

Don't ask what it means, but rather how it is used.

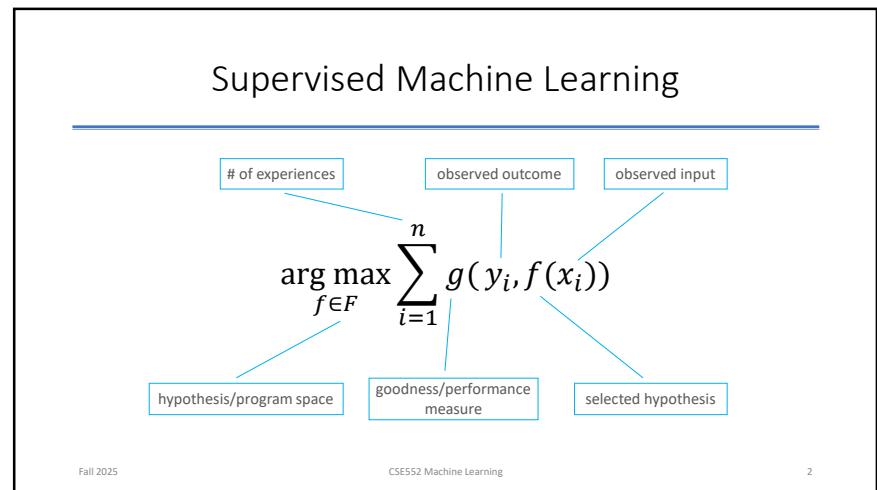
- L. Wittgenstein

CSE 552

Machine Learning

Spring 2025
Ensemble Learning
© 2013-2025 Yakup Genc

1



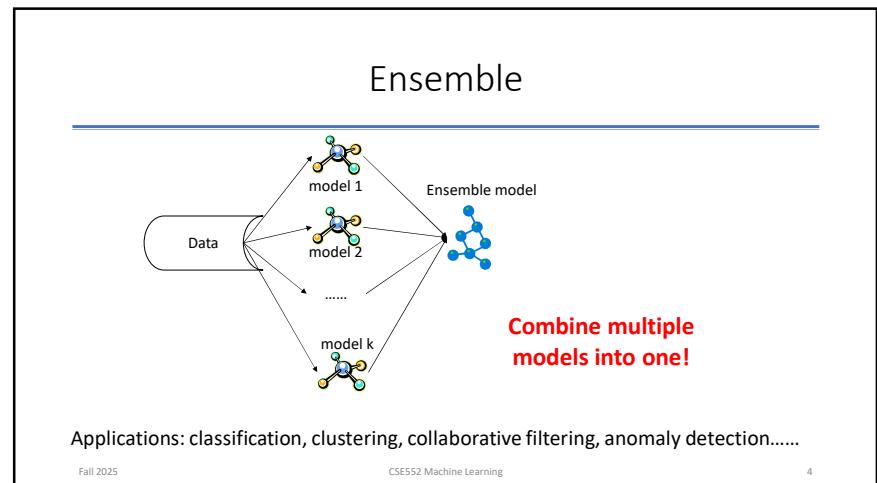
2

Slide Credits

- Gao et al.: Jing Gao, Wei Fan, Jiawei Han,
<http://ews.uiuc.edu/~jinggao3/sdm10ensemble.htm>
- Holloway: Todd Holloway, 2007.

Fall 2025 CSE552 Machine Learning 3

3



4

Stories of Success



Fall 2025

- **Million-dollar prize**

- Improve the baseline movie recommendation approach of Netflix by 10% in accuracy
- The top submissions all combine several teams and algorithms as an ensemble

- **Data mining competitions**

- Classification problems
- Winning teams employ an ensemble of classifiers

CSE552 Machine Learning

5

Netflix Prize

- **Supervised learning task**

- Training data is a set of users and ratings (1,2,3,4,5 stars) those users have given to movies.
- Construct a classifier that given a user and an unrated movie, correctly classifies that movie as either 1, 2, 3, 4, or 5 stars
- \$1 million prize for a 10% improvement over Netflix's current movie recommender ($MSE = 0.9514$)

- **Competition**

- Began October 2006
- At first, single-model methods are developed, and performances are improved
- However, improvements slowed down
- Later, individuals and teams merged their results, and significant improvements are observed

CSE552 Machine Learning

6

Netflix Prize

Just three weeks after it began, at least 40 teams had bested the Netflix classifier.

Top teams showed about 5% improvement.

Team Name	Best Score	% Improvement
No Grand Prize candidates yet	—	—
Grand Prize - RMSE <= 0.8563	—	—
How low can he go?	0.9046	4.92
ML@UToronto A	0.9046	4.92
ssoton	0.9099	4.47
wwwzgutting.com	0.9103	4.32
The Thought Gang	0.9113	4.21
MPG Reget	0.9110	4.16
slimonth	0.9145	3.98
Boss_The_Clown	0.9177	3.54
EllipticChaos	0.9179	3.52
datracker	0.9183	3.48
Foresier	0.9214	3.15
bsdfish	0.9229	3.00
Three Blind Mice	0.9234	2.94
Boschimack	0.9238	2.90
Remco	0.9252	2.75
karmatics	0.9301	2.24
Chapeltor	0.9314	2.10
Fimod	0.9325	1.99
mthrox	0.9328	1.96

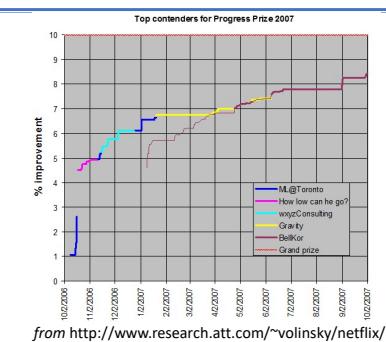
Fall 2025

CSE552 Machine Learning

7

Netflix Prize

However,
improvement
slowed...



CSE552 Machine Learning

8

8

Netflix Prize

Today, the top team has posted a 8.5% improvement.

Ensemble methods are the best performers...

No Progress Prize candidates yet			
Progress Prize - RMSE <= 0.8625			
1	BellKor	0.8705	8.50
Progress Prize 2007 - RMSE = 0.8712 - Winning Team: Korfell			
2	Korfell	0.8712	8.43
3	When Gravity and Dinosaurs Unite	0.8717	8.38
4	Gravify	0.8743	8.10
5	bashfi	0.8746	8.07
6	Dinosaur Planet	0.8753	8.00
7	ML@UToronto_A	0.8787	7.64
8	Arak Paterek	0.8789	7.52
9	NIPS Reedit	0.8808	7.42
10	Just a guy in a garage	0.8834	7.15
11	Ensemble Experts	0.8841	7.07
12	mathematical capital	0.8844	7.04
13	HowLowCanHeGo2	0.8847	7.01
14	The Thought Gang	0.8849	6.99
15	Reel Ingenuity	0.8855	6.93
16	studentfemale	0.8859	6.88
17	NIPS Submission	0.8861	6.86
18	Three Blind Mice	0.8869	6.78
19	TrainOnTest	0.8869	6.78
20	GeoffDean	0.8869	6.76
21	Rootless	0.8872	6.75
22	Paul Harrison	0.8872	6.75
23	ATTEAM	0.8873	6.74
24	wwwconsulting.com	0.8874	6.73
25	ICMLsubmission	0.8875	6.72

Fall 2025 CSE552 Machine Learning

9

Netflix Prize

Rookies

“Thanks to Paul Harrison’s collaboration, a simple mix of our solutions improved our result from 6.31 to 6.75”

No Progress Prize candidates yet			
Progress Prize - RMSE <= 0.8625			
1	BellKor	0.8705	8.50
Progress Prize 2007 - RMSE = 0.8712 - Winning Team: Korfell			
2	Korfell	0.8712	8.43
3	When Gravity and Dinosaurs Unite	0.8717	8.38
4	Gravify	0.8743	8.10
5	bashfi	0.8746	8.07
6	Dinosaur Planet	0.8753	8.00
7	ML@UToronto_A	0.8787	7.64
8	Arak Paterek	0.8789	7.52
9	NIPS Reedit	0.8808	7.42
10	Just a guy in a garage	0.8834	7.15
11	Ensemble Experts	0.8841	7.07
12	mathematical capital	0.8844	7.04
13	HowLowCanHeGo2	0.8847	7.01
14	The Thought Gang	0.8849	6.99
15	Reel Ingenuity	0.8855	6.93
16	studentfemale	0.8859	6.88
17	NIPS Submission	0.8861	6.86
18	Three Blind Mice	0.8869	6.78
19	TrainOnTest	0.8869	6.78
20	GeoffDean	0.8869	6.76
21	Rootless	0.8872	6.75
22	Paul Harrison	0.8872	6.75
23	ATTEAM	0.8873	6.74
24	wwwconsulting.com	0.8874	6.73
25	ICMLsubmission	0.8875	6.72

Fall 2025 CSE552 Machine Learning

10

Netflix Prize

Arek Paterek

“My approach is to combine the results of many methods (also two-way interactions between them) using linear regression on the test set. The best method in my ensemble is regularized SVD with biases, post processed with kernel ridge regression”

http://rainbow.mimuw.edu.pl/~ap/ap_kdd.pdf

No Progress Prize candidates yet			
Progress Prize - RMSE <= 0.8625			
1	BellKor	0.8705	8.50
Progress Prize 2007 - RMSE = 0.8712 - Winning Team: Korfell			
2	Korfell	0.8712	8.43
3	When Gravity and Dinosaurs Unite	0.8717	8.38
4	Gravify	0.8743	8.10
5	bashfi	0.8746	8.07
6	Dinosaur Planet	0.8753	8.00
7	ML@UToronto_A	0.8779	7.64
8	Arak Paterek	0.8789	7.52
9	NIPS Reedit	0.8808	7.42
10	Just a guy in a garage	0.8834	7.15
11	Ensemble Experts	0.8841	7.07
12	mathematical capital	0.8844	7.04
13	HowLowCanHeGo2	0.8847	7.01
14	The Thought Gang	0.8849	6.99
15	Reel Ingenuity	0.8855	6.93
16	studentfemale	0.8859	6.88
17	NIPS Submission	0.8861	6.86
18	Three Blind Mice	0.8869	6.78
19	TrainOnTest	0.8869	6.78
20	GeoffDean	0.8869	6.76
21	Rootless	0.8872	6.75
22	Paul Harrison	0.8872	6.75
23	ATTEAM	0.8873	6.74
24	wwwconsulting.com	0.8874	6.73
25	ICMLsubmission	0.8875	6.72

Fall 2025 CSE552 Machine Learning

11

Netflix Prize

U of Toronto

“When the predictions of multiple RBM models and multiple SVD models are linearly combined, we achieve an error rate that is well over 6% better than the score of Netflix’s own system.”

