



# Classification:

# Analyzing Sentiment

CS 229: Machine Learning

Emily Fox

Stanford University

January 24, 2024

©2024 Emily Fox

1

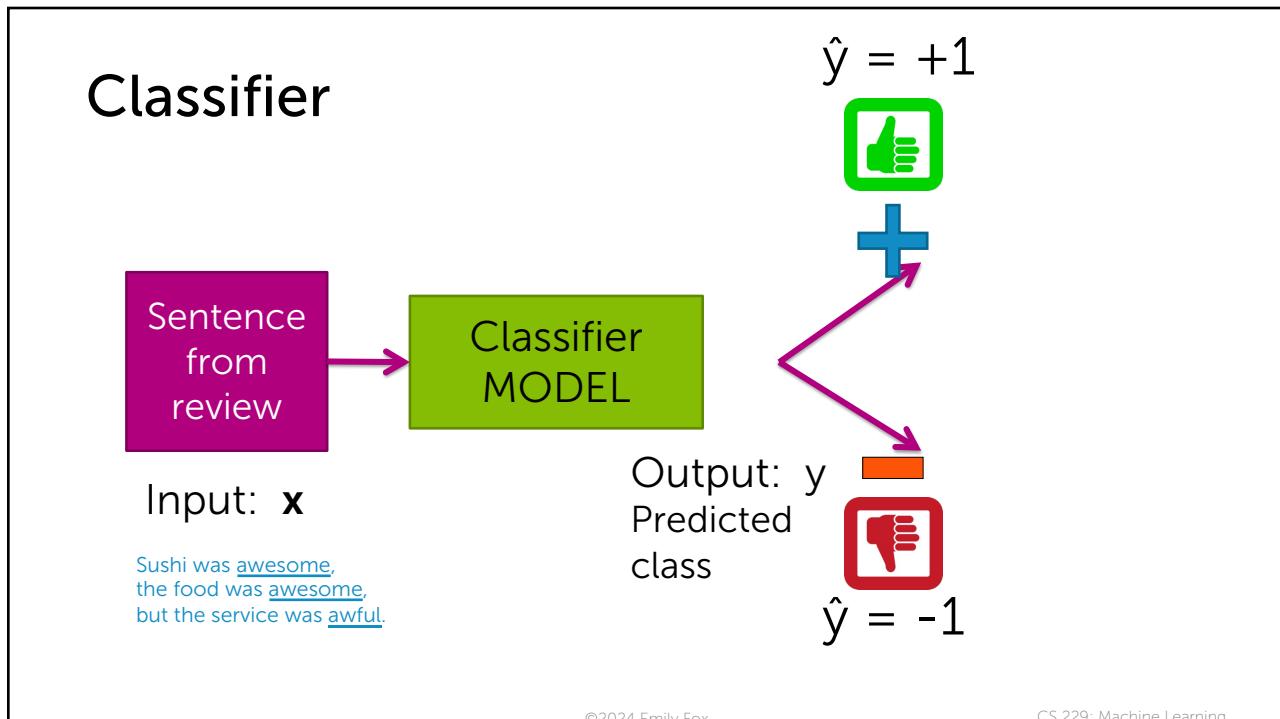
# Classification

©2024 Emily Fox

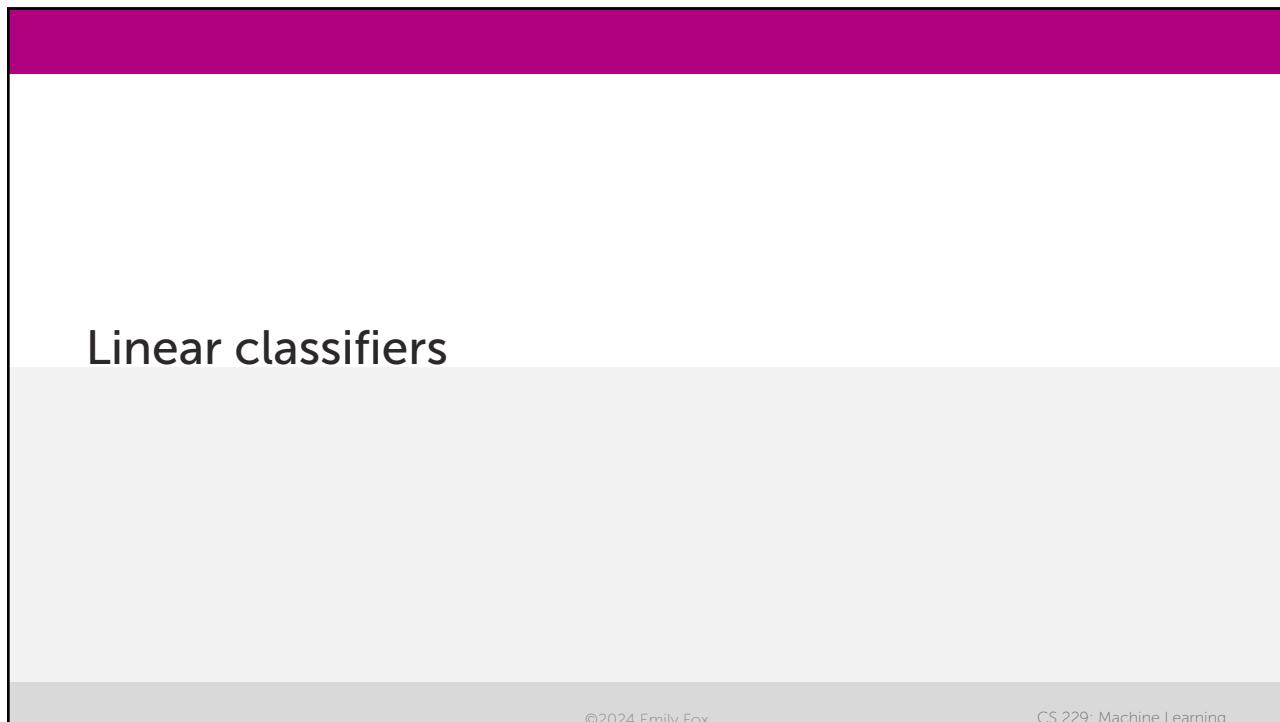
CS 229: Machine Learning

2

1



3



4

Feature	Coefficient
...	...

↓

### Simple linear classifier

Score( $\mathbf{x}$ ) = weighted sum of features of sentence

If Score ( $\mathbf{x}$ ) > 0:  
 $\hat{y} = +1$

Else:  
 $\hat{y} = -1$

Sentence from review

Input:  $\mathbf{x}$



©2024 Emily Fox

CS 229: Machine Learning

5

## A simple example: Word counts

Feature	Coefficient
good	1.0
great	1.2
awesome	1.7
bad	-1.0
terrible	-2.1
awful	-3.3
restaurant, the, we, where, ...	0.0
...	...

Input  $\mathbf{x}_i$ :

Sushi was great,  
the food was awesome,  
but the service was terrible.

*word counts*

$$\begin{aligned} \text{Score}(\mathbf{x}_i) &= 1.2(1) + 1.7(1) - 2.1(1) \\ &= 0.8 > 0 \end{aligned}$$

$\Rightarrow \hat{y}_i = +1$

positive review

Called a linear classifier, because score is weighted sum of features.

©2024 Emily Fox

CS 229: Machine Learning

6

## More generically...

Model:  $\hat{y}_i = \text{sign}(\text{Score}(\mathbf{x}_i))$

$\text{Sign}(\text{pos } \#) = +1$   
 $\text{Sign}(\text{neg } \#) = -1$

$$\text{Score}(\mathbf{x}_i) = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i)$$

$$= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{h}(\mathbf{x}_i)$$

feature 1 =  $h_0(\mathbf{x})$  ... e.g., 1

feature 2 =  $h_1(\mathbf{x})$  ... e.g.,  $\mathbf{x}[1] = \#\text{awesome}$

feature 3 =  $h_2(\mathbf{x})$  ... e.g.,  $\mathbf{x}[2] = \#\text{awful}$

or,  $\log(\mathbf{x}[7]) \mathbf{x}[2] = \log(\#\text{bad}) \times \#\text{awful}$

or, tf-idf("awful")

...

feature  $D+1 = h_D(\mathbf{x})$  ... some other function of  $\mathbf{x}[1], \dots, \mathbf{x}[d]$

