



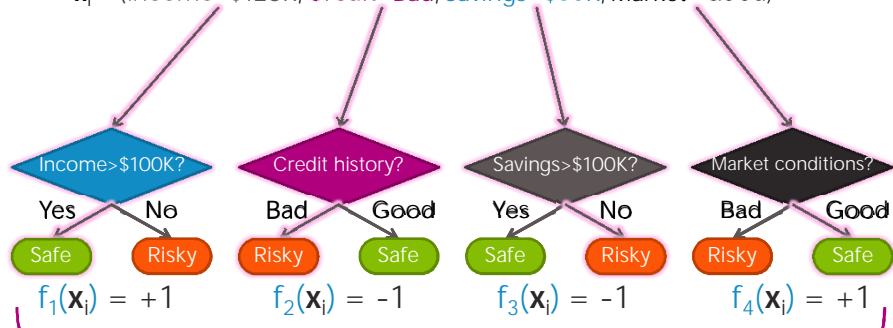
Boosting

CSE 446: Machine Learning
Emily Fox
University of Washington
February 1, 2017

©2017 Emily Fox

Ensemble methods: Each classifier “votes” on prediction

$\mathbf{x}_i = (\text{Income}=\$120\text{K}, \text{Credit}=\text{Bad}, \text{Savings}=\$50\text{K}, \text{Market}=\text{Good})$



Combine?

Ensemble model

Learn coefficients

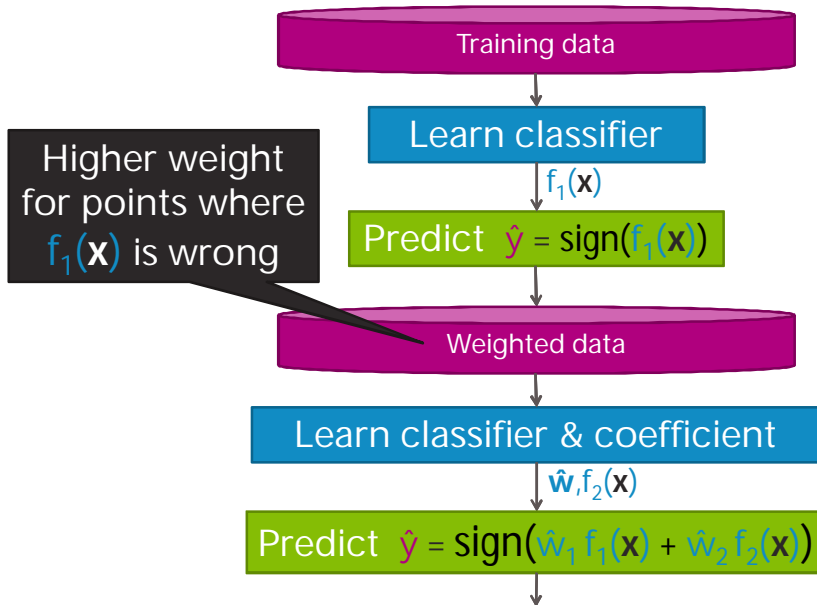
$$F(\mathbf{x}_i) = \text{sign}(w_1 f_1(\mathbf{x}_i) + w_2 f_2(\mathbf{x}_i) + w_3 f_3(\mathbf{x}_i) + w_4 f_4(\mathbf{x}_i))$$

2

©2017 Emily Fox

CSE 446: Machine Learning

Boosting = Greedy learning ensembles from data



3

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost: learning ensemble

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

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

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

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

- Final model predicts by:

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

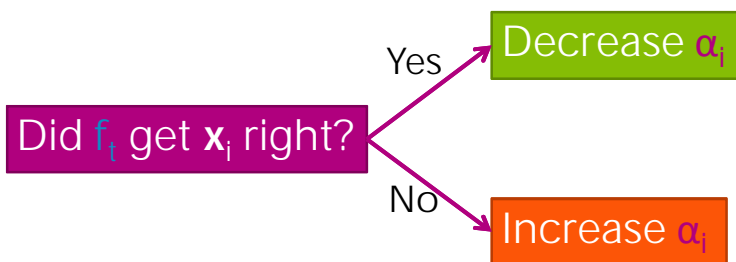
4

©2017 Emily Fox

CSF 446: Machine Learning

Recompute weights α_j

AdaBoost: Updating weights α_j based on where classifier $f_t(\mathbf{x})$ makes mistakes



AdaBoost: Formula for updating weights α_i

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

	$f_t(\mathbf{x}_i) = y_i$?	\hat{w}_t	Multiply α_i by	Implication
Did f_t get \mathbf{x}_i right?	Yes	2.3	$e^{-2.3} = 0.1$	decrease importance of \mathbf{x}_i, y_i
	yes	0	$e^0 = 1$	no change
	no	2.3	$e^{2.3} = 9.98$	increase importance of \mathbf{x}_i, y_i
	no	0	$e^0 = 1$	no change

7

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost: learning ensemble

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

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

- Learn $f_t(\mathbf{x})$ with data weights α_i

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

- Compute coefficient \hat{w}_t

- Recompute weights α_i

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

- Final model predicts by:

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

8

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost: Normalizing weights α_i

If x_i often mistake,
weight α_i gets very
large

If x_i often correct,
weight α_i gets very
small

Can cause numerical instability
after many iterations

Normalize weights to
add up to 1 after every iteration

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

9

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost: learning ensemble

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

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

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

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

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

- Final model predicts by:

