

Recap of what we've done

- 0) ML flow
- 1) Linear regression (Regression)
- 2) Naive Bayes (Classification)

Differentiate → Regression vs Classification } Supervised Learning

Student	Math	ML	DB	Salary
1				5000
2				8000
⋮				
1000				100 k
new sample → 1001	81	79	51	?

Student	Math	ML	DB	Work at Microsoft?
1				YES
2				No
⋮				
1000				YES
1001	81	79	01	?

Decision Trees

- Can be used for classification AND regression
- But mainly famous because of classification.

Day	Outlook	Temp.	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cold	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Entropy, Information Gain,
Gini Coefficient

D14 Rain Hot Low Weak ?

Entropy -

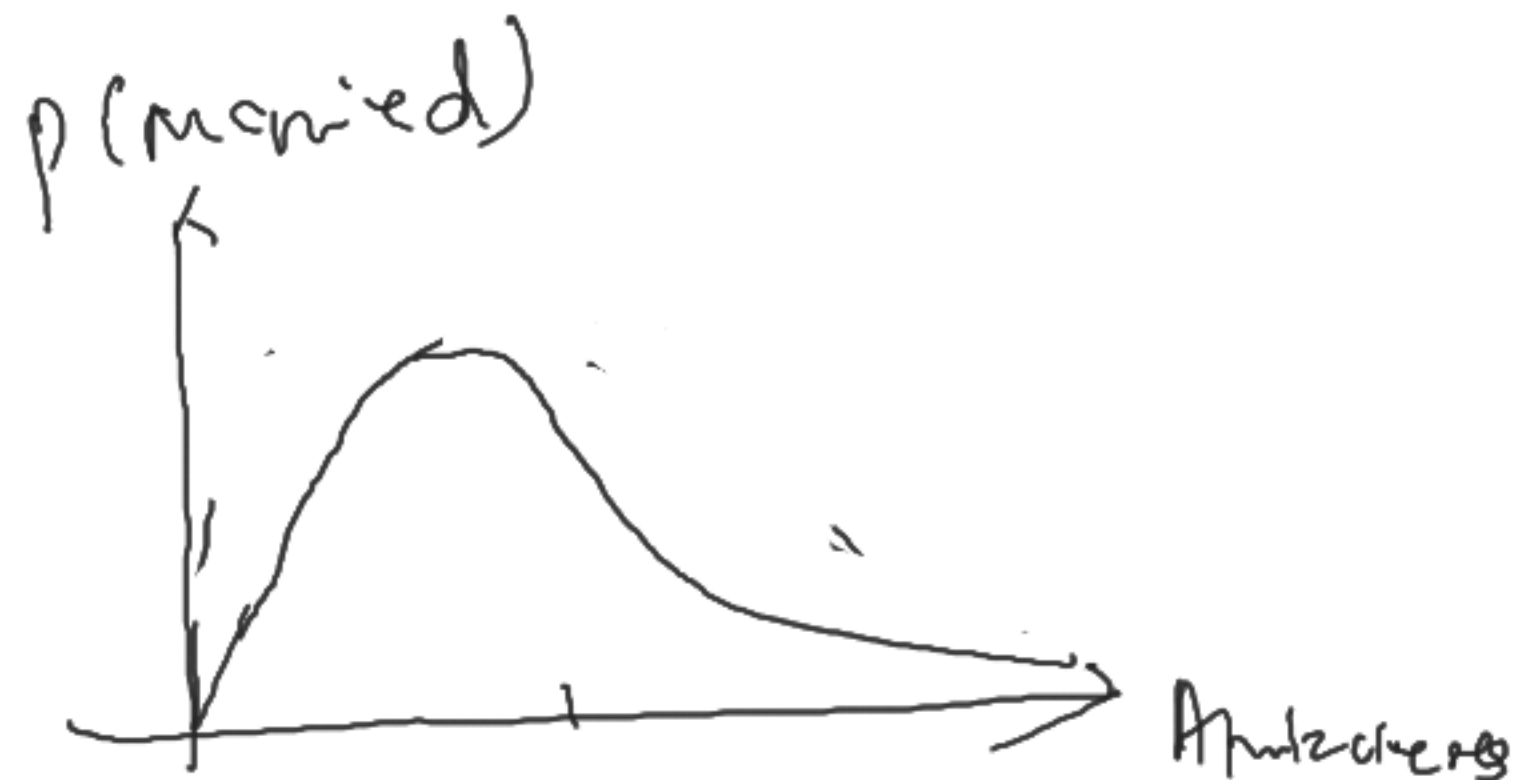
$$H(X) = - \sum P(X) \log_2 P(X)$$



Information Gain IG

$$IG(X, Y) = H(Y) - H(Y|X)$$

Measuring the amount
of surprise



6. We will use the dataset below to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. (For this problem, you can write your answers using \log_2 , but it may be helpful to note that $\log_2 3 \approx 1.6$.)

$$H(x) = - \sum_{k=1}^K p(x) \log_2 [p(x)]$$

Weather = Sunny, Cloudy, Rainy

$$H(\text{weather}) = - \sum_{k=1}^K p(\text{weather}) \log_2 p(\text{weather})$$

	GPA	Studied	Passed
1	L	F	F
2	L	T	T
3	M	F	F
4	M	T	T
5	H	F	T
6	H	T	T

$$P(\text{Pass} = T) = 4/6 = 2/3$$

$$P(\text{Pass} = F) = 2/6 = 1/3$$

$$H(\text{Pass}) = - \sum_{k=1}^K p(\text{Pass}) \log_2 p(\text{Pass})$$

$$= - \left[\overbrace{P(\text{Pass} = T) \log_2 P(\text{Pass} = T)}^{\text{Pass}} + \overbrace{P(\text{Pass} = F) \log_2 P(\text{Pass} = F)}^{\text{Fail}} \right]$$

