



Department of Computational and Data Sciences



AI: Module 01 Week 01

Deepak Subramani

Assistant Professor

Dept. of Computational and Data Science

Indian Institute of Science Bengaluru



Week 01

- Part 1

- 1 ✓ • Decision Tree – Structure, Prediction, Attributes
- 2 ✓ • Decision Boundary – Fundamental Concept in AI for Classification
- 3 ✓ • Cart Training Algorithm

- Part 2

- 4 ✓ • Metrics of performance – Regression and Classification
- 5 ✓ • Development-Testing Paradigm
- ✓ • Overfit vs Underfit – Regularization to prevent overfit
- ✓ • Hyperparameter Selection through k-fold CV

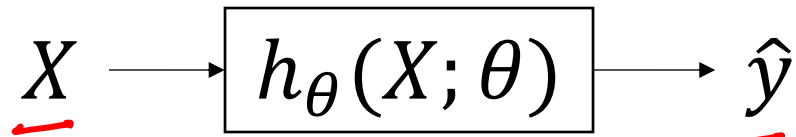
ML: Mental Model

Data that can be collected



Quantity that must be predicted to make money

Data that can be collected



Machine's Prediction

↗ Mathematical Fn.

θ : parameters of the machine ML Model

↙ Machine Learning

Gradient Boosted Decision Tree

Table

	C	T	SR	R
→	.	.	.	⊗
→				
→				

→ yes on no

1000 x 3

The Machine Learning Workflow

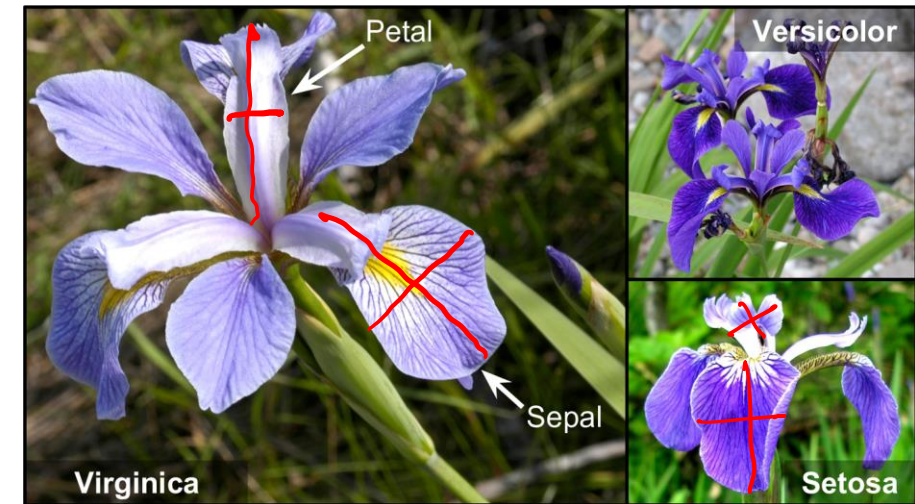
1. Frame the ML problem by looking at the business need
 - a. Identify subproblems (One/more of the 5 tasks a computer can do)
 - b. Establish a current baseline (What is currently done?)
 - c. Define success
2. Gather the data and do Data Munging/Wrangling + Baselines
 - a. Explore the data
 - b. Clean data and prepare for the downstream ML models
 - c. Establish a data, domain and SoTA baseline
3. Explore different models, improve them through Cross Validation and perhaps new model design
4. Form an ensemble of multiple models and solutions
5. Present your solution
 - a. Say a story with the data
6. Deploy

Decision Trees

- Decision Trees are data driven models for classification and regression
- They are a versatile ML Algorithm capable of fitting complex data
- They are trained by a greedy optimization algorithm called CART
- Plan of action:
 1. See commands for training DT in sklearn
 2. Predict with DT and see how they make a decision
 3. Understand the math of the optimization algorithm
 4. Discuss pros and cons

Iris Dataset

Sepal Length	Sepal Width	Petal Length	Petal Width	Species
				.
				.
				.
				.
				.
				.



150

3

• The problem ~~we did so far~~ is a classification problem

• Species = $h_{\theta}(SL, SW, PL, PW; \theta)$

DT

θ : parameters

h_{θ} : DT



Training a DT in sklearn

- from sklearn.datasets import load_iris ✓
- from sklearn.tree import DecisionTreeClassifier
- ⇒ • iris = load_iris()
- ⇒ • X = iris.data[:,2:] #Make 2 features in the data
- y = iris.target 0,1,2
- ⇒ • tree_clf = DecisionTreeClassifier(max_depth=2)
- tree_clf.fit(X,y)