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

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

10

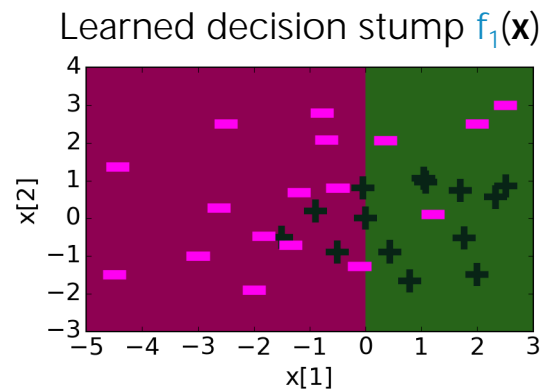
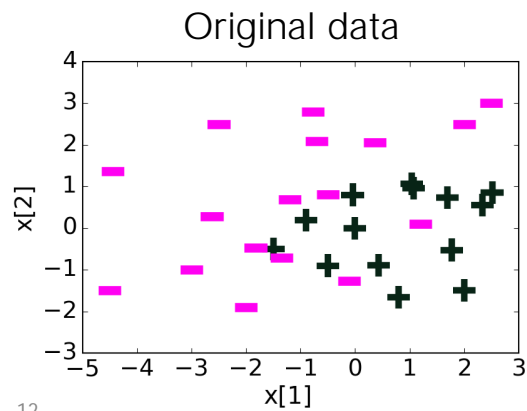
©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost example: A visualization

t=1: Just learn a classifier on original data

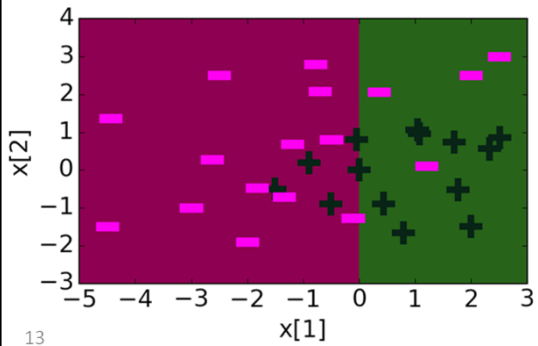
all pts have same weight *standard learning*



