

INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 29

HASSAN ASHTIANI

ENSEMBLE LEARNING

SIMPLE RECIPE TO IMPROVE THE PERFORMANCE:

1. TRAIN MULTIPLE CLASSIFIERS
2. AGGREGATE THEIR DECISIONS
 - E.G., VOTING

THE RESULT CAN BE BETTER THAN
THE INDIVIDUAL CLASSIFIERS!

-  REDUCING VARIANCE?
-  REDUCING BIAS?



THE NETFLIX CHALLENGE

- GIVEN SOME USER RATINGS FOR VARIOUS FILMS,
 - PREDICT USER RATINGS FOR OTHER FILMS
- COLLABORATIVE FILTERING

THE NETFLIX CHALLENGE

- 2006: COMPETITION BEGAN, 1M USD FOR IMPROVING 10% OVER NETFLIX'S OWN METHOD
- 2007: 8.43% IMPROVEMENT (BELLKOR WON 50k)
- 2008: NO INDIVIDUAL TEAM BETTER THAN 9.43%
 - BELLKOR+BIGCHAOS MERGED... >9.43% IMPROVEMENT!
- 2009: TOP THREE MERGE! BELLKOR+BIGCHAOS+PRAGMATIC >10%
 - NEW TEAM IN THE LAST MONTH: GRANDPRIZETEAM (9.46%)
 - ANYONE COULD JOIN, AND SHARE THE PRIZE BASED ON THE IMPROVEMENT
 - ENSEMBLE: GRANDPRIZETEAM+VANDERLAY INDUSTRIES (>10%) ←
BELLKOR+BIGCHAOS+PRAGMATIC (10.09%) AND ENSEMBLE (10.10%)
 - THEY BOTH GET THE EXACT SAME ACCURACY ON THE PRIVATE TEST SET!
 - BELLKOR+BIGCHAOS+PRAGMATIC SUBMITTED 20MINS EARLIER....
- 2007: LINKAGE ATTACKS WITH IMDB
- 2010: PRIVACY CONCERNS AND CLASS-ACTION LAWSUITS...COMPETITION WAS CANCELED

AGGREGATION

- WHEN DOES COMBINING MODELS HELP?
 - THE BASE LEARNERS SHOULD BE ACCURATE ↙
 - THE BASE LEARNERS SHOULD BE DIVERSE (LESS CORRELATED) ↙
- EXAMPLE FOR CLASSIFICATION

1	✓	✓	✗	✓	✓
2	✓	✓	✓	✗	✓
3	✗	✓	✓	✓	✓
	✓	✓	✓	✓	✓

AGGREGATION

- EXAMPLE FOR REGRESSION

1	1.1		
2	1.5		
3	0.9		

median/mean/...

- AGGREGATION CAN REDUCE THE VARIANCE
 - HELPS TACKLING OVERFITTING
- HOW TO DIVERSIFY THE BASE LEARNERS?

THE ENSEMBLE

HOW TO CREATE A DIVERSE ENSEMBLE OF LEARNERS?

- DIFFERENT CLASSIFIERS (NEURAL NETS, LINEAR, NEAREST NEIGHBOR,...)
- DIFFERENT HYPER-PARAMETERS
 - WEIGHT INITIALIZATION IN NEURAL NETWORKS (RANDOM SEED)
 - NETWORK ARCHITECTURES
- DIFFERENT TRAINING SUBSETS
- DIFFERENT FEATURE SUBSETS

BAGGING

- USING NON-OVERLAPPING TRAINING SUBSETS CREATES TRULY INDEPENDENT/DIVERSE CLASSIFIERS
 - I.I.D. ASSUMPTION!
- BUT CAN BE WASTEFUL
 - EACH CLASSIFIER IS TRAINED USING ONLY A SMALL TRAIN SET...
- BAGGING (BOOTSTRAP AGGREGATING)
 - RANDOM SAMPLING WITH REPLACEMENT!