↑ ↑
i/p o/p

↑↑
key word
argument
hyper parameters

$$\hat{y} = h_{\theta}(PL, PW; \theta)$$

Petal Length
Petal Width

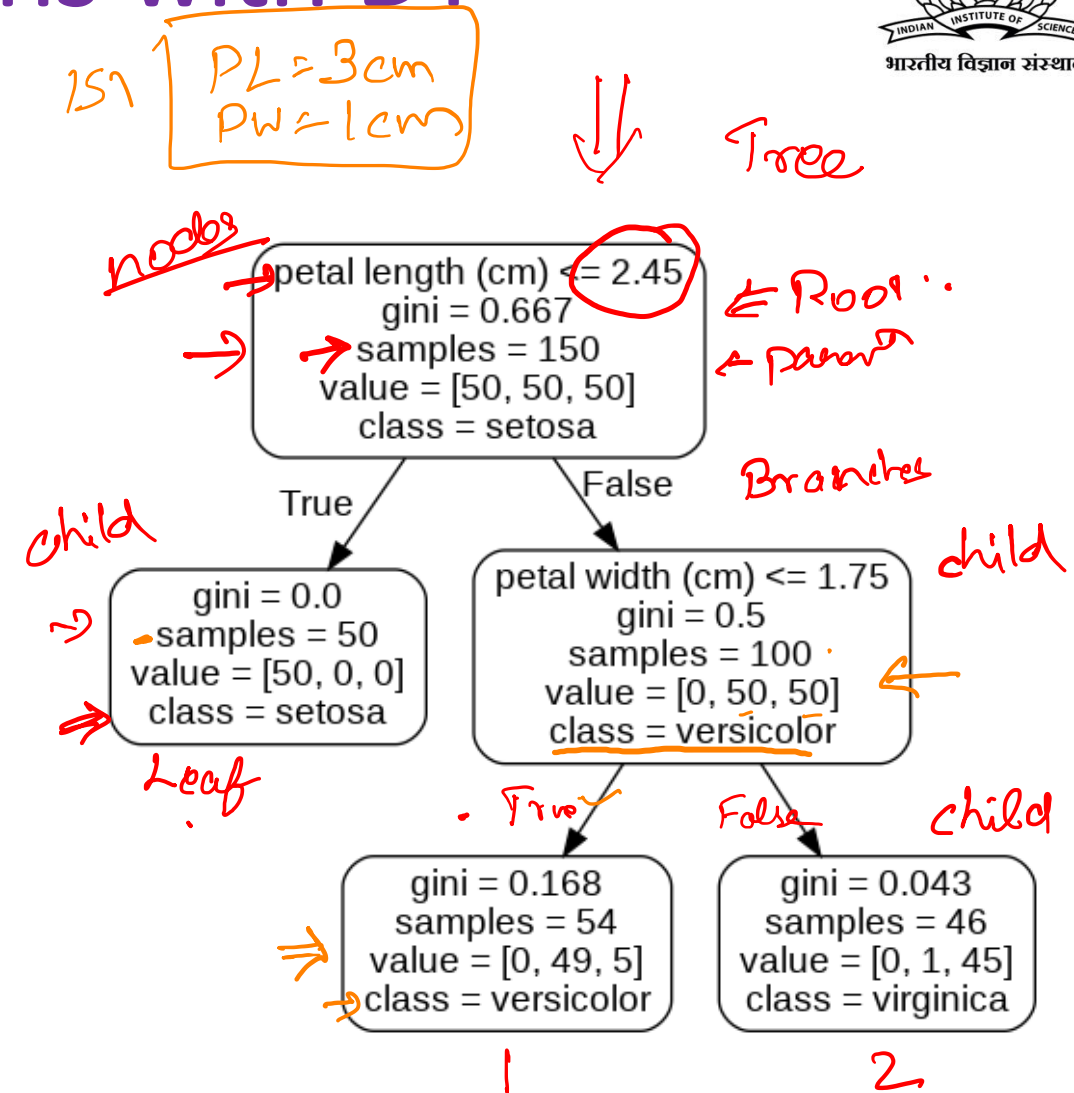
$$\hat{y} = h_{\theta}(X; \theta)$$

↑
tgt i/p
o/p

- 0 Setosa
- 1 Versicolour
- 2 Virginica

Making Predictions with DT

- How does DT classify a new data point?
- Start from the top, and move down asking the question at each node and following the answer
- Is petal length (cm) ≤ 2.45 ?
 - True – move left; False – move right
- Keep moving until you reach a leaf node
 - Leaf node – no children
- The class of the leaf node a datapoint ends up in is its class!
- That simple!



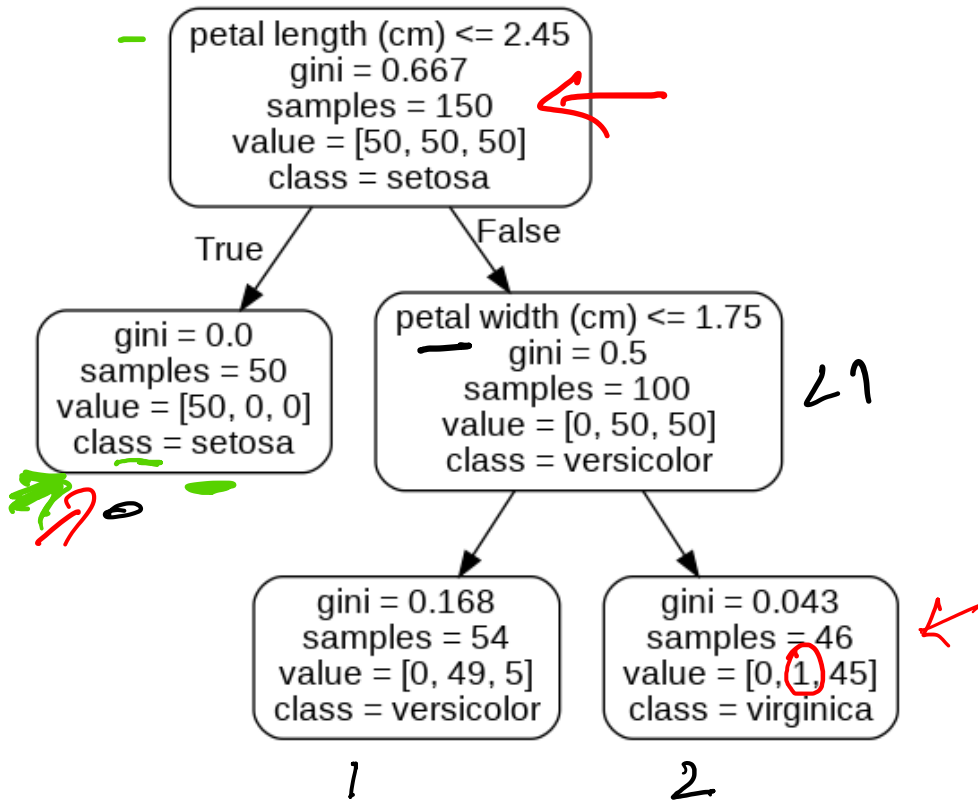
Understanding DT Attributes

- Gini attribute:
 - Measures impurity at a node
 - $G_i = 1 - \sum_{k=1}^n p_{i,k}^2$
 - $p_{i,k}$ is the fraction of k th class in i th node. n =total classes
- Samples attribute:
 - Number of training instances for which this question was applied
- Value attribute:
 - How many training instances of each class this node applies to
- Class attribute:
 - The label which will be applied to the instance if we stop there
- Probability of a class at a particular node
 - Value/samples