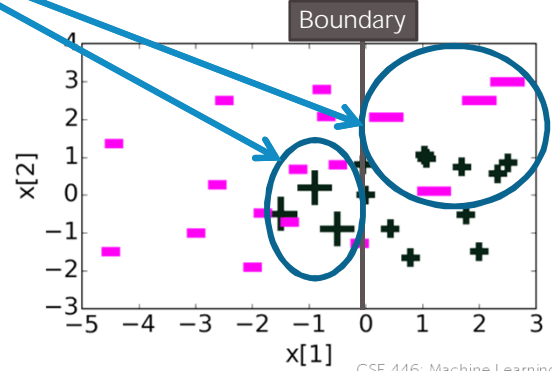
Updating weights α_i

Increase weight α_i of misclassified points

Learned decision stump $f_1(x)$



New data weights α_i

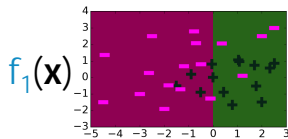


13

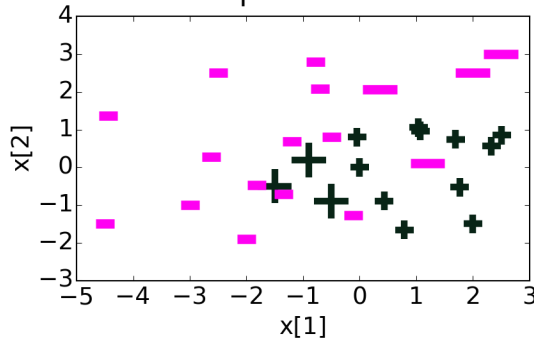
©2017 Emily Fox

CSE 446: Machine Learning

t=2: Learn classifier on weighted data

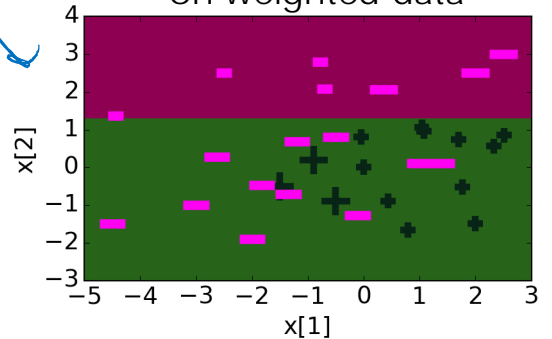


Weighted data: using α_i chosen in previous iteration



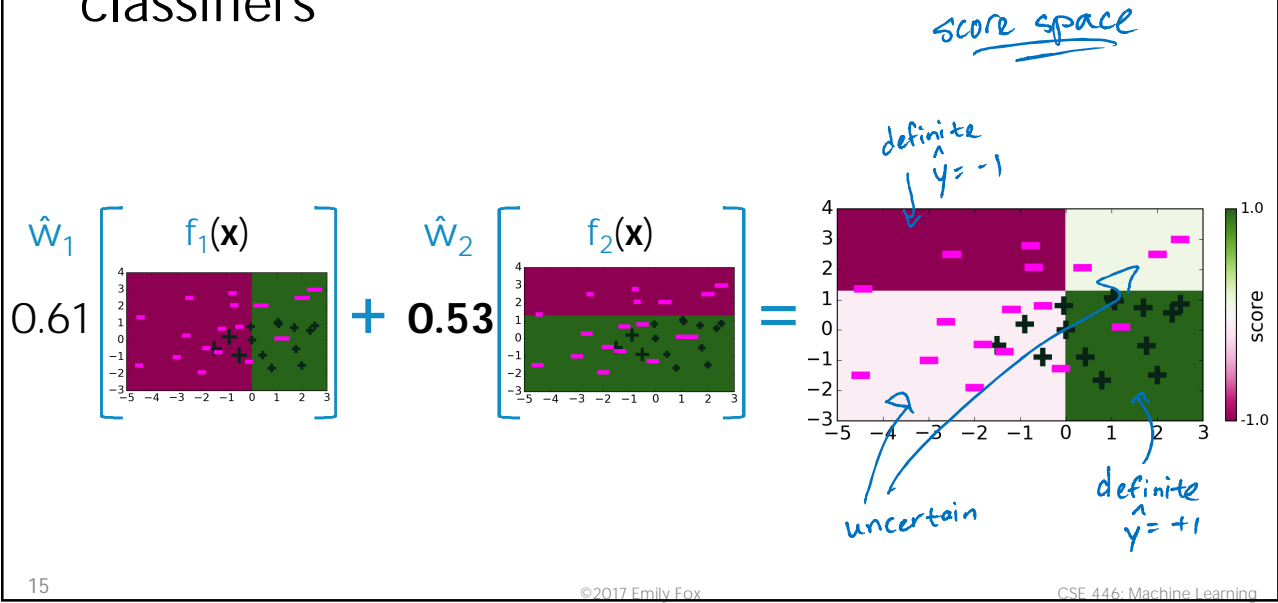
better split for these weights

Learned decision stump $f_2(x)$ on weighted data



©2017 Emily Fox

Ensemble becomes weighted sum of learned classifiers

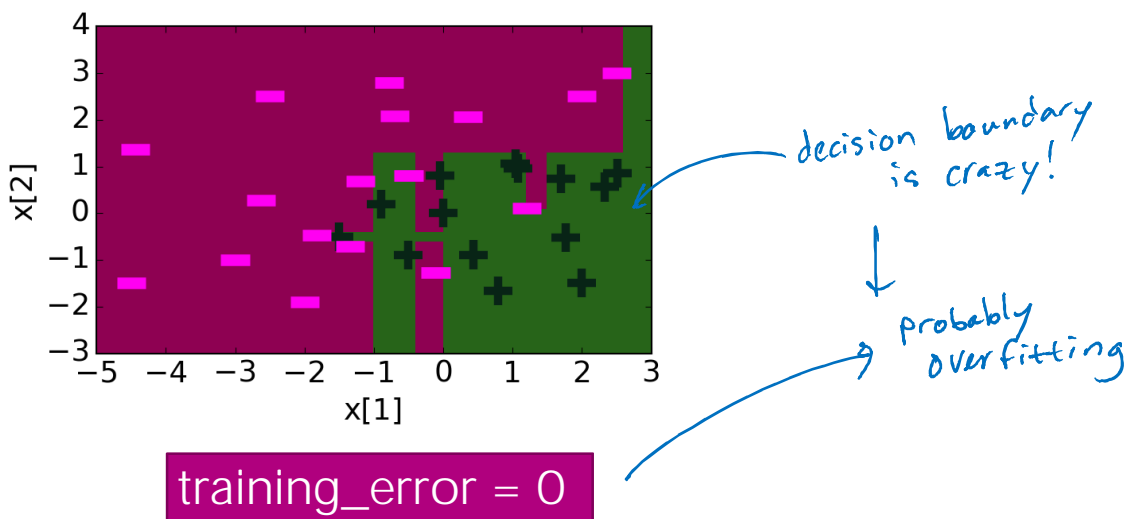


15

©2017 Emily Fox

CSF 446: Machine Learning

Decision boundary of ensemble classifier after 30 iterations



16

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost example: Boosted decision stumps step-by-step

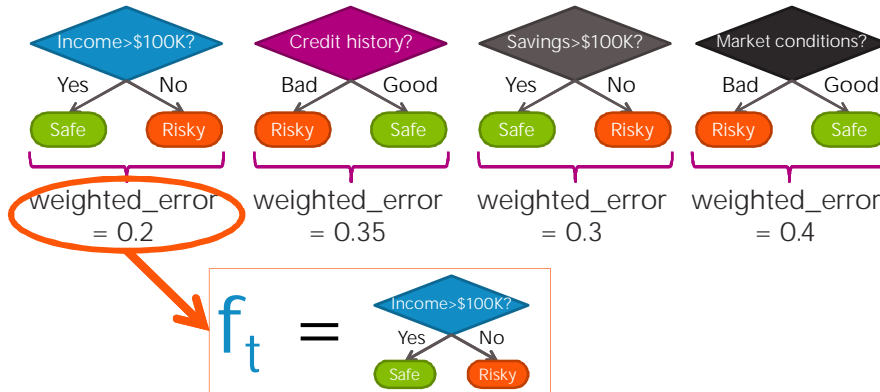
Boosted decision stumps

- Start same weight for all points: $\alpha_i = 1/N$
- For $t = 1, \dots, T$
 - Learn $f_t(\mathbf{x})$: pick decision stump with lowest weighted training error according to α_i
 - Compute coefficient \hat{w}_t
 - Recompute weights α_i
 - Normalize weights α_i
- Final model predicts by:

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

Finding best next decision stump $f_t(\mathbf{x})$

Consider splitting on each feature:



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

19

©2017 Emily Fox

CSF 446: Machine Learning

Boosted decision stumps

- Start same weight for all points: $\alpha_i = 1/N$
- For $t = 1, \dots, T$
 - Learn $f_t(\mathbf{x})$: pick decision stump with lowest weighted training error according to α_i
 - Compute coefficient \hat{w}_t
 - Recompute weights α_i
 - Normalize weights α_i
- Final model predicts by:

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

20

©2017 Emily Fox

CSF 446: Machine Learning

Updating weights α_j



$$\alpha_j \leftarrow \begin{cases} \alpha_j e^{-\hat{w}_t} = \alpha_j e^{-0.69} = \alpha_j / 2, & \text{if } f_t(\mathbf{x}_i) = y_i \\ \alpha_j e^{\hat{w}_t} = \alpha_j e^{0.69} = 2 \alpha_j, & \text{if } f_t(\mathbf{x}_i) \neq y_i \end{cases}$$