©2024 Emily Fox

CS 229: Machine Learning

7

## Decision boundaries

©2024 Emily Fox

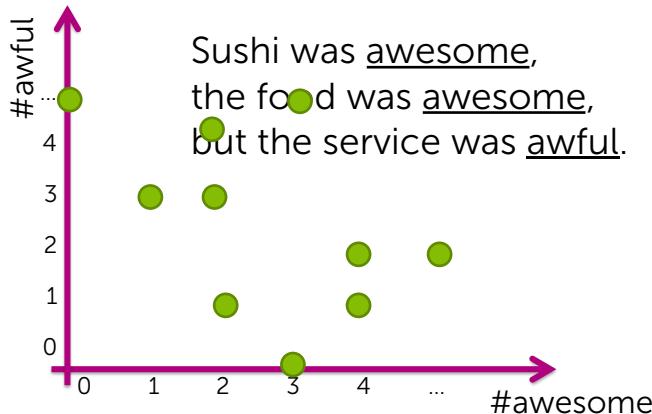
CS 229: Machine Learning

8

## Suppose only two words had non-zero coefficient

Input	Coefficient	Value
	$w_0$	0.0
#awesome	$w_1$	1.0
#awful	$w_2$	-1.5

$$\rightarrow \text{Score}(x) = 1.0 \# \text{awesome} - 1.5 \# \text{awful}$$



©2024 Emily Fox

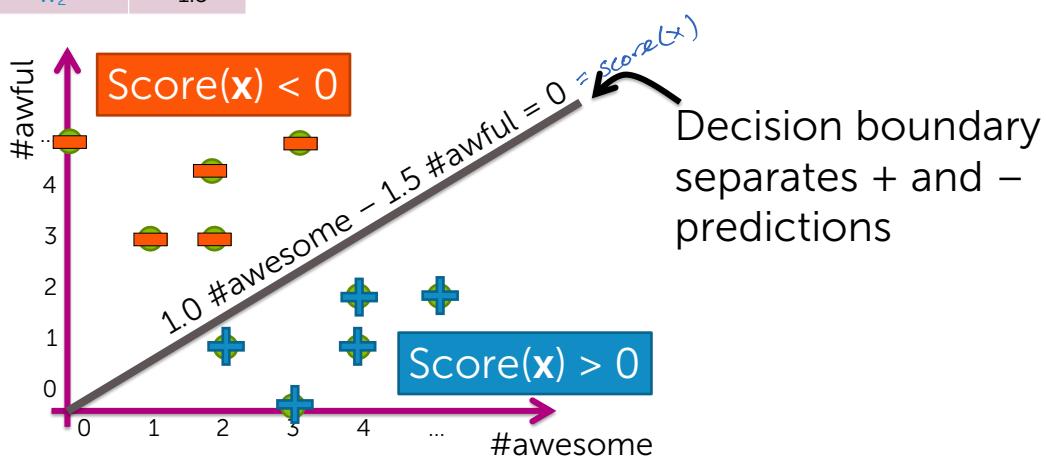
CS 229: Machine Learning

9

## Decision boundary example

Input	Coefficient	Value
	$w_0$	0.0
#awesome	$w_1$	1.0
#awful	$w_2$	-1.5

$$\rightarrow \text{Score}(x) = 1.0 \# \text{awesome} - 1.5 \# \text{awful}$$

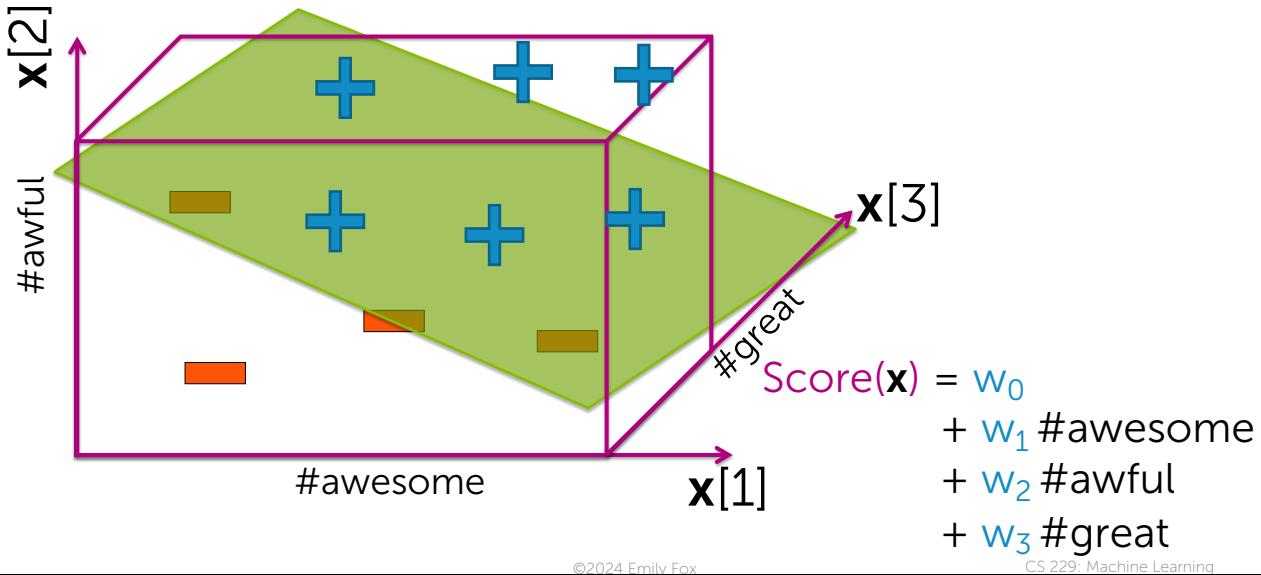


©2024 Emily Fox

CS 229: Machine Learning

10

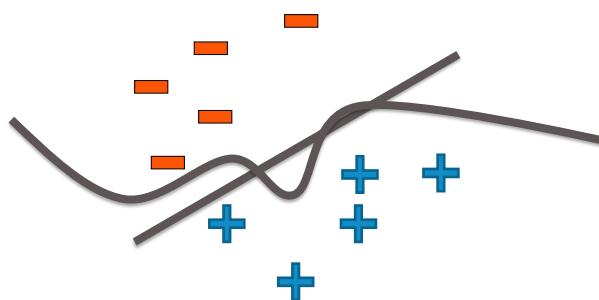
## For more inputs (linear features)...



11

## For general features...

For more general classifiers (not just linear features)  
**→ more complicated shapes**



©2024 Emily Fox

CS 229: Machine Learning

12



# Linear classifiers: Logistic regression

CS 229: Machine Learning

Emily Fox

Stanford University

January 24, 2024

©2024 Emily Fox

13

Are you sure about the prediction?  
**Class probability**

©2024 Emily Fox

CS 229: Machine Learning

14

## How confident is your prediction?

- Thus far, we've outputted a prediction **+1** or **-1**
- But, how sure are you about the prediction?

*"The sushi & everything else were awesome!"*

Definite **+1**

*"The sushi was good, the service was OK"*

Not sure

©2024 Emily Fox

CS 229: Machine Learning

15

## Using probabilities in classification

©2024 Emily Fox

CS 229: Machine Learning

16

## How confident is your prediction?

"The sushi & everything else were awesome!"

Definite +1

$$P(y=+1|x=\text{"The sushi & everything else were awesome!"}) = 0.99$$

"The sushi was good, the service was OK"

Not sure

$$P(y=+1|x=\text{"The sushi was good, the service was OK"}) = 0.55$$

Many classifiers provide a degree of certainty:

$$\text{Output label} \rightarrow P(y|x) \leftarrow \text{Input sentence}$$

Extremely useful in practice

©2024 Emily Fox

CS 229: Machine Learning

17

## Goal: Learn conditional probabilities from data

Training data:  $N$  observations  $(\mathbf{x}_i, y_i)$

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
2	1	+1
0	2	-1
3	3	-1
4	1	+1
...	...	...

Optimize **quality metric**  
on training data

Find best model  $\hat{P}$   
by finding best  $\hat{w}$

Useful for predicting  $\hat{y}$

©2024 Emily Fox

CS 229: Machine Learning

18

**Predict most likely class**

$\hat{P}(y|x)$  = estimate of class probabilities

Sentence from review → If  $\hat{P}(y=+1|x) > 0.5$ :  
 $\hat{y} = +1$   
 Else:  
 $\hat{y} = -1$

Input:  $x$

Estimating  $\hat{P}(y|x)$  improves **interpretability**:

- Predict  $\hat{y} = +1$  **and** tell me how sure you are

19

©2024 Emily Fox

CS 229: Machine Learning

20

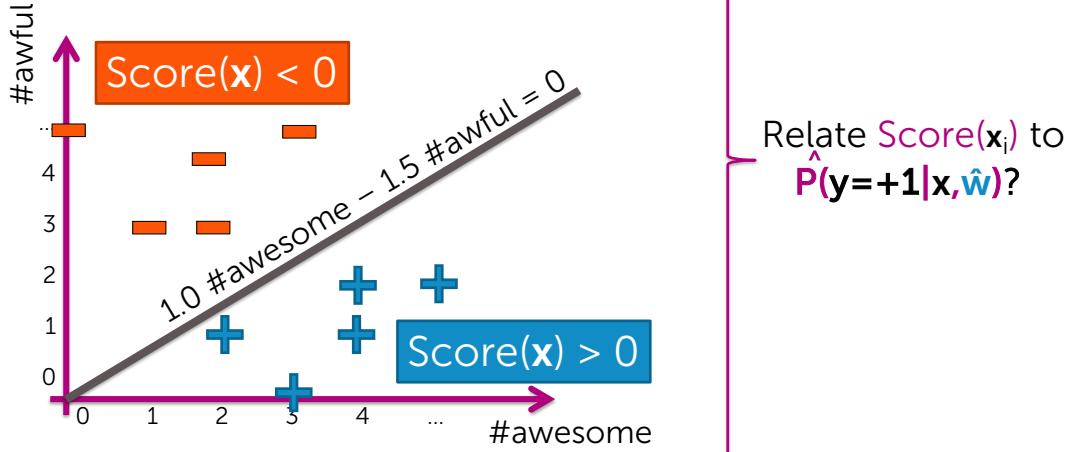
©2024 Emily Fox

CS 229: Machine Learning

Predicting class probabilities with logistic regression

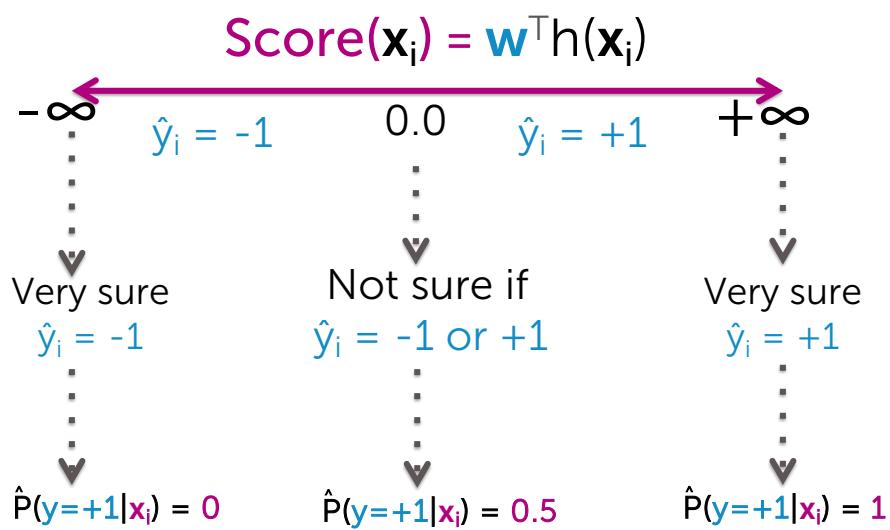
Thus far, we focused on decision boundaries

$$\begin{aligned} \text{Score}(\mathbf{x}_i) &= w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) \\ &= \mathbf{w}^\top \mathbf{h}(\mathbf{x}_i) \end{aligned}$$



21

## Interpreting $\text{Score}(\mathbf{x}_i)$

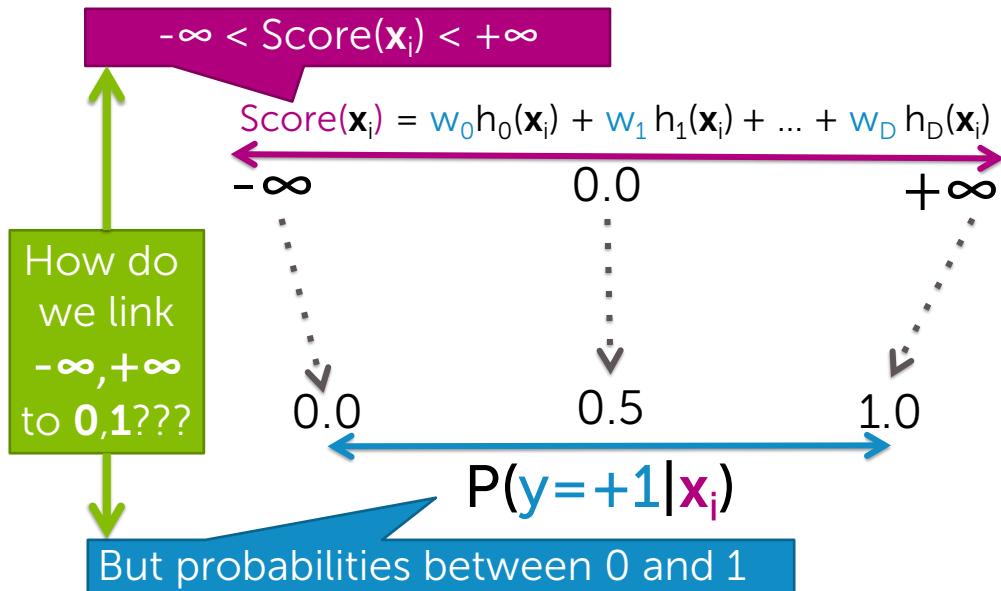


©2024 Emily Fox

CS 229: Machine Learning

22

## Why not just use regression to build classifier?



©2024 Emily Fox

CS 229: Machine Learning

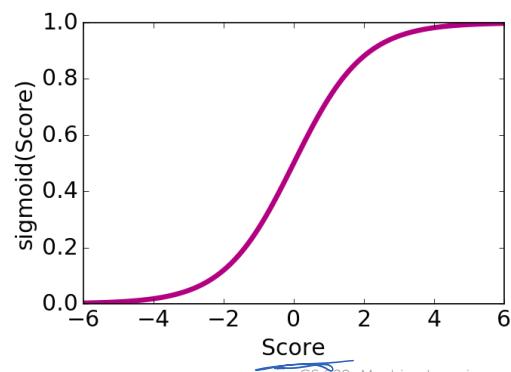
23

## Logistic function (sigmoid, logit)

$$\star \text{sigmoid}(\text{Score}) = \frac{1}{1 + e^{-\text{Score}}}$$

"link fcn"  
used in  
logistic  
regression

Score	$-\infty$	-2	0.0	+2	$+\infty$
$\text{sigmoid}(\text{Score})$	$\frac{1}{1+e^{\infty}}$		$\frac{1}{1+e^0}$		$\frac{1}{1+e^{-\infty}}$