Decision Boundaries

in
Feature Space

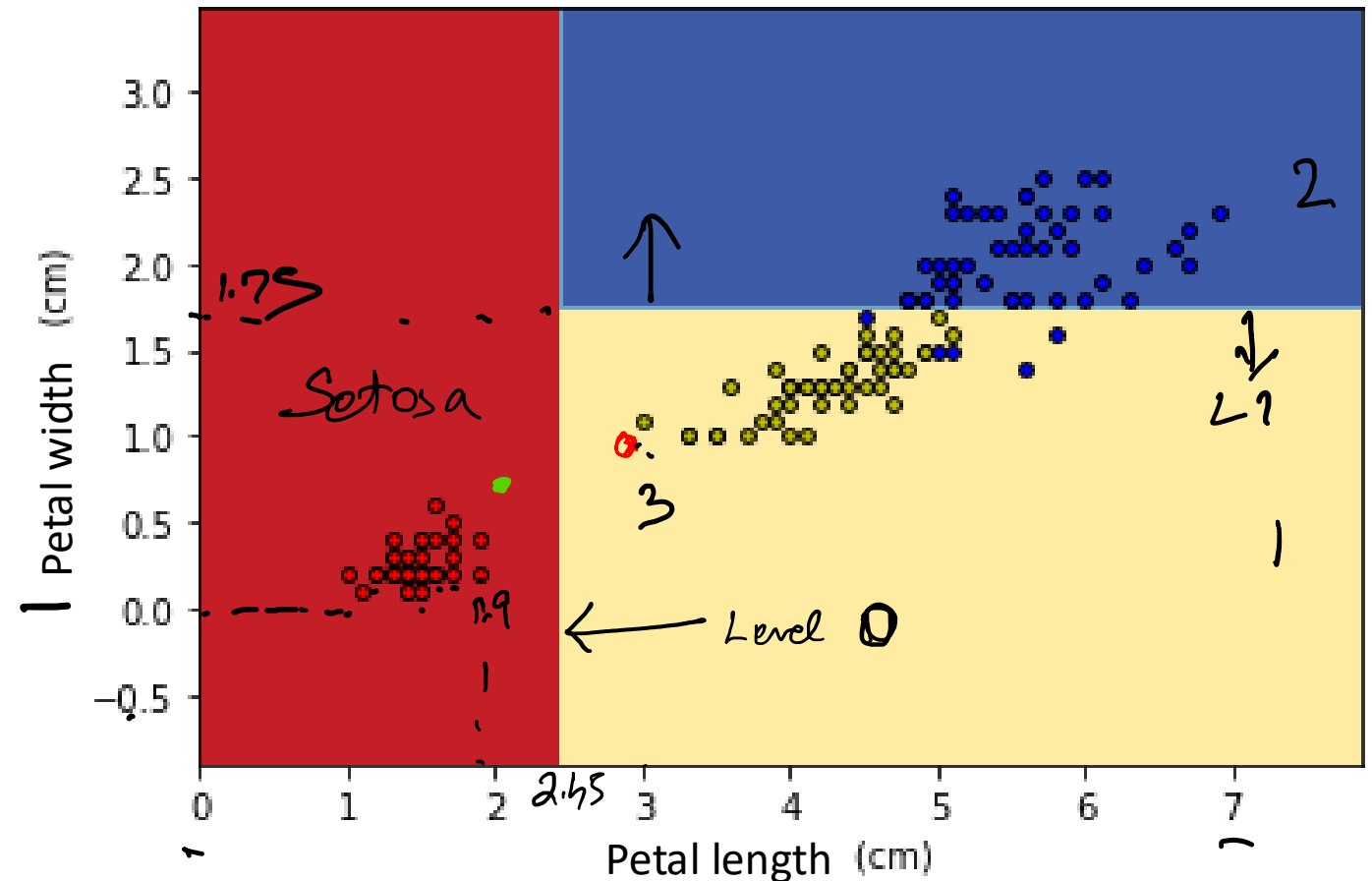
11:00 AM



46 +
50

96

350
700



0 0.01 0.02 ...

CART Training Algorithm

1. At every node, split the training set into two subsets based on one single feature x_k and a threshold t_k on it
2. The pair (x_k, t_k) is chosen by what results in the purest subset measured by Gini Coefficient or Entropy

Loss fn:
$$J(x_k, t_k) = \frac{m_{left}}{m} G_{left} + \frac{m_{right}}{m} G_{right}$$

$G_{l/r}$ is the impurity and $m_{l/r}$ is the number of instances

3. Continue the same at the two child nodes
4. Stop when no further split reduces impurity or maximum depth is reached



Bag: 150 0 1 2
50 50 50
R G B

$$p(\text{picking class 0 and labeling it wrong}) = \frac{1}{3} \left(1 - \frac{1}{3}\right) = P_0(1 - P_0)$$

$$p(\text{picking class 1 and labeling it wrong}) = \frac{1}{3} \left(1 - \frac{1}{3}\right) = P_1(1 - P_1) + P_2(1 - P_2)$$

$$p(\text{class 2}) =$$

$$\begin{aligned} &= P_0(1 - P_0) + P_1(1 - P_1) + P_2(1 - P_2) \\ &= P_0 + P_1 + P_2 - (P_0^2 + P_1^2 + P_2^2) \\ &= 1 - \sum P_i^2 \end{aligned}$$



Class :

$$G = P_0(1-P_0) + P_1(1-P_1) + P_2(1-P_2)$$
$$E = -\left[P_0 \ln P_0 + P_1 \ln P_1 + P_2 \ln P_2\right]$$