Credit	Income	y	\hat{y}	Previous weight α	New weight α
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

21

©2017 Emily Fox

CSF 446: Machine Learning

Boosting convergence & overfitting

Boosting question revisited

"Can a set of weak learners be combined to create a stronger learner?" Kearns and Valiant (1988)



Yes! Schapire (1990)



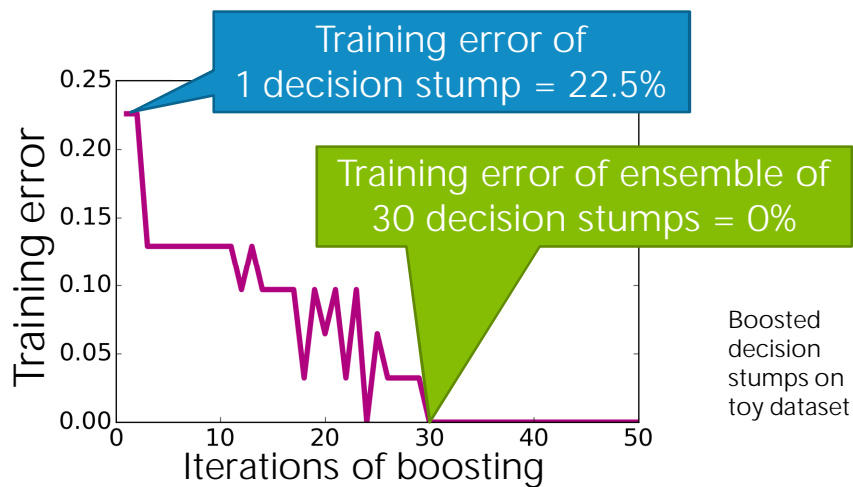
Boosting

23

©2017 Emily Fox

CSF 446: Machine Learning

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



24

©2017 Emily Fox

CSF 446: Machine Learning

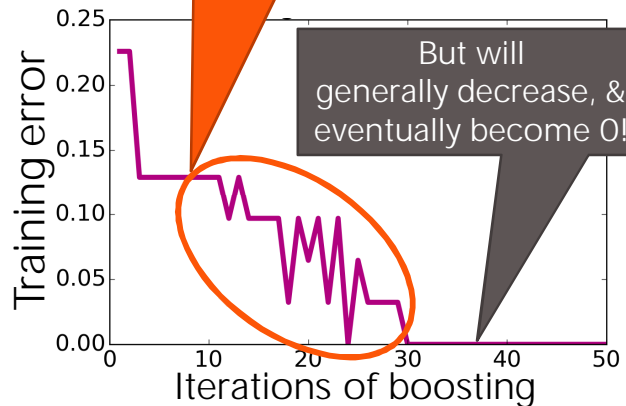
AdaBoost Theorem

Under some technical conditions...



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

May oscillate a bit



25

©2017 Emily Fox

CSF 446: Machine Learning

Condition of 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 $\text{weighted_error}(f_t) < 0.5$

Not always possible

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



Nonetheless, boosting often yields great training error

26

©2017 Emily Fox

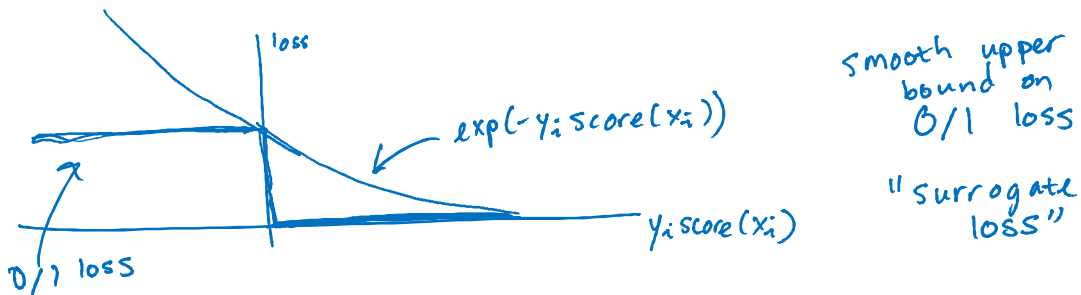
CSF 446: Machine Learning

AdaBoost Theorem more formally

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}[F(x_i) \neq y_i] \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i \text{score}(x_i))$$

Where $\text{score}(x) = \sum_t \hat{w}_t f_t(x)$; $F(x) = \text{sign}(\text{score}(x))$



27

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost Theorem more formally

Training error of final classifier is bounded by:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}[F(x_i) \neq y_i] \leq \frac{1}{N} \sum_{i=1}^N \exp(-y_i \text{score}(x_i)) = \prod_{t=1}^T Z_t$$

$$Z_t = \sum_{i=1}^N D_t(i) \exp(-\hat{w}_t y_i f_t(x_i))$$

Where $\text{score}(x) = \sum_t \hat{w}_t f_t(x)$; $F(x) = \text{sign}(\text{score}(x))$

Pf: HW " ... "telescopic sums"
 if $Z_t < 1 \forall t \Rightarrow$ as $T \rightarrow \infty$, training error $\rightarrow 0$

28

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost Theorem more formally

If we minimize $\prod_{t=1}^T Z_t$, we minimize our training error

We can tighten this bound greedily by choosing \hat{w}_t, f_t on each iteration to minimize:

$$Z_t = \sum_{i=1}^N D_t(i) \exp(-\hat{w}_t y_i f_t(x_i))$$

take der. +
set = 0

For boolean target function, this is accomplished by [Freund & Schapire '97]:

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

weighted class. error

29

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost Theorem more formally

If each classifier is (at least slightly) **better than random**

$$\epsilon_t < 0.5$$

AdaBoost will achieve **zero training error** (exponentially fast):

$$\frac{1}{N} \sum_{i=1}^N \mathbb{I}[F(x_i) \neq y_i] \leq \prod_{t=1}^T Z_t \leq \exp \left(-2 \sum_{t=1}^T (1/2 - \epsilon_t)^2 \right)$$

upper bound
for $\epsilon_t < 1/2$

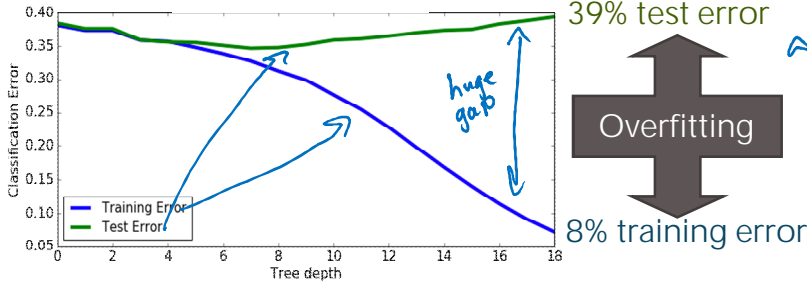
as ϵ_t moves
away from $1/2$,
exponent gets
bigger

30

©2017 Emily Fox

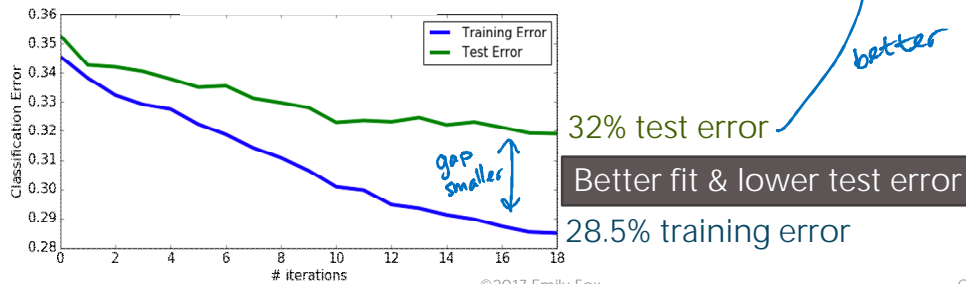
CSF 446: Machine Learning

Decision trees on loan data



39% test error
Overfitting
8% training error

Boosted decision stumps on loan data



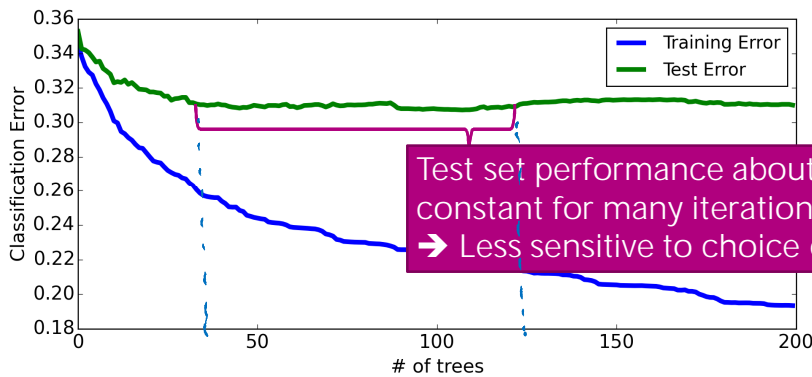
32% test error
Better fit & lower test error
28.5% training error