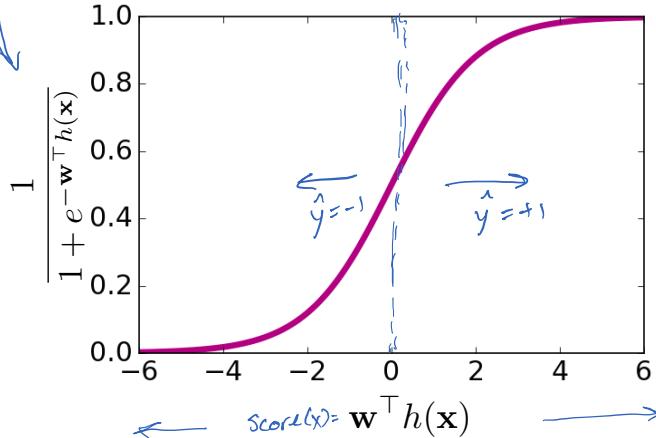
©2024 Emily Fox

CS 229: Machine Learning

24

## Understanding the logistic regression model

$$P(y=+1|x_i, \mathbf{w}) = \text{sigmoid}(\text{Score}(\mathbf{x}_i)) = \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}}$$



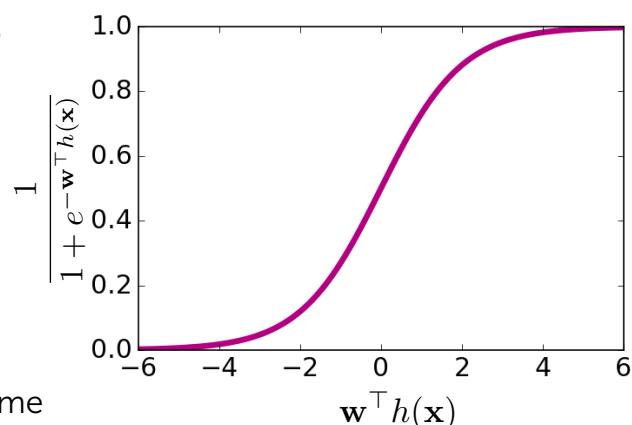
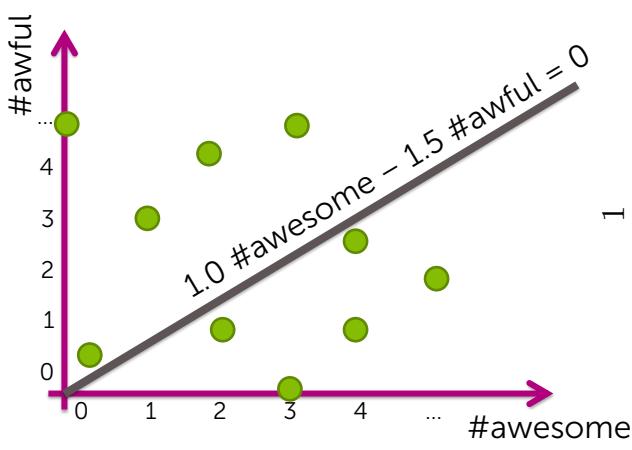
Score( $\mathbf{x}_i$ )	$P(y=+1 \mathbf{x}_i, \mathbf{w})$
0	0.5
-2	0.12 < 0.5 $\Rightarrow \hat{y} = -1$
2	0.88 > 0.5 $\Rightarrow \hat{y} = +1$
4	0.98

©2024 Emily Fox

CS 229: Machine Learning

25

## Logistic regression → Linear decision boundary

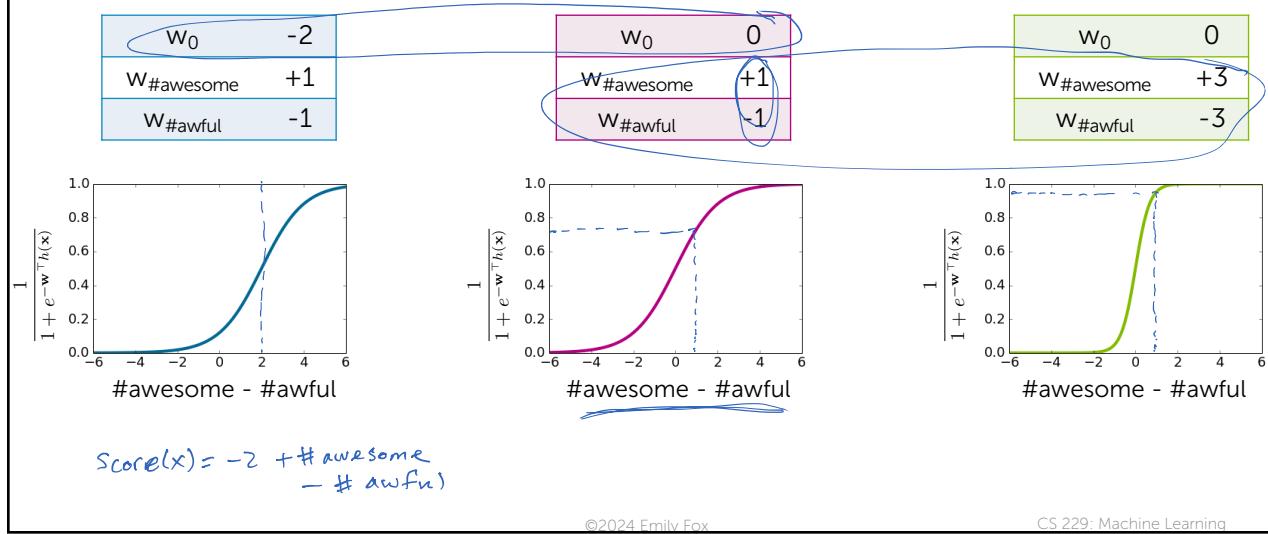


©2024 Emily Fox

CS 229: Machine Learning

26

## Effect of coefficients on logistic regression model



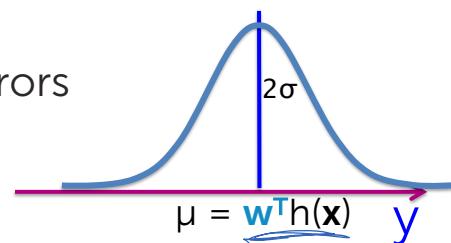
27

## Compare and contrast regression models

- Linear regression with Gaussian errors

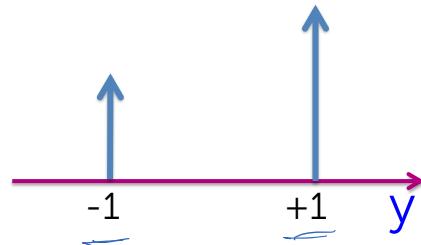
$$y_i = \mathbf{w}^T \mathbf{h}(\mathbf{x}_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$\rightarrow p(y|\mathbf{x}, \mathbf{w}) = N(y; \mathbf{w}^T \mathbf{h}(\mathbf{x}), \sigma^2)$$



- Logistic regression

$$P(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{h}(\mathbf{x})}} & y = +1 \\ \frac{e^{-\mathbf{w}^T \mathbf{h}(\mathbf{x})}}{1 + e^{-\mathbf{w}^T \mathbf{h}(\mathbf{x})}} & y = -1 \end{cases}$$



©2024 Emily Fox

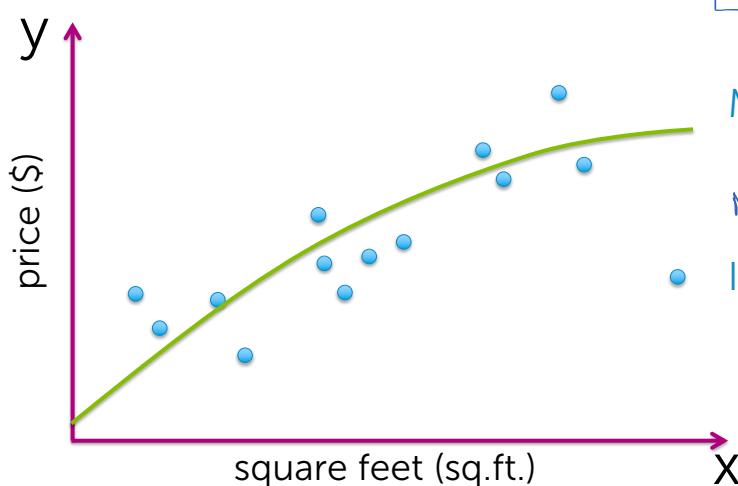
CS 229: Machine Learning

28

Loss function for logistic regression:  
 (Negative log-) Likelihood for maximum likelihood estimation (MLE)

29

## Recall: Gaussian linear regression model



$$N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

parameters:  $\mu$ ,  $\sigma^2$

Model for  $\epsilon_i$ :

$$\text{So far : } E[\epsilon_i] = 0$$

Now, stronger assumption:  
 $\epsilon_i \sim N(0, \sigma^2)$

- Implied distribution on  $y_i$ :

$$y_i = h^\top(x_i) w + \epsilon_i$$

$\underbrace{h^\top(x_i) w}_{\text{"given" const}}$

$$\underline{y_i | x_i; w} \sim N(h^\top(x_i) w, \sigma^2)$$

30

## Recall: Maximum likelihood estimation

$$\arg \max_w p(\mathbf{y} | \mathbf{X}; w) \stackrel{\text{E}_i \text{ iid}}{=} \arg \max_w \prod_{i=1}^N p(y_i | x_i; w) = \arg \max_w \ln \left\{ \prod_{i=1}^N p(y_i | x_i; w) \right\}$$