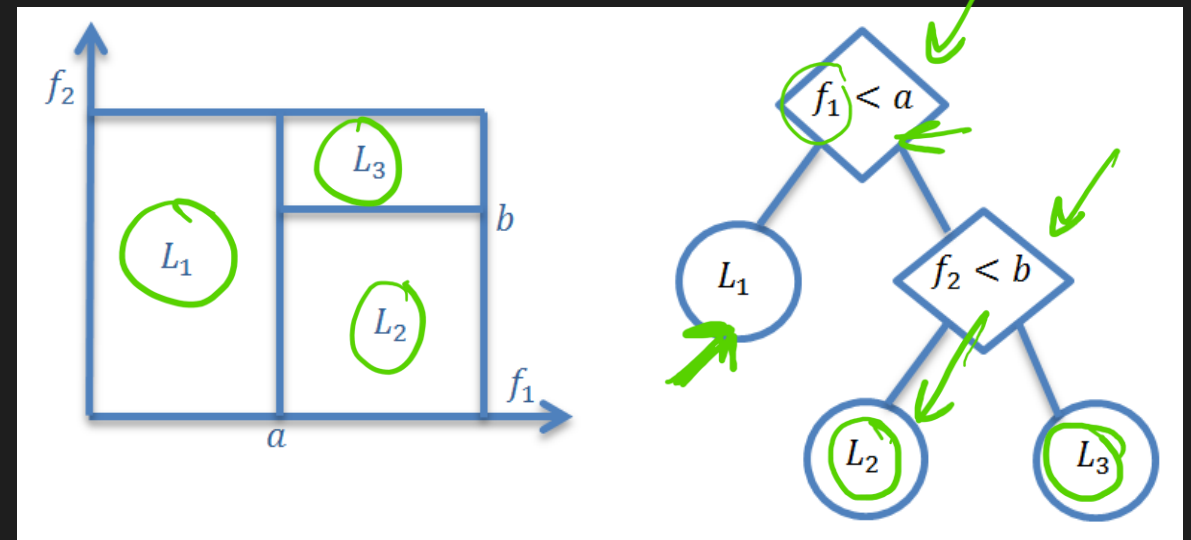
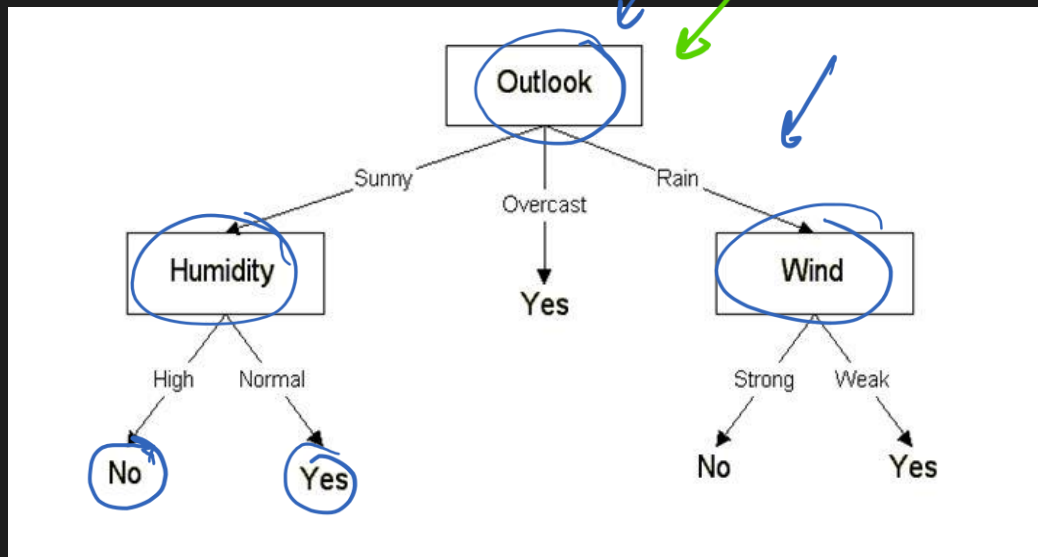
RANDOM SUBSPACE METHOD

- TRAIN EACH CLASSIFIER USING A RANDOM SUBSET OF FEATURES
- SO EACH CLASSIFIER OPERATES IN A RANDOM SUBSPACE
- ALSO CALLED FEATURE BAGGING, OR ATTRIBUTE BAGGING
- ARE THE CLASSIFIERS DIVERSE?
 - THERE IS CORRELATION BETWEEN THE FEATURES
 - THERE IS ONLY SO MUCH YOU CAN LEARN FROM A DATA POINT