31

©2017 Emily Fox

CSF 446: Machine Learning

Boosting tends to be robust to overfitting



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

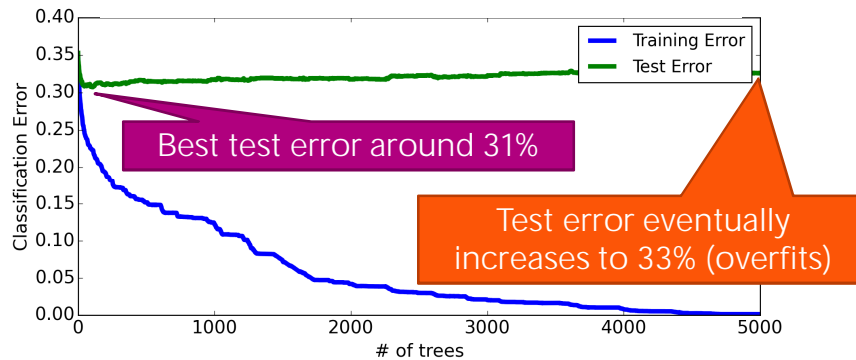
any of these values of T would be fine

32

©2017 Emily Fox

CSF 446: Machine Learning

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

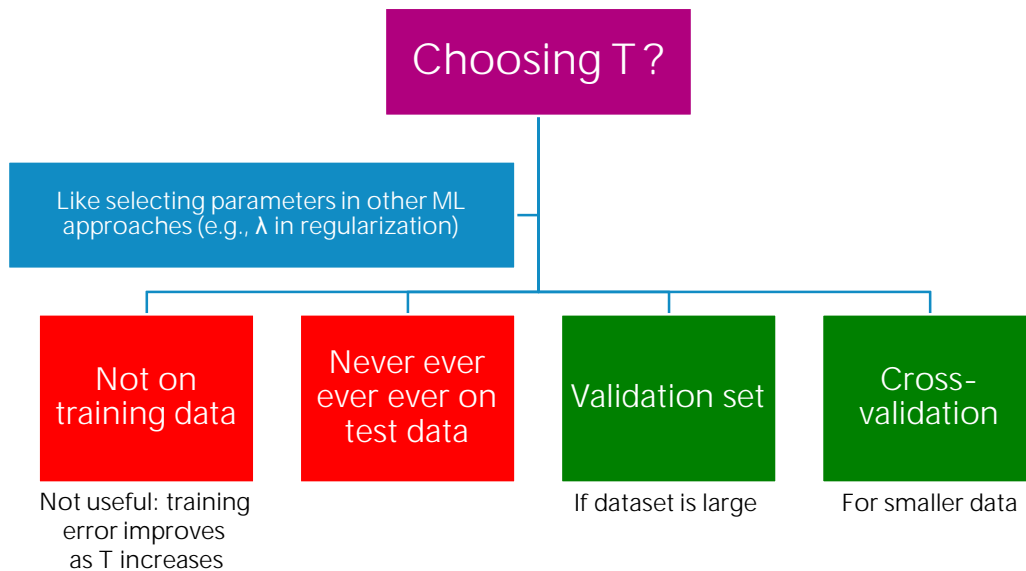


33

©2017 Emily Fox

CSF 446: Machine Learning

How do we decide when to stop boosting?



34

©2017 Emily Fox

CSF 446: Machine Learning

AdaBoost vs logistic regression

What is logistic regression minimizing?

Logistic regression assumes:

$$P(y = +1 | x) = \frac{1}{1 + \exp(-\text{score}(x_i))}$$

$$\text{score}(x_i) = \sum_j w_j h_j(x_i)$$

And tries to maximize data likelihood:

$$P(\mathcal{D} | \mathbf{w}, \mathbf{x}) = \prod_{i=1}^N \frac{1}{1 + \exp(-y_i \text{score}(x_i))}$$

$$\Rightarrow \min_{\mathbf{w}} -\ln P(\mathcal{D} | \mathbf{w}, \mathbf{x})$$

Equivalent to minimizing log loss

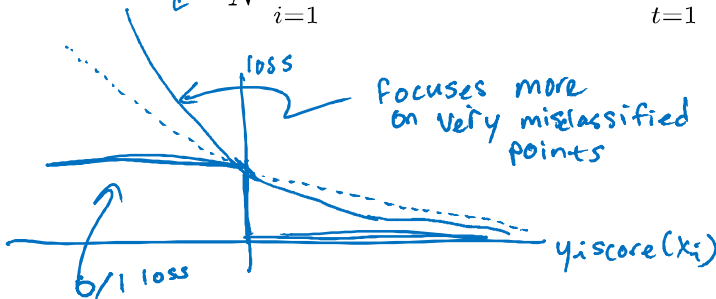
$$\sum_{i=1}^N \ln(1 + \exp(-y_i \text{score}(x_i)))$$

Logistic regression vs boosting objectives

Logistic regression minimizes: $\sum_{i=1}^N \ln(1 + \exp(-y_i \text{score}(x_i)))$

Boosting minimizes similar loss function: $\text{score}(x) = \sum_t \hat{w}_t f_t(x)$

$$\frac{1}{N} \sum_{i=1}^N \exp(-y_i \text{score}(x_i)) = \prod_{t=1}^T Z_t$$



Both smooth approximations of 0/1 loss!

37

©2017 Emily Fox

CSF 446: Machine Learning

Logistic regression vs boosting overview

Logistic regression:

- Minimize loss fn

$$\sum_{i=1}^N \ln(1 + \exp(-y_i \text{score}(x_i)))$$

- Define

$$\text{score}(x) = \sum_j \hat{w}_j h_j(x)$$

where features $h_j(\mathbf{x})$ are predefined

- Weights \hat{w}_j learned in joint optimization

Boosting:

- Minimize loss fn

$$\sum_{i=1}^N \exp(-y_i \text{score}(x_i))$$

loss fn (handwritten blue arrow pointing to the equation)

- Define

$$\text{score}(x) = \sum_t \hat{w}_t f_t(x);$$

where $f_t(\mathbf{x})$ defined dynamically to fit data (not a linear classifier)

- Weights \hat{w}_t learned incrementally

38

©2017 Emily Fox

CSF 446: Machine Learning

Summary of boosting

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 learn ensembles, most important:

Random forests

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

Impact of boosting (spoiler alert... HUGE IMPACT)

Amongst most useful ML methods ever created

Extremely useful in computer vision

- Standard approach for face detection, for example

Used by **most winners** of ML competitions (Kaggle, KDD Cup,...)

- 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

41

©2017 Emily Fox

CSF 446: Machine Learning

What you can do now...

- Identify notion ensemble classifiers
- Formalize ensembles as the weighted combination of simpler classifiers
- Outline the boosting framework – sequentially learn classifiers on weighted data
- Describe the AdaBoost algorithm
 - Learn each classifier on weighted data
 - Compute coefficient of classifier
 - Recompute data weights
 - Normalize weights
- Implement AdaBoost to create an ensemble of decision stumps
- Discuss convergence properties of AdaBoost & how to pick the maximum number of iterations T
- Derive why AdaBoost leads to zero training error exponentially fast
- Compare and contrast what AdaBoost and logistic regression are minimizing

42

©2017 Emily Fox

CSF 446: Machine Learning

Instance-Based Learning:

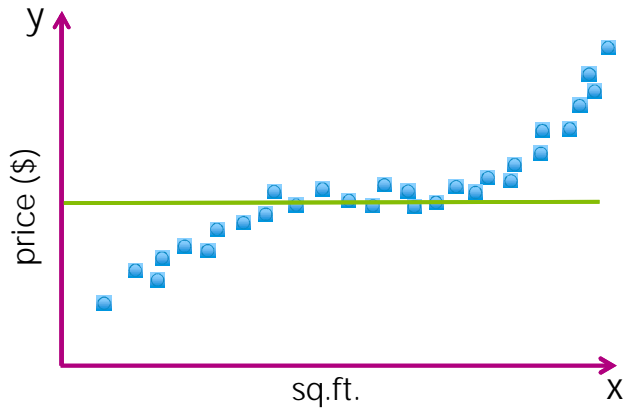
Nearest neighbor and kernel regression and classification