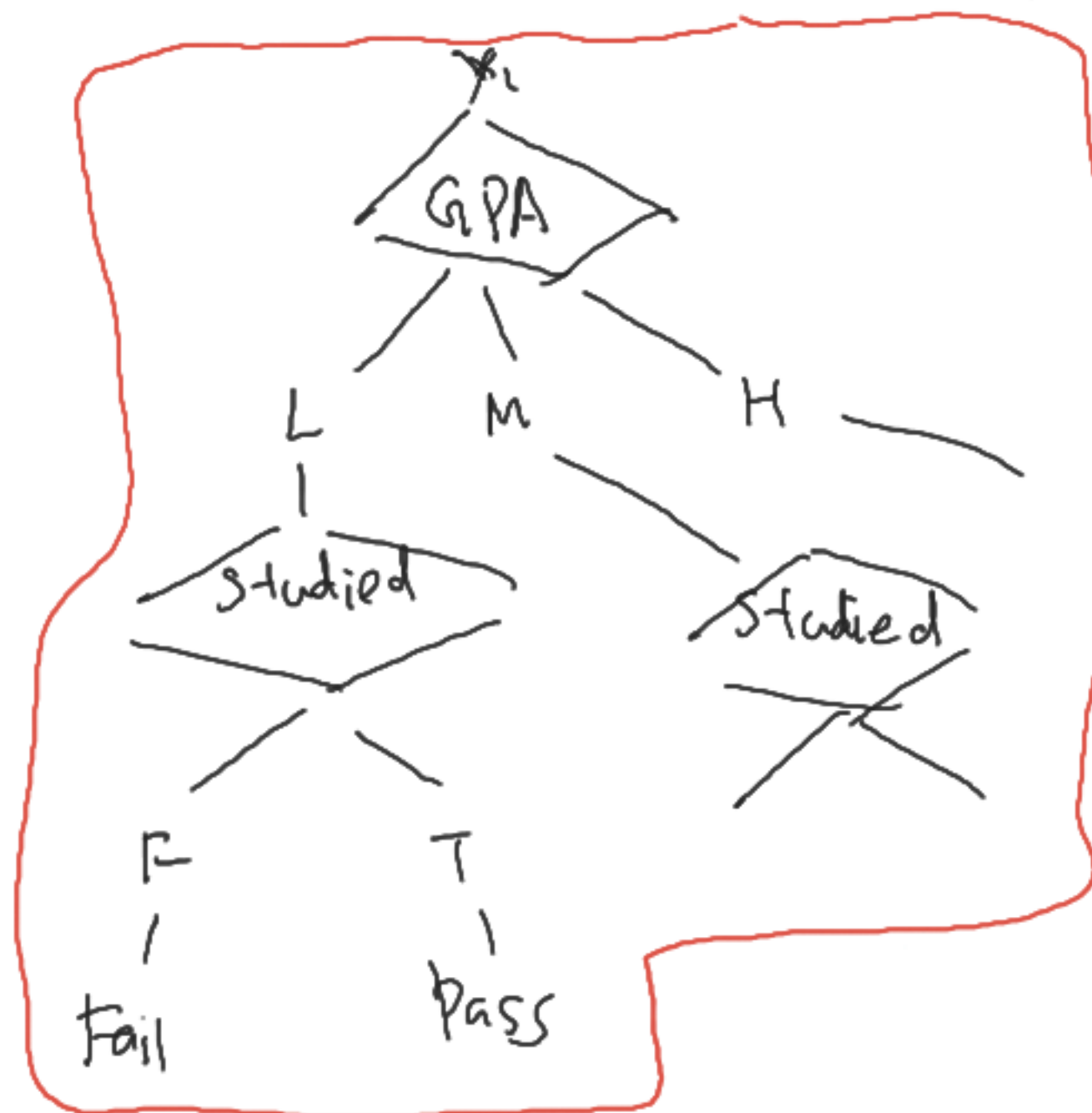
$$= - \left[\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] \approx 0.92$$

6. We will use the dataset below to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. (For