↑ strictly inc. fun

Maximize log-likelihood wrt  $w$

$$\begin{aligned} \arg \max_w \ln p(\mathbf{y} | \mathbf{X}; w) &= \ln \left\{ \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(y_i - h(x_i)^\top w)^2}{2\sigma^2}} \right\} \\ &= \arg \max_w N \ln \left( \frac{1}{\sigma \sqrt{2\pi}} \right) - \sum_{i=1}^N \frac{(y_i - h^\top(x_i) w)^2}{2\sigma^2} \\ &= \arg \min_w \sum_{i=1}^N (y_i - h^\top(x_i) w)^2 \quad \leftarrow \text{MLE for Gauss model} = \min \text{RSS! (LS line)} \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

31

## Finding best coefficients

x[1] = #awesome	x[2] = #awful	y = sentiment
2	1	+1
0	2	-1
3	3	-1
4	1	+1
1	1	+1
2	4	-1
0	3	-1
0	1	-1
2	1	+1

©2024 Emily Fox

CS 229: Machine Learning

32

## Finding best coefficients

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
0	2	-1
3	3	-1
2	4	-1
0	3	-1
0	1	-1
2	4	-1
0	3	-1
0	1	-1

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
2	1	+1
4	1	+1
1	1	+1
2	1	+1
1	1	+1
2	1	+1

©2024 Emily Fox

CS 229: Machine Learning

33

## Finding best coefficients

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
0	2	-1
3	3	-1
2	4	-1
0	3	-1
0	1	-1

$$P(y=+1|x_i, \hat{w}) = 0.0$$

$x[1] = \#\text{awesome}$	$x[2] = \#\text{awful}$	$y = \text{sentiment}$
2	1	+1
4	1	+1
1	1	+1
2	1	+1

$$P(y=+1|x_i, \hat{w}) = 1.0$$

Want  $\hat{w}$  that makes

©2024 Emily Fox

CS 229: Machine Learning

34

## Learn logistic regression model with maximum likelihood estimation (MLE)

Data point	x[1]	x[2]	y	Choose $w$ to maximize
$\mathbf{x}_1, y_1$	2	1	+1	$P(y=+1 \mathbf{x}[1]=2, \mathbf{x}[2]=1, w)$
$\mathbf{x}_2, y_2$	0	2	-1	$P(y=-1 \mathbf{x}[1]=0, \mathbf{x}[2]=2, w)$
$\mathbf{x}_3, y_3$	3	3	-1	$P(y=-1 \mathbf{x}[1]=3, \mathbf{x}[2]=3, w)$
$\mathbf{x}_4, y_4$	4	1	+1	$P(y=+1 \mathbf{x}[1]=4, \mathbf{x}[2]=1, w)$

$$\ell(w) = \prod_{i=1}^N P(y_i | \mathbf{x}_i, w)$$

*← pick  $w$  to make this large*

©2024 Emily Fox

CS 229: Machine Learning

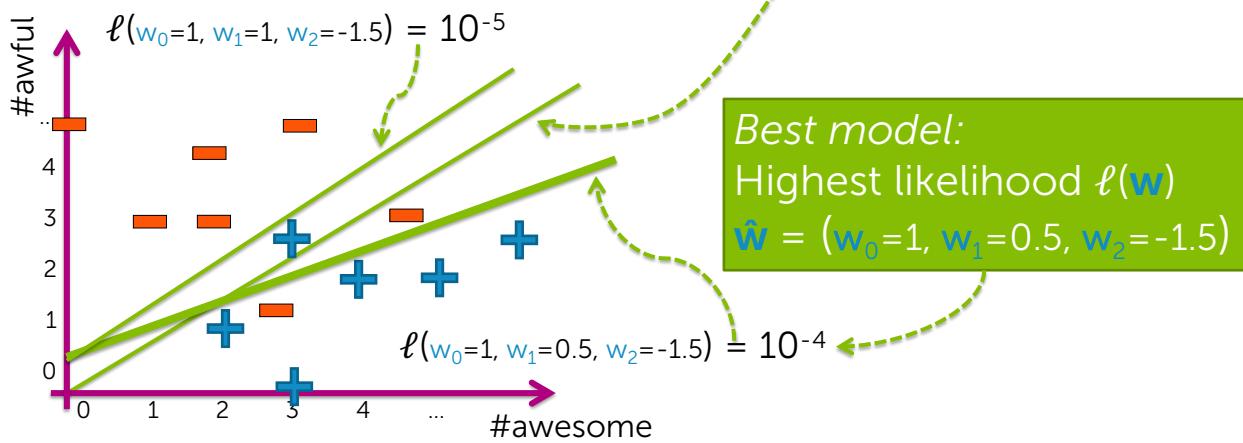
35

## Find “best” classifier

Maximize likelihood over all possible  $w_0, w_1, w_2$

$$\ell(w) = \prod_{i=1}^N P(y_i | \mathbf{x}_i, w)$$

$$\ell(w_0=0, w_1=1, w_2=-1.5) = 10^{-6}$$



36

If max likelihood

as

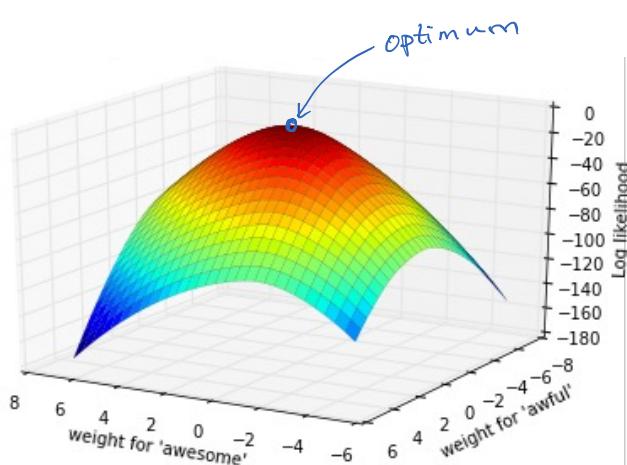
## Gradient descent for logistic regression

©2024 Emily Fox

CS 229: Machine Learning

37

## Maximizing likelihood



Maximize function over all possible  $w_0, w_1, w_2$

$$\max_{w_0, w_1, w_2} \prod_{i=1}^N P(y_i | \mathbf{x}_i, \mathbf{w})$$

$\ell(w_0, w_1, w_2)$  is a function of 3 variables

©2024 Emily Fox

CS 229: Machine Learning

38

19

## Our optimization objective

- Can compute gradient, but no closed-form solution to:

$$\nabla \ell(\mathbf{w}) = 0$$

- Use **gradient descent**

- As with MLE for Gaussians, rewrite objective as:

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} l(\mathbf{w}) = \arg \max_{\mathbf{w}} \ell l(\mathbf{w})$$

↑ log-likelihood

©2024 Emily Fox

CS 229: Machine Learning

39

## Gradient of logistic log-likelihood

$$\frac{\partial \ell(\mathbf{w})}{\partial \mathbf{w}_j} = \sum_{i=1}^N h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}) \right)$$

$\mathbb{1}[y_i = +1]$   
 $= \begin{cases} 1 & \text{if } y_i = +1 \\ 0 & \text{otherwise} \end{cases}$

Sum over data points      Feature value      Difference between truth and prediction

If  $h_j(\mathbf{x}_i) = 1$ :

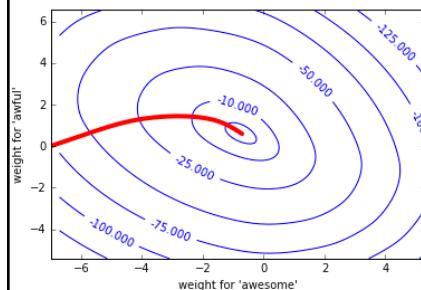
	$P(y = +1   \mathbf{x}_i, \mathbf{w}) \approx 1$	$P(y = +1   \mathbf{x}_i, \mathbf{w}) \approx 0$
$y_i = +1$	$\Delta_i \approx (1 - 1) = 0$ → don't change anything	$\Delta_i \approx 1 \Rightarrow \text{increase } w_j$ $\Rightarrow P(y = +1   \mathbf{x}_i, \mathbf{w}) \text{ increase}$
$y_i = -1$	$\Delta_i \approx -1 \Rightarrow w_j \text{ decreases}$ $\Rightarrow P(y = +1   \mathbf{x}_i, \mathbf{w}) \text{ decreases}$	$\Delta_i \approx 0$ → don't change anything

©2024 Emily Fox

CS 229: Machine Learning

40

# Gradient ascent for logistic regression



init  $\mathbf{w}^{(1)} = \mathbf{0}$  (or randomly, or smartly),  $t=1$

**while**  $\|\nabla \ell(\mathbf{w}^{(t)})\| > \epsilon$

Difference between  
truth and prediction

**for**  $j=0, \dots, D$

$$\text{partial}[j] = \sum_{i=1}^N h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}) \right)$$

$$\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} + \eta \text{partial}[j]$$

$$t \leftarrow t + 1$$

$$\frac{\partial \ell(\mathbf{w}^{(t)})}{\partial \mathbf{w}_j}$$

©2024 Emily Fox

CS 229: Machine Learning

41

## Summary for linear classifiers and logistic regression (example of linear classifier)

©2024 Emily Fox

CS 229: Machine Learning

42

## What you can do now...

- Describe decision boundaries and linear classifiers
- Use class probability to express degree of confidence in prediction
- Define a logistic regression model
- Interpret logistic regression outputs as class probabilities
- Describe impact of coefficient values on logistic regression output
- Measure quality of a classifier using the likelihood function
- Learn a logistic regression model with gradient descent for negative log-likelihood loss

©2024 Emily Fox

CS 229: Machine Learning

43



# Linear classifiers: Overfitting

CS 229: Machine Learning  
Emily Fox  
Stanford University  
January 24, 2024

©2024 Emily Fox

44

# Overfitting in classification

©2024 Emily Fox

CS 229: Machine Learning

45

More complex models tend to have less bias...

Sentiment classifier using  
single words can do OK, but...

Never classifies correctly:  
“The sushi was not good.”