RANDOM FORESTS

- COMBINES THE IDEAS OF BAGGING AND RANDOM SUBSPACE METHODS
- USES DECISION TREES AS BASE CLASSIFIERS

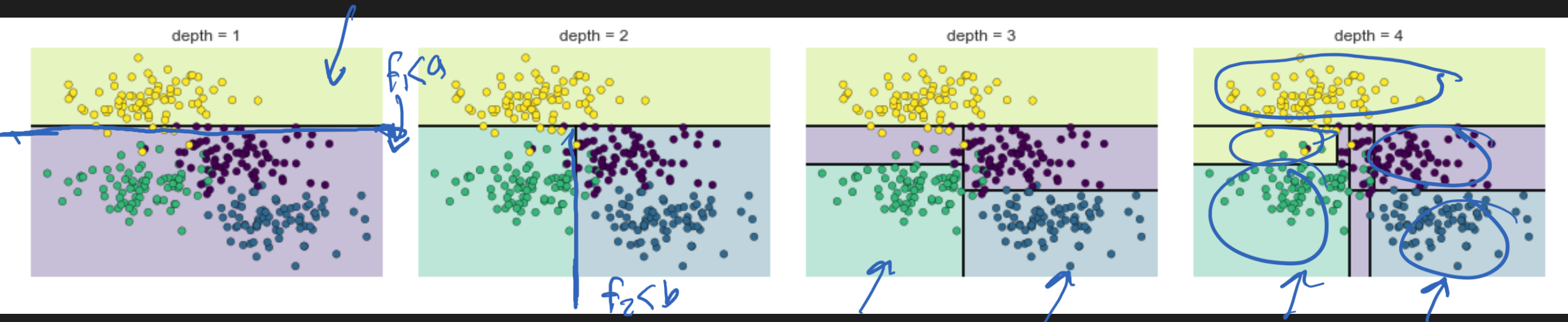
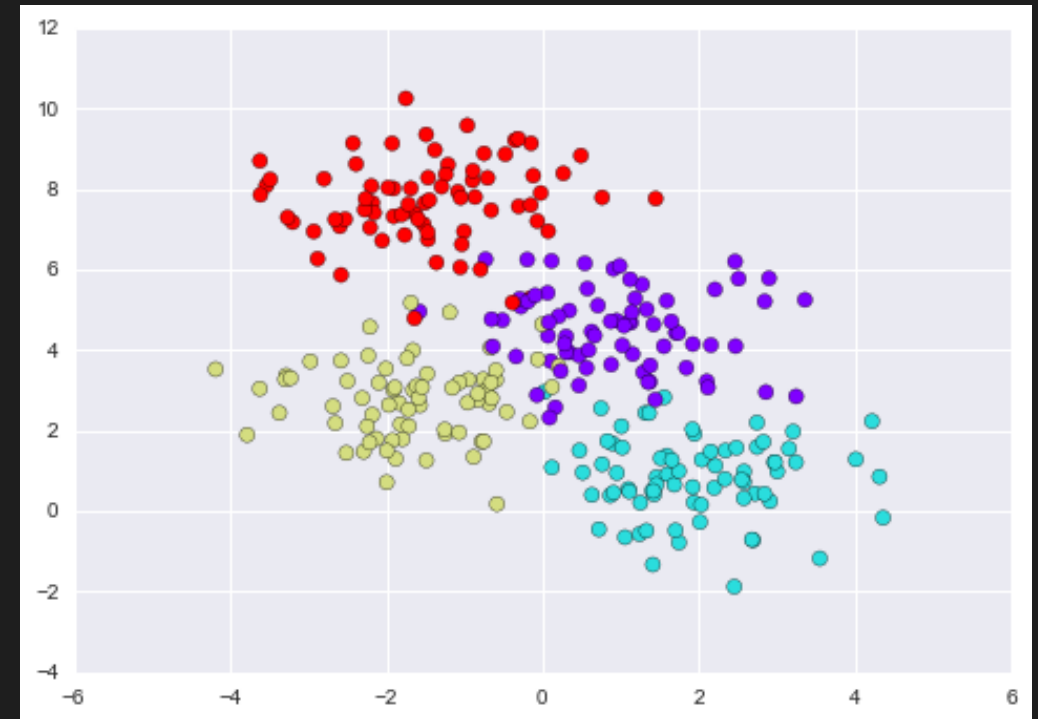
DECISION TREE



- CAN HANDLE CATEGORICAL FEATURES
- DEEPER TREES CAN OVERFIT EASILY
- HOW DO YOU “TRAIN” DECISION TREES?

