

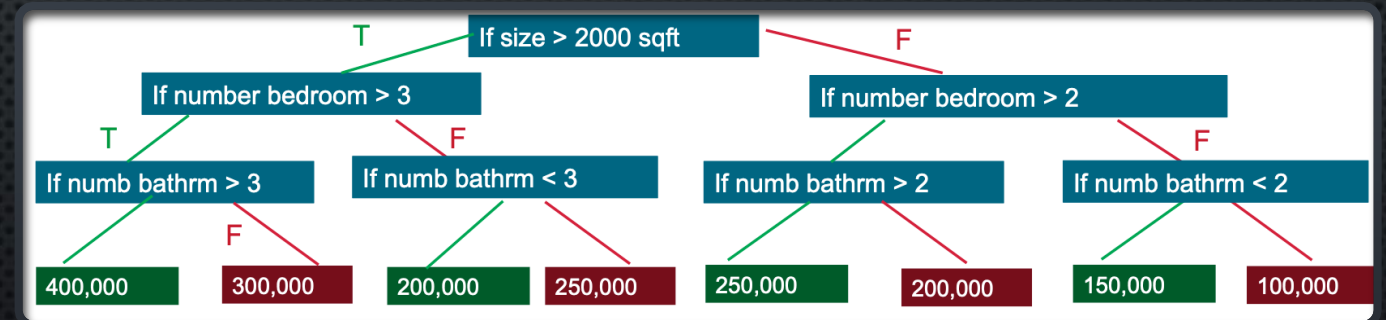
XGBOOST

GROUP 4

NICK KARTSCHOKE, ALEX MORA, CHRIS NG, JAKOB THUNEN

WHAT IS XGBOOST?

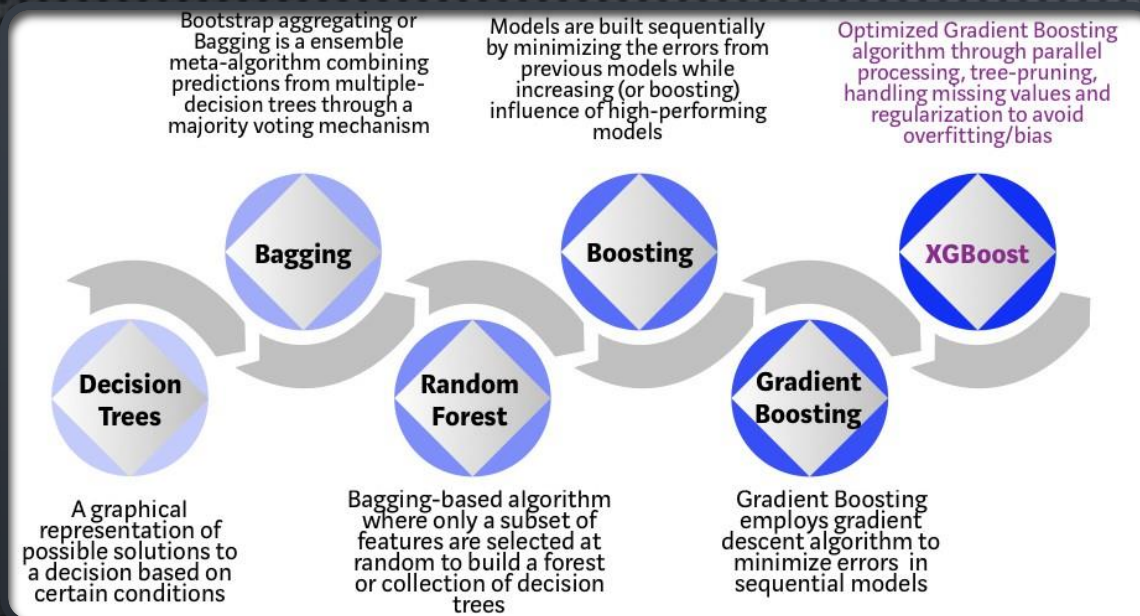
XGBOOST, SHORT FOR "EXTREME GRADIENT BOOSTING", IS A MACHINE LEARNING ALGORITHM BASED ON COLLECTIONS OF RANDOM DECISION TREES, LIKE THE ONE PICTURED HERE.



HOW DOES IT WORK?

XGBOOST IS SIMILAR TO OTHER RANDOM DECISION TREE-BASED ALGORITHMS LIKE RANDOM FOREST IN THAT IT USES A COLLECTION OF RANDOM DECISION TREES TO GENERATE ITS PREDICTIVE MODEL. UNLIKE A RANDOM FOREST, HOWEVER, IT "BOOSTS" THE FOREST BY COMBINING SMALLER, WEAKER TREES INTO LARGER, STRONGER ONES.

BAGGING AND BOOSTING



"BAGGING" IN A RANDOM FOREST ALGORITHM REFERS TO CHOOSING RANDOM SAMPLES OF DATA TO CONSTRUCT DECISION TREES, MINIMIZING VARIANCE AND OVER-FITTING.

