

BOOSTING

The boosting question

Simple (weak) classifiers are good!

Logistic regression
w/ simple features

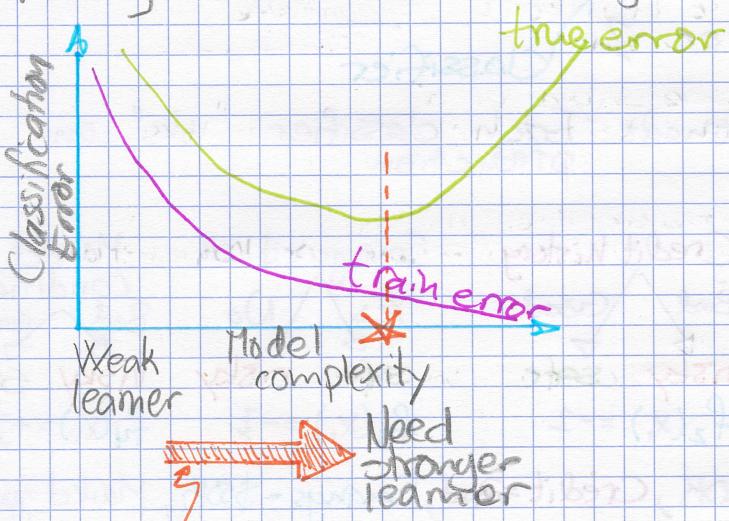
Shallow decision trees

Decision stumps

Low variance, learning is fast!

But high bias

Finding a classifier that's just right



Option 1: add more features or depth

Option 2: ???

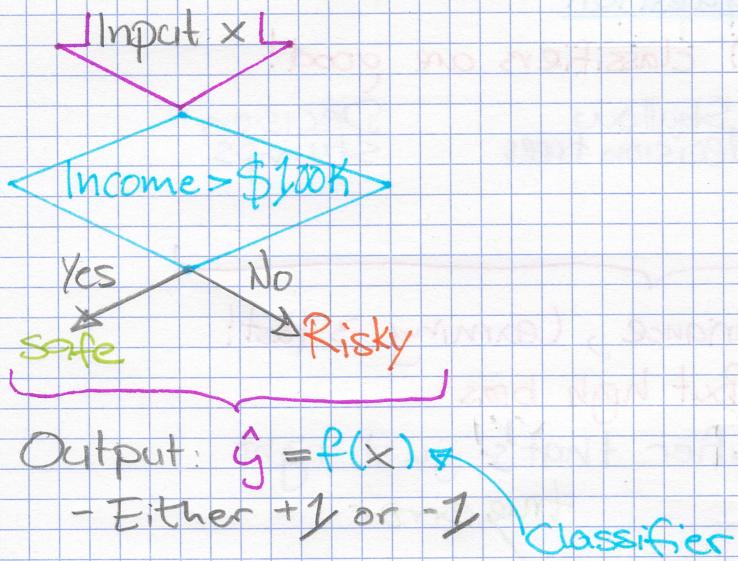
"Can a set of weak learners be combined to create a stronger learner?"

Yes! Schapire (1990)

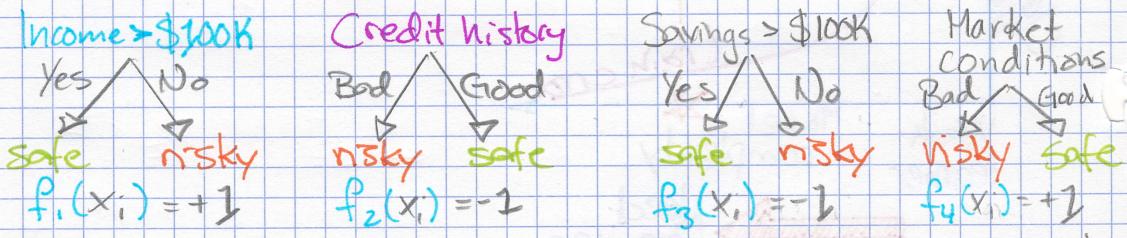
Boosting

{ Amazing impact • simple approach
Widely used in industry
Wins most Kaggle competition