Reg :

$$\sum (\bar{y} - y_j)^2$$

MSE

$$\hat{y} = \bar{y}$$

Example of a node split

50, 0, 0

$$\begin{aligned} p_0 &= 1 \\ p_1 &= 0 \\ p_2 &= 0 \end{aligned}$$

• PL ≤ 2.45

- True (left) – Samples 50 Value = [50, 0, 0]
- False (right) – Samples 100 Value = [0, 50, 50]

• G_left = 0

• G_right = 0.5

• J = $\frac{50}{150} \times 0 + \frac{100}{150} \times 0.5 = \frac{2}{3} \times \frac{1}{2} = 0.33$

$$G_{\text{left}} = 1 - (1^2 + 0^2 + 0^2) = 0$$

$$G_{\text{right}} = 1 - (0^2 + 0.5^2 + 0.5^2) = 1 - 0.5 = 0.5$$

$$\begin{aligned} p_0 &= 0 \\ p_1 &= 0.5 \\ p_2 &= 0.5 \end{aligned}$$

• PL ≤ 3.5

- True (left) – Samples 55 Value = [50, 5, 0]
- False (right) – Samples 95 Value = [0, 45, 50]

• G_left =

• G_right =

• J =

G_left

$$p_0 = \frac{50}{55} \quad p_1 = \frac{5}{55} \quad p_2 = \frac{0}{55}$$

$$G_{\text{left}} = 1 - \left(\left(\frac{50}{55} \right)^2 + \left(\frac{5}{55} \right)^2 + \left(\frac{0}{55} \right)^2 \right)$$

1

What is Gini Impurity?

- The Gini impurity measures the probability that a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in the subset.
- Worked out example
 - Let us consider class A, B and C.
 - we randomly pick a point and label it according to the distribution.
 - Let the label distribution be (0.3,0.5,0.2).
 - probability of picking up a data point randomly with label A is 0.3. The probability of labeling it B or C is $1-0.3=0.7$.
 - probability of picking up a data point randomly with label B is 0.5. The probability of labeling it A or C is $1-0.5=0.5$.
 - probability of picking up a data point randomly with label C is 0.2. The probability of labeling it A or B is $1-0.2=0.8$.
- The total probability of misclassification of a randomly selected data point is
- $$p(mis) = 0.3 * (1 - 0.3) + 0.5 * (1 - 0.5) + 0.2 * (1 - 0.2) = 0.3 + 0.5 + 0.2 - (0.3^2 + 0.5^2 + 0.2^2) = 1 - (0.3^2 + 0.5^2 + 0.2^2) = 1 - \sum_{k=1}^K p_k^2 = 1 - 0.38 = 0.62$$
- It ranges from 0 to 0.5 (for binary classification) or 0 to $1 - 1/J$ (for multi-class problems).
- 0 represents a pure node (all elements belong to the same class). Higher values indicate more mixed classes.

Week 01

- Part 1
 - Decision Tree – Structure, Prediction, Attributes
 - Decision Boundary – Fundamental Concept in AI for Classification
 - ✓ • Cart Training Algorithm
 - ~~Collection of Decision Trees – Random Forest~~
- Part 2
 - ✓ • Metrics of performance – Regression and Classification
 - Development-Testing Paradigm
 - Overfit vs Underfit – Regularization to prevent overfit
 - Hyperparameter Selection through k-fold CV

Evaluation Metrics – Classification

Confusion Matrix

P
N

→ GT+ ↓ ↓

	CLS+	True Positive	False Positive	54
	CLS-	False Negative	True Negative	96
		50	100	

- Recall = $TP / (TP + FN)$ – Also called Sensitivity in Statistics
- Precision = $TP / (TP + FP)$ $49/54$
- F1 Score = Harmonic mean of Recall and Precision
- False Negative Rate = $FN / (TP + FN)$
- Specificity = $TN / (FP + TN)$
- False Positive Rate = $FP / (FP + TN)$

150

Example:

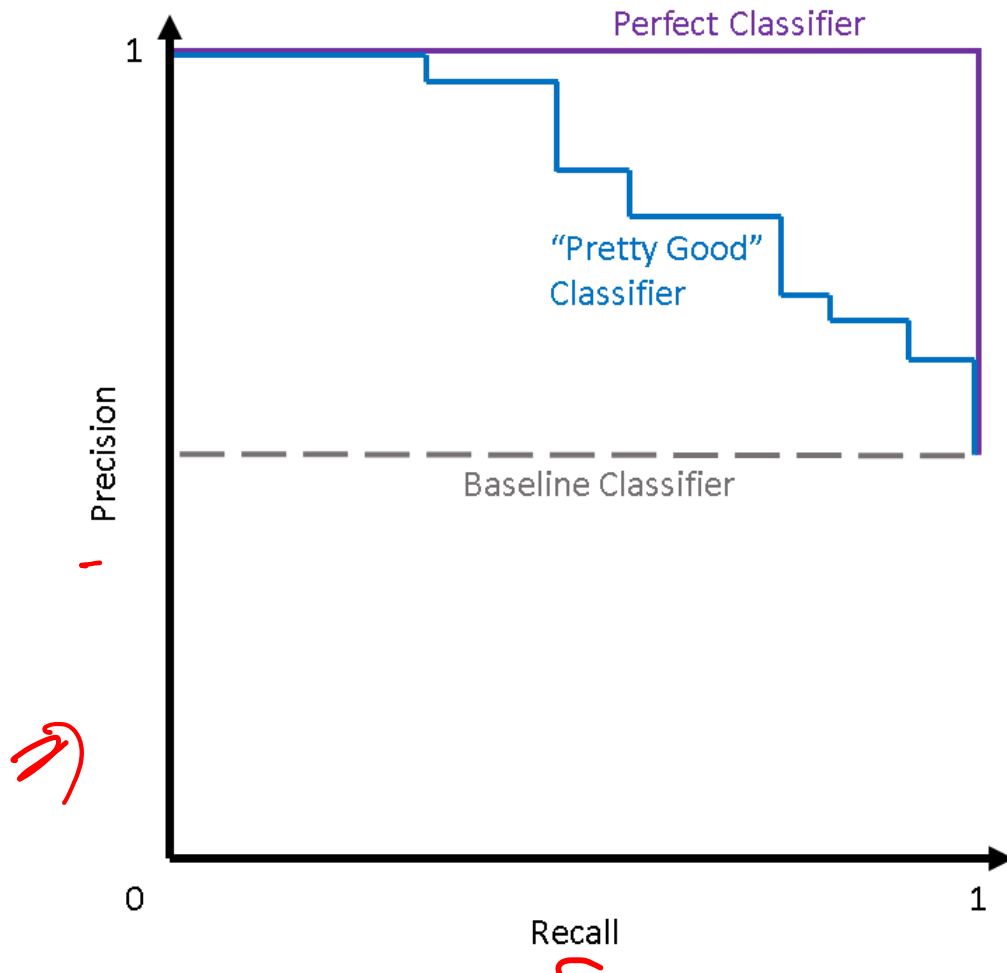
+ : Versicolor 50
- : Not Versicolor 100

