

16- Tree Based Method - Decision Tree Learning

MITRA

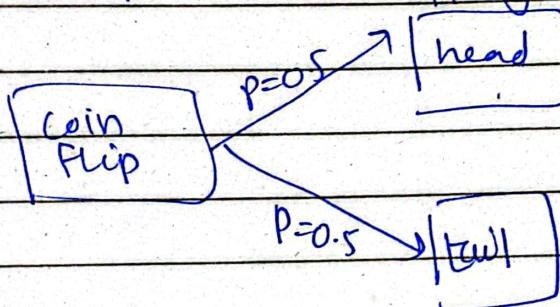
DATE: 11/11/2023

Three Main methods

- Decision Tree
- Random forest
- Boosted Trees

2- Decision Tree History

The general term "decision tree" can refer to a flowchart mapping out outcomes:



check the slides for more information-

003 - Decision Tree - Terminology

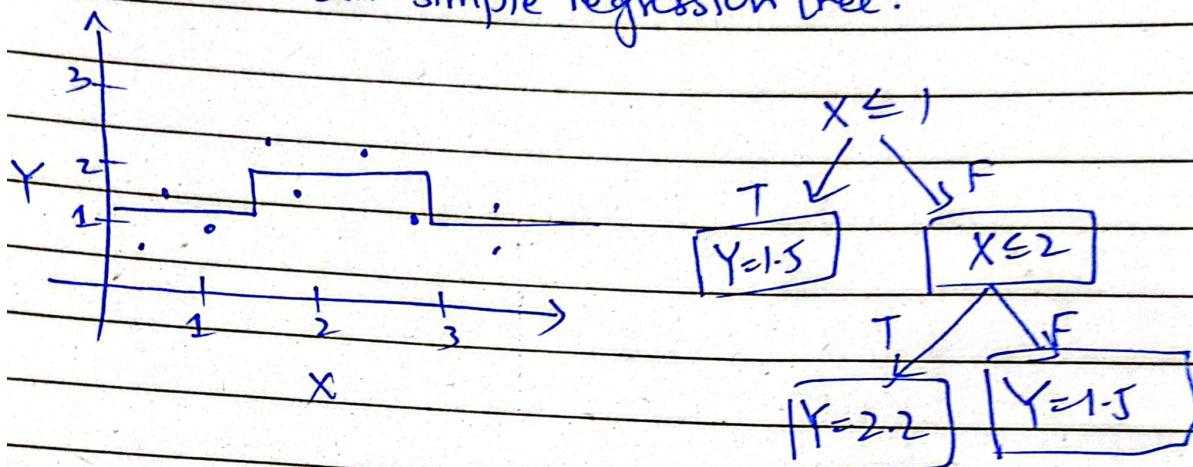
Decision Tree Basics

M T W T F S

DATE: _____

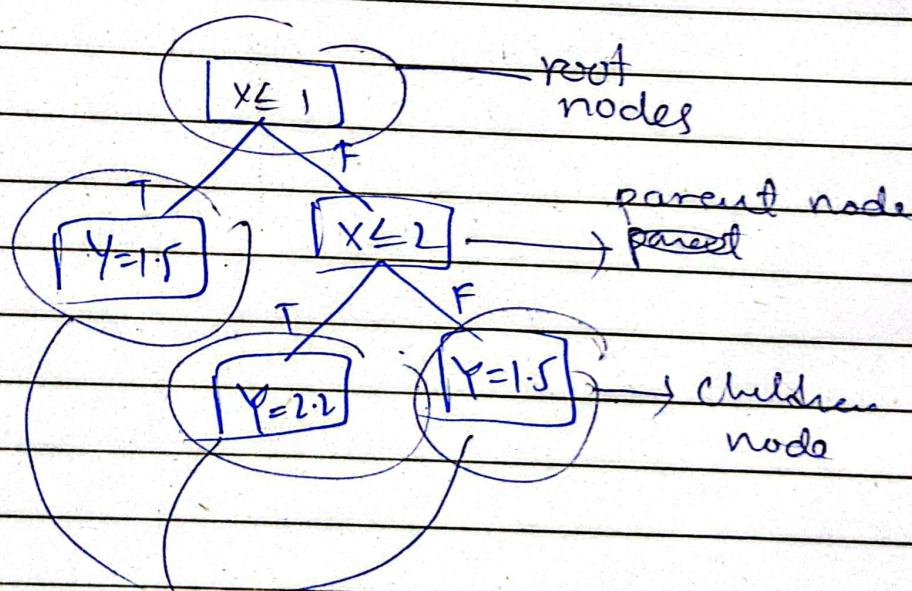
- To begin understanding a decision tree, we first need to review some terminology about the dataset decision tree components.