Ensemble classifier



Ensemble methods: Each classifier "votes" on prediction



$$x_i = (\text{income} = \$120K, \text{Credit} = \text{Bad}, \text{Savings} = \$50K, \text{Market} = \text{good})$$

Combine?

$$F(x_i) = \text{Sign}(w_1 f_1(x_i) + w_2 f_2(x_i) + w_3 f_3(x_i) + w_4 f_4(x_i))$$

↑
↑
↑
↑

Ensemble model learn coefficients

Prediction with ensemble

$$= \text{Sign}(2(+1) + 1.5(-1) + 1.5(-1) + 0.5(+1)) =$$

$$= \text{Sign}(-0.5) \Rightarrow \hat{y}_i = -1$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{matrix} \begin{matrix} 2 \\ 1.5 \\ 1.5 \\ 0.5 \end{matrix}$$

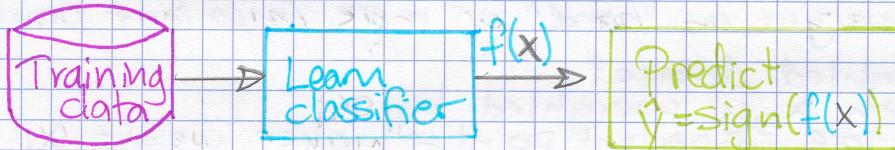
Ensemble classifier in general

- Goal:
 - Predict output y
 - Either $+1$ or -1
 - From input x
- Learn ensemble model:
 - Classifiers: $f_1(x), f_2(x), \dots, f_T(x)$
 - Coefficients: $\hat{\omega}_1, \hat{\omega}_2, \dots, \hat{\omega}_T$
- Prediction:

$$\hat{y} = \text{sign}\left(\sum_{t=1}^T \hat{\omega}_t f_t(x)\right)$$

Boosting

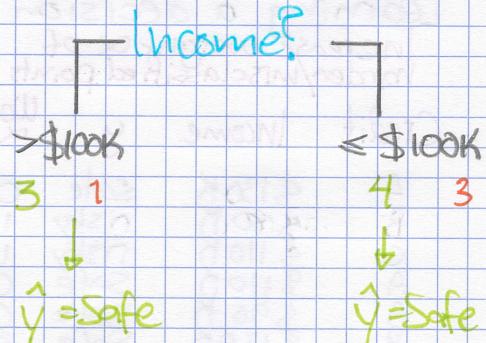
Training a classifier



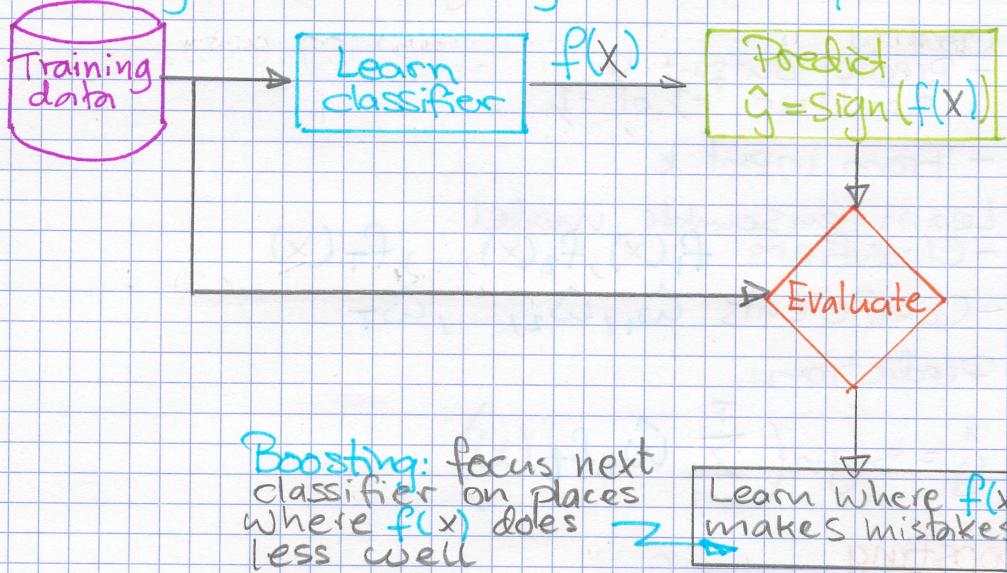
Learning decision stump

Credit Income y

A	\$130K	Safe
B	\$80K	risky
C	\$110K	risky
D	\$110K	Safe
E	\$90K	Safe
F	\$120K	Safe
G	\$30K	risky
H	\$60K	risky
I	\$95K	Safe
J	\$60K	Safe
K	\$98K	Safe



Boosting = Focus learning on "hard" points



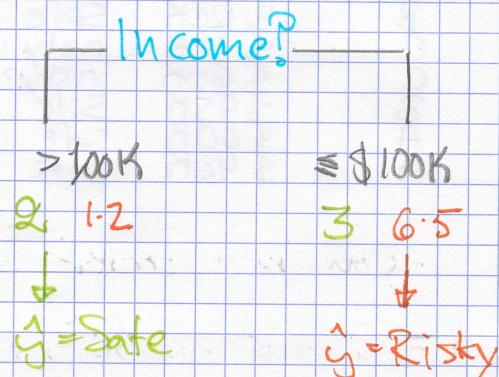
Learning on weighted data:
More weight on "hard" or more important points

- Weighted dataset:
 - Each x_i, y_i weighted by α_i
 - More important point = higher weight α_i
- Learning:
 - Data point i counts as α_i data points
 - Eg., $\alpha_i=2 \rightarrow$ count point twice

Learning a decision stump on weighted data

Increase weight α of hard/misclassified points ↗

Credit	Income	y	Weight
A	\$130K	safe	0.5
B	\$80K	risky	1.5
C	\$110K	risky	1.2
A	\$110K	safe	0.8
A	\$90K	safe	0.6
B	\$120K	safe	0.7
C	\$30K	risky	3
C	\$60K	risky	2
B	\$95K	safe	0.8
A	\$65K	safe	0.7
A	\$98K	safe	0.9



Learning from weighted data in general

- Usually, learning from weighted data
 - Data point i counts as α_i data points
- E.g., gradient ascent for logistic regression:
Sum over data points

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \sum_{i=1}^N \alpha_i h_j(x_i) (1[y_i = +1] - P(y = +1 | x_i, w^{(t)}))$$

Weigh each point by α_i

Boosting = Greedy learning ensembles from data

Training data

↓
Learn classifier

Higher weight for points where $f_i(x)$ is wrong

Predict $\hat{y} = \text{sign}(f_i(x))$

↓
Weighted data

↓
Learn classifier & coefficient

Predict $\hat{y} = \text{sign}(\hat{w}_1 f_1(x) + \hat{w}_2 f_2(x))$

AdaBoost overview

AdaBoost: learning ensemble
[Freund & Schapire 1999]

- Start same weight for all points: $\alpha_i = 1/N$

For $t = 1, T$

- Learn $f_t(x)$ with data weights α_i
- Compute coefficient \hat{w}_t
- Recompute weights α_i

Problem 1: how much do I trust f_t ?

Problem 2: weigh mistakes more?

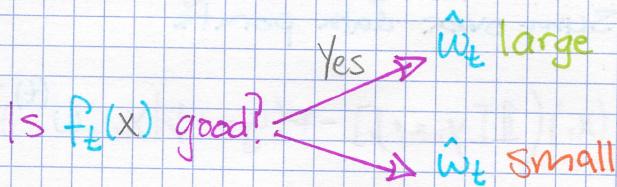
- Final model predicts by:

$$\hat{y} = \text{sign} \left(\sum_{t=1}^T \hat{w}_t f_t(x) \right)$$

Weighted errors

Computing coefficient \hat{w}_t

AdaBoost: Computing coefficient \hat{w}_t of classifier $f_t(x)$



- $f_t(x)$ is good $\rightarrow f_t$ has low training error
- Measuring error in weighted data?
 - Just weighted # of misclassified points