<http://www.cs.toronto.edu/~rsalakh/papers/bmcf.pdf>

No Progress Prize candidates yet			
Progress Prize - RMSE <= 0.8625			
1	BellKor	0.8705	8.50
Progress Prize 2007 - RMSE = 0.8712 - Winning Team: Korfell			
2	Korfell	0.8712	8.43
3	When Gravity and Dinosaurs Unite	0.8717	8.38
4	Gravify	0.8743	8.10
5	bashfi	0.8746	8.07
6	Dinosaur Planet	0.8753	8.00
7	ML@UToronto_A	0.8779	7.64
8	Arak Paterek	0.8789	7.52
9	NIPS Reedit	0.8808	7.42
10	Just a guy in a garage	0.8834	7.15
11	Ensemble Experts	0.8841	7.07
12	mathematical capital	0.8844	7.04
13	HowLowCanHeGo2	0.8847	7.01
14	The Thought Gang	0.8849	6.99
15	Reel Ingenuity	0.8855	6.93
16	studentfemale	0.8859	6.88
17	NIPS Submission	0.8861	6.86
18	Three Blind Mice	0.8869	6.78
19	TrainOnTest	0.8869	6.78
20	GeoffDean	0.8869	6.76
21	Rootless	0.8872	6.75
22	Paul Harrison	0.8872	6.75
23	ATTEAM	0.8873	6.74
24	wwwconsulting.com	0.8874	6.73
25	ICMLsubmission	0.8875	6.72

Fall 2025 CSE552 Machine Learning

12

Netflix Prize

Gravity

Table 5: Best results of single approaches and their combinations

Method/Combination	RMSE
MF	0.9190
NB	0.9013
CL	0.9006
NB + CL	0.9275
MF + CL	0.9137
MF + NB	0.9089
MF + NB + CL	0.9089

home.mit.bme.hu/~gtakacs/download/gravity.pdf

Fall 2025 CSE552 Machine Learning

13

Netflix Prize

When Gravity and Dinosaurs Unite

"Our common team blends the result of team Gravity and team Dinosaur Planet."

Might have guessed from the name...

Method/Combination	RMSE
BellKor	0.8705
KorBell	0.8712
When Gravity and Dinosaurs Unite	0.8713
Gravity	0.8743
Dinosaur Planet	0.8746
ML@UToronto A	0.8753
AxesPaterek	0.8787
NIPS Reedit	0.8808
Just a guy in a garage	0.8834
Ensemble Experts	0.8841
mathematical capital	0.8844
HowLowCanHeGo2	0.8847
The Thought Gang	0.8849
Reel Ingenuity	0.8855
studentfemale	0.8859
NIPS Submission	0.8861
Three Blind Mice	0.8869
TrainOnTest	0.8869
GeoffDean	0.8869
Rookee	0.8872
Paul Harrison	0.8872
ATTEAM	0.8873
wyzconsulting.com	0.8874
ICMLsubmission	0.8875

Fall 2025 CSE552 Machine Learning

14

Netflix Prize

BellKor / KorBell

And, yes, the top team which is from AT&T...

"Our final solution (RMSE=0.8712) consists of blending 107 individual results."

Fall 2025 CSE552 Machine Learning

15

Leaderboard

"Our final solution (RMSE=0.8712) consists of blending 107 individual results."

"Predictive accuracy is substantially improved when blending multiple predictors. Our experience is that most efforts should be concentrated in deriving substantially different approaches, rather than refining a single technique."

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
1	BellKor's Pragmatic Chaos	0.8567	10.05	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.05	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Vanderlei Instituto and Vanderlei United	0.8588	9.54	2009-07-10 21:23:31
5	PragmaticTheo	0.8591	9.81	2009-07-10 00:33:20
6	PragmaticTheo	0.8594	9.77	2009-06-24 12:06:56
7	BellKor n BioChass	0.8601	9.70	2009-05-13 05:14:09
10	Roohas	0.8623	9.47	2009-04-07 23:35:59
11	Odear Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-25 17:19:11
13	sanalling	0.8642	9.27	2009-07-15 14:53:22
14		0.8642	9.26	2009-04-22 14:11:12
20	azemill	0.8668	9.00	2009-03-21 16:20:50
	Progress Prize 2007 - RMSE = 0.8723 - Winning Team: KorBell			
	Cinematch score - RMSE = 0.9525			

Fall 2025 CSE552 Machine Learning

16

Motivations

- Motivations of ensemble methods

- Ensemble model improves accuracy and robustness over single model methods
- Applications:
 - distributed computing
 - privacy-preserving applications
 - large-scale data with reusable models
 - multiple sources of data
- Efficiency: a complex problem can be decomposed into multiple sub-problems that are easier to understand and solve (divide-and-conquer approach)

Fall 2025

CSE552 Machine Learning

17

Why Ensemble Works?

- Intuition

- combining diverse, independent opinions in human decision-making as a protective mechanism (e.g. stock portfolio)

- Uncorrelated error reduction

- Suppose we have 5 completely independent classifiers for majority voting
- If accuracy is 70% for each
 - $10 \cdot (.7^3)(.3^2) + 5 \cdot (.7^4)(.3) + (.7^5)$
 - 83.7% majority vote accuracy
- 101 such classifiers
 - 99.9% majority vote accuracy

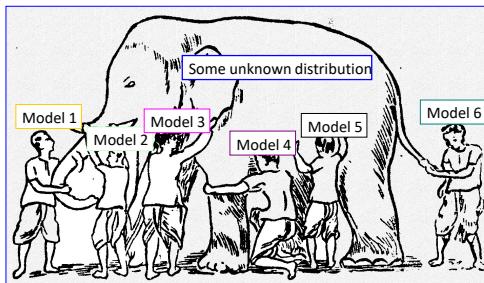
Fall 2025

CSE552 Machine Learning

18

Why Ensemble Works?

Why Ensemble Works?



Fall 2025

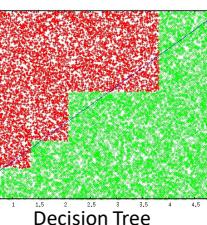
CSE552 Machine Learning

19

Why Ensemble Works?

- Overcome limitations of single hypothesis

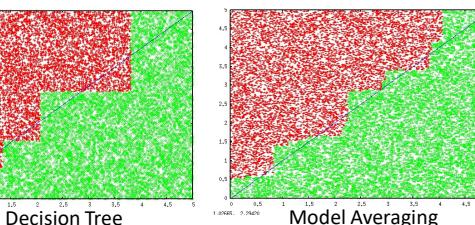
- The target function may not be implementable with individual classifiers, but may be approximated by model averaging



Fall 2025

CSE552 Machine Learning

20



19

20

Intuitions

- Utility of combining diverse, independent opinions in human decision-making
 - Protective Mechanism (e.g. stock portfolio diversity)
- Violation of Ockham's Razor**
 - Identifying the best model requires identifying the proper "model complexity"

See Domingos, P. Occam's two razors: the sharp and the blunt. KDD. 1998.

Fall 2025

CSE552 Machine Learning

21

Intuitions

Majority vote

Suppose we have 5 completely **independent** classifiers...

- If accuracy is 70% for each
 - $10(.7^3)(.3^2)+5(.7^4)(.3)+(.7^5)$
 - 83.7% majority vote accuracy**
- 101 such classifiers
 - 99.9% majority vote accuracy**

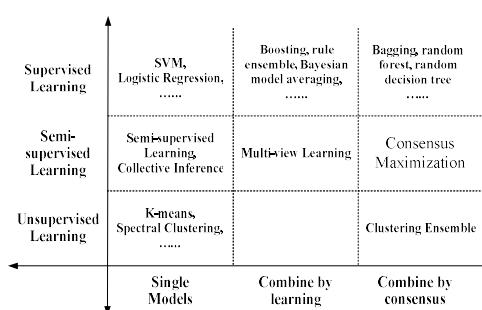
Fall 2025

CSE552 Machine Learning

22

21

Summary

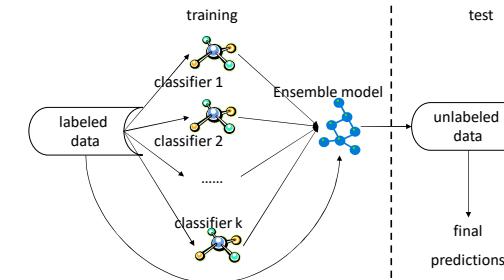


Fall 2025

CSE552 Machine Learning

23

Ensemble of Classifiers—Learn to Combine



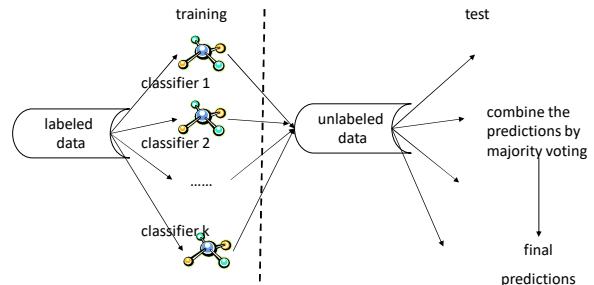
Fall 2025

CSE552 Machine Learning

24

23

Ensemble of Classifiers—Consensus



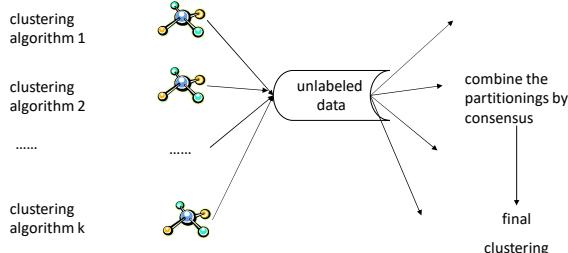
Algorithms: bagging, random forest, random decision tree, model averaging of probabilities...