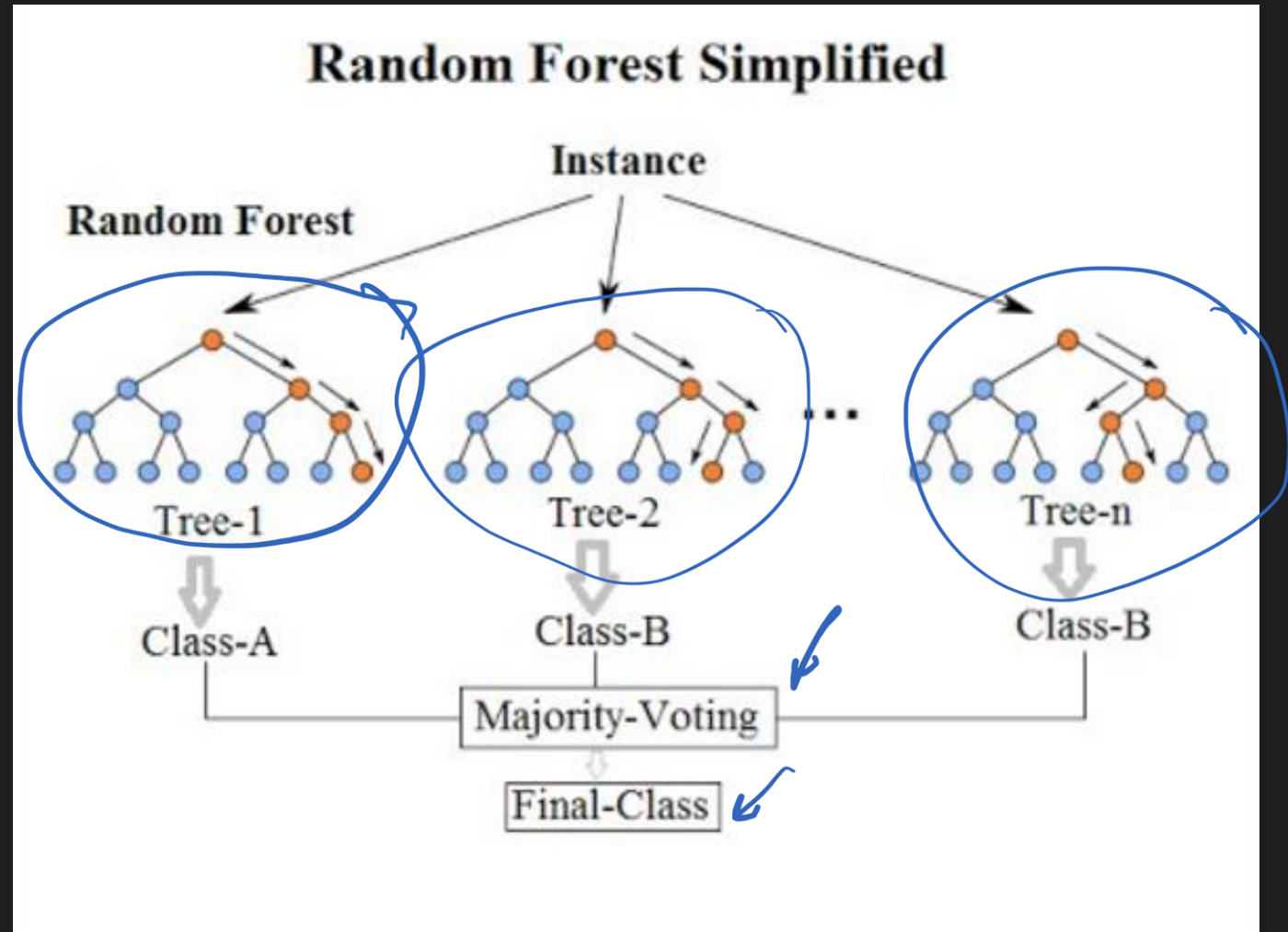
DECISION TREES

- A POSSIBLE APPROACH: SELECT THE TREE NODES RANDOMLY!
- LABEL LEAVES BY DOING MAJORITY VOTE IN THE TRAINING DATA



RANDOM FORESTS

- CREATE MANY DEEP RANDOM TREES
- USE RANDOM SUBSETS OF DATA FOR EACH TREE TO DETERMINE THE LABEL OF LEAVES
- FOR A TEST POINT, TAKE MAJORITY VOTE BETWEEN THE TREES