Weighted classification error

Data point		"Sushi was great" "Food was OK"	Learned classifier $\hat{y} = +, +$	Correct!	Weight of Correct
(Sushi was great, +, $\alpha=1.2$)					1.2
(Food was OK, -, $\alpha=0.5$)				Mistake! +	0.5
Hide label					Weight of mistakes

Weighted classification error

- Total weight of mistakes:

$$= \sum_{i=1}^N \alpha_i \cdot \mathbb{I}(y_i \neq \hat{y}_i)$$

- Total weight of all points:

$$= \sum_{i=1}^N \alpha_i$$

- Weighted error measures fraction of weight of mistakes

$$\text{Weighted_error} = \frac{\text{total weight of mistakes}}{\text{total weight of all data points}}$$

- Best possible value is 0.0 \rightarrow worst 1.0

Computing coefficient of each ensemble component

AdaBoost: Formula for computing coefficient \hat{w}_t of classifier $f_t(x)$

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$

	Weighted_error(f_t) on training data	$\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)}$	\hat{w}
Is $f_t(x)$ good?	Yes 0.01	$\frac{1 - 0.01}{0.01} = 99$	$\frac{1}{2} \ln 99 \approx 2.3$
No	$\begin{cases} 0.5 \\ 0.99 \end{cases}$	$\frac{1 - 0.5}{0.5} = 1$ $\frac{1 - 0.99}{0.99} = 0.01$	0 -2.3

⇒ Terrible classifier, but $1 - f_t$ is awesome!!

AdaBoost: learning ensemble

- Start same weight for all points: $\alpha_i = 1/N$

- For $t=1, \dots, T$
 - Learn $f_t(x)$ with data weights α_i

- Compute coefficient \hat{w}_t

- Recompute weights α_i

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$

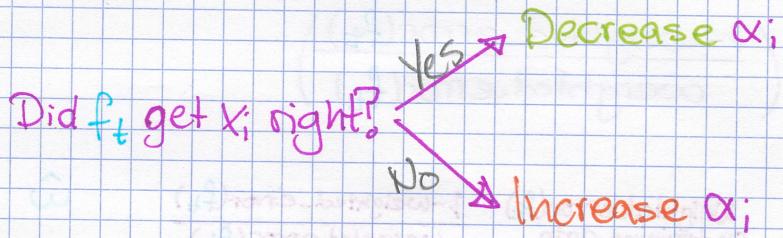
$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(x_i) = y_i \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(x_i) \neq y_i \end{cases}$$

- Final model predicts by

$$\hat{y} = \text{Sign} \left(\sum_{t=1}^T \hat{w}_t f_t(x) \right)$$

Recompute weights α_i

AdaBoost: Updating weights α_i based on where classifier $f_t(x)$ makes mistakes



AdaBoost: Formula for updating weights α_i :

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(x_i) = y_i \leftarrow \text{correct} \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(x_i) \neq y_i \leftarrow \text{mistake} \end{cases}$$

Did f_t get x_i right?	$f_t(x_i) = y_i$	\hat{w}_t	Multiply α_i by $e^{-\hat{w}_t}$	Implication: Decrease importance of x_i, y_i
Yes	Correct	2.3	$e^{-2.3} = 0.1$	
No	Correct	0	$e^0 = 1$	Keep importance the same
	Mistake	9.3	$e^{9.3} = 9.98$	Increase importance of x_i, y_i
	Mistake	0	$e^0 = 1$	Keep importance the same

Normalizing weights

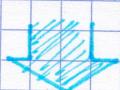
AdaBoost: Normalizing weights α_i :