More complex model:  
consider pairs of words (bigrams)

Word	Weight
good	+1.5
not good	-2.1

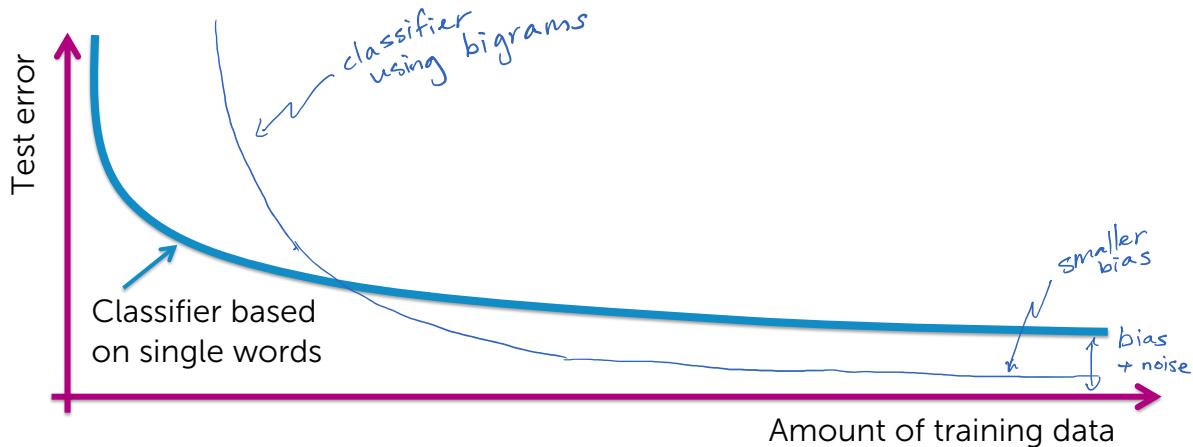
Less bias →  
potentially more accurate,  
needs more data to learn

©2024 Emily Fox

CS 229: Machine Learning

46

Models with less bias tend to need more data to learn well, but do better with sufficient data



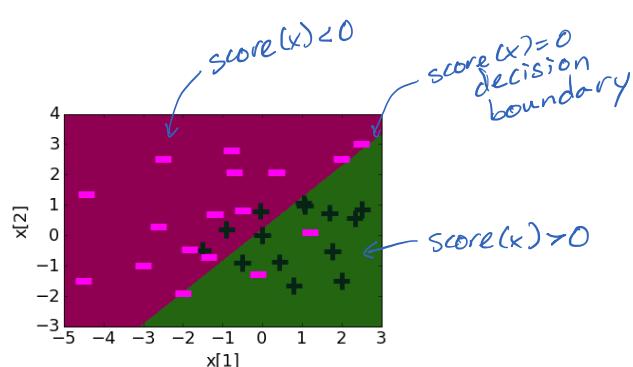
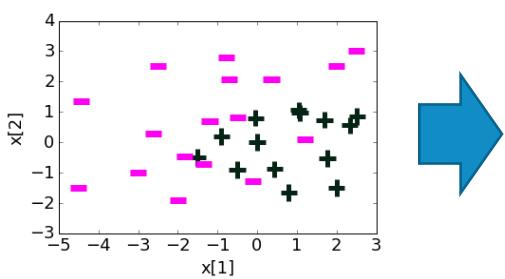
©2024 Emily Fox

CS 229: Machine Learning

47

## Learned decision boundary

Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	0.23
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	1.12
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-1.07



©2024 Emily Fox

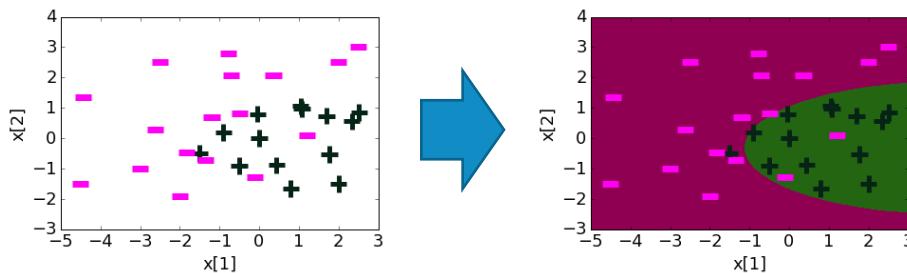
CS 229: Machine Learning

48

## Quadratic features (in 2d)

Note: we are not including cross terms for simplicity

Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	1.68
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	1.39
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-0.59
$h_3(\mathbf{x})$	$(\mathbf{x}[1])^2$	-0.17
$h_4(\mathbf{x})$	$(\mathbf{x}[2])^2$	-0.96



©2024 Emily Fox

CS 229: Machine Learning

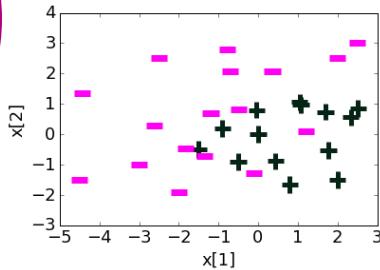
49

## Degree 6 features (in 2d)

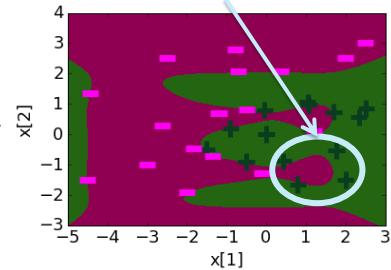
Note: we are not including cross terms for simplicity

Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	21.6
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	5.3
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-42.7
$h_3(\mathbf{x})$	$(\mathbf{x}[1])^2$	-15.9
$h_4(\mathbf{x})$	$(\mathbf{x}[2])^2$	-48.6
$h_5(\mathbf{x})$	$(\mathbf{x}[1])^3$	-11.0
$h_6(\mathbf{x})$	$(\mathbf{x}[2])^3$	67.0
$h_7(\mathbf{x})$	$(\mathbf{x}[1])^4$	1.5
$h_8(\mathbf{x})$	$(\mathbf{x}[2])^4$	48.0
$h_9(\mathbf{x})$	$(\mathbf{x}[1])^5$	4.4
$h_{10}(\mathbf{x})$	$(\mathbf{x}[2])^5$	-14.2
$h_{11}(\mathbf{x})$	$(\mathbf{x}[1])^6$	0.8
$h_{12}(\mathbf{x})$	$(\mathbf{x}[2])^6$	-8.6

Coefficient values getting large



Score( $\mathbf{x}$ ) < 0



©2024 Emily Fox

CS 229: Machine Learning

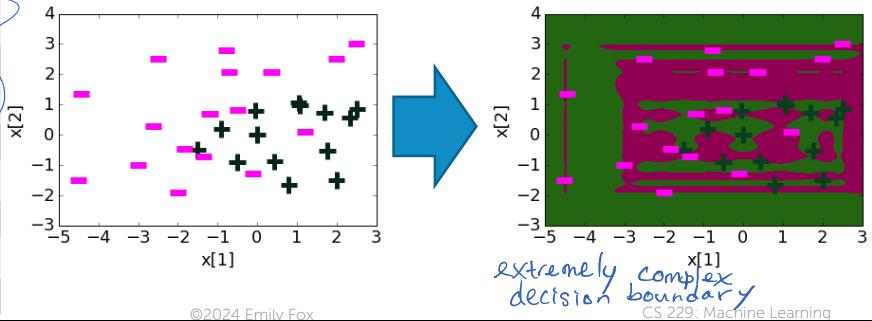
50

## Degree 20 features (in 2d)

Note: we are not including cross terms for simplicity

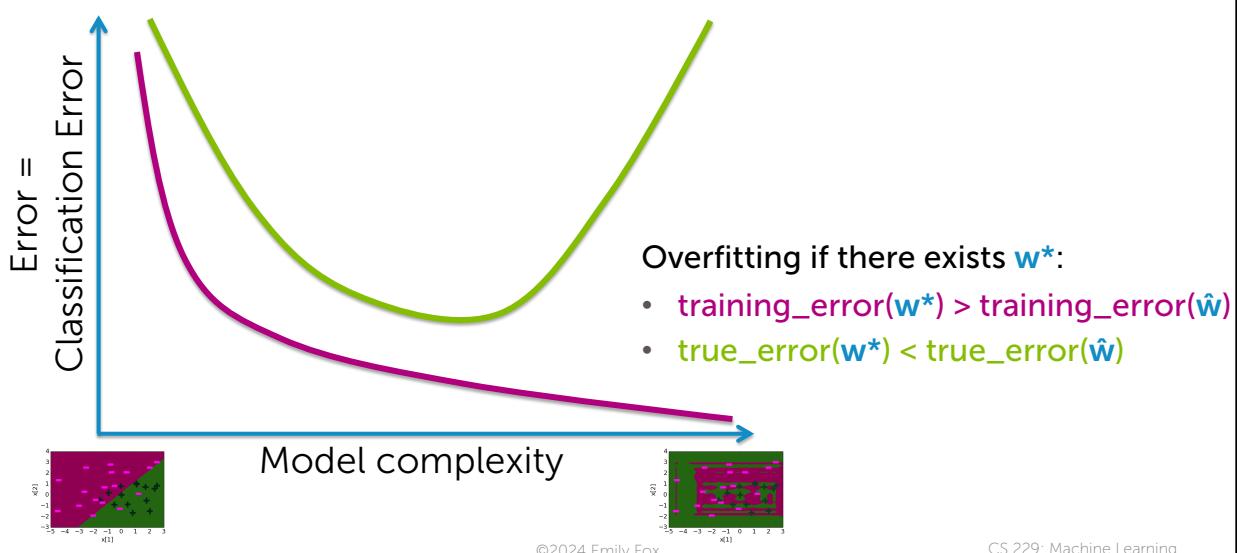
Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	8.7
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	5.1
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	78.7
...	...	...
$h_{11}(\mathbf{x})$	$(\mathbf{x}[1])^6$	-7.5
$h_{12}(\mathbf{x})$	$(\mathbf{x}[2])^6$	3803
$h_{13}(\mathbf{x})$	$(\mathbf{x}[1])^7$	21.1
$h_{14}(\mathbf{x})$	$(\mathbf{x}[2])^7$	-2406
...	...	...
$h_{37}(\mathbf{x})$	$(\mathbf{x}[1])^{19}$	$-2 \cdot 10^{-6}$
$h_{38}(\mathbf{x})$	$(\mathbf{x}[2])^{19}$	-0.15
$h_{39}(\mathbf{x})$	$(\mathbf{x}[1])^{20}$	$-2 \cdot 10^{-8}$
$h_{40}(\mathbf{x})$	$(\mathbf{x}[2])^{20}$	0.03

Often, overfitting associated with very large estimated coefficients  $\hat{\mathbf{w}}$



51

## Overfitting in classification



52

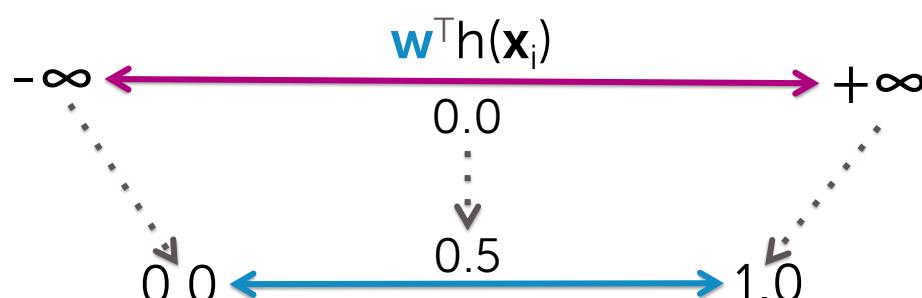
Overfitting in classifiers →  
Overconfident predictions