- Recall our simple regression tree:



- splitting

notes nodes



leaf (Terminal)
nodes

Stops Splitting

pruning:-

stop
start the mode
data from further
Splitting

4 - Decision Tree - Understanding Gini impurity

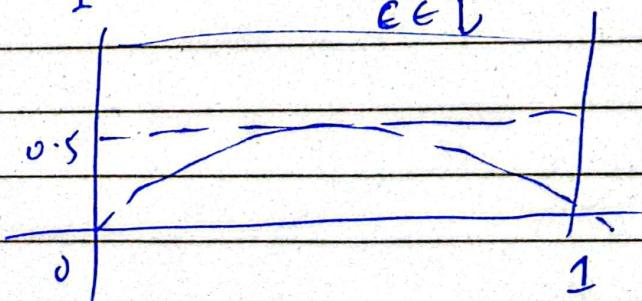
MONDAY

DATE:

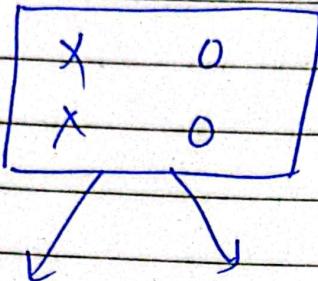
- Before we explore how splitting criterion is used in constructing decision trees, let's explore the common information measurement for decision for decision trees, gini impurity.
- Gini impurity is a mathematical measurement of how pure "the information in a data set is.
- In regards to classification, we can think of this as a measurement of class uniformity.
- Let see how this relates to the simplest case of two classes-
- Gini impurity for classification:
 - for a set of classes C for a given dataset Q , P_c is probability of class C .

$$P_c = \frac{1}{N_Q} \sum_{x \in Q} I(y_{\text{class}}=c) G(Q)$$
$$= \sum_{C \in C} P_c(1-P_c)$$

$$G(Q) = \sum_{C \in C} P_c(1-P_c)$$



$$\cdot G(Q) = \sum_{C \in C} P_C (1 - P_C)$$



class X
 $(2/4)(1 - 2/4) = 0.25$

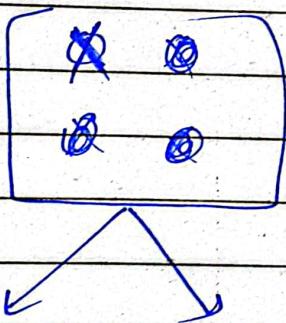
class 0
 $(2/4)(1 - 2/4) = 0.25$

"Maximum" Impurity possible

$$G(Q) = \sum_{C \in C} P_C (1 - P_C)$$

gini impurity = $0.25 + 0.25 = 0.5$

• Data is more "pure" (less impurity)



class X
 $(1/4)(1 - 1/4) = 0.1875$

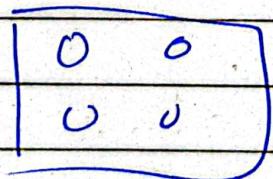
$(3/4)(1 - 3/4) = 0.1875$

Gini impurity

$0.1875 + 0.1875 = 0.375$

Dataset is completely "pure" (no impurity)

$$G(Q) = \sum_{C \in C} P_C (1 - P_C)$$



class X
 $(0/4)(1 - 0) = 0$

class 0

$(4/4)(1 - 4/4) = 0$
 $0 + 0 = 0$

Gini impurity

- If the goal of a decision tree is to separate out classes, we can use gini impurity to decide on data split value -
- We want to minimize the gini impurity at leaf nodes -
- Minimized impurity at leaf nodes mean we are separating classes effectively!

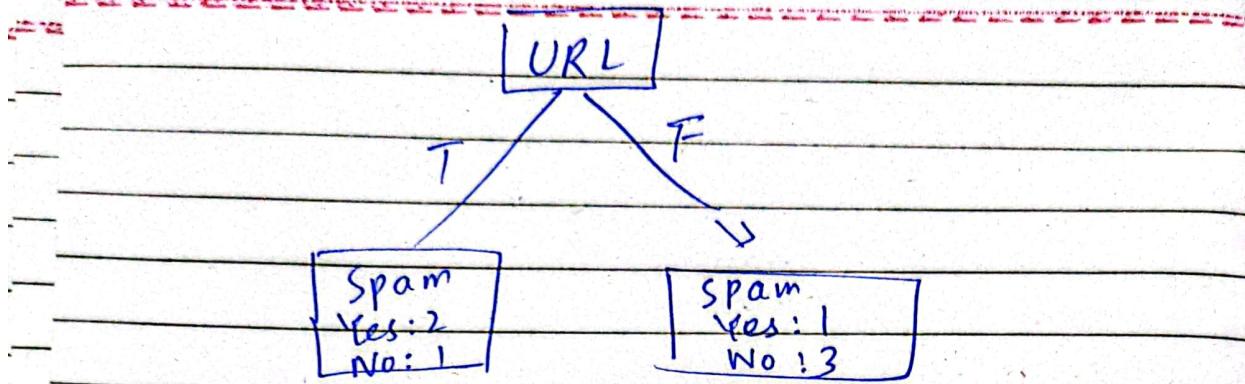
5. Constructing Decision Tree with Gini Impurity

