## Recall:

- Bagging
- Random Forest
- Boosting          —o    One group of ensemble mainly aims to **REDUCE VARIANCE!**
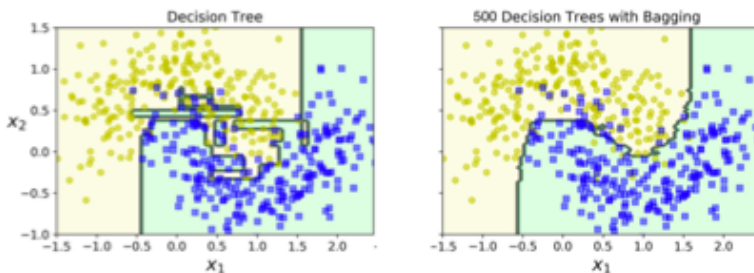- Stacking               Another one aims to **REDUCE BIAS!**

## Bagging

- meaning → **Bootstrap Aggregation**

- **Training:**
  o Make **Bootstrapped subsets** of the training data
  o Train a **separate model on each subset.**
  o Given that samples come from same dataset, **not truly indep.**

> NOTE:
> {  Creating diversity in the data, introduces diversity in the models

- **Prediction:**
  o **combination of predictions of all models**

- **Final Step:**
  o typically use **HARD VOTING** to combine predictions for classifications.
  o typically use **AVERAGING** for regression.

- **Scenario:**
  o I have a classifier which has **HIGH VARIANCE PROB.**
  o It has **low BIAS.**
  o I have **computational power available**

### Try Bagging!!!

An example of Decision Tree (left) vs. a bagging ensemble of 500 trees (right)



We can see how we provide a great way to **AVOID OVERFITTING**

- **Validation:**
  o **Out-of-Bag evaluation**
  o We know that in **BAGGING** we use **Bootstrapped data** → **Sampling with Replacement**
    - Not all datapoints will be in a given bag!
    - ≠ datapoint are incl./exclud. in ≠ bags.
  o On **avg 37%** are **Out-of-bag samples** if bags are the same size as dataset
  o Out of Bag samples **will be used for validation!**
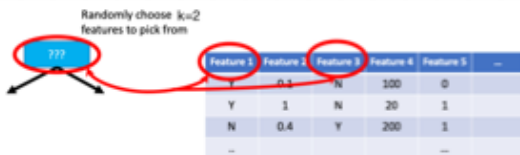  o ∀ data point in original set we **combine pred** where it wasn't sampled

- Validation Error is the AVE error across these combined preds.
- the OOB prediction error is an unbiased estimation of the test error.
- this method can REPLACE cross-validation.

- **ADVANTAGES:**
  - Easy to implement
  - Training in parallel
  - Pred is ave or vote
  - Reduces variance
  - Out-of-bag validation

# Random Forest

- Essentially a BAGGING DECISION TREE with MODIFIED SPLITTING CRITERION
- **PROCESS:**
  - MAKE Bootstrapped DATA SETS FROM TRAINING SET
  - $\forall$ DATASET, TRAIN A FULL DECISION TREE WITH MODIFIED SPLIT STRATEGY:
    - Before finding each split:
      * Randomly choose k features
      * Only consider those feature or split
  - Output of random forest is the AGGREGATED output of all trees



Random Forest are nice because!

- Only few Hyper-parameters:
  - m = # of bootstrapped datasets.
  - k = # of randomly selected features.
- Based on decision trees which avoid hard process.
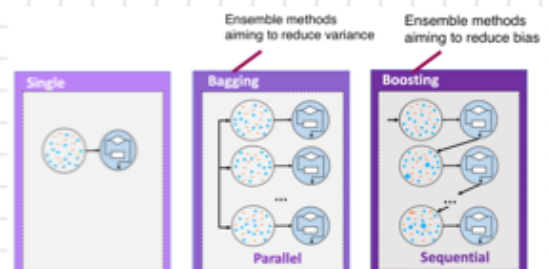- Average depth estimates feature importance!

# Extra Trees

- MEANING → Extremely Randomized Trees
- Unlike bagging & random forest, they fit each D.T. on WHOLE DATASET. ~~Bootstrap~~
- Unlike random forest, their algo picks a random cut-point at each of the k features.
- Like random forest, their algo randomly samples k features to find the good split.

* We know that Bagging & associated methods → REDUCE VARIANCE
* What about reducing Bias though?

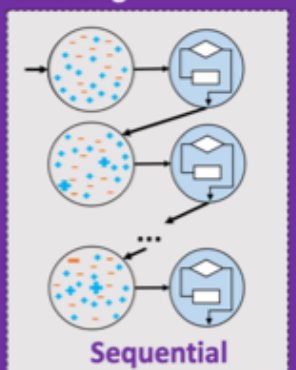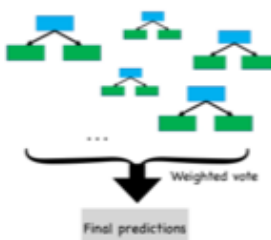# Use Boosting !!!

- A BIT BETTER THAN RANDOM GUESS

- **SEQUENCE** OF (WEAK LEARNERS) COMBINED INTO A STRONG LEARNER

- **STEPS:**
  - **T**RAIN A WEAK LEARNER
  - GET TRAINING ERROR
  - TRAIN NEXT WEAK LEARNER BASED ON THAT ERROR
  - ⋮
  - FINAL PREDICTIONS: A COMBINATION OF ALL TRAINED WEAK LEARNERS.