CSE 446: Machine Learning
Emily Fox
University of Washington
February 1, 2017

©2017 Emily Fox

Fit globally vs. fit locally

Parametric models of $f(x)$

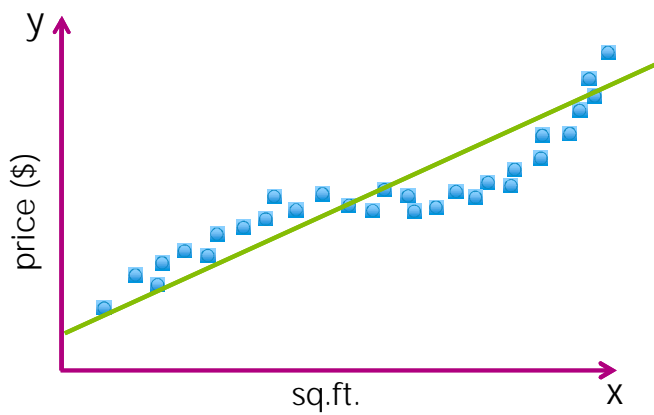


45

©2017 Emily Fox

CSF 446: Machine Learning

Parametric models of $f(x)$

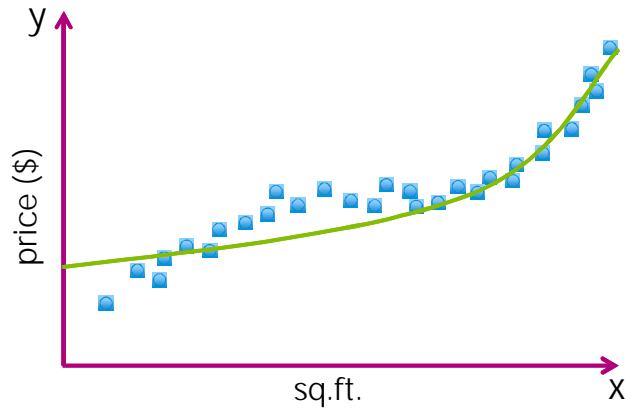


46

©2017 Emily Fox

CSF 446: Machine Learning

Parametric models of $f(x)$

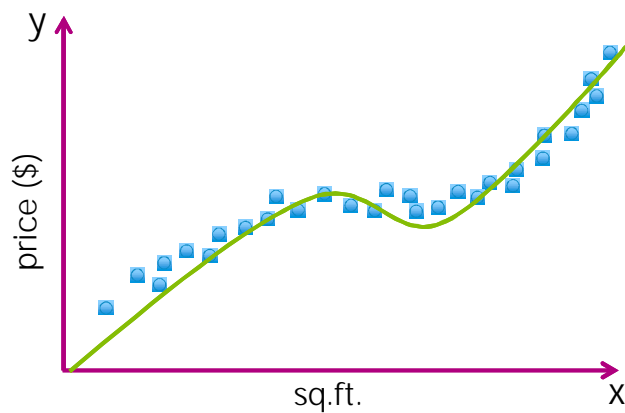


47

©2017 Emily Fox

CSF 446: Machine Learning

Parametric models of $f(x)$

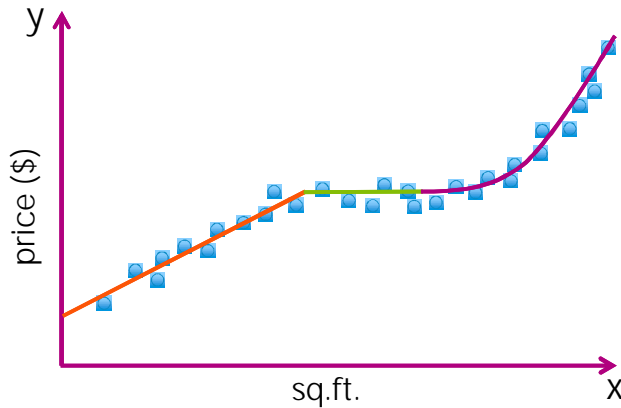


48

©2017 Emily Fox

CSF 446: Machine Learning

$f(x)$ is not really a polynomial



49

©2017 Emily Fox

CSF 446: Machine Learning

What alternative do we have?

If we:

- Want to allow flexibility in $f(\mathbf{x})$ having **local structure**
- Don't want to infer "structural breaks"

What's a simple option we have?

- Assuming we have **plenty of data...**

50

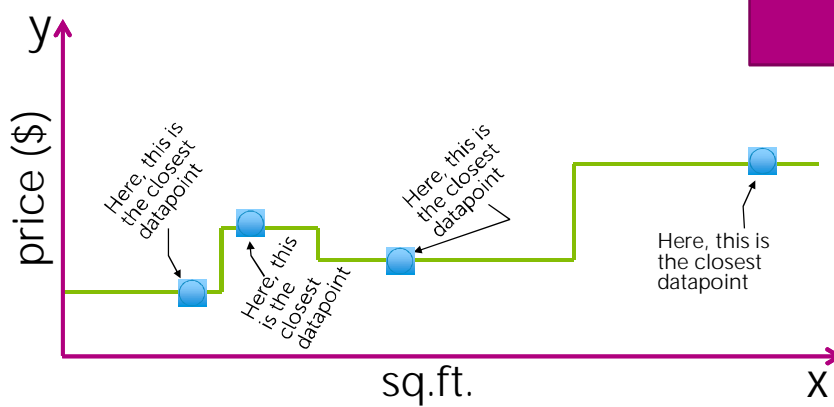
©2017 Emily Fox

CSF 446: Machine Learning

Simplest approach: Nearest neighbor regression

Fit locally to each data point

Predicted value = "closest" y_i



What people do naturally...