"BOOSTING" REFERS TO THE COMBINATION OF WEAKER TREES TO CREATE A STRONGER ONE, AS IN GRADIENT-BOOSTED DECISION TREE MODELS. THIS REDUCES BIAS AND UNDER-FITTING.

XGBOOST UTILIZES BOTH OF THESE TECHNIQUES TO GET THE BEST OF BOTH ALGORITHMS, LEADING IT TO BE THE MODEL OF CHOICE FOR NUMEROUS ML COMPETITIONS.

WHAT ELSE DOES IT DO?

UNLIKE GBDTs, XGBOOST TREES ARE BUILT IN PARALLEL TO EACH OTHER INSTEAD OF ONE AT A TIME. THE MODEL IS ANALYZED PER LAYER INSTEAD OF PER TREE, ENSURING THE OPTIMAL QUALITY OF EACH SPLIT IN THE DECISION TREES.

BUILDING THE DECISION TREES IN PARALLEL ALSO INCREASES EFFICIENCY AND SCALABILITY BY ALLOWING FOR PARALLEL PROCESSING. THIS EFFECT IS ESPECIALLY PROFOUND WHEN USING GPUS WITH THOUSANDS OF PROCESSING CORES TO CREATE MODELS.

ADVANTAGES/DISADVANTAGES

ADVANTAGES

- TREE PRUNING
- HANDLING OF MISSING/SPARSITY-AWARE: XGBOOST CONTAINS A SPARSITY-AWARE SPLIT FINDING ALGORITHM AND HANDLES DIFFERENT TYPES OF SPARSITY PATTERNS IN THE DATA
- BUILT-IN CROSS VALIDATION
- GOOD JOB LIMITING OVERFITTING TREES
 - USES LASSO AND RIDGE REGULARIZATION

DISADVANTAGES

- HAS A LOT OF HYPER PARAMETERS, THUS HARD TO TUNE
- CAN BE OVERFITTED IF PARAMETERS ARE NOT IN TUNE
- SENSITIVE TO OUTLIERS
- PARTIALLY A “BLACK BOX”
 - HARD TO INTERPRET THE INNER WORKINGS

PROCESSING STEPS

- IN ORDER FOR XGBOOST TO BE ABLE TO USE THE DESIRED DATASET, A TRANSFORMATION IS NEEDED.
 - DMATRIX FORMAT
 - TRANSFORMS A NUMPY ARRAY TO DMATRIX FORMAT
- WITH XGBOOST (SPARSE AWARE), WILL CREATE A DEFAULT DIRECTION FOR MISSING VALUE NODES. AT THE PREDICTION TIME, IF THE PATH GOES THROUGH A MISSING VALUE, THE DEFAULT PATH IS FOLLOWED. IT IS GOOD PRACTICE TO CLEAN MISSING DATA, BUT XGBOOST DOES ALLOW FOR MISSING VARIABLES.
- TREE BASED CLASSIFIER - THE NEED TO STANDARDIZE FEATURES IS NOT REQUIRED, BUT IF THERE IS A WANT FOR HIGH PERFORMANCE, IT IS ENCOURAGED.

```
D_train = xgb.DMatrix(X_train, label=Y_train)
D_test = xgb.DMatrix(X_test, label=Y_test)
```

HYPERPARAMETERS

- GENERAL PARAMETERS – GUIDE THE OVERALL FUNCTIONING OF THE MODEL AND RELATE TO WHICH BOOSTER IS BEING USED
- BOOSTER PARAMETERS – GUIDE EACH BOOSTING STEP
- TASK PARAMETERS – USED TO DEFINE THE OPTIMIZATION OBJECTIVE BY SPECIFYING THE LEARNING TASK AND THE CORRESPONDING LEARNING OBJECTIVE

COMMON HYPERPARAMETERS

- BOOSTER — DEFAULT IS GBTREE, COULD ALSO BE GBLINEAR OR DART
- OBJECTIVE — DETERMINES THE LOSS FUNCTION TO BE USED I.E. REG:LINEAR OR REG:LOGISTIC
- COLSAMPLE_BYTREE — PERCENTAGE OF FEATURES USED PER TREE (HIGH VALUE CAN LEAD TO OVERFITTING)
- LEARNING_RATE (ETA) - $[0-1]$ A HIGHER LEARNING RATE MEANS FASTER COMPUTATION, A LOWER LEARNING RATE MEANS YOU CAN GET CLOSER TO THE PERFECT MODEL FIT
- MAX_DEPTH — DETERMINES HOW DEEP EACH TREE CAN GROW WHILE BOOSTING
- ALPHA — $[0-\infty]$ THE HIGHER THE ALPHA, THE MORE CONSERVATIVE THE MODEL/THE MORE REGULARIZATION
- N_ESTIMATORS — NUMBER OF TREES YOU WANT TO BUILD