- VARIOUS BOOSTING METHODS EXIST → AdaBoost, GradientBoosting...

**Boosting**



Sequential

# AdaBoost



Weighted vote

Final predictions

AFTER MAKING EACH TREE, THE ERRORS:
- ○ INFLUENCE WEIGHTS OF NEXT TRAINING.
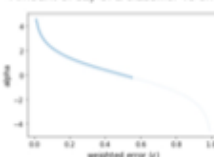- ○ DETERMINE WEIGHT OF TREE IN FINAL.

# AdaBoost
## Training

Start with equal weights for all the training samples, e.g., $w_i^{(1)} = \frac{1}{N}$

For $c = 1, \ldots, M$ iterations (trees):

- ○ Train a model $h_c(x)$ based on the weighted training samples
- ○ Calculate the total error on the training set: $\epsilon_c = \sum_{i, h_c(x_i) \neq y_i} w_i^{(c)}$
- ○ Based on the total error, calculate the model's amount of say: $\alpha_c = ln(\frac{1 - \epsilon_c}{\epsilon_c})$
- ○ Update the samples' weights:
  - ○ Increase for misclassified samples: $w_i^{(c)} = w_i^{(c)} e^{\alpha_c}$
  - ○ Unchanged for correct classifications: $w_i^{(c)} = w_i^{(c)}$
- ○ Normalize all weights so they sum up to 1: $w_i^{(c+1)} = \frac{w_i^{(c)}}{\sum_{j=1}^{n} w_j^{(c)}}$

Amount of say of a classifier vs error



weighted error (ε)

Some implementations add a hyper-parameter (learning rate) here to avoid overfitting

Some implementations also decrease the weights of correctly classified samples and then normalize
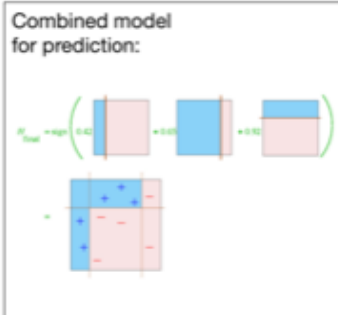
**two types of weights:**
- $w_i$ = TRAINING SAMPLE
- $\alpha_c$ = ∀ MODEL

**PREDICTION:**
$$SIGN \left( \sum_{c=1}^{M} \alpha_c h_c(x) \right)$$
+1 OR -1

**WEAK LEARNERS = STUMPS**

## Example AdaBoost
with decision stumps as base classifier

Training in three iterations:



Iteration 1:   Iteration 2:   Iteration 3:

Combined model for prediction:



We said "Train a model $h_c(x)$ based on the weighted training samples"
But how do we incorporate the weights into the training process?

Think about it!



WE CAN INCORPORATE THE WEIGHTS INTO THE ASSESSMENT OF SPLITS.

A DATAPOINT IS NO MORE WORTH 1 UNIT, IT COUNTS AS ITS WEIGHT $w_i$ IN IMPURITY CALCULATION.

⎛ → THEIR WEIGHTED AVG = $P_L \phi(t_L) + P_R \phi(t_R)$

$$P_K = \sum_{i: y_i = K} w_i$$

# GRADIENT BOOSTING

· Commonly used for both REGRESSION & CLASSIFICATION.



## Gradient Boost:
### small steps on the right direction

- Sequential algorithm like AdaBoost
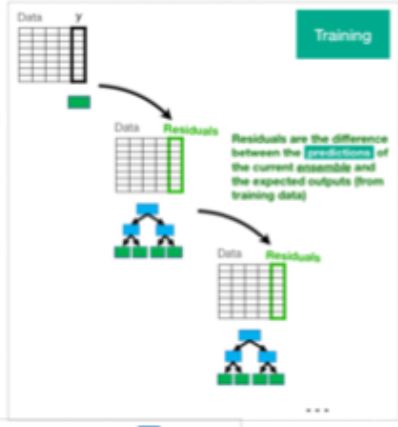- Use of tress (max-depth=3) as weak learner is common

- **Training:**
  - Start with a single leaf (initial prediction, e.g., the average of all labels)
  - use the current **ensemble** to predict the training data and compute the residuals
  - Fit a new tree to the residuals

Side note: Stochastic gradient boosting is a variation of the method that only uses a subset of data in each iteration

Residuals are the difference between the predictions of the current ensemble and the expected outputs (from training data)

- **Prediction:** Predictions = ◼ + (Learning rate) × 🌳 + (Learning rate) × 🌳 ...

learning rate (often 0.1) must be used to avoid overfitting

# BAGGING VS BOOSTING

| BAGGING | BOOSTING |
|---|---|
| Usually HELPS DECREASING VARIANCE | HELPS DECREASING BIAS |
| Not much EFFECT WITH STABLE models | DOESN'T SOLVE OVERFITTING |
| Easy to PARALLELIZE | SOMETIMES HURTS PERFORMANCE |
| | CANNOT BE FULLY PARALLELIZED |

Bottom Line: IT DEPENDS!!

# STACKING

· Instead of VOTING FOR AGGR. → TRAIN META-LEARNER.

ONE OF OUR CHOICE

SHOULD GET PREDS OF INDIVIDUAL MODELS AS INPUT & EXPECTED LABELS FOR DATA SAMPLES AS OUTPUT.