Fall 2025

CSE552 Machine Learning

25

Clustering Ensemble—Consensus



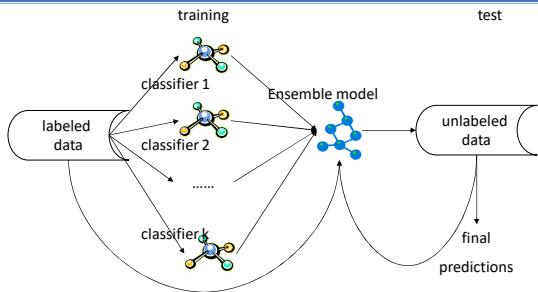
Algorithms: direct approach, object-based, cluster-based, object-cluster-based approaches, generative models

Fall 2025

CSE552 Machine Learning

26

Semi-Supervised Ensemble-Learn to Combine



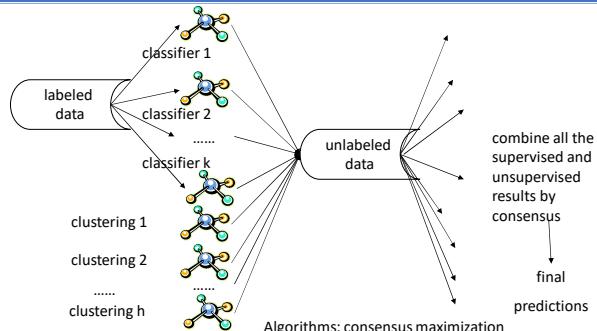
learn the combination from both labeled and unlabeled data
Algorithms: multi-view learning

Fall 2025

CSE552 Machine Learning

27

Semi-supervised Ensemble—Consensus



CSE552 Machine Learning

28

27

28

Pros and Cons

	Combine by learning	Combine by consensus
Pros	<p>Get useful feedbacks from labeled data</p> <p>Can potentially improve accuracy</p>	<p>Do not need labeled data</p> <p>Can improve the generalization performance</p>
Cons	<p>Need to keep the labeled data to train the ensemble</p> <p>May overfit the labeled data</p> <p>Cannot work when no labels are available</p>	<p>No feedbacks from the labeled data</p> <p>Require the assumption that consensus is better</p>

Fall 2025

CSE552 Machine Learning

29

Supervised Ensemble Methods

- Problem

- Given a data set $D=\{x_1, x_2, \dots, x_n\}$ and their corresponding labels $L=\{l_1, l_2, \dots, l_n\}$

- An ensemble approach computes:

- A set of classifiers $\{f_1, f_2, \dots, f_k\}$, each of which maps data to a class label: $f_i(x) =$
- A combination of classifiers f^* which minimizes generalization error: $f^*(x) = w_1f_1(x) + w_2f_2(x) + \dots + w_kf_k(x)$

Fall 2025

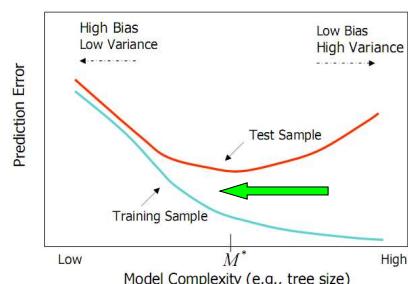
CSE552 Machine Learning

30

Bias and Variance

- Ensemble methods

- Combine learners to reduce variance



Fall 2025

CSE552 Machine Learning

31

Generating Base Classifiers

- Sampling training examples

- Train k classifiers on k subsets drawn from the training set

- Using different learning models

- Use all the training examples, but apply different learning algorithms

- Sampling features

- Train k classifiers on k subsets of features drawn from the feature space

- Learning “randomly”

- Introduce randomness into learning procedures

Fall 2025

CSE552 Machine Learning

32

31

32

BAGGING

Fall 2025

CSE552 Machine Learning

33

33

Main types of Ensemble methods

Combine multiple models together

- Committees – Bagging
 - Regression: Take an average of the predictions made by each model
 - Classification: Make classification by voting over a collection of classifiers
- Boosting – Adaboost
 - Train multiple models in sequence
- Decision trees
 - Different models are responsible for making predictions in different regions of input space

Fall 2025

CSE552 Machine Learning

34

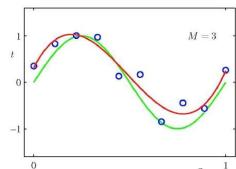
34

Committees [1]

Construct a committee ----
 average the predictions of a set of individual models

For example, regression problem:

- $T(x) = \sin 2\pi x + \text{noise}$
- $N = 10$
- Polynomial, order 3



Fall 2025

CSE552 Machine Learning

35

35

Polynom Fitting

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

where λ is the regularization parameter

Fall 2025

CSE552 Machine Learning

36

36

Committees [2]

train multiple polynomials

-- one fitting function learned from each i.i.d. subset of the whole

$T(x) = \sin 2\pi x + \text{noise}$

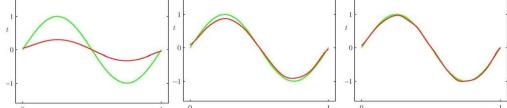
$L = 100$ data sets $\rightarrow 100$ fits (only 20 shown)

Each dataset: $N = 25$ data points

Regularisation factor λ

24 Gaussian basis functions (+bias)

average the resulting functions



Fall 2025

CSE552 Machine Learning

37

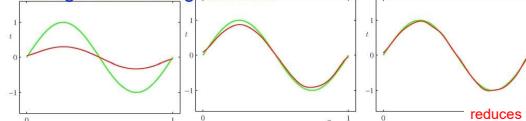
Committees [3]

train multiple polynomials

-- one fitting function learned from each i.i.d. subset of the whole data

$\ln \lambda = 2.6$, $\ln \lambda = -0.31$, $\ln \lambda = -2.4$

average the resulting functions



Fall 2025

CSE552 Machine Learning

38

Error Difference [1]

- true regression function $h(x)$,
- output of each model learned from one of the M data sets,
 $y_m(x) = h(x) + \epsilon_m(x)$ "bootstrap"
- average error made by the models acting individually

$$E_{AV} = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_x [\epsilon_m(x)^2]$$

- Committee prediction

$$y_{COM}(x) = \frac{1}{M} \sum_{m=1}^M y_m(x)$$

- the expected error from the committee

$$E_{COM} = \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{m=1}^M y_m(x) - h(x) \right)^2 \right] = \mathbb{E}_x \left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(x) \right)^2 \right]$$

Fall 2025

CSE552 Machine Learning

39

Error Difference [2]

- If we assume that the errors have zero mean and are uncorrelated, so that

$$\begin{aligned} \mathbb{E}_x [\epsilon_m(x)] &= 0 \\ \mathbb{E}_x [\epsilon_m(x)\epsilon_l(x)] &= 0, \quad m \neq l \end{aligned}$$

- then we obtain

$$E_{COM} = \frac{1}{M} E_{AV}$$

The average error of a model can be reduced by a factor of M simply by averaging M versions of the model.

But, the assumption may not be true. The errors are typically highly correlated, but still

$$E_{COM} \leq E_{AV}$$

39

40

40

Bias and Variance

- Ensemble methods**
 - Combine learners to reduce variance

Fall 2025 CSE552 Machine Learning 41

41

Main types of Ensemble methods

- Combine multiple models together
 - Committees – Bagging**
 - Regression:** Take an average of the predictions made by each model
 - Classification:** Make classification by voting over a collection of classifiers
 - Boosting – Adaboost**
 - Train multiple models in sequence
 - Decision trees**
 - Different models are responsible for making predictions in different regions of input space

Fall 2025 CSE552 Machine Learning 42

42

Bagging

- Bootstrap**
 - Sampling (random) with replacement
 - Contains around 63.2% original records in each sample, the rest duplicates
- Bootstrap Aggregation**
 - Train a classifier on each bootstrap sample
 - Use majority voting to determine the class label of ensemble classifier
- Note that we only have 1 single dataset

*[Breiman96]

Fall 2025 CSE552 Machine Learning 43

Bagging

Fall 2025 CSE552 Machine Learning 44

43

44

The Bagging Algorithm

BAGGING

Training phase

1. Initialize the parameters
 - $\mathcal{D} = \emptyset$, the ensemble.
 - L , the number of classifiers to train.
 2. For $k = 1, \dots, L$
 - Take a bootstrap sample S_k from \mathbf{Z} .
 - Build a classifier D_k using S_k as the training set.
 - Add the classifier to the current ensemble, $\mathcal{D} = \mathcal{D} \cup D_k$.
 3. Return \mathcal{D} .
- Classification phase**
4. Run D_1, \dots, D_L on the input \mathbf{x} .
 5. The class with the maximum number of votes is chosen as the label for \mathbf{x} .