↑ CLS

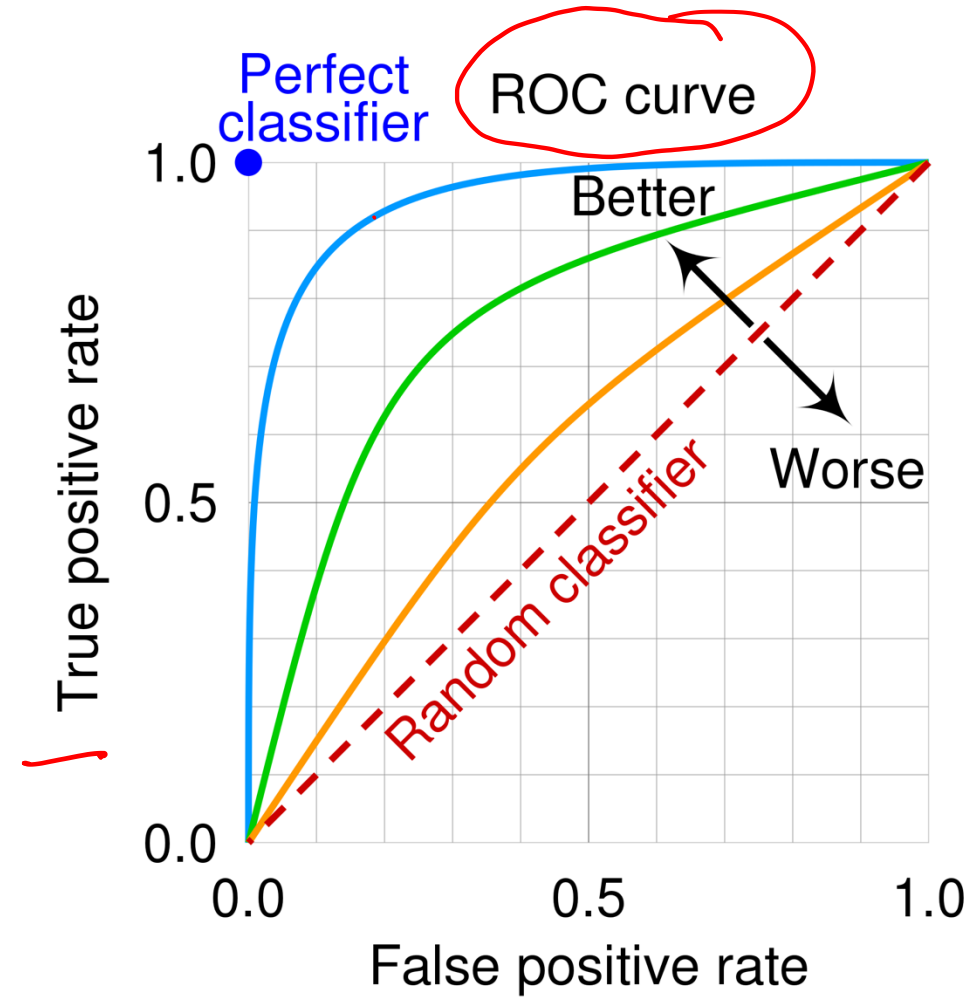
	CLS+	TP=49	FP=5	54
	CLS-	FN=1	TN=96	96
		50	100	

