

CS 307

SPRING 2024

DALPIAZ

ENSEMBLES

MAIN IDEA

- FIT MANY MODELS TO DATA
- MAKE PREDICTION w/ EACH
- "AVERAGE" PREDICTIONS

"BASE LEARNER"

WHY?

$$X_1, X_2, \dots, X_n \sim F$$

$$E[X_i] = \mu$$

$$V[X_i] = \sigma^2$$

$$\bar{X} = \frac{1}{n} \sum X_i$$

$$E[\bar{X}] = \mu$$

$$V[\bar{X}] = \sigma^2/n + [\text{Covariance Term}]$$

IF X_i NOT IND

RANDOM FOREST

- FIT MANY DEEP TREES IN PARALLEL
- BOOTSTRAP DATA FOR EACH TREE
- RANDOM SPLITS
- "BIG" FILE SIZE
- SLOWER THAN SINGLE TREE TO TRAIN
- FEATURE IMPORTANCE
- LOW TUNING

BOOSTING

- SHALLOW TREES IN SEQUENCE
- LOTS OF TUNING
- WIN KAGGLE!
 - XGBOOST
 - LIGHTGBM
 - CATBOOST

BOOTSTRAP

"A BOOTSTRAP RESAMPLE"

SAMPLE THE ROWS
w / REPLACEMENT