Fall 2025

CSE552 Machine Learning

45

Bagging

- Works for unstable base classifiers
 - Neural Networks
 - Decision Trees
 - KNN
- Sampling with replacement...
 - Can emphasize replaced samples

Fall 2025

CSE552 Machine Learning

46

Bagging

Original Data:

x	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
y	1	1	1	-1	-1	-1	-1	1	1	1
H_1							H_1			
						H_1				
H_2						H_2				
							H_2			
H_3						H_3				
							H_3			
H_4							H_4			
								H_4		
H_5								H_5		
									H_5	

Combine predictions by majority voting

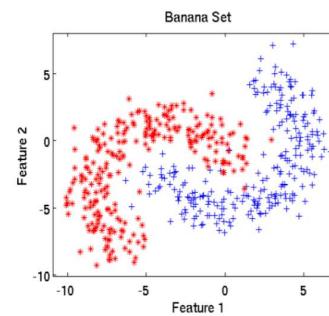
H_c		H_c		H_c		H_c
-------	--	-------	--	-------	--	-------

Fall 2025

CSE552 Machine Learning

47

Bagging – Example 2



Training data...

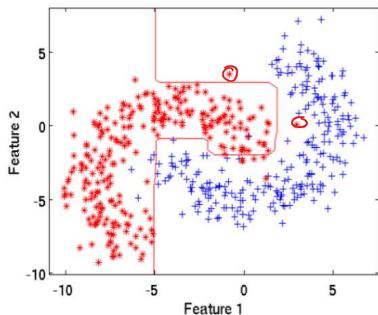
Base classifiers:
Decision trees

Fall 2025

CSE552 Machine Learning

48

Bagging – Example 2



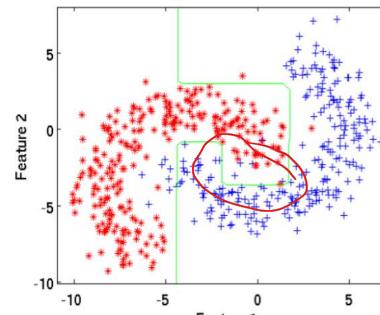
Decision boundaries produced by first tree...

Fall 2025

CSE552 Machine Learning

49

Bagging – Example 2



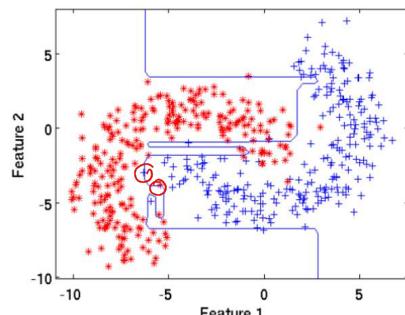
Decision boundaries produced by second tree...

Fall 2025

CSE552 Machine Learning

50

Bagging – Example 2



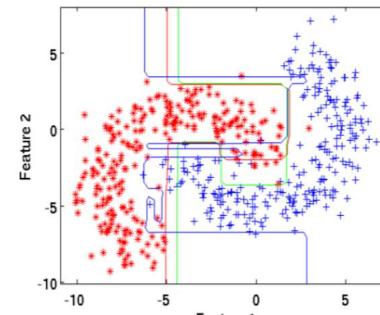
Decision boundaries produced by third tree...

Fall 2025

CSE552 Machine Learning

51

Bagging – Example 2



Decision boundaries produced by all trees...

Fall 2025

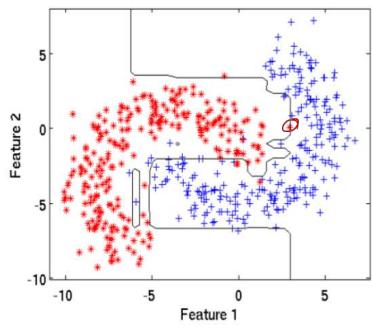
CSE552 Machine Learning

52

51

13

Bagging – Example 2



Final decision
boundaries produced
by bagging all trees...

Fall 2025

CSE552 Machine Learning

53

When Bagging Works

- Expected error is the expected discrepancy between the estimated $\hat{f}(x) = E[\hat{f}(x)]$ and true function $f(x)$

$$E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

- Bias is squared discrepancy between averaged estimated and true function

$$(E[\hat{f}(x)] - E[f(x)])^2$$

- Variance is expected divergence of the estimated function vs. its average value

$$E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

Fall 2025

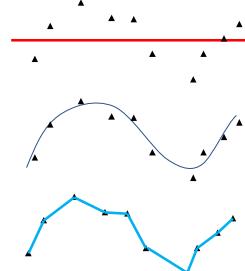
CSE552 Machine Learning

54

53

When Bagging Works

- Under-fitting:
 - High bias (models are not accurate)
 - Small variance (smaller influence of examples in the training set)
- Over-fitting:
 - Small bias (models flexible enough to fit well to training data)
 - Large variance (models depend very much on the training set)



Fall 2025

CSE552 Machine Learning

55

Averaging Decreases Variance

- Example
 - Assume we measure a random variable x with a $N(\mu, \sigma^2)$ distribution
 - If only one measurement x_1 is done,
 - The expected mean of the measurement is μ
 - Variance is $\text{Var}(x_1) = \sigma^2$
- If random variable x is measured K times (x_1, x_2, \dots, x_K) and the value is estimated as: $(x_1, x_2, \dots, x_K)/K$,
 - Mean of the estimate is still μ
 - But, variance is smaller:

$$\frac{[\text{Var}(x_1) + \text{Var}(x_2) + \dots + \text{Var}(x_K)]}{K^2} = \frac{K\sigma^2}{K^2} = \frac{\sigma^2}{K}$$
- Observe: Bagging is a kind of averaging!

Fall 2025

CSE552 Machine Learning

56

55

When Bagging Works

- Main property of Bagging (no proof will be given)
 - Bagging decreases variance of the base model without changing the bias!!!
 - Why? Averaging!
- Bagging typically helps
 - When applied with an over-fitted base model
 - High dependency on actual training data
- It does not help much
 - High bias. When the base model is robust to the changes in the training data (due to sampling)

Fall 2025

CSE552 Machine Learning

57

57

Boosting

Fall 2025

CSE552 Machine Learning

58

58

Main types of Ensemble methods

- Combine multiple models together
- Committees – Bagging
 - *Regression*: Take an average of the predictions made by each model
 - *Classification*: Make classification by voting over a collection of classifiers
 - **Boosting – Adaboost**
 - Train multiple models in sequence
 - Decision trees
 - Different models are responsible for making predictions in different regions of input space

Fall 2025

CSE552 Machine Learning

59

59

Boosting

• Principles

- Boost a set of weak learners to a strong learner
- Make records currently misclassified more important

• Example

- Record 4 is hard to classify
- Its weight is increased, therefore it is more likely to be chosen again in subsequent rounds

Original Data	1	2	3	4	5	6	7	8	9	10
Boosting (Round 1)	7	3	2	8	7	9	4	10	6	3
Boosting (Round 2)	5	4	9	4	2	5	1	7	4	2
Boosting (Round 3)	4	4	8	10	4	5	4	6	3	4

*[FrSc97]

Fall 2025

CSE552 Machine Learning

60

60

Boosting

- AdaBoost

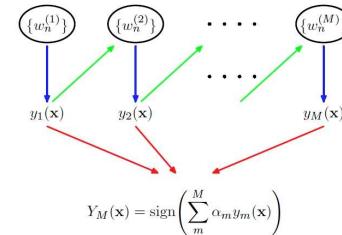
- Initially, set uniform weights on all the records
- At each round
 - Create a bootstrap sample based on the weights
 - Train a classifier on the sample and apply it on the original training set
 - Records that are wrongly classified will have their weights increased
 - Records that are classified correctly will have their weights decreased
 - If the error rate is higher than 50%, start over
- Final prediction is weighted average of all the classifiers with weight representing the training accuracy

Fall 2025

CSE552 Machine Learning

61

Boosting



Fall 2025

CSE552 Machine Learning

62

Boosting

- Determine the weight

- For classifier i , its error is

$$\varepsilon_i = \frac{\sum_{j=1}^N w_j \delta(C_i(x_j) \neq y_j)}{\sum_{j=1}^N w_j}$$

- The classifier's importance is represented as:

$$\alpha_i = \frac{1}{2} \ln \left(\frac{1 - \varepsilon_i}{\varepsilon_i} \right)$$

- The weight of each record is updated as:

$$w_j^{(i+1)} = \frac{w_j^{(i)} \exp(-\alpha_i y_j C_i(x_j))}{Z^{(i)}}$$

- Final combination:

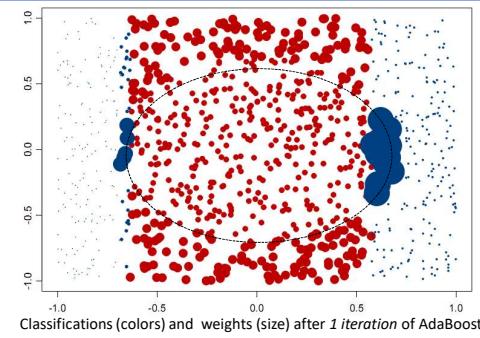
$$C^*(x) = \arg \max_y \sum_{i=1}^K \alpha_i \delta(C_i(x) = y)$$