$$Acc = \frac{49 + 96}{150}$$

PR Curve and ROC AUC

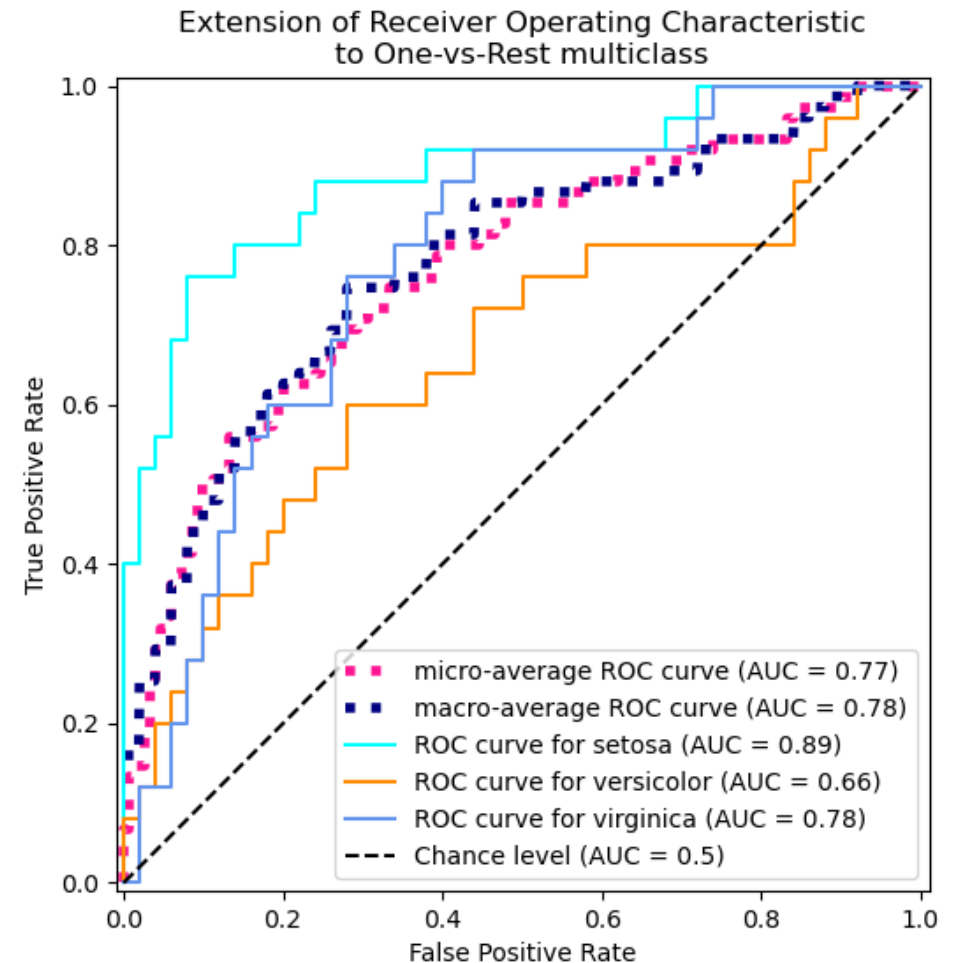


Baseline – predict all instances as +ve class



ROC AUC for Multi Class

- **Micro-average**
 - Ignores class distinctions and gives equal weight to each instance. It's a global strategy that calculates the performance metric based on the total counts of true positives, false positives, and false negatives across all classes.
- **Macro-average**
 - Gives equal weight to each class, regardless of the number of instances. It calculates each class's performance metric (e.g., precision, recall) and then takes the arithmetic mean across all classes.



Evaluation Metrics – Regression

- Root Mean Square Error
- Mean Absolute Error
- ✓ • Relative Errors *MAPE*
- $R^2 = 1 - \text{MSE} / \text{Variance}$

$$R^2 = 1 - \frac{4}{3} < 0$$

Development-Testing Paradigm

- Consider a data set with m rows and n columns
- Development Set
 - Used to train a ML model
- Training a model involves
 - Finding parameters or growing trees
 - Uses data and an optimization algorithm working on some loss
- Development involves training and hyper-parameter tuning
 - K-Fold Cross Validation is used
- Testing involves using a developed model on totally unseen data and evaluating an expected real world performance

	X_1	X_2	X_3	y	
Dev. Set					80%
Test Set					20%

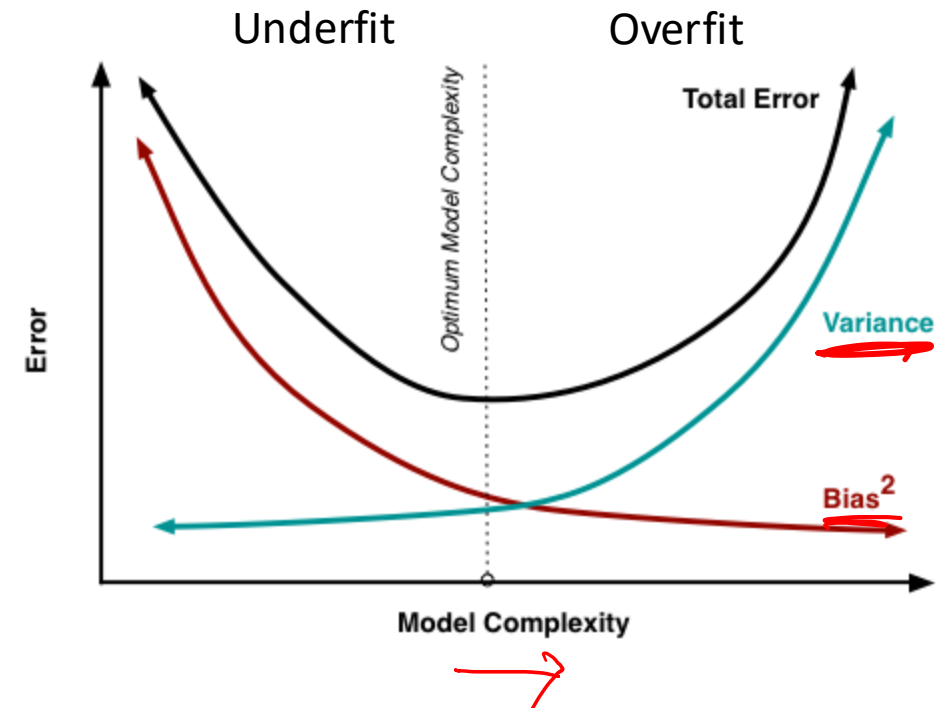
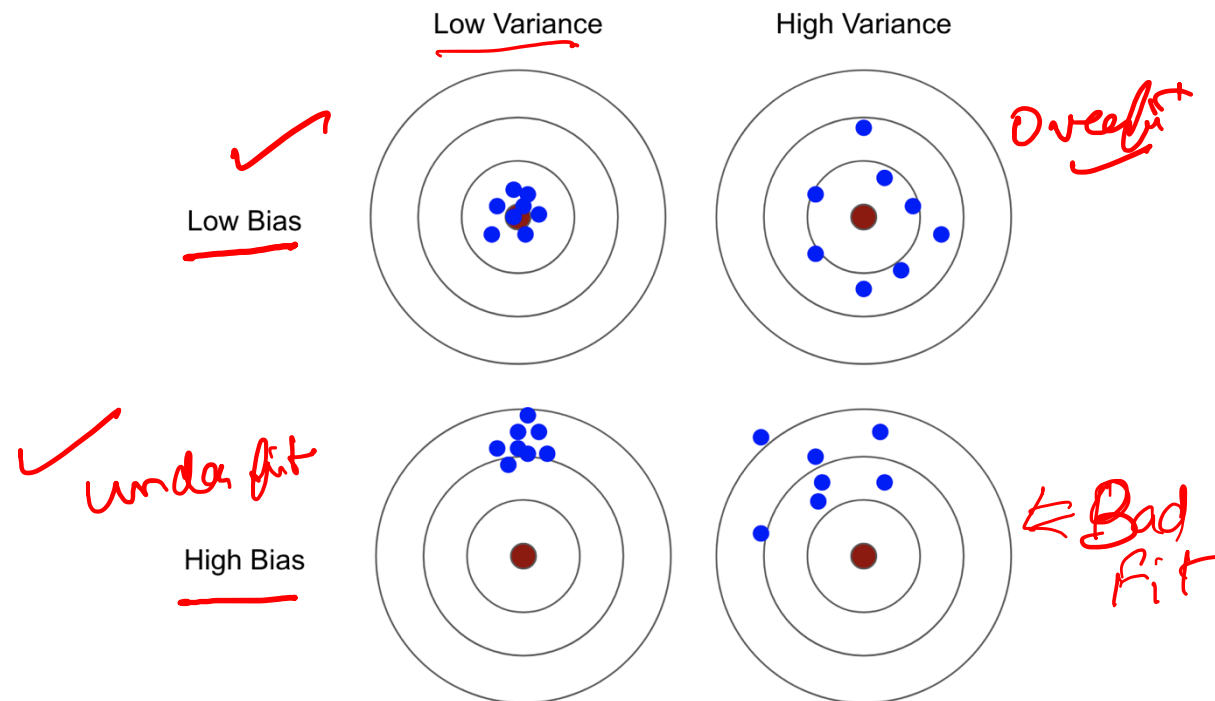
150