- When first constructing a tree, we need to decide what features will be used as the root node -
- we can use gini impurity to compare the information contained within features for the training data -

$$p_c = \frac{1}{N_Q} \sum_{x \in Q} I(y_{class}=c) \quad G(Q) = \sum_{c \in C} p_c(1-p_c)$$

- let take a look at this data set:

X-URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
Yes	No
Yes	No



Treat Yes Spam and No Spam as C classes

- left leaf Node

$$\left(\frac{2}{3}\right)\left(1 - \frac{2}{3}\right) + \left(\frac{1}{3}\right)\left(1 - \frac{1}{3}\right)$$

$$\text{left } Gini = 0.44$$

On

- Right leaf Node

$$\left(\frac{1}{4}\right)\left(1 - \frac{1}{4}\right) + \left(\frac{3}{4}\right)\left(1 - \frac{3}{4}\right)$$

$$\text{Right leaf node} = 0.375$$

- Weight average of both:

$$\text{left leaf } Gini = 0.44$$

$$\text{Right leaf } Gini = 0.375$$

$$\begin{aligned} \text{Total Email} &= (2+1) + (1+3) = 7 \\ \text{Email} &= (2+1) + (1+3) = 7 \end{aligned}$$

$$\text{left Email} = 3$$

$$\text{Right Email} = 4$$

$$\left(\frac{3}{7}\right)^* 0.44 + \left(\frac{4}{7}\right)^* 0.375 = 0.403$$

But if we have multiple features -

We still have issue to consider :

- Multiple features
- Continuous Features
- Multi-class categorical features

we can incorporate the gini impurity to each of these to solve for best ~~possible~~ root nodes and best split parameters for leaves -

part II

- We explored how to create gini impurity for a binary categorical feature (only consisting of two category)

Now let us explore the following :

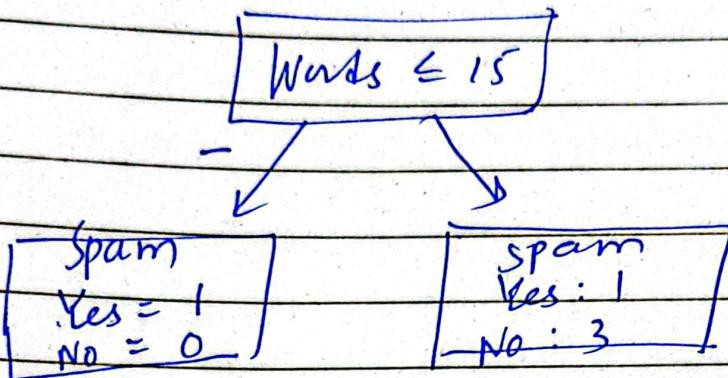
- continuous numeric features
- Multi-categorical features ($N > 2$)
- choosing a root node feature

X - words in Email	Y Sign
15	10
20	40
35	20
45	50
	30

M T W T F S

DATE: _____

words ≤ N



$$G(Q) = \left(\frac{4}{5} \right) (0+0) + \left(\frac{1}{5} \right) \left(\frac{1}{4} \right) (1-\frac{1}{4}) + \left(\frac{1}{4} \right) \left(\frac{3}{4} \right)$$

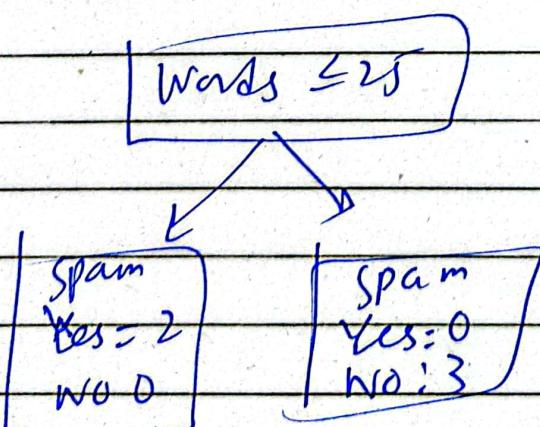
$$= 0.3$$

15 — 0.3

25 — 0 Gini Impurity

35 — 0.26

45 — 0.4



$$G(Q) = 0$$

- Let's explore gini impurity for a features that is multiclass.