©2024 Emily Fox

CS 229: Machine Learning

53

## Logistic regression model



$$P(y=+1|x_i, w) = \text{sigmoid}(w^T h(x_i))$$

©2024 Emily Fox

CS 229: Machine Learning

54

## The subtle (negative) consequence of overfitting in logistic regression

Overfitting → Large coefficient values



$\hat{\mathbf{w}}^T \mathbf{h}(\mathbf{x}_i)$  is very positive (or very negative) → sigmoid( $\hat{\mathbf{w}}^T \mathbf{h}(\mathbf{x}_i)$ ) goes to 1 (or to 0)



Model becomes extremely overconfident of predictions

©2024 Emily Fox

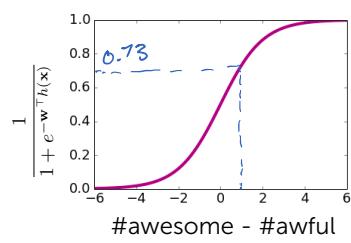
CS 229: Machine Learning

55

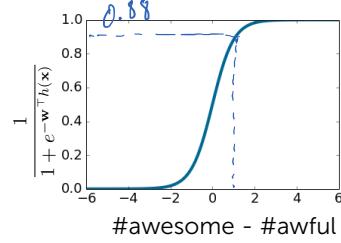
## Effect of coefficients on logistic regression model

Input  $\mathbf{x}$ : #awesome=2, #awful=1 ⭐

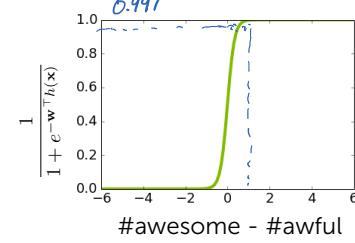
$w_0$	0
$w_{\text{awesome}}$	+1
$w_{\text{awful}}$	-1



$w_0$	0
$w_{\text{awesome}}$	+2
$w_{\text{awful}}$	-2



$w_0$	0
$w_{\text{awesome}}$	+6
$w_{\text{awful}}$	-6



©2024 Emily Fox

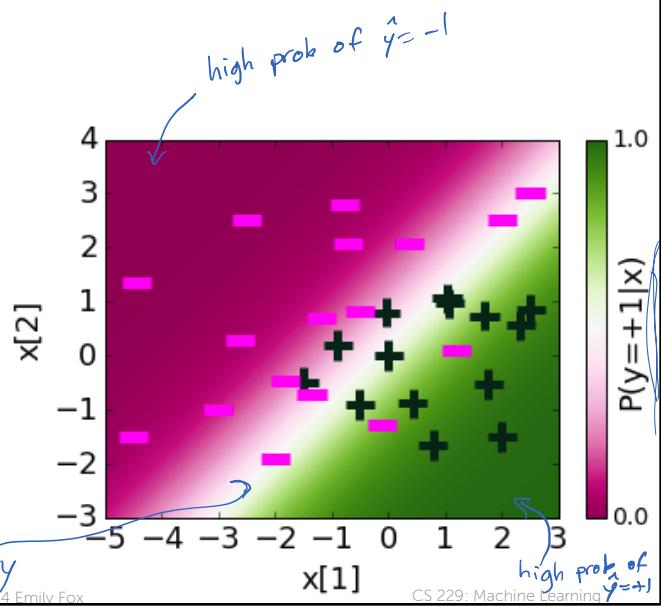
CS 229: Machine Learning

56

## Learned probabilities

Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	0.23
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	1.12
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-1.07

$$P(y = +1 \mid \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}}$$



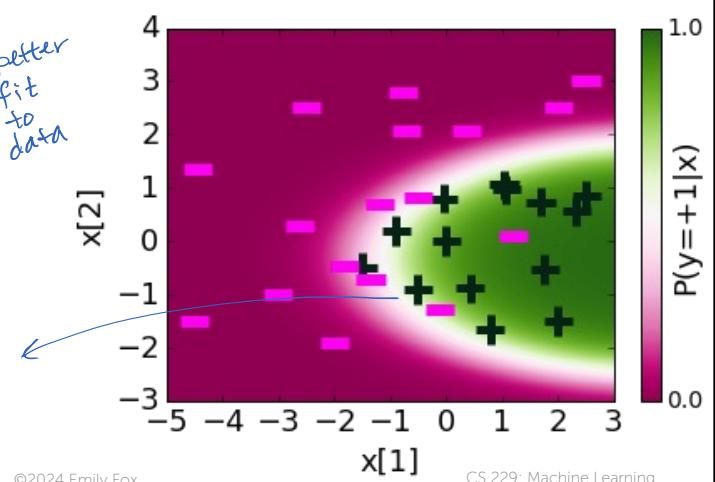
57

## Quadratic features: Learned probabilities

Feature	Value	Coefficient learned
$h_0(\mathbf{x})$	1	1.68
$h_1(\mathbf{x})$	$\mathbf{x}[1]$	1.39
$h_2(\mathbf{x})$	$\mathbf{x}[2]$	-0.58
$h_3(\mathbf{x})$	$(\mathbf{x}[1])^2$	-0.17
$h_4(\mathbf{x})$	$(\mathbf{x}[2])^2$	-0.96

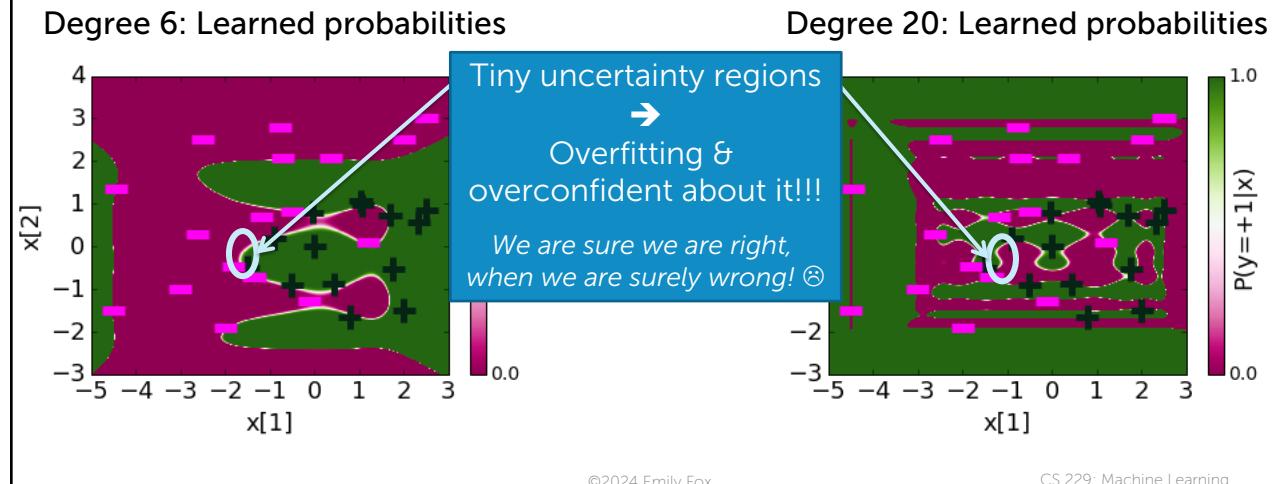
$$P(y = +1 \mid \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}}$$

uncertainty region narrows



58

# Overfitting → Overconfident predictions

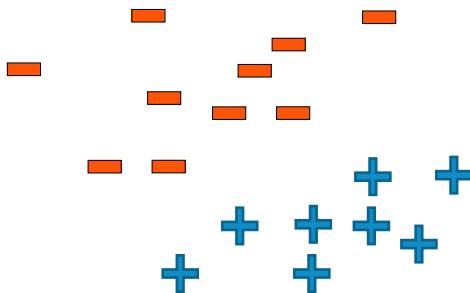


59

## Overfitting in logistic regression: Another perspective

60

## Linearly-separable data



Note 1: If you are using D features, linear separability happens in a D-dimensional space

Note 2: If you have enough features, data are (almost) always linearly separable

Data are linearly separable if:

- There exist coefficients  $\hat{w}$  such that:
  - For all positive training data  
 $score(x) = \hat{w}^T h(x) > 0$
  - For all negative training data  
 $score(x) = \hat{w}^T h(x) < 0$