this problem, you can write your answers using \log_2 , but it may be helpful to note that $\log_2 3 \approx 1.6$.)

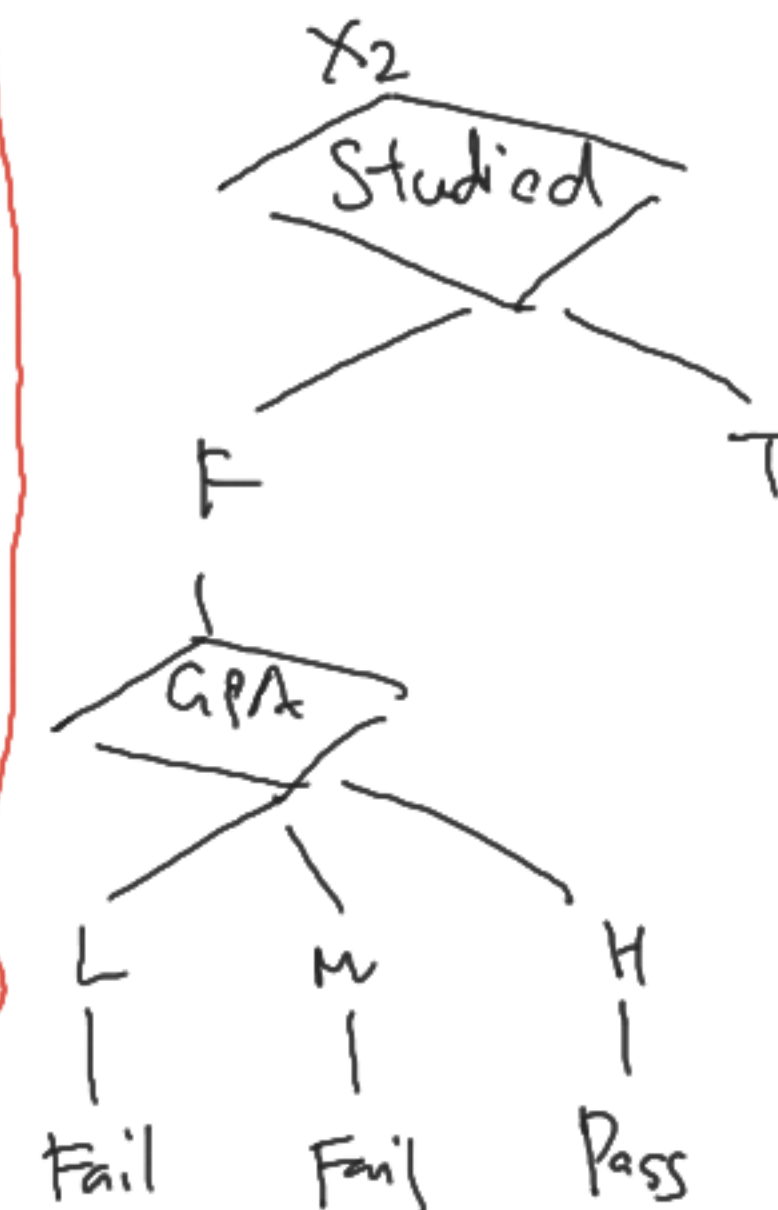
x_1	x_2	y
GPA	Studied	Passed
L	F	F
L	T	T
M	F	F
M	T	T
H	F	T
H	T	T