Bias-Variance Tradeoff

- $\text{Error} = \text{Bias} + \text{Variance} + \text{Irreducible Error}$
- **Bias:** Error from erroneous assumptions in the model
 - High bias \rightarrow underfitting
 - Model has limited flexibility to learn
 - Analogy: Overly simplistic assumptions about people make you a biased person
- **Variance:** Error from sensitivity of model to small perturbations in training data
 - High variance \rightarrow overfitting
 - Model has too much flexibility to learn
 - If you have a lot of data and over complex model, you can make each parameter of the model “by-heart” one data point
 - When a new data point comes, then the model will change wildly to accommodate the new point

Bias/Variance Equation

- $Error = E[(y - \hat{y})^2] \rightarrow \text{RMSE}$
- $Error = (E[\hat{y}] - y)^2 + E[(\hat{y} - E[\hat{y}])^2] + \sigma_e^2$
- Error = Bias + Variance + Irreducible Error



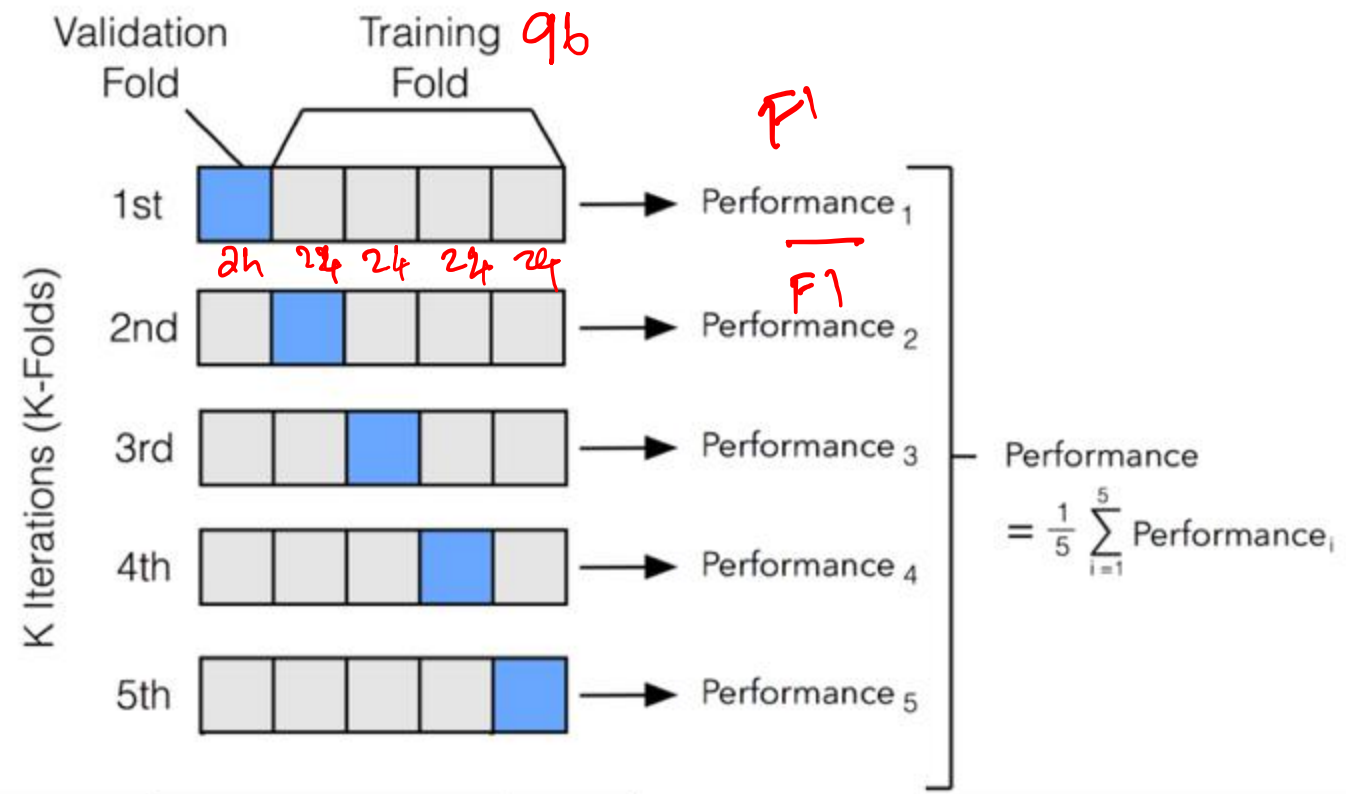


Regularization

- DT are what is called as non-parametric models
 - They don't have a pre-defined number of parameters
 - On the other hand, linear models are parametric, with predefined number of parameters
- Thus DT can overfit any complex training data
- To avoid overfitting, we restrict a DT's degrees of freedom
- Use the following hyperparameters to regularize and avoid overfitting:
 - max_depth: the maximum number of levels
 - min_samples_split: the minimum number of samples a node must have before it can be split,
 - min_samples_leaf: the minimum number of samples a leaf node must have,
 - min_weight_fraction_leaf: same as min_samples_leaf but expressed as a fraction of the total number of weighted instances,
 - max_leaf_nodes: the maximum number of leaf nodes, and
 - max_features: the maximum number of features that are evaluated for splitting at each node.
 - Increasing min_ hyperparameters or reducing max_ hyperparameters will regularize the model.
- Can also build a tree without restrictions and prune

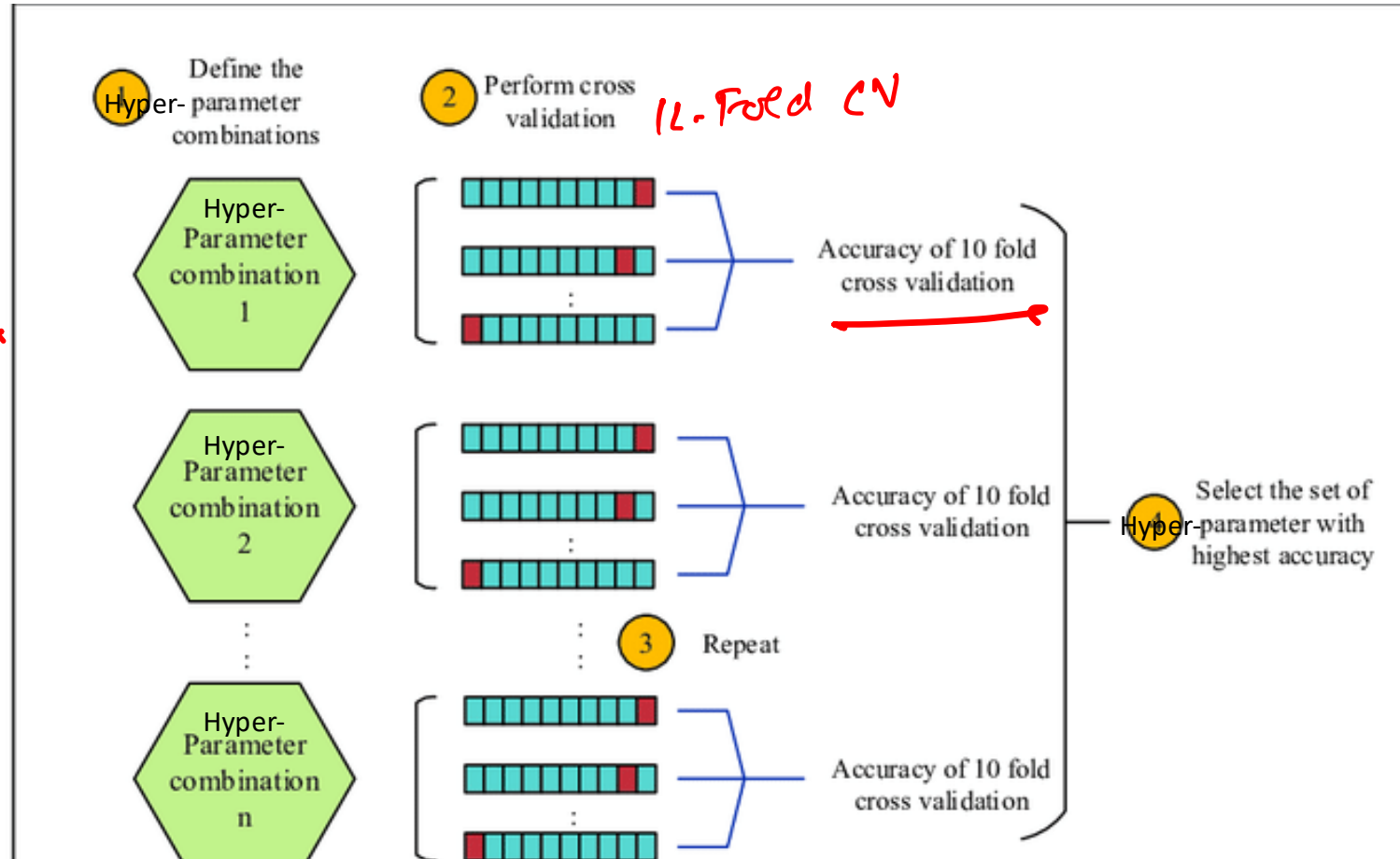


K-Fold Cross Validation



120
 $K=5$

Hyper-parameter Tuning



2 →
3 →
10 →

Audience Poll

1. Ensemble Learning can be used only with Decision Trees
 - True, False
2. Feature importance can be calculated ONLY for tree-based models
 - True, False
3. Permutation feature importance measures the importance of each feature by using an adversarial approach of permuting that feature and evaluating the model
 - True, False