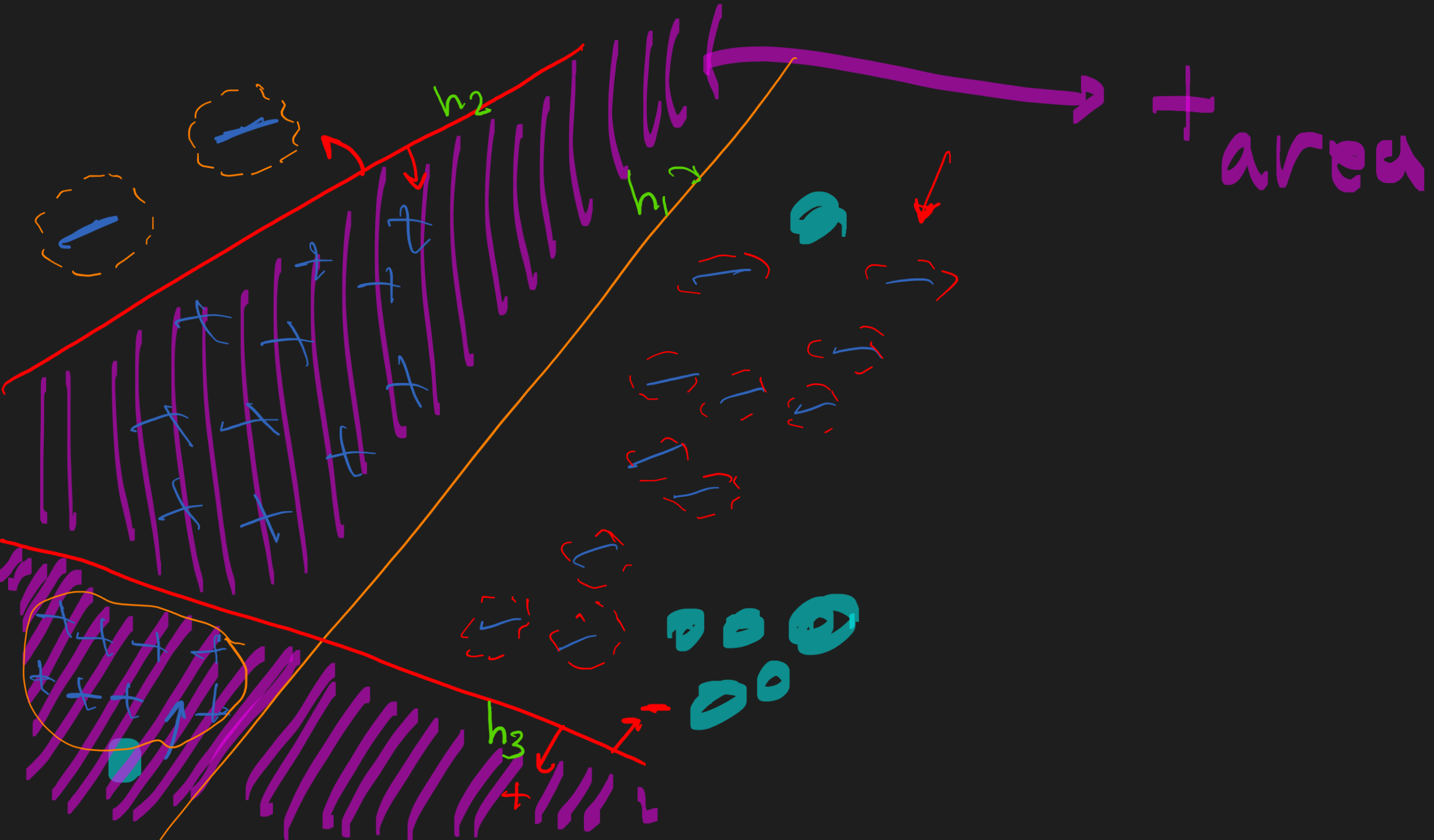
BOOSTING

- UP TO NOW WE PICKED THE BASE CLASSIFIERS INDEPENDENTLY
- THE GOAL WAS TO REDUCE VARIANCE
- BUT CAN WE COMBINE CLASSIFIERS TO REDUCE BIAS?
 - A “STRONGER CLASSIFIER” OUT OF “WEAK CLASSIFIERS”?

