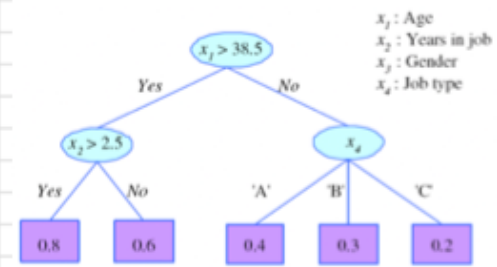


Lecture 15 - 24/10/24

Feature Importance

- BASAL FEATURES ARE MORE GLOBALLY IMPORTANT.
- SOME FEATURES MAY NOT BE USED (FEATURE EXTRACTION)



IN THIS CASE GENDER IS USELESS

Q.A.

- Are D.T. SENSITIVE TO ROTATION?
- How ABOUT SCALING?
- How ABOUT SHIFTING?
- * Are D.T. SENSITIVE TO VARIATION IN DATA?
- Yes, THEY ARE, SLIGHT CHANGES COULD LEAD TO TOTALLY DIFFERENT TREES.

Dealing with Missing Features.

- We MAY HAVE MISSING VALUES FOR SOME FEATURES IN SOME TRAINING SAMPLE.
- ↳ Imagine EACH FEATURE HAS 5% CHANCE OF BEING MISSING.
- ↳ LEADS TO A PROB OF 92.3% IF 50 FEATURES WERE SAMPLED.

Solution → USE SURROGATES!!!

SURROGATE SPLIT: REPLACEMENT FEATURE TO BE USED FOR SPLITTING WHEN FEATURE IS MISSING.

STOP CONDITIONS:

- * A NODE IS FULLY PURE.
- * NO SPLIT REDUCES THE IMPURITY.
- * MIN # OF DATA POINT IN LEAF.
- * MAX # OF NODES IN THE TREE.
- * MAX DEPTH FOR THE TREE.

Variance in Decision Trees

- Decision trees CAN BE GROWN TO COMPLETE PURITY ON TRAINING DATASET
- IT TENDS TO MEMORIZE TRAINING SET, POOR GENERALIZATION.
- HIGH VARIANCE → OVERFITTING
- * BUT THERE IS A REMEDY → REGULARIZATION (SHRINKAGE)
 - ↳ Early Stopping
 - ↳ Pruning

Early Stopping

- RECALL THE STOPPING CONDITIONS ABOVE!
- STOP TREE FROM GROWING TOO LARGE BASED ON SOME STOPPING CRITERIA:
 - MIN # OF DATA POINTS IN A LEAF
 - MAX # OF NODES IN THE TREE
 - MAX DEPTH FOR THE TREE

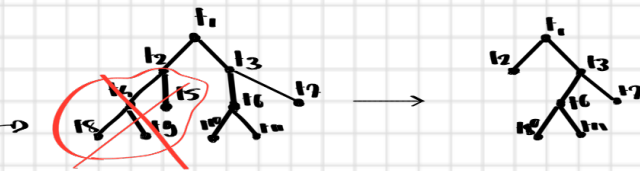
TREE TOO RESTRICTED = HIGH BIAS

Pruning → Post-pruning

- Grow a **LARGE TREE**, THEN **PRUNE IT BACK**

✓ SUBTREE:

- **COMPARE** SUBTREE TO A **LEAF**
- **COMPARE** TREE WITH SUBTREE IN TERMS OF **COST-COMPLEXITY**
- IF **LEAF (WITHOUT SUBTREE)** IS **AS GOOD AS SUBTREE** → **DISCARD SUBTREE, KEEP LEAF.**



PROS VS CONS → OF DECISION TREES

PROS

- **WHITE BOXES** → **VERY EASY TO EXPLAIN.**
- CAN HANDLE BOTH **CATEGORICAL & NUMERICAL.**
- CAN **MATCH ANY DISTRIBUTION.**
- **INDEPENDENT OF SCALING, CENTERING.**
- **INHERENT FEATURE SELECTION.**
- **CHEAP**

CONS

- **NOT AS GOOD IN PREDICTION** AS OTHER STUFF
- **VERY DECISION BOUNDARY**
- **SENSITIVE TO ROTATION!!** (VARIATION IN DATA)
- **BIAS/VARIANCE PROBLEM**

ENSEMBLE METHODS

THE LAW OF LARGE NUMBERS

- IF WE SAMPLE **INDEPENDENT & IDENTICALLY DISTRIBUTED** RANDOM VARIABLES:

$$\frac{1}{m} \sum_{i=1}^m x_i \rightarrow \bar{x} \quad \text{AS } m \rightarrow \infty$$

WHERE \bar{x} IS THE **TRUE EXPECTED VALUE** OF THE DIST.