$$H(\text{Pass}) = 0.92$$

Q) How do we expand the tree / Choose the root node?



Two features : GPA (x_1), Studied (x_2)
 one target : Passed (y)



Do we choose GPA or Studied as the root node?

We have to calculate Information Gain

- $IG(\text{Pass}, \text{GPA})$
- $IG(\text{Pass}, \text{Studied})$

6. We will use the dataset below to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. (For this problem, you can write your answers using \log_2 , but it may be helpful to note that $\log_2 3 \approx 1.6$.)

We know this already



$$IG(\text{Pass}, \text{GPA}) = H(\text{Pass}) - H(\text{Pass} | \text{GPA})$$

$$H(\text{Pass} | \text{GPA}) = - \sum P(\text{GPA}) \sum P(\text{Pass} | \text{GPA}) \log_2 P(\text{Pass} | \text{GPA})$$

Pass = TRUE (Pass)

$$P(\text{Pass} = \underline{T} | \text{GPA} = \underline{L}) = \underline{1/2}$$

$$P(\text{Pass} = T | \text{GPA} = M) = 1/2$$

$$P(\text{Pass} = T | \text{GPA} = H) = 1$$

Pass = FALSE (Fail)

$$P(\text{Pass} = F | \text{GPA} = L) = 1/2$$

$$P(\text{Pass} = F | \text{GPA} = M) = 1/2$$

$$P(\text{Pass} = F | \text{GPA} = H) = 0$$

$$P(\text{GPA} = L) = 2/6$$

$$P(\text{GPA} = M) = 2/6$$

$$P(\text{GPA} = H) = 2/6$$

	GPA	Studied	Passed
1	L	F	F
2	L	T	T
3	M	F	F
4	M	T	T
5	H	F	T
6	H	T	T

6. We will use the dataset below to learn a decision tree which predicts if people pass machine learning (Yes or No), based on their previous GPA (High, Medium, or Low) and whether or not they studied. (For this problem, you can write your answers using \log_2 , but it may be helpful to note that $\log_2 3 \approx 1.6$.)