Fall 2025

CSE552 Machine Learning

63

Example



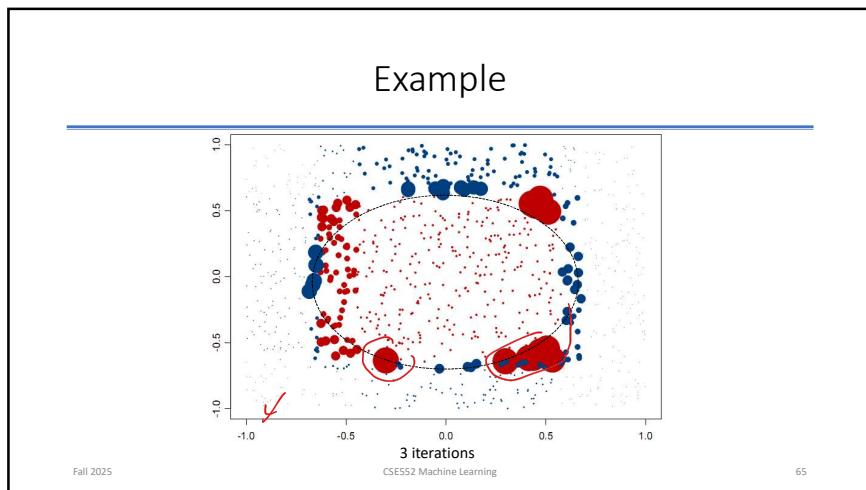
Fall 2025

CSE552 Machine Learning

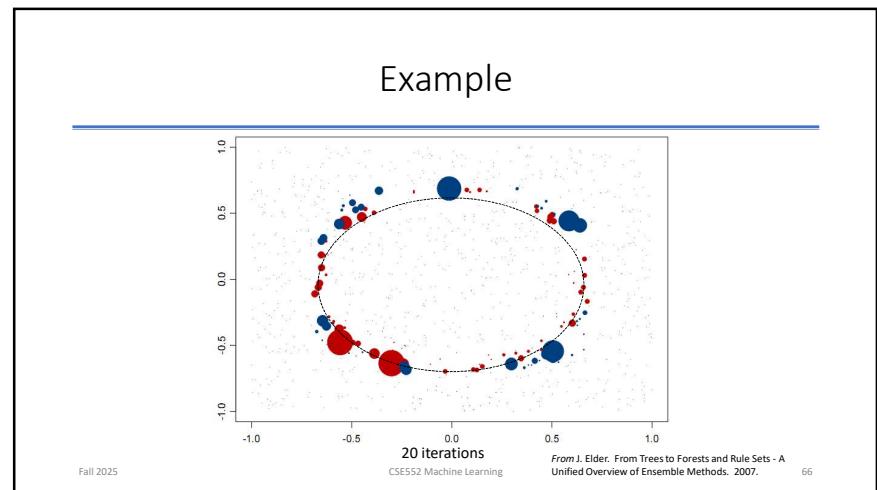
64

63

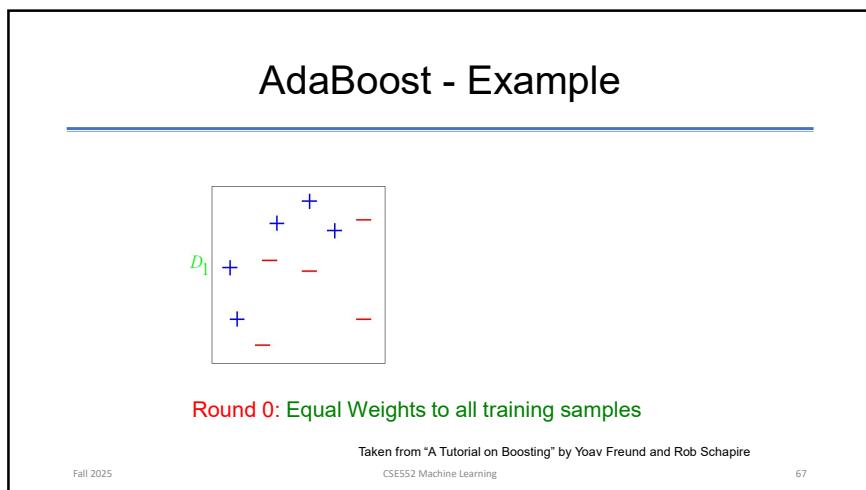
64



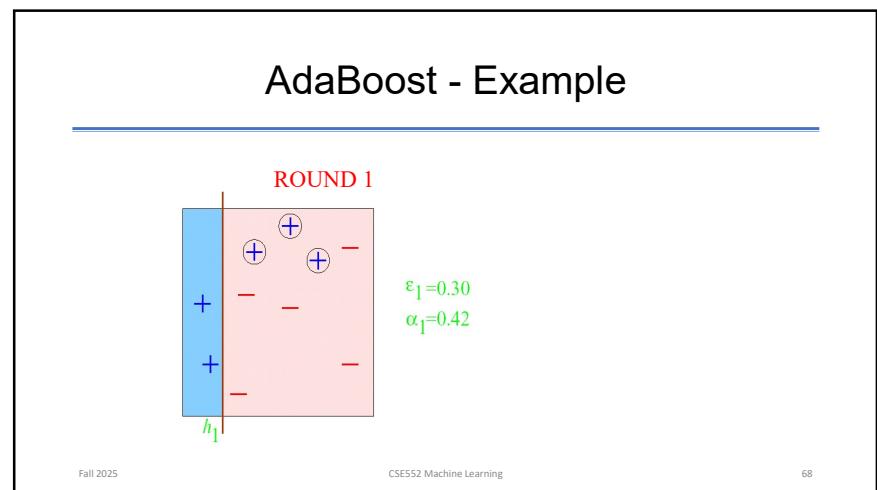
65



66

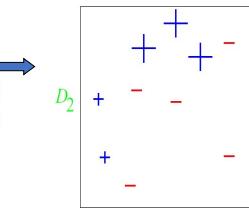
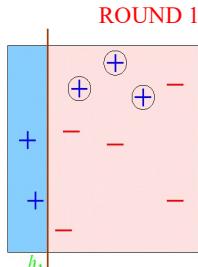


67



68

AdaBoost - Example

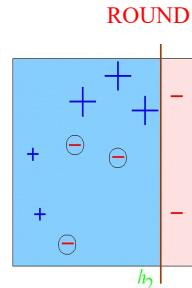


Fall 2025

CSE552 Machine Learning

69

AdaBoost - Example

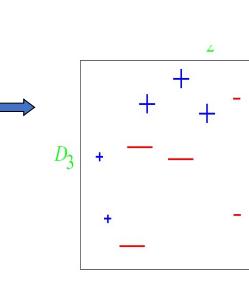
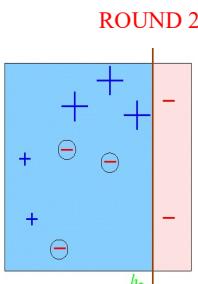

 $\epsilon_2 = 0.21$
 $\alpha_2 = 0.65$

Fall 2025

CSE552 Machine Learning

70

AdaBoost - Example

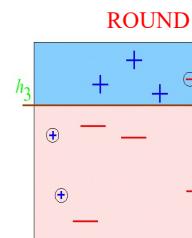


Fall 2025

CSE552 Machine Learning

71

AdaBoost - Example


 $\epsilon_3 = 0.14$
 $\alpha_3 = 0.92$

Fall 2025

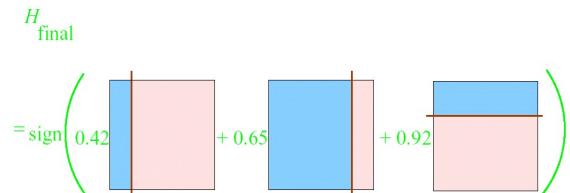
CSE552 Machine Learning

72

71

18

AdaBoost - Example



Fall 2025

CSE552 Machine Learning

73

73

Boosting – Theoretical Foundations

- PAC: Probably Approximately Correct framework
 - $\epsilon - \delta$ solution
- PAC learning:
 - Learning with the pre-specified accuracy ϵ and confidence δ
 - The probability that the misclassification error is larger than ϵ is smaller than δ
$$P(ME(c) > \epsilon) \leq \delta$$
- Accuracy(ϵ): Percent of correctly classified samples in test
- Confidence(δ): The probability that in one experiment some accuracy will be achieved

Fall 2025

CSE552 Machine Learning

74

74

PAC Learnability

Strong (PAC) learnability

- There exists a learning algorithm that **efficiently** learns the classification with a pre-specified **accuracy** and **confidence**

Strong (PAC) learner:

- A learning algorithm P that given an arbitrary
 - classification error $\epsilon (< \frac{1}{2})$, and
 - confidence $\delta (< \frac{1}{2})$
- Outputs a classifier
 - With a classification accuracy $> (1 - \epsilon)$
 - A confidence probability $> (1 - \delta)$
 - And runs in time polynomial in $\frac{1}{\delta}, \frac{1}{\epsilon}$
 - Implies: number of samples N is polynomial in $\frac{1}{\delta}, \frac{1}{\epsilon}$

Fall 2025

CSE552 Machine Learning

75

75

Weak Learner

Weak learner:

- A learning algorithm (learner) W
 - Providing classification accuracy $> (1 - \epsilon_0)$
 - With probability $> (1 - \delta_0)$
- For some **fixed and uncontrollable**
 - Classification error $\epsilon_0 (< \frac{1}{2})$
 - Confidence $\delta_0 (< \frac{1}{2})$
- And this on an arbitrary distribution of data entries

Fall 2025

CSE552 Machine Learning

76

76

Weak Learnability = Strong (PAC) Learnability

- Assume there exists a weak learner
 - it is better than a random guess (50%) with confidence higher than 50% on any data distribution
- Question:
 - Is problem also PAC-learnable?
 - Can we generate an algorithm P that achieves an arbitrary $\epsilon - \delta$ accuracy?
- Why is this important?
- Usual classification methods (decision trees, neural nets), have specified, but uncontrollable performances
- Can we improve performance to achieve pre-specified accuracy (confidence)?

Fall 2025

CSE552 Machine Learning

77

Weak=Strong Learnability!!!

- Proof due to R. Schapire
 - An arbitrary $(\epsilon - \delta)$ improvement is possible
- Idea: combine multiple weak learners together
 - Weak learner W with confidence δ_0 and maximal error ϵ_0
 - It is possible:
 - To improve (boost) the confidence
 - To improve (boost) the accuracy

by training different weak learners on slightly different datasets

Fall 2025

CSE552 Machine Learning

78

Boosting Accuracy

- Training
 - Sample randomly from the distribution of examples
 - Train hypothesis H_1 on the sample
 - Evaluate accuracy of H_1 on the distribution
 - Sample randomly such that for the half of samples H_1 provides correct, and for another half, incorrect results; Train hypothesis H_2 .
 - Train H_3 on samples from the distribution where H_1 and H_2 classify differently
- Test
 - For each example, decide according to the majority vote of H_1 , H_2 and H_3

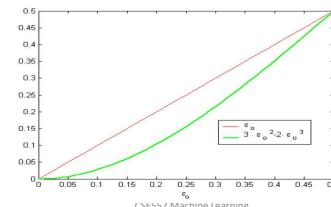
Fall 2025

CSE552 Machine Learning

79

Theorem

- If each hypothesis has an error ϵ_0 , the final classifier has error $< g(\epsilon_0) = 3\epsilon_0^2 - 2\epsilon_0^3$
- Accuracy improved !!!!
- Apply recursively to get to the target accuracy !!!