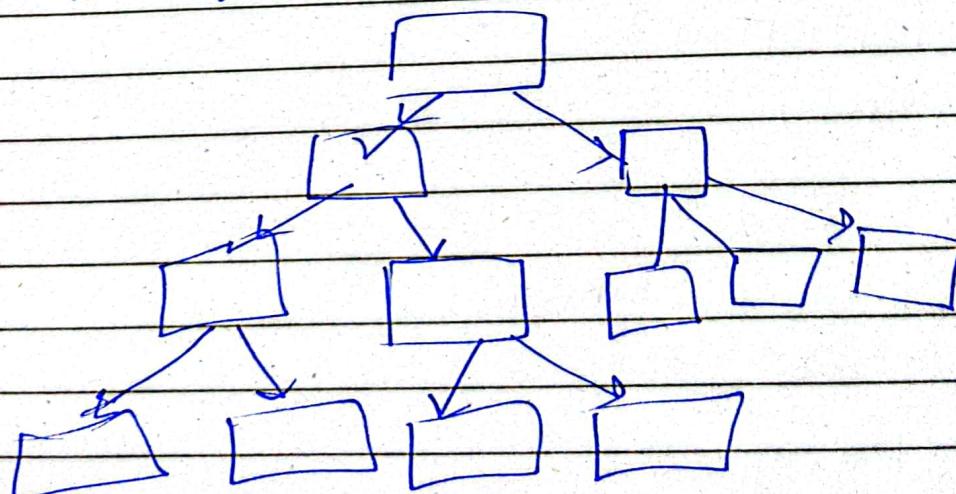
X-Sender	Y-Spam	Senders = Abe
Abe	Yes	
Bob	Yes	
Clay	No	
abe	No	
bob	No	

Sender == Abe or bob Sender == bob

Sender == Clay or bob Sender == Clay

Sender == abc or cd

A large overfitted tree:



M T W T F S

DATE: _____

- Add minimum gini impurity decrease
- We can also mandate a max depth!
prune after after
certain split

007 Coding Decision Tree

- part One the data

```
df = pd.read_csv("penguinsize")  
df.head()
```

```
df["species"].unique()
```

```
df.isnull().sum()  
df.info()
```

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

```
df.head()
```

```
df["island"].unique()  
df["sex"].unique()
```

```
df[df["sex"]=="F"]
```

describe()

```
df[df["species"]=="Gentoo"].groupby("sex")
```

`df.at[336, 'sex'] = 'FEMALE'`

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

`sns.catplot(x='species', y='culmen-length',
data=df, kind='box',
col='sex')`

`X = pd.get_dummies(df.drop(['species'], axis=1))
y = df['species']` dropfirst=True

`X_train, X-test, y_train, y-test = train-test-split
(X, y, test_size=)`

part two

`from sklearn.tree import DecisionTreeClassifier`

`model = DecisionTreeClassifier()`

`base_pred = model.fit(X_train, X-test)`

`base_pred = model.predict(X-test)`

`y-classification-report(y-test, y-pred)`

`confusion-matrix(y-test, y-pred)`

`model.feature-importance`

`nb1()`

`1('sex')`

X. columns

Index

```
pd.DataFrame(index=X.columns, data=model.F
              columns=['feature1in']) .sort_values
              ('feature')
              ('inpr')
```

```
from sklearn.tree import plot_tree
plt.figure(figsize=(12,8), dpi=200)
plot_tree(model);
```

```
feature_name = X.columns,
filled = True -
```

```
def report_model(model):
    model_preds = model.predict(X-test)
    print(classification_report(y-test, model))
    print('\n')
    plt.figure(figsize=(12,8), dpi=200)
    plot_tree(model, feature=X.columns,
              filled=True);
```

```
report_model(model)
```

```
pruned_tree = DecisionTreeClassifier(
    max_depth=2)
```

```
pruned_tree.fit(X-train) f-train)
```

```
class report_model(model)
```

QUESTION

DATE:

maxleaf tree = Decision tree efficient { max-leaf-
mean leaf tree · fit() nodes = 3
~~Decision tree C~~

entropy tree = Decision tree class / criterium = 1 Entropy