GPA	Studied	Passed
L	F	F
L	T	T
M	F	F
M	T	T
H	F	T
H	T	T

$$H(\text{Pass} | \text{GPA}) =$$

$$= \frac{2}{6} \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right]$$

$$+ \frac{2}{6} \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right]$$

$$+ \frac{2}{6} \left[\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{0}{2} \log_2 \left(0 \right) \right]$$

$$H(\text{Pass} | \text{GPA}) \approx 0.66$$

$$IG(\text{Pass}, \text{GPA}) = H(\text{Pass}) - H(\text{Pass} | \text{GPA})$$

$$= 0.92 - 0.66$$

$$= 0.26$$

0.26

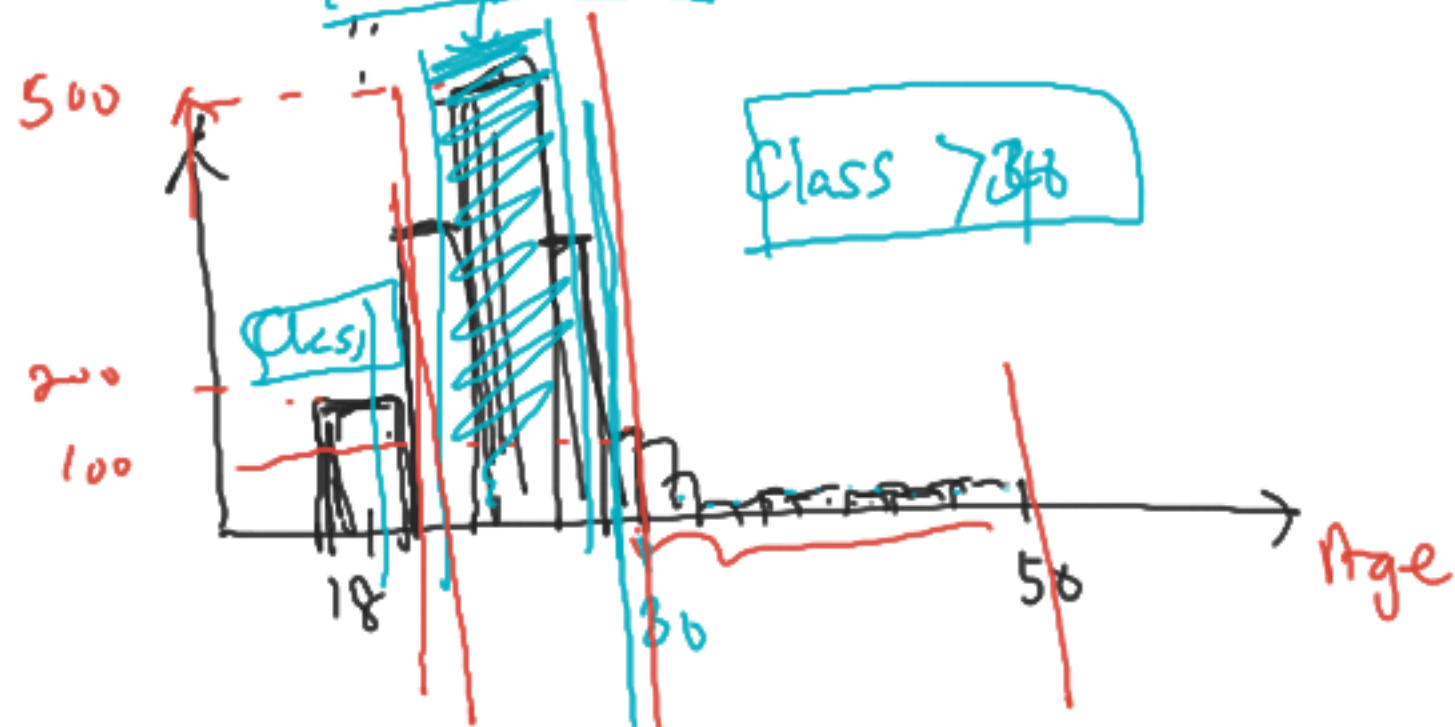
$$H(\text{Pass} | \text{GPA}) = - \sum P(\text{GPA}) \sum P(\text{Pass} | \text{GPA}) \log_2 P(\text{Pass} | \text{GPA})$$