Real estate agent assesses value by finding sale of most similar house



53

©2017 Emily Fox

CSF 446: Machine Learning

1-NN regression more formally

Dataset of (house, \$) pairs: $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$

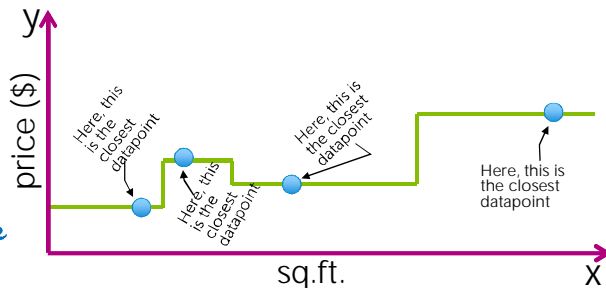
Query point: \mathbf{x}_q *big lime green house \$?*

1. Find "closest" \mathbf{x}_i in dataset

$x_{NN} \leftarrow \min_i \text{distance}(x_i, x_q)$

2. Predict

$\hat{y}_q = y_{NN}$ *sales price of big pink house*

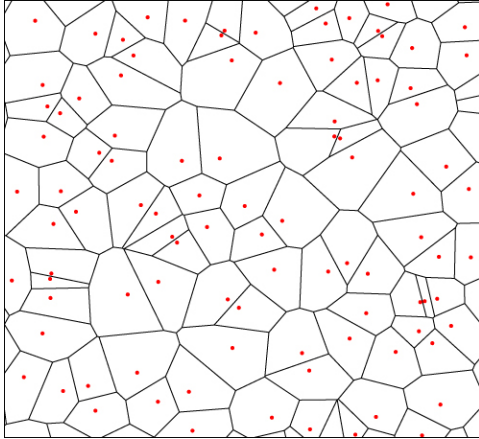


54

©2017 Emily Fox

CSF 446: Machine Learning

Visualizing 1-NN in multiple dimensions



Voronoi tessellation (or diagram):

- Divide space into N regions, each containing 1 datapoint
- Defined such that any \mathbf{x} in region is "closest" to region's datapoint

Don't explicitly form!

55

©2017 Emily Fox

CSF 446: Machine Learning

Distance metrics: Defining notion of "closest"

In 1D, just Euclidean distance:

$$\text{distance}(x_j, x_q) = |x_j - x_q|$$

In multiple dimensions:

- can define many interesting distance functions
- most straightforwardly, might want to weight different dimensions differently

56

©2017 Emily Fox

CSF 446: Machine Learning

Weighting housing inputs

Some inputs are more relevant than others



bedrooms
bathrooms
sq.ft. living
sq.ft. lot
floors
year built
year renovated
waterfront



57

©2017 Emily Fox

CSF 446: Machine Learning

Scaled Euclidean distance

Formally, this is achieved via

$$\text{distance}(\mathbf{x}_j, \mathbf{x}_q) = \sqrt{a_1(\mathbf{x}_j[1] - \mathbf{x}_q[1])^2 + \dots + a_d(\mathbf{x}_j[d] - \mathbf{x}_q[d])^2}$$

weight on each input
(defining relative importance)

Other example distance metrics:

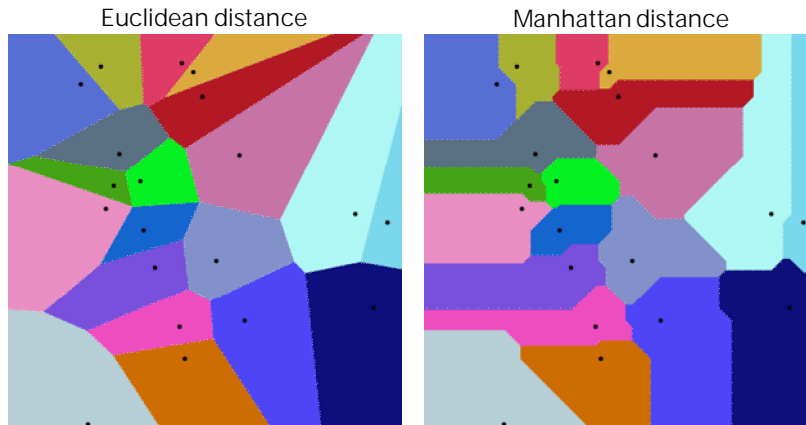
- Mahalanobis, rank-based, correlation-based, cosine similarity, Manhattan, Hamming, ...

58

©2017 Emily Fox

CSF 446: Machine Learning

Different distance metrics lead to different predictive surfaces

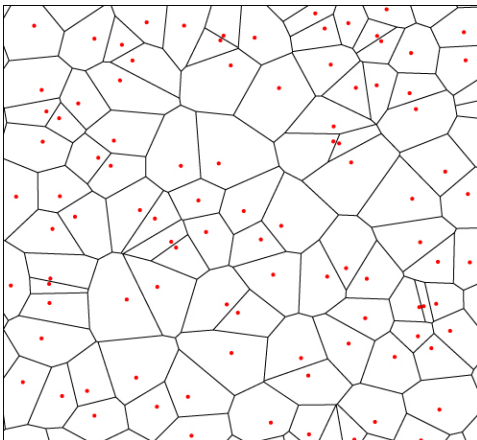


59

©2017 Emily Fox

CSF 446: Machine Learning

Can 1-NN be used for classification?



Yes!!

Just predict class of neighbor

60

©2017 Emily Fox

CSF 446: Machine Learning