If x_i often mistake, weights α_i gets very large



If x_i often correct, weights α_i gets very small

Can cause numerical instability after many iterations



COUT'D

Normalize weights to add up to 1 after every iteration

$$\alpha_i \leftarrow \frac{\alpha_i}{\sum_{j=1}^N \alpha_j}$$

AdaBoost: learning ensemble

- Start same weight for all points $\alpha_i = 1/N$

- For $t=1, \dots, T$

- Learn $f_t(x)$ with data weights α_i

- Compute coefficient \hat{w}_t

- Recompute weights α_i

- Normalize weights α_i

- Final model predicts by

$$\hat{y} = \text{Sign} \left(\sum_{t=1}^T \hat{w}_t f_t(x) \right)$$

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right)$$

$$\alpha_i \leftarrow \begin{cases} \alpha_i e^{-\hat{w}_t}, & \text{if } f_t(x_i) = y_i; \\ \alpha_i e^{\hat{w}_t}, & \text{if } f_t(x_i) \neq y_i \end{cases}$$

$$\alpha_i \leftarrow \frac{\alpha_i}{\sum_{j=1}^N \alpha_j}$$

Boosted decision stumps

Finding best next decision stump $f_t(x)$

Consider splitting on each feature:

Income > \$100K Credit history Savings > \$100K

Yes ↗ No ↘

Bad ↗ Good ↘

Yes ↗ No ↘

Safe ↗ Risky ↘

Risky ↗ Safe ↘

Safe ↗ Risky ↘

weighted error
= 0.2

weighted error
= 0.35

weighted error
= 0.3

Market
Conditions

Bad ↗ Good

Risky ↗ Safe

weighted error
= 0.4

$$f_t = \begin{cases} \text{Income} > \$100k \\ \text{Yes} \quad \text{No} \\ \text{Safe} \quad \text{Risky} \end{cases}$$

$$\hat{w}_t = \frac{1}{2} \ln \left(\frac{1 - \text{weighted_error}(f_t)}{\text{weighted_error}(f_t)} \right) = 0.69$$

Updating weights α_i

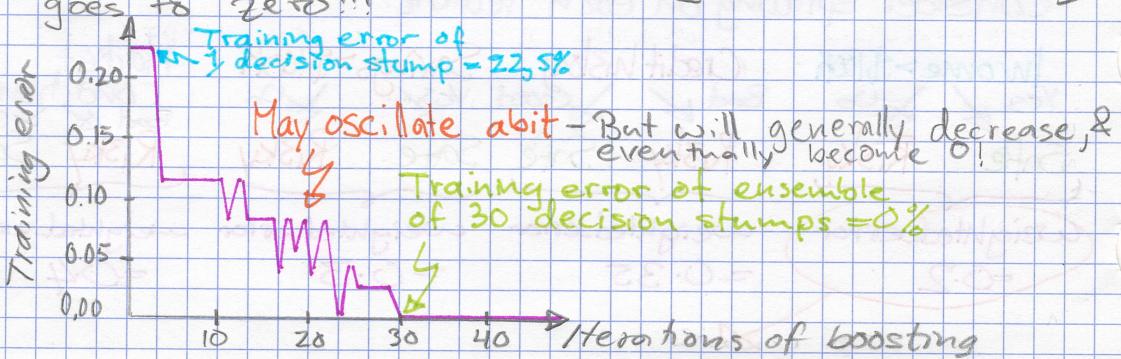
Credit	Income	y	\hat{y}	Previous weight: α_i	New weight: α_i'
A	\$130K	Safe	Safe	0.5	$0.5/2 = 0.25$
B	\$80K	Risky	Risky	1.5	0.75
C	\$110K	Risky	Safe	1.5	$2 * 1.5 = 3$
A	\$110K	Safe	Safe	2	1
A	\$90K	Safe	Risky	1	2
B	\$120K	Safe	Safe	2.5	1.25
C	\$30K	Risky	Risky	3	1.5
C	\$60K	Risky	Risky	2	1
B	\$95K	Safe	Risky	0.5	1
A	\$60K	Safe	Risky	1	2
A	\$98K	Safe	Risky	0.5	1

$$\text{Income} > \$100k \quad \begin{cases} \text{Yes} \quad \text{Safe} \\ \text{No} \quad \text{Risky} \end{cases} \quad \begin{cases} \alpha_i e^{-0.69} = \alpha_i / 2, f_t(x_i) = y_i \\ \alpha_i e^{0.69} = 2\alpha_i, f_t(x_i) \neq y_i \end{cases}$$

The Boosting Theorem

Boosting convergence & overfitting

After some iterations, training error of boosting goes to zero!!!