Decision Trees

GPA	Age	Studied	Pass
L	32	T	
H	18	T	
M		F	
M	50	F	



IG (Pass, Age > 40)



Young middle old

- Number of samples for each 'category'
- Information gain of each category

Problems with DT

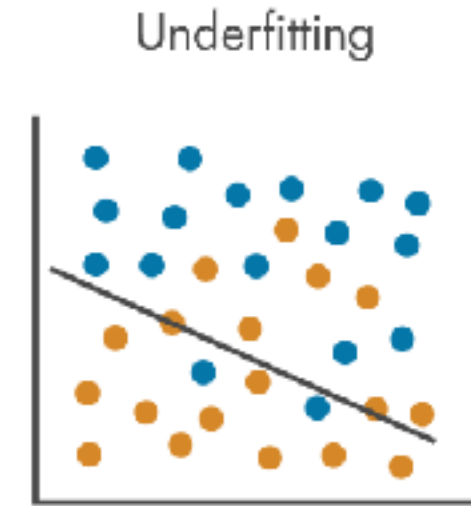
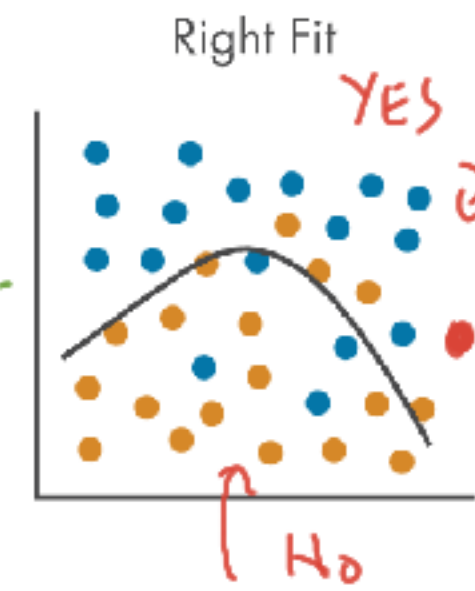
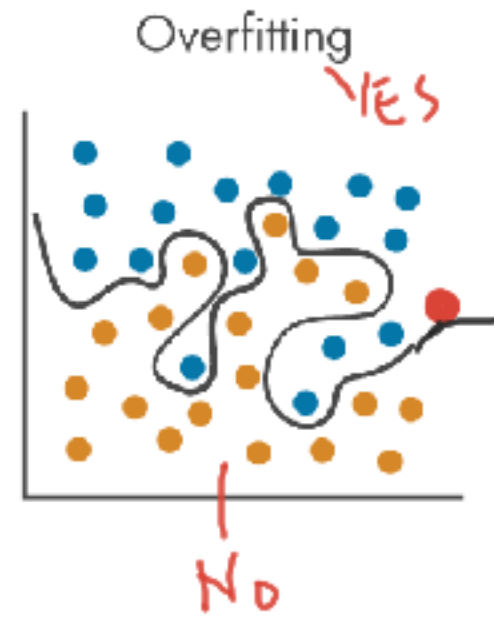
1) Overfitting
↓

