-test-split(x, y, test_size=0.15,  
random_state=101)
```

```
from sklearn.model_selection import GridSearchCV
```

```
from sklearn
```

```
n_estimators = [64, 100, 128, 100]
```

```
max_features = [2, 3, 4]
```

```
bootstrap = [True, False]
```

```
oob_score = [True, False],
```

```
params = grid = {
```

```
}
```

```
rfc = randomForestClassifier()
```

```
grid = GridSearchCV(rfc, param_grid)
```

```
grid.best_params_
```

rfc = RandomF C ( max\_features=2,  
n\_estimators=200,  
oob\_score=True)

rfc.fit(x-train, y-train)

rfc.oob\_score → True or False

rfc.oob\_score\_

predict = rfc.predict(x-test)

classification-report()

confusion\_matrix()

error = []

missclassification = []

for n in range(1, 200):

rfc = RandomForestClassifier(n= )

rfc.fit(x-train, y-train)

err = 1 - accuracy(y-test, y-preds)

n\_missed = np.sum(preds != y-test)

errors.append(err)

missclassification.append(n\_missed)

plt.plot(range(1, 200), errors)

008

# Coding Regression with Random Forest - part one

M T W T F S

DATE:

NLIDW

## Discussion about Dataset

def

## 1 Coding Regression with part 1

```
df = pd.read_csv("rock-rebound")
```

```
df.columns = ['signal', 'Density']
```

```
sns.scatterplot(x='signal', y='Density', data=  
=df)
```

```
x = df['signal'].values.reshape(-1, 1) error  
if not  
done
```

```
y = df['density']
```

```
X_train, y_test, u
```

```
Lr = model = LinearRegression()
```

```
Lr_model.fit(X_train, y_train)
```

```
Lr_preds = Lr_model.predict(X_test)
```

```
mean_absolute_error(y_test, Lr_preds)
```

plots scatter plot

```
signal_range = np.arange(0, 100)
```

```
signal_pred = Lr_model.predict(.signal_range,
```

```
reshape(-1, 1)
```

# → 10 Regression - part 2

NAME \_\_\_\_\_

DATE: \_\_\_\_\_

```
def runmodel(model,X-train,y-train,X-test,y-test)  
    model.fit(X-train,y-train)  
    pred = model.predict(X-test)  
    rmse = np.sqrt(meansq(y-test,pred))  
    mae =  
    print(f'MAE: {mae}')  
    print(f'RMSE: {rmse}')
```

II plot Result

Signal range = np.arange(0,100)

sigtpred = model.predict(sig.range  
(-1,1))

plt.figure()

sns.scatterplot(x=

plt.plot

```
from sklearn.pipeline import make_pipeline
```

```
from sklearn.preprocessing import PolynomialFe
```

```
pipe = make_pipeline(PolynomialFeatures(degree=2))
```

```
LinearRegression())
```

```
run_model(pipe,X-train,y-train,X-test,y-test)
```

# 11 - Building regression with Random Forest Regressor

M T W T F S

DATE:

# 1

from sklearn.neighbors import KNeighborsRegressor

K\_values = [1, 5, 10]

for n in K\_values:

model = KNeighborsModel(n\_neighbors=n)

runModel(model, X\_train, y\_train, Xtest)

model

= DecisionTree

SVR from sklearn.svm import SVR

SVR = SVR()

paramgrids = {  
 'C': [0.01, 0.1, 1, 5, 10, 100, 1000],  
 'gamma': ['auto', 'scale']}

grid = GridSearchCV(p\_SVR, paramgrids)

from sklearn.ensemble import RandomForestRegressor

RandomForestRegressor(n\_estimators=10)