Fall 2025

CSE552 Machine Learning

80

Theoretical Boosting Algorithm

- Similar to boosting the accuracy we can boost the confidence at some restricted accuracy cost
- The key result:** we can improve both the accuracy and confidence
- Problems with the theoretical algorithm
 - A good (better than 50 %) classifier on all data problems
 - We cannot properly sample from data-distribution
 - Method requires large training set
- Solution to the sampling problem:
 - Boosting by sampling
 - AdaBoost algorithm and variants

Fall 2025

CSE552 Machine Learning

81

81

Boosting

- Explanation
 - Among the classifiers of the form:
$$f(x) = \sum_{i=1}^K \alpha_i C_i(x)$$
- The classifier's importance is represented as:

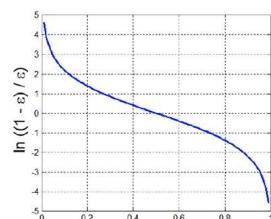
$$\alpha_i = \frac{1}{2} \ln \left(\frac{1 - \epsilon_i}{\epsilon_i} \right)$$
- We seek to minimize the exponential loss function:

$$\sum_{j=1}^N e^{-y_j f(x_j)}$$

Fall 2025

CSE552 Machine Learning

83



83

Boosting - AdaBoost

- In practice overfitting rarely occurs (Bishop)

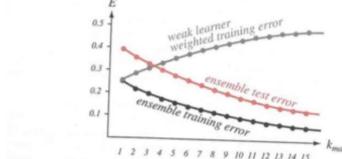


FIGURE 9.7. AdaBoost applied to a weak learning system can reduce the training error E exponentially as the number of component classifiers, k_{max} , is increased. Because AdaBoost "focuses on" difficult training patterns, the training error of each successive component classifier (measured on its own weighted training set) is generally larger than that of the previous component classifier (shown in gray). Nevertheless, so long as the component classifiers perform better than chance (e.g., have error less than 0.5 on a two-category problem), the weighted ensemble decision of Eq. 36 ensures that the training error will decrease, as given by Eq. 37. It is often found that the test error decreases in boosted systems as well, as shown in red.

Fall 2025

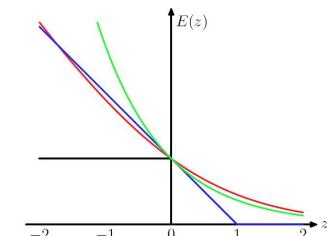
CSE552 Machine Learning

82

82

AdaBoost – Loss Function

- Sequential minimization of the exponential error function
 - Penalizes large negative values heavily
 - Not robust to outliers!
 - No probabilistic interpretation
 - Only for 2 classes



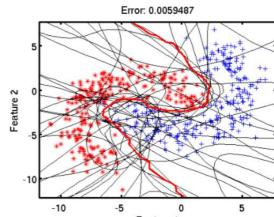
Fall 2025

CSE552 Machine Learning

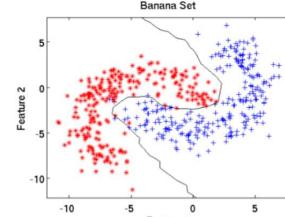
84

84

Boosting Example Again



AdaBoost using 20 neural nets
[bpnnc] default settings



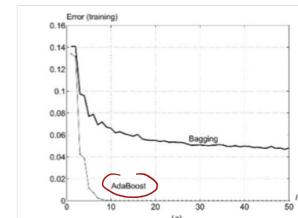
Final output of AdaBoost with 20
neural nets

Fall 2025

CSE552 Machine Learning

85

Boosting vs Bagging Example



(a)

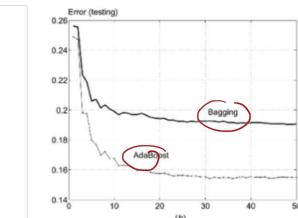


Fig. 7.11 Training and testing error of bagging and AdaBoost for the rotated check-board example.

Krogh 95

Fall 2025

CSE552 Machine Learning

86

Boosting vs Bagging

- Committee/Bagging
 - Base classifiers are trained in parallel on samples of data set
- Boosting
 - Base classifiers are trained in sequence by using weighted form of the data set
 - Weights for each data point depend on the performance of the previous classifiers
 - Misclassified points are given greater weight when used to train the next classifier in the sequence
 - Boosting tends to achieve greater accuracy, but also risks overfitting the model to misclassified data

Fall 2025

CSE552 Machine Learning

87

Boosting and Computer Vision

88

88

Example: Face Detection

- Frontal faces are a good example of a class where global appearance models + a sliding window detection approach fit well:
 - Regular 2D structure
 - Center of face almost shaped like a “patch”/window



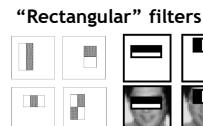
- Let see how the Viola-Jones face detector works using AdaBoost

Fall 2025

CSE552 Machine Learning

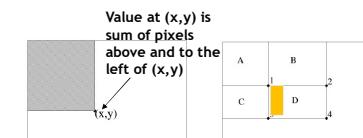
89

Features



“Rectangular” filters

Feature output is difference between adjacent regions



Efficiently computable with integral image: any sum can be computed in constant time

Avoid scaling images → scale features directly for same cost

Viola & Jones, CVPR 2001

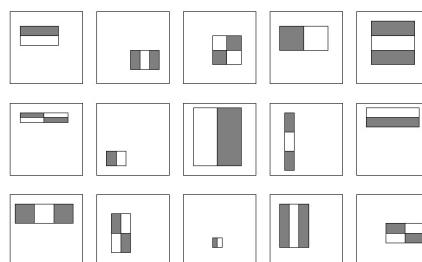
Fall 2025

CSE552 Machine Learning

$$\begin{aligned} D &= 1+4-(2+3) \\ &= 4+(A+B+C+D)-(A+C+A+B) \\ &= D \end{aligned}$$

90

Large Library of Filters



Considering all possible filter parameters: position, scale, and type:

180,000+ possible features associated with each 24 x 24 window

Use AdaBoost both to select the informative features and to form the classifier

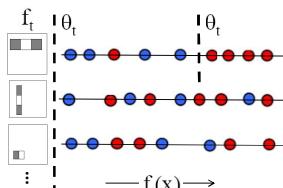
Fall 2025

CSE552 Machine Learning

91

AdaBoost for Feature+Classifier Selection

Want to select the single rectangle feature and threshold that best separates **positive** (faces) and **negative** (non-faces) training examples, in terms of *weighted error*.



Outputs of a possible rectangle feature on faces and non-faces.

Fall 2025

CSE552 Machine Learning

92

Resulting weak classifier:

$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ -1 & \text{otherwise} \end{cases}$$

For next round, reweight the examples according to errors, choose another filter/threshold combo.

AdaBoost Algorithm

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{t,0} = \frac{1}{m}, \frac{1}{n}$ for $y_i = 0, 1$ respectively, where m, n are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:
 - Normalize the weights,

$$w_{t,j} \leftarrow \frac{w_{t-1,j}}{\sum_{j=1}^m w_{t-1,j}}$$
 so that $\{w_t\}$ is a probability distribution.
 - For each feature j , train a classifier h_j which is restricted to use a single feature. The error is evaluated with respect to w_t , $\epsilon_j = \sum_i w_t h_j(x_i) - y_i$.
 - Choose the classifier h_t with the lowest error ϵ_t .
 - Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{\epsilon_t}$$
 where $\epsilon_t = 0$ if example x_i is classified correctly, $\epsilon_t = 1$ otherwise, and $\beta_t = \frac{e^{-\epsilon_t}}{1 - e^{-\epsilon_t}}$.
- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

Start with uniform weights on training examples



For T rounds

Evaluate weighted error for each feature, pick best.

Re-weight the examples:

Incorrectly classified \rightarrow more weight
Correctly classified \rightarrow less weight

Final classifier is combination of the weak ones, weighted according to error they had.

Freund & Schapire 1995

93

Fall 2025

CSE552 Machine Learning

AdaBoost for Efficient Feature Selection

- Image Features = Weak Classifiers
- For each round of boosting:
 - Evaluate each rectangle filter on each example
 - Sort examples by filters values
 - Select best threshold for each filter (min error)
 - Sorted list can be quickly scanned for the optimal threshold
 - Select best filter/threshold combination
 - Weight on this feature is a simple function of error rate
 - Reweight examples

Viola & Jones, CVPR 2001

Fall 2025

CSE552 Machine Learning

94

94

Problem

- Even if the filters are fast to compute, each new image has a lot of possible windows to search
- How to make the detection more efficient?