2) overcome by 2
methods

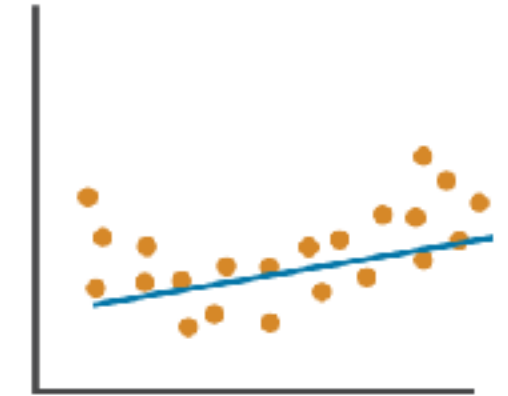
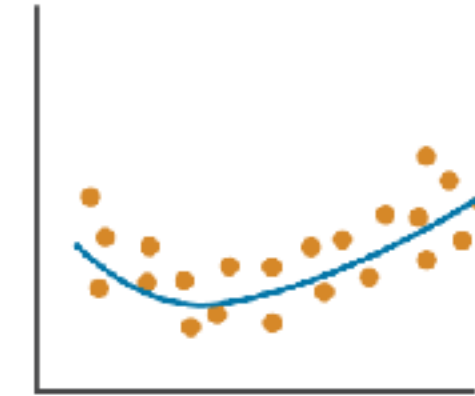
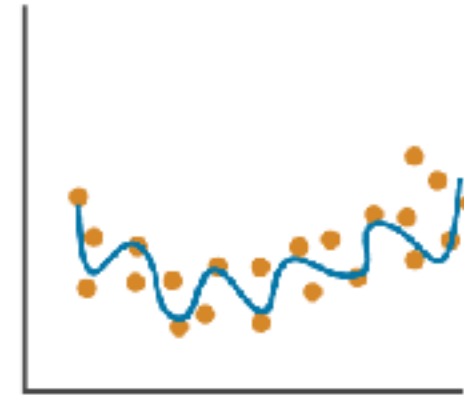
a) Pre-Pruning

b) Post-Pruning

Classification



Regression



<https://medium.datadriveninvestor.com/decision-tree-algorithm-with-hands-on-example-e6c2afb40d38>

<https://www.kaggle.com/code/prashant111/decision-tree-classifier-tutorial#14.-Decision-Tree-Classifier-with-criterion-entropy->

Ensemble Learning → The foundation of Random Forest

DT + SVM + NB

Why ensemble learning?

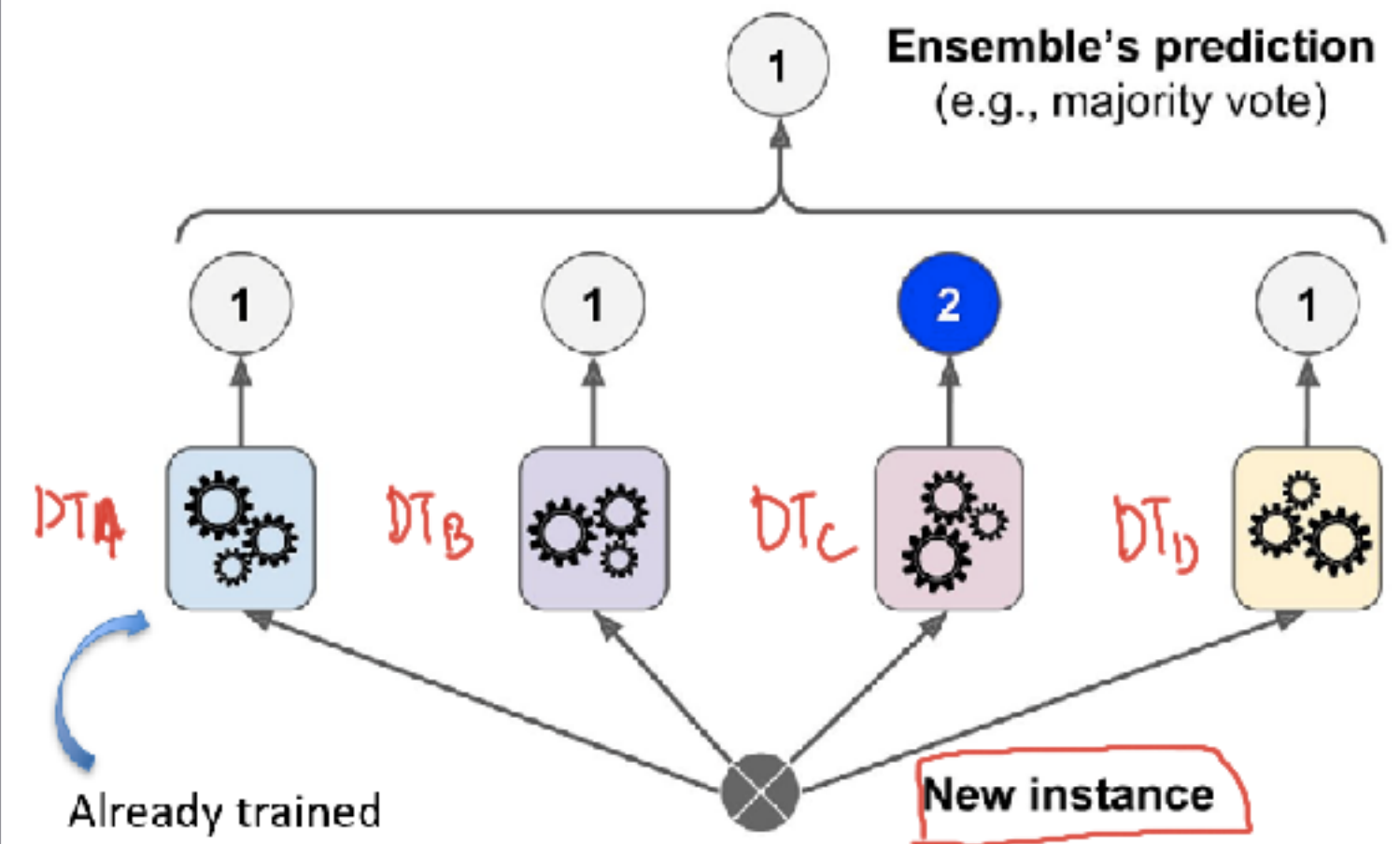
RF = DT + DT + DT + ...