BOOSTING

A GREEDIER APPROACH

- PICK THE BASE CLASSIFIERS ONE-BY-ONE (INCREMENTALLY)
- EACH NEW CLASSIFIER (CALLED A WEAK LEARNER) TRIES TO ADDRESS THE SHORTCOMINGS OF THE PREVIOUS ONES
- THE COMBINATION OF “WEAK LEARNERS” CAN BE A “STRONG LEARNER”



TRAINING ON A WEIGHTED DATA SET

- REGULAR TRAINING

- $\text{MIN}_{\theta} \frac{1}{n} \sum_{i=1}^n l(y^i, h_{\theta}(x^i))$

- WEIGHTED TRAINING

- $\text{MIN}_{\theta} \sum_{i=1}^n \underbrace{D_i}_{\text{weight}} \cdot l(y^i, h_{\theta}(x^i))$

- WE CAN PUT MORE EMPHASIS ON SOME OF THE TRAINING POINTS

BOOSTING

1. INITIALIZE THE WEIGHTS OF ALL TRAINING POINTS TO BE EQUAL
2. DO FOR A NUMBER OF ITERATIONS:

- TRAIN A WEAK LEARNER FOR THE WEIGHTS (FROM THE BASE CLASS)
- STORE THE ACCURACY OF THIS WEAK LEARNER (α_i) $1 - \epsilon_i = \alpha_i$
- SEE WHERE THE LEARNER MAKES MISTAKES
- INCREASE THE WEIGHTS OF THOSE MISCLASSIFIED POINTS (D_j)
- SO THAT THEY ARE CLASSIFIED CORRECTLY IN THE NEXT ROUNDS

THE FINAL CLASSIFIER IS A WEIGHTED MAJORITY OF ALL WEAK CLASSIFIERS WHERE THE WEIGHTS ARE PROPORTIONAL α_i

AdaBoost

$$y_i = \{-1, +1\}, \quad y_i = h(x_i)$$

input:

training set $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

weak learner WL

number of rounds T

initialize $\mathbf{D}^{(1)} = (\frac{1}{m}, \dots, \frac{1}{m})$.

for $t = 1, \dots, T$:

invoke weak learner $h_t = \text{WL}(\mathbf{D}^{(t)}, S)$

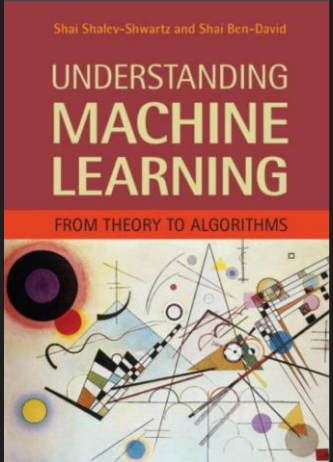
compute $\epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]}$

let $w_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$



update $D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-w_t y_i h_t(\mathbf{x}_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-w_t y_j h_t(\mathbf{x}_j))}$ for all $i = 1, \dots, m$

output the hypothesis $h_s(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^T w_t h_t(\mathbf{x}) \right)$.

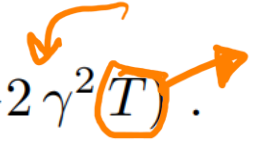
$$\mathbb{1}_{[y_i \neq h_t(\mathbf{x}_i)]} = \begin{cases} 1 & y_i \neq h_t(\mathbf{x}_i) \\ 0 & y_i = h_t(\mathbf{x}_i) \end{cases}$$



BOOSTING THEORY

- IF ALL THE INTERMEDIATE WEAK LEARNERS ARE BETTER THAN RANDOM (E.G., ERROR < 49% FOR BINARY CLASSIFICATION)

- THEN THE TRAINING ERROR OF THE COMBINED MODEL CONVERGES QUICKLY TO 0!


THEOREM 10.2 *Let S be a training set and assume that at each iteration of AdaBoost, the weak learner returns a hypothesis for which $\epsilon_t \leq 1/2 - \gamma$. Then, the training error of the output hypothesis of AdaBoost is at most*

$$L_S(h_S) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{[h_S(\mathbf{x}_i) \neq y_i]} \leq \exp(-2\gamma^2 T)$$


BOOSTING THEORY

- SO THE GOAL IS NOT REDUCING THE VARIANCE ANYMORE
- THE GOAL IS REDUCING THE BIAS!
- WHAT ABOUT THE TEST ERROR?