IMAGINE A SCENARIO WHERE:

$$\text{Coin } M_{\text{oss}} \begin{cases} \text{HEAD} & P(\text{HEAD}) = 0.51 \\ \text{TAIL} & P(\text{TAIL}) = 0.49 \end{cases}$$



WE WANNA
DECIDE
BASED ON
majority!

WEAK CLASSIFIERS

- ASSUME **NOT-DO-GOOD CLASSIFIERS**
- **INDIVIDUALLY ONLY CORRECT 51% OF THE TIME**
- WE USE **MAJORITY VOTE METHOD.**

# classifiers	# True pred.	Probability of correct prediction	Probability of correct prediction
1	1	0.51	0.51
	0	0.49	
3	3	0.51^3	
	2	$3 * 0.49 * 0.51^2$	0.515
	1	$3 * 0.51 * 0.49^2$	
	0	0.49^3	

with only 1 classifier

with 3 classifiers
(at least 2 out of 3 is correct)

- With nearly 1000 classifiers, probability of correct prediction comes close to 0.75

with 1000 classifiers

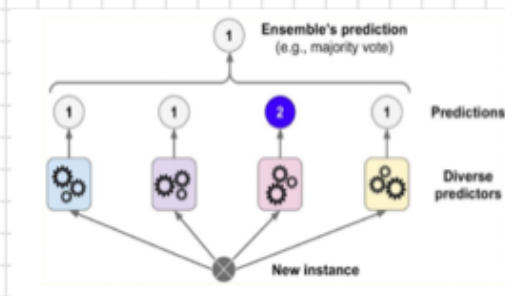
Voting Methods

THERE ARE **TWO MAIN VOTING METHODS**:

- **HARD VOTING** (majority rule)
- **SOFT VOTING**

HARD VOTING

- Each classifier makes its separate prediction on the test data
- Majority rule → The final pred is the class that gets most votes!



SOFT VOTING

- Requires that all classifiers can estimate posterior class probs $p(k|x)$
- Final prediction is the class with highest probs, **AVG** over all classifiers
- Often performs better than Hard Voting

Example:

- We have three classifiers and two classes:

Classifier 1 estimates a posterior probability 91% for A
Classifier 2 estimates a posterior probability 49% for A
Classifier 3 estimates a posterior probability 49% for A

Hard Voting:

1 vote for A
2 votes for B → (B)
3 votes for B

Soft Voting:

Prob of A is:

$$\frac{(91 + 49 + 49)}{3} = 63\% \rightarrow (A)$$

Now onto Ensemble Methods

- A group of separate classifiers or regressors:
 - They form an ensemble (committee)
- An ensemble can perform better than any of its individual base-learners.
- Any type of ML model can be used as base-learner
- Can be used both for classification & regression.

Popular Ensembles:

BAGGING Random Forest BOOSTING STACKING

- An ensemble mainly aims to reduce variance (Law of Large #)
- Another aims to reduce bias.

REDUCING VARIANCE

- AS SAID AN ENSEMBLE PERFORMS BETTER THAN INDIVIDUALS.

BUT ONLY IF BASE-LEARNERS ARE COMPLETELY **INDEPENDENT**.

- ERRORS SHOULD BE UNCORRELATED
- IDENTICAL LEARNERS DO NOT INCREASE PERFORMANCE.

We need diversity to get closer to the independence condition!

Introducing diversity:

* DIVERSITY IN PREDICTORS

- HETEROGENEOUS ENSEMBLE (^{type} ≠ CLASSIFIERS) TRAINED ON THE SAME DATA

* DIVERSITY IN TRAINING DATA

- HOMOGENEOUS ENSEMBLE (^{type} SAME CLASSIFIERS) TRAINED ON ≠ DATA

Diversity in the predictor

Use a heterogeneous ensemble (different types of classifiers) trained on the same dataset.



Bootstrapping



This is a way of varying (diversifying) the training data for each predictor

It is difficult to have many datasets...