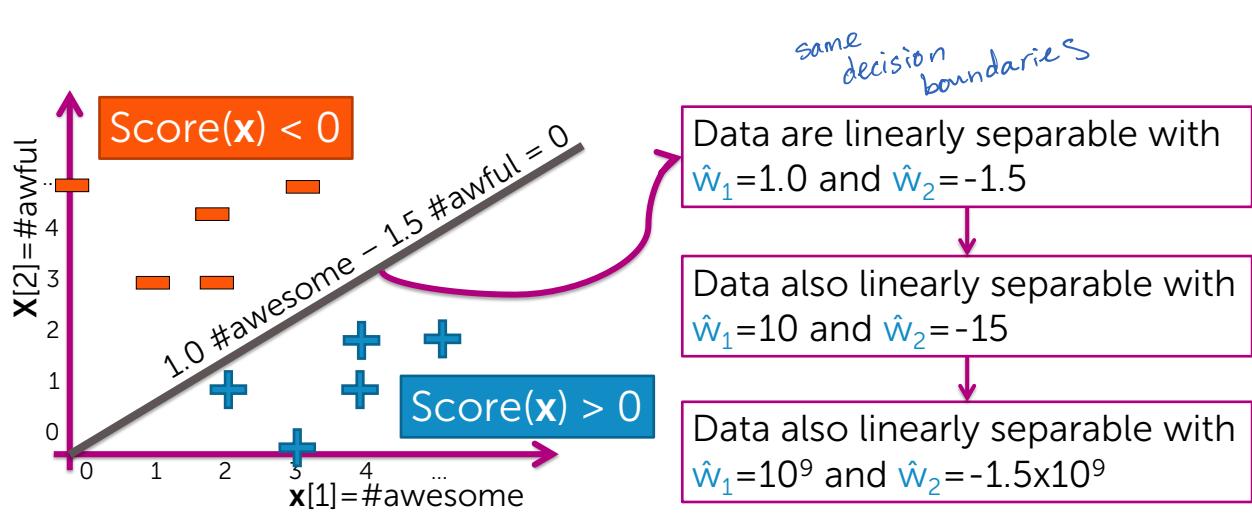
$$\text{training\_error}(\hat{w}) = 0$$

©2024 Emily Fox

CS 229: Machine Learning

61

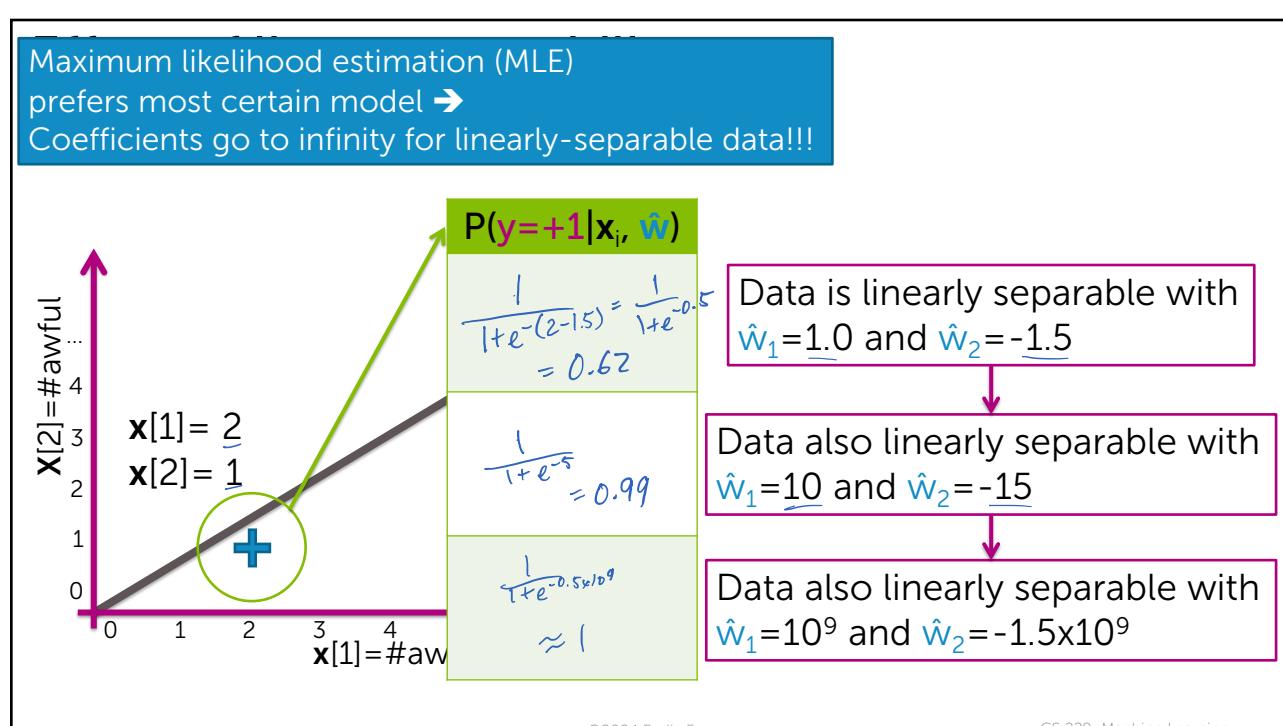
## Effect of linear separability on coefficients



©2024 Emily Fox

CS 229: Machine Learning

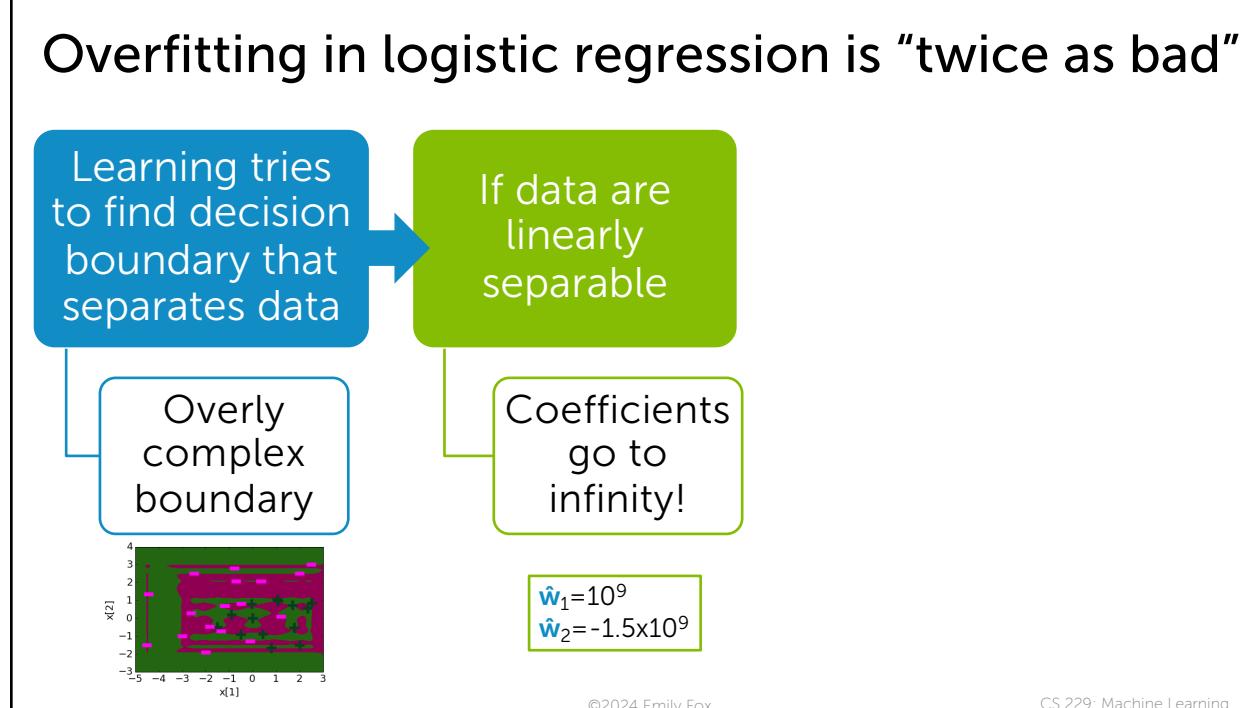
62



©2024 Emily Fox

CS 229: Machine Learning

63



©2024 Emily Fox

CS 229: Machine Learning

64

## Penalizing large coefficients to mitigate overfitting

©2024 Emily Fox

CS 229: Machine Learning

65

## Desired total cost format

Want to balance:

- i. How well function fits data
- ii. Magnitude of coefficients

$$\text{Total quality} = \text{measure of fit} - \text{measure of magnitude of coefficients}$$

want to balance

↑  
(data likelihood)  
large # = good fit to training data

↑  
large # = overfit

©2024 Emily Fox

CS 229: Machine Learning

66

## L<sub>2</sub> regularized logistic regression

©2024 Emily Fox

CS 229: Machine Learning

67

## Consider resulting objective

Select  $\hat{\mathbf{w}}$  to maximize:

$$\ell(\mathbf{w}) - \lambda \|\mathbf{w}\|_2^2$$

↑ tuning parameter = balance of fit and magnitude

L<sub>2</sub> regularized  
logistic regression

Pick  $\lambda$  using:

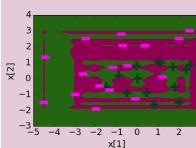
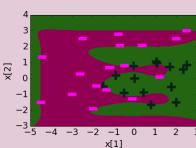
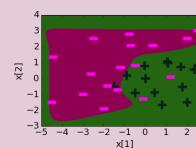
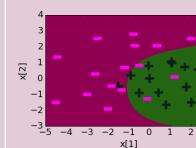
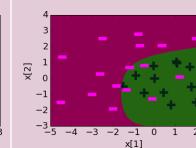
- Validation set (for large datasets)
- Cross-validation (for smaller datasets)  
(as in ridge/lasso regression)

©2024 Emily Fox

CS 229: Machine Learning

68

## Degree 20 features, effect of regularization penalty $\lambda$

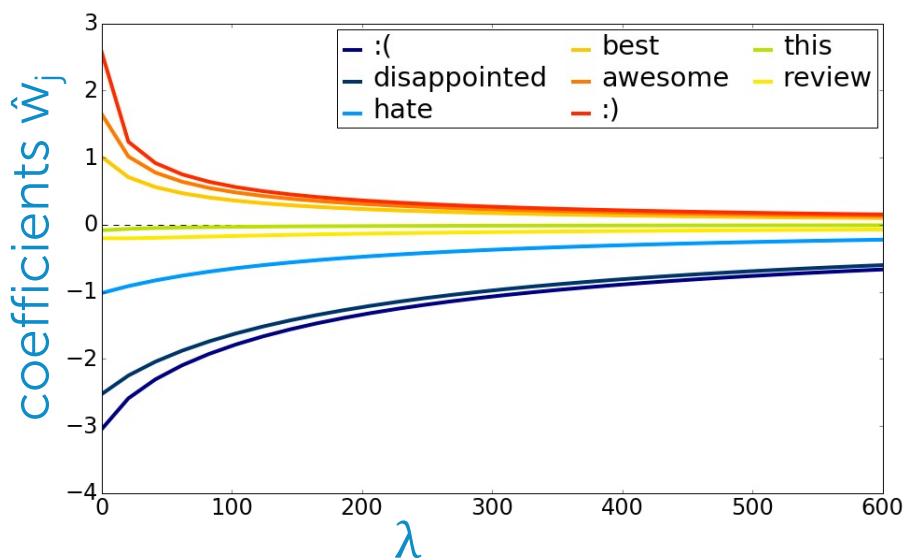
Regularizatio n	$\lambda = 0$	$\lambda = 0.00001$	$\lambda = 0.001$	$\lambda = 1$	$\lambda = 10$
Range of coefficients	-3170 to 3803	-8.04 to 12.14	-0.70 to 1.25	-0.13 to 0.57	-0.05 to 0.22
Decision boundary					

©2024 Emily Fox

CS 229: Machine Learning

69

## Coefficient path

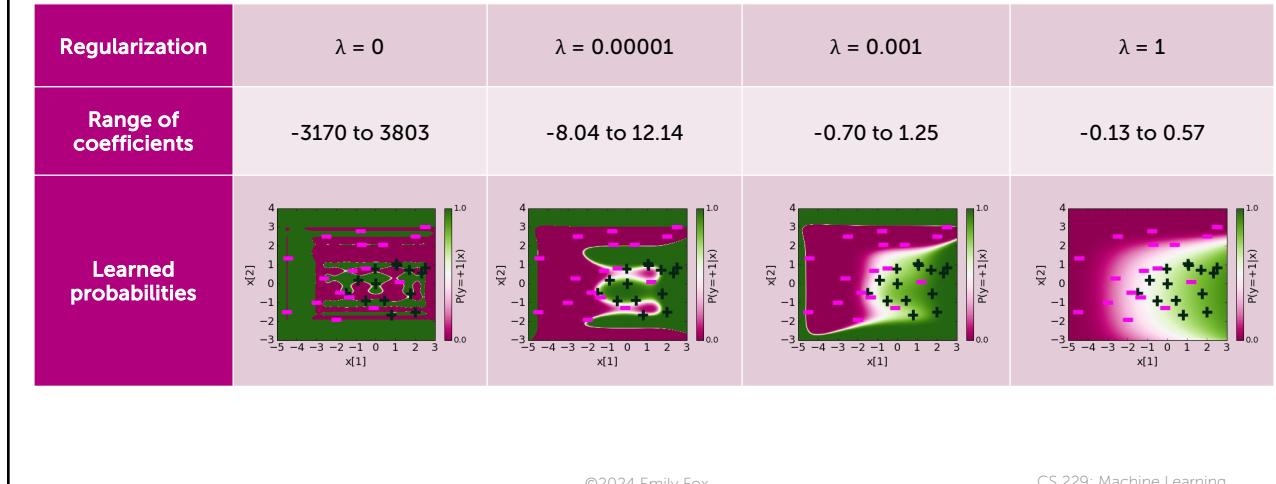


©2024 Emily Fox

CS 229: Machine Learning

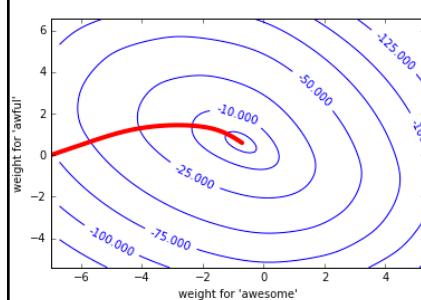
70

## Degree 20 features: regularization reduces “overconfidence”



71

## Gradient ascent for $L_2$ regularized logistic regression



```

init  $w^{(1)} = 0$  (or randomly, or smartly),  $t = 1$ 
while not converged:
  for  $j = 0, \dots, D$ 
     $\text{partial}[j] = \sum_{i=1}^N h_j(x_i) (\mathbb{1}[y_i = +1] - P(y = +1 | x_i, w^{(t)}))$ 
     $w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta (\text{partial}[j] - 2\lambda \underline{w_j^{(t)}})$ 
     $t \leftarrow t + 1$ 
      only change