AdaBoost Theorem

Under some technical conditions...

Training error of boosted classifier $\rightarrow 0$ as $T \rightarrow \infty$

Condition = At every t , can find a weak learner with weighted error(f_t) < 0.5

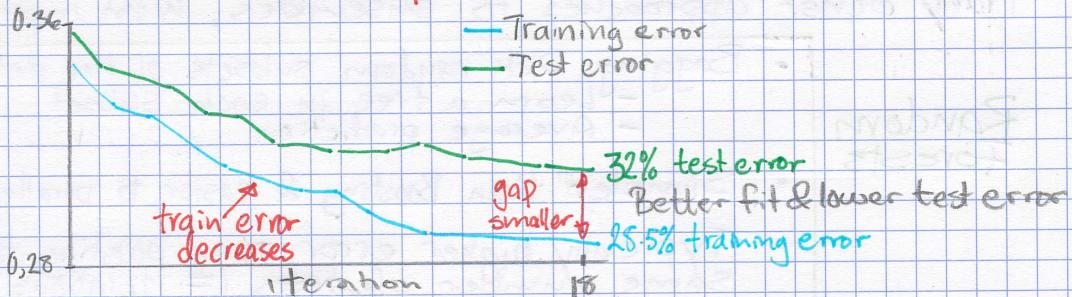
Not always possible!

Extreme example:
No classifier can separate a +1 on top of -1

Nonetheless, boosting often yields great training error

Overfitting in Boosting

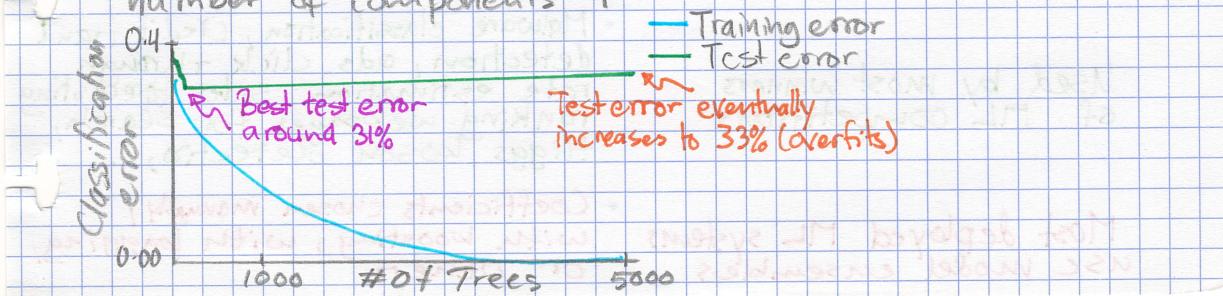
Boosted decision stumps on loan data



Boosting tends to be robust to overfitting

"Test set performance about constant for many iterations
→ Less sensitive to choice of T "

But boosting will eventually overfit, so must choose max number of components T



How do we decide when to stop boosting

Choosing T?

like selecting parameters in other ML approaches (e.g., λ in regularization)

Not on training data

Not useful: training error improves as T increases

Never ever ever ever on test data

Validation set

If dataset is large

Cross-validation

For smaller data

Variants of Boosting and related algorithms

There are hundreds of variants of boosting; most important

Gradient boosting

- Like AdaBoost, but useful beyond basic classification

Many other approaches to ensembles, most important

Random forests

- Bagging: Pick random subsets of the data
 - Learn a tree in each subset
 - Average prediction
- Simpler than Boosting & easier to parallelize
- Typically higher error than boosting for same number of trees (# iterations T)

Impact of Boosting

Amongst most useful ML methods ever created

Extremely useful in computer vision

- Standard approach for face detection

Used by most winners of ML competitions

- Malware classification, credit fraud detection, ads click-through rate estimation, sales forecasting, ranking webpages for search, Higgs boson detection, ...

Most deployed ML systems use model ensembles

- Coefficients chosen manually, with boosting, with bagging, or others