Example: 3 classifiers, each is correct with 60% accuracy

- Ensemble: combination of this 3 → take the majority vote
- What is the accuracy of this ensemble?
- Answer: 3 votes - ensemble is wrong if there are at least 2 wrong predictions
- Case 1: Exactly 2 wrong predictions: $3 \cdot 0.4 \cdot 0.4 \cdot 0.6 = 0.288$
- Case 2: All 3 are wrong: $0.4^3 = 0.064$
- Total: 0.352 → ensemble is correct with 64.8% accuracy
- Another example: 1000 classifiers with 51% accuracy each → 75% accuracy

Law of the large numbers

Ensemble learning



Common practices

Classifiers should be independent: D_{TA} is not dependant on D_{TB}

- Correlated classifiers don't work well together (why?)
 - E.g., 2 classifiers: both make mistakes at the same data points
 - Their combination will not be better at all

- Ensemble (completely) different classifiers

RF uses mainly
↑ bagging

SVM + NB + DT

✓

SVM + SVM + NB

✗

- Vary the training data (bagging and pasting)
- Vary the feature sets (random patches & random subspaces)

Bagging and pasting

Handwritten notes:

- d_1, \dots, d_n (circled)
- Sampling
- $\sim \mathcal{N}(\mu, \sigma)$
- $- d_i \sim \mathcal{N}(\mu, \sigma)$
- $\sim \mathcal{N}(\mu, \sigma)$

- Use combination of same type of classifiers (e.g., all are SVMs)
- Randomly sample a subset of training data to train each classifier
- **Bagging** (= bootstrap aggregating): choose with replacement (can sample same data point for the same classifier multiple times)
- **Pasting**: choose without replacement (**cannot** sample the same data point for same classifier multiple times)

Handwritten note: "Samples are perurbed"

Random Forest

Sample randomly with replacement $N=800$

Class	Studied	Pass
1	F	T
2	T	F
3	T	T
4	T	T
5	F	F

Train data (red box)

Test data (green box)