```

©2024 Emily Fox

CS 229: Machine Learning

72

## $L_1$ regularized logistic regression

©2024 Emily Fox

CS 229: Machine Learning

73

## Sparse logistic regression with $L_1$ penalty

Select  $\hat{\mathbf{w}}$  to maximize:

$$\ell(\mathbf{w}) - \lambda \|\mathbf{w}\|_1$$

↑ tuning parameter = balance of fit and magnitude

$L_1$  regularized  
logistic regression

Pick  $\lambda$  using:

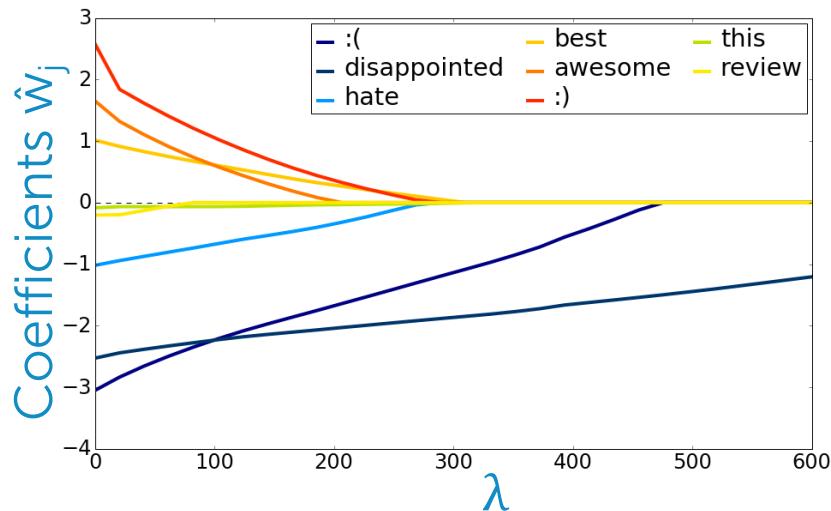
- Validation set (for large datasets)
- Cross-validation (for smaller datasets)  
(as in ridge/lasso regression)

©2024 Emily Fox

CS 229: Machine Learning

74

## Coefficient path – $L_1$ penalty



©2024 Emily Fox

CS 229: Machine Learning

75

Summary of overfitting in logistic regression

©2024 Emily Fox

CS 229: Machine Learning

76

## What you can do now...

- Describe symptoms and effects of overfitting in classification
  - Identify when overfitting is happening
  - Relate large learned coefficients to overfitting
  - Describe the impact of overfitting on decision boundaries and predicted probabilities of linear classifiers
- Use regularization to mitigate overfitting
  - Motivate the form of L2 regularized logistic regression quality metric
  - Describe the use of L1 regularization to obtain sparse logistic regression solutions
  - Describe what happens to estimated coefficients as tuning parameter  $\lambda$  is varied
  - Estimate L2 regularized logistic regression coefficients using gradient ascent
  - Interpret coefficient path plot

©2024 Emily Fox

CS 229: Machine Learning

77



# Linear classifiers: More fundamentals

CS 229: Machine Learning  
Emily Fox  
Stanford University  
January 24, 2024

©2024 Emily Fox

78

39

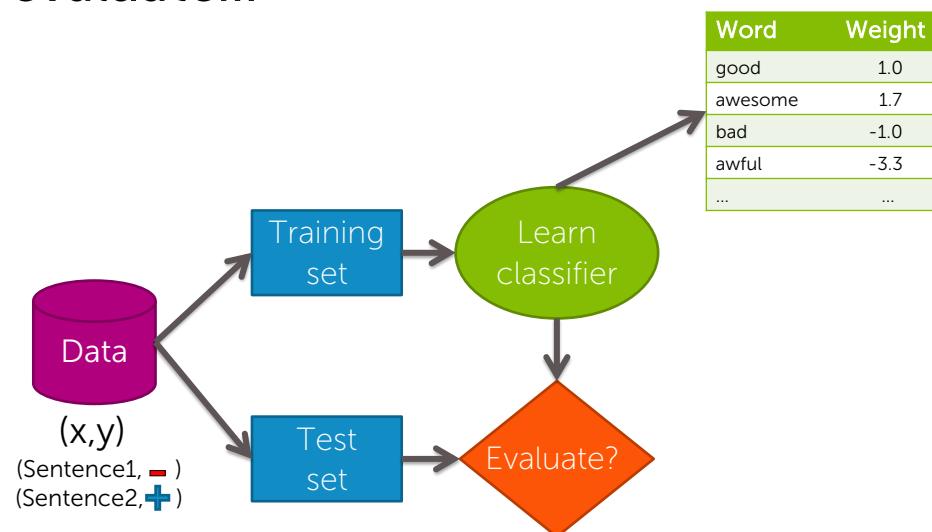
## Evaluating a classifier

©2024 Emily Fox

CS 229: Machine Learning

79

Training a classifier = Learning the weights  
Then evaluate...

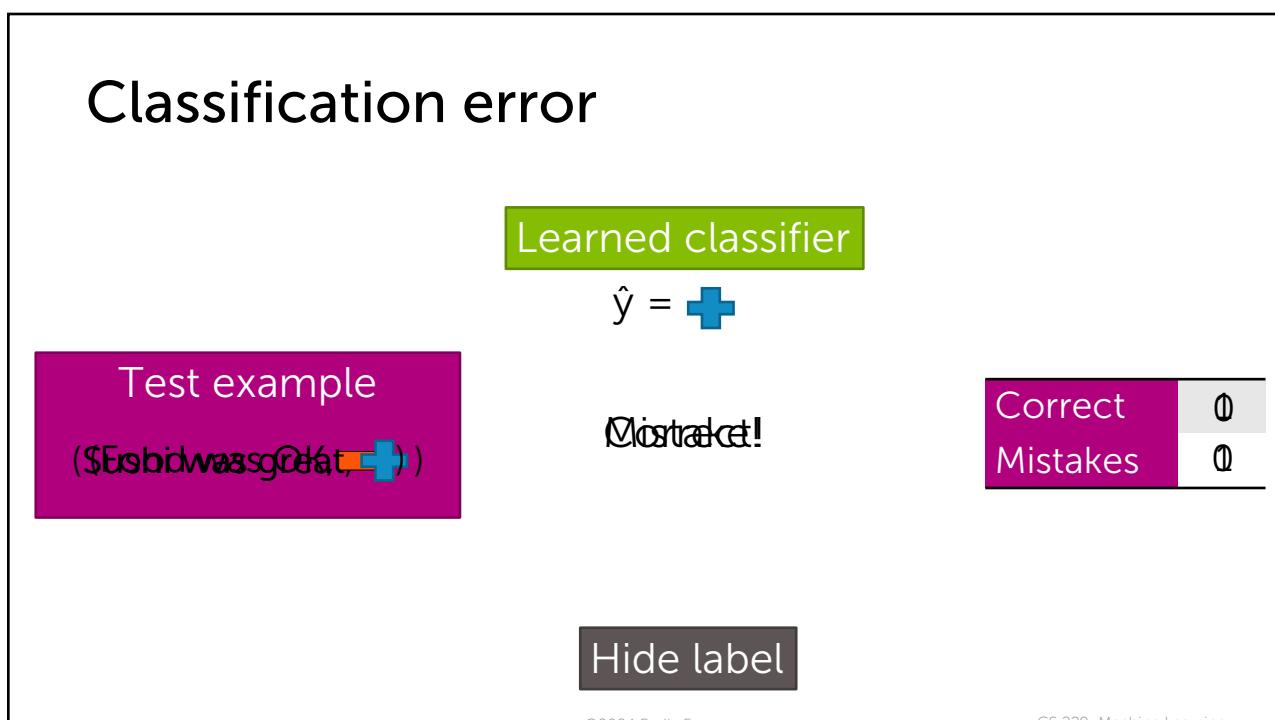


©2024 Emily Fox

CS 229: Machine Learning

80

## Classification error



©2024 Emily Fox

CS 229: Machine Learning

81

## Classification error & accuracy

- Error measures fraction of mistakes

$$\text{error} = \frac{\# \text{ of mistakes}}{\text{total } \# \text{ of sentences}}$$

← test set

- Best possible value is 0.0

$\text{error} = 1 - \text{accuracy}$

- Often, measure **accuracy**

- Fraction of correct predictions

$$\text{accuracy} = \frac{\# \text{ correct}}{\text{total } \# \text{ of sentences}}$$

- Best possible value is 1.0

©2024 Emily Fox

CS 229: Machine Learning

82

## What's a good accuracy?

©2024 Emily Fox

CS 229: Machine Learning

83

## What if you ignore the sentence, and just guess?

- For binary classification:
  - Half the time, you'll get it right! (on average)  
→ accuracy =  $0.5$
- For  $k$  classes, accuracy =  $1/k$ 
  - $\frac{1}{3}$  for 3 classes,  $\frac{1}{4}$  for 4 classes,...

At the very, very, very least,  
you should healthily beat random...  
Otherwise, it's (usually) pointless...

©2024 Emily Fox

CS 229: Machine Learning

84

## Is a classifier with 90% accuracy good? Depends...

2010 data shows:  
“90% emails sent are spam!”

Predicting every email is spam  
gets you 90% accuracy!!!

Majority class prediction

Amazing performance when  
there is class imbalance  
(but silly approach)

- One class is more common than others
- Beats random (if you know the majority class)

©2024 Emily Fox

CS 229: Machine Learning

85

## So, always be digging in and asking the hard questions about reported accuracies

- Is there **class imbalance**?
- How does it compare to a **simple, baseline approach**?
  - Random guessing
  - Majority class
  - ...
- Most importantly:  
***What accuracy does my application need?***
  - What is good enough for my user’s experience?
  - What is the impact of the mistakes we make?

©2024 Emily Fox

CS 229: Machine Learning

86

## False positives, false negatives, and confusion matrices

©2024 Emily Fox

CS 229: Machine Learning

87

## Types