Fall 2025

CSE552 Machine Learning

95

Cascading Classifiers for Detection

For efficiency, apply less accurate but faster classifiers first to immediately discard windows that clearly appear to be negative; e.g.,

- Filter for promising regions with an initial inexpensive classifier
- Build a chain of classifiers, choosing cheap ones with low false negative rates early in the chain

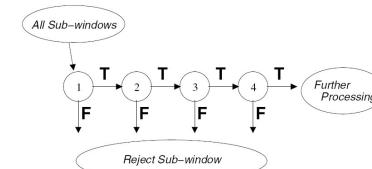


Figure from Viola & Jones CVPR 2001

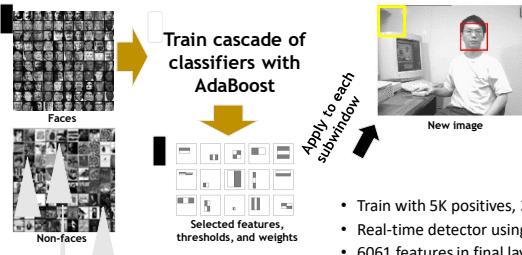
Fall 2025

CSE552 Machine Learning

96

96

Viola-Jones Face Detector: Summary

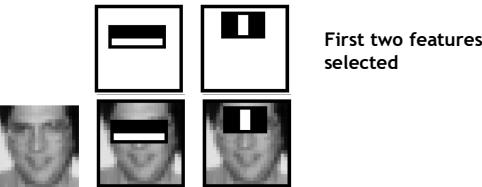


Fall 2025

CSE552 Machine Learning

97

Viola-Jones Face Detector: Results

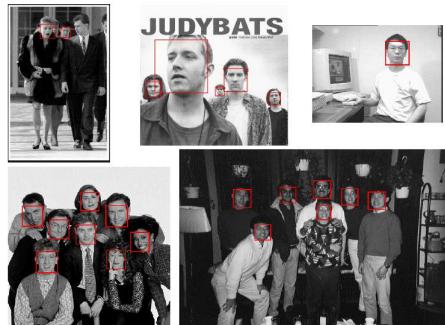


Fall 2025

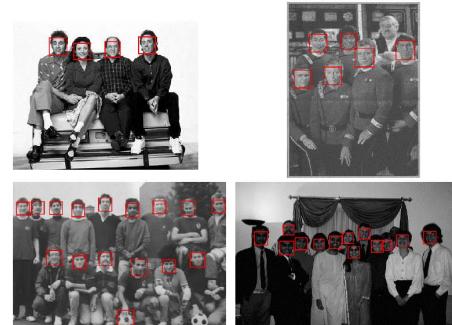
CSE552 Machine Learning

98

Viola-Jones Face Detector: Results



Viola-Jones Face Detector: Results



99

100

Viola-Jones Face Detector: Results



Fall 2025

CSE552 Machine Learning

101

101

Viola-Jones Face Detector: Results



Fall 2025

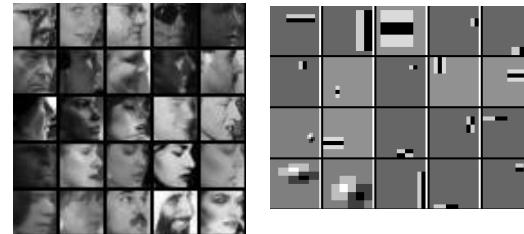
CSE552 Machine Learning

103

103

Detecting Profile Faces?

Detecting profile faces requires training separate detector with profile examples.



Fall 2025

CSE552 Machine Learning

102

102

Limitations

- High computational complexity
 - For example: 250,000 locations x 30 orientations x 4 scales = 30,000,000 evaluations!
 - If training binary detectors independently, means cost increases linearly with number of classes
- With so many windows, false positive rate better be low

Fall 2025

CSE552 Machine Learning

104

104

RANDOM FORESTS

Fall 2025

CSE552 Machine Learning

105

105

Main types of Ensemble methods

Combine multiple models together

- Committees – Bagging
 - *Regression*: Take an average of the predictions made by each model
 - *Classification*: Make classification by voting over a collection of classifiers
- Boosting – Adaboost
 - Train multiple models in sequence
- Decision trees
 - Different models are responsible for making predictions in different regions of input space

Fall 2025

CSE552 Machine Learning

106

106

Random Forests

Algorithm

- 
- Choose T — number of trees to grow
 - Choose $m \leq M$ (M is the number of total features) — number of features used to calculate the best split at each node (typically 20%)
 - For each tree
 - Choose a training set by choosing N times (N is the number of training examples) with replacement from the training set
 - For each node, randomly choose m features and calculate the best split
 - Fully grown and not pruned
 - Use majority voting among all the trees (or average for regression)
[Breiman'01]

Fall 2025

CSE552 Machine Learning

107

107

Random Forests

The common element in all of these procedures is that for the k th tree, a random vector Θ_k is generated, independent of the past random vectors $\Theta_1, \dots, \Theta_{k-1}$ but with the same distribution; and a tree is grown using the training set and Θ_k , resulting in a classifier $h(\mathbf{x}, \Theta_k)$ where \mathbf{x} is an input vector. For instance, in bagging the random vector Θ is generated as the counts in N boxes resulting from N darts thrown at random at the boxes, where N is number of examples in the training set. In random split selection Θ consists of a number of independent random integers between 1 and K . The nature and dimensionality of Θ depends on its use in tree construction.

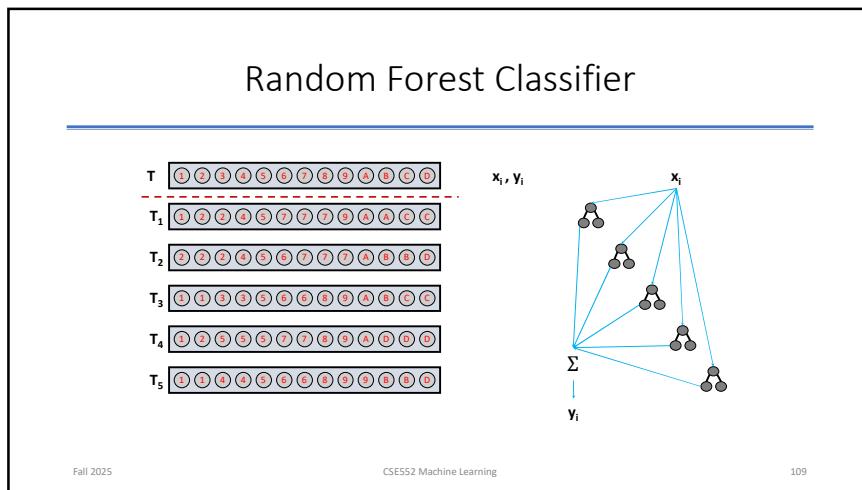
Definition 1.1. A random forest is a classifier consisting of a collection of tree-structured classifiers $\{h(\mathbf{x}, \Theta_k), k = 1, \dots\}$ where the $\{\Theta_k\}$ are independent identically distributed random vectors and each tree casts a unit vote for the most popular class at input \mathbf{x} .

Fall 2025

CSE552 Machine Learning

108

108



109

- ## Random Forests
- **Discussions**
 - Bagging + random features
 - Improve accuracy
 - Incorporate more diversity and reduce variances
 - Improve efficiency
 - Searching among subsets of features is much faster than searching among the complete set

Fall 2025

CSE552 Machine Learning

110

110

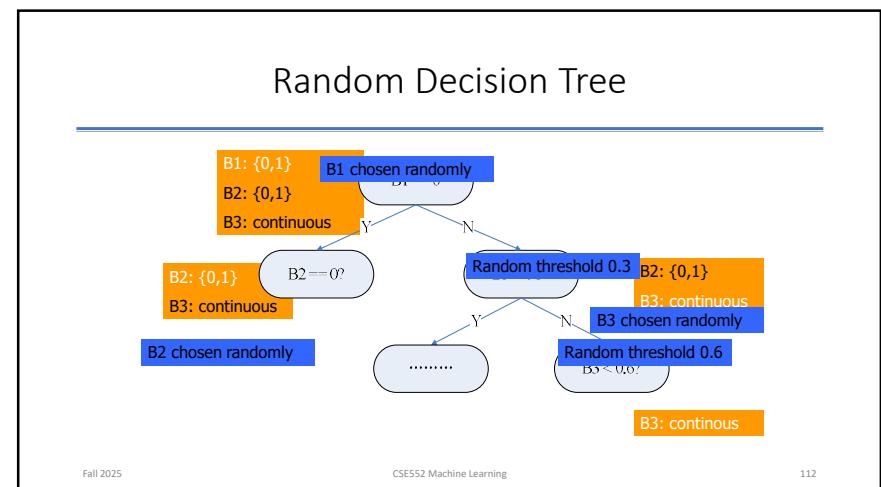
- ## Random Decision Tree
- **Single-model learning algorithms**
 - Fix structure of the model, minimize some form of errors, or maximize data likelihood (eg., Logistic regression, Naive Bayes, etc.)
 - Use some “free-form” functions to match the data given some “preference criteria” such as information gain, gini index and MDL. (eg., Decision Tree, Rule-based Classifiers, etc.)
 - **Such methods will make mistakes if**
 - Data is insufficient
 - Structure of the model or the preference criteria is inappropriate for the problem
 - **Learning as Encoding**
 - Make no assumption about the true model, neither parametric form nor free form
 - Do not prefer one base model over the other, just average them
- [FWM+03]

Fall 2025

CSE552 Machine Learning

111

111



Fall 2025

CSE552 Machine Learning

112

112

Random Decision Tree

Potential Advantages

- Training can be very efficient – particularly true for very large datasets
 - No cross-validation based estimation of parameters for some parametric methods
- Natural multi-class probability
- Imposes very little about the structures of the model

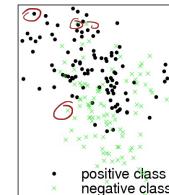
Fall 2025

CSE552 Machine Learning

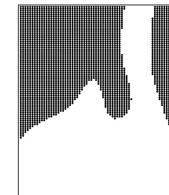
113

Optimal Decision Boundary

Figure 3.5: Gaussian mixture training samples and optimal boundary.



training samples



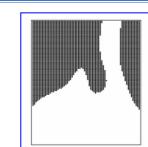
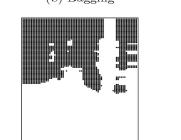
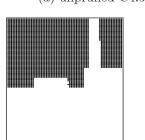
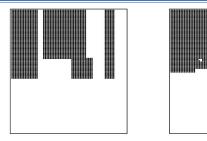
optimal boundary

Fall 2025

CSE552 Machine Learning

114

Optimal Decision Boundary



RDT looks
like the optimal
boundary

Fall 2025

CSE552 Machine Learning

115

Another Example

- Number of Trees: 500
- Iteration: 400
- Accuracy estimation: 10 fold cross validation
- Maximal terminal node size for PRF: 20

Fall 2025

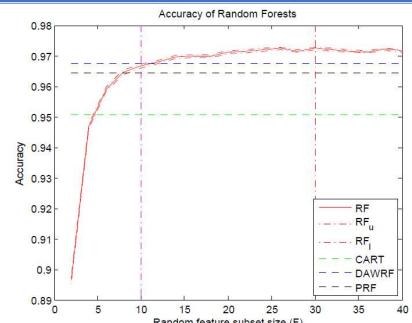
CSE552 Machine Learning

116

115

116

Another Example

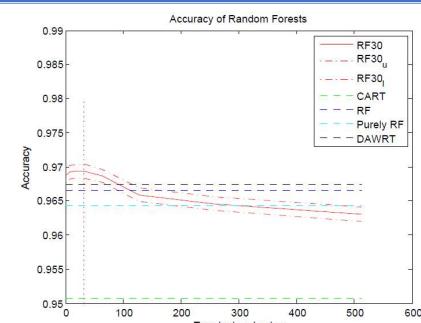


Fall 2025

CSE552 Machine Learning

117

Another Example



Fall 2025

CSE552 Machine Learning

118

Another Example

	CART	PRF	DAWRF	RF	RF-best
mean	0.9508	0.9643	0.9674	0.9666	0.9724
sd	0.0112	0.0102	0.0096	0.0114	0.0086
F	NA	1	NA	10	30

Fall 2025

CSE552 Machine Learning

119

Good

- Accuracy is as good as Adaboost and sometimes better
 - Randomness induces vastly more between-tree differences... (mutations in genetic algorithms)
 - Bootstrapping (even weights) alone does not ...
- Relatively robust to outliers and noise
- Fast Computation
- Gives a wealth of important insights (e.g. Estimate of error, variable importance, proximity)
- Simple

Fall 2025

CSE552 Machine Learning

120

119

120

Class Vote Proportions and Margin

- At the end of an RF run, for every record, the proportion of votes for each class represents the probability of class membership
- The **Margin** – the proportion of votes for the true class minus the maximum proportion of votes for the other classes
- The larger the margin, the higher the confidence of classification

Fall 2025

CSE552 Machine Learning

121

Out-of-bag Error

Assume a method for constructing a classifier from any training set. Given a specific training set T , form bootstrap training sets T_k , construct classifiers $h(x, T_k)$ and let these vote to form the bagged predictor. For each y, x in the training set, aggregate the votes only over those classifiers for which T_k does not contain y, x . Call this the out-of-bag classifier.

Breiman's studies provides empirical evidence showing that the **out-of-bag estimate is as accurate as using a test set of the same size as the training set**.
 → Using the out-of-bag error estimate removes the need for a set aside test set.

Fall 2025

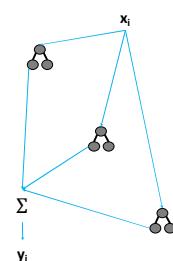
CSE552 Machine Learning

122

Out-of-bag Error

T	(1 2 3 4 5 6 7 8 9 0)
T_1	(1 2 3 4 5 6 7 8 9 0)
T_2	(1 2 3 4 5 6 7 8 9 0)
T_3	(1 2 3 4 5 6 7 8 9 0)
T_4	(1 2 3 4 5 6 7 8 9 0)
T_5	(1 2 3 4 5 6 7 8 9 0)

(D)

 x_i, y_i 

Fall 2025

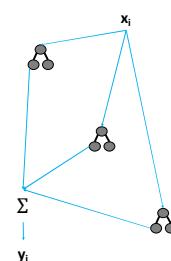
CSE552 Machine Learning

123

Out-of-bag Error

T	(1 2 3 4 5 6 7 8 9 0)
T_1	(1 2 3 4 5 6 7 8 9 0)
T_2	(1 2 3 4 5 6 7 8 9 0)
T_3	(1 2 3 4 5 6 7 8 9 0)
T_4	(1 2 3 4 5 6 7 8 9 0)
T_5	(1 2 3 4 5 6 7 8 9 0)

(D)

 x_i, y_i 

Fall 2025

CSE552 Machine Learning

124

123

124

Variable Importance

- The concept of margin allows new “unbiased” definition of variable importance
- To estimate the importance of the m^{th} variable:
 - Take the OOB (out of bag) cases for the k^{th} tree, assume that we already know the margin for those cases M_0
 - Randomly permute all values of the variable m
 - Apply the k^{th} tree to the OOB cases with the permuted values
 - Compute the new margin M
 - Compute the difference $M_0 - M$
- The variable importance is defined as the average lowering of the margin across all OOB cases and all trees in the RF
- This procedure is fundamentally different from the intrinsic variable importance scores computed by CART –the latter are always based on the LEARN data and are subject to the overfitting issues

Fall 2025

CSE552 Machine Learning

125

Proximity Measure

- RF introduces a novel way to define proximity between two observations:
 - Initialize proximities to zeroes
 - For any given tree, apply the tree to all cases
 - If case i and case j both end up in the same node, increase proximity $\text{prox}(ij)$ between i and j by one
 - Accumulate over all trees in RF and normalize by twice the number of trees in RF
- The resulting matrix of size $N \times N$ provides intrinsic measure of proximity
 - The measure is invariant to monotone transformations
 - The measure is clearly defined for any type of independent variables, including categorical

Fall 2025

CSE552 Machine Learning

126

Using Proximities

- Having the full intrinsic proximity matrix opens new horizons
 - Informative data views using metric scaling
 - Missing value imputation
 - Outlier detection
- Unfortunately, things get out of control when dataset size exceeds 5,000 observations (25,000,000+ cells are needed)
- RF switches to “compressed” form of the proximity matrix to handle large datasets –for any case, only M closest cases are recorded. M is usually less than 100.

Fall 2025

CSE552 Machine Learning

127

Scaling Coordinates

- The values $1 - \text{prox}(ij)$ can be treated as Euclidean distances in a high dimensional space
- The theory of metric scaling solves the problem of finding the most representative projections of the underlying data “cloud” onto low dimensional space using the data proximities
- The theory is similar in spirit to the principal components analysis and discriminant analysis
- The solution is given in the form of ordered “scaling coordinates”
- Looking at the scatter plots of the top scaling coordinates provides informative views of the data

Fall 2025

CSE552 Machine Learning

128

127

128

Outlier Detection

- Outliers are defined as cases having small proximities to all other cases belonging to the same target class
- The following algorithm is used:
 - For a case n , compute the sum of the squares of $\text{prox}(nk)$ for all k in the same class as n
 - Take the inverse –it will be large if the case is “far away” from the rest
 - Standardize using the median and standard deviation
 - Look at the cases with the largest values –those are potential outliers
- Generally, a value above 10 is reason to suspect the case of being an outlier

Fall 2025

CSE552 Machine Learning

129

Missing Value Imputation

- RF offers two ways of missing value imputation
- The Cheap Way – conventional median imputation for continuous variables and mode imputation for categorical variables
- A Better Way:
 - Suppose case n has x coordinate missing
 - 1. Do the Cheap Way imputation for starters
 - 2. Grow a full size RF
 - 3. We can now re-estimate the missing value by a weighted average over all cases k with non-missing x using weights $\text{prox}(nk)$
 - 4. Repeat steps 2 and 3 several times to ensure convergence

Fall 2025

CSE552 Machine Learning

130

Bad?

- Why maximal tree (no pruning)?
- Optimal random feature subset size(F)?
- Bootstrap sample?

Fall 2025

CSE552 Machine Learning

131

Why Maximal Tree?

- Lin and Jeon (2006), JASA
 - Breiman’s classifier can be viewed as adaptively weighted k-potential nearest neighbors methods in regression.
 - Terminal node size should be made to increase with the sample size.
- Biau et al (2008), JMLR
 - Using stopping rule is not necessary in some cases.
 - Empirical evidence...

Fall 2025

CSE552 Machine Learning

132

131

132

Why are Single Random Trees Predictive?

- A single tree in an RF forest can be predictive because it is a form of nearest neighbor classifier
- We know that nearest neighbor prediction is completely model free
- A good nearest neighbor system should achieve an error rate about twice that of the theoretically best model
 - If the best possible model is correct for 90% of all cases it is wrong for 10%
 - Nearest neighbor should then be correct for 80% and wrong for 20%
- Follow the path down a big tree is ...

Fall 2025

CSE552 Machine Learning

133

133

Optimal Random Feature Subset Size(F)?

- Many empirical studies

Fall 2025

CSE552 Machine Learning

134

134

Bootstrap Sample?

- Bootstrap sample is not essential for prediction
- Using bootstrap sample provides useful information
- But we can get same information by cross validation

Fall 2025

CSE552 Machine Learning

135

135

Supervised and Unsupervised RFs

- A Random Forest can be supervised or unsupervised
- Supervised:
 - In a supervised Random Forest, groupings for the training data are input to the algorithm
 - Estimated classification error is computed using out-of-bag data

Fall 2025

CSE552 Machine Learning

136

136

The Road to Clustering

- RF to solve unsupervised learning problems, in particular, clustering problems and missing value imputation in the general sense
- RF generates a synthetic target variable in order to proceed with a regular run:
 - Give class label 1 to the original data
 - Create a copy of the data such that each variable is sampled independently from the values available in the original dataset
 - Give class label 2 to the copy of the data
 - RF on two class → dissimilarity measure for clustering...
 - Note that the second copy has marginal distributions identical to the first copy, whereas the possible dependency among predictors is completely destroyed
 - A necessary drawback is that the resulting dataset is twice as large as the original

Fall 2025

CSE552 Machine Learning

137

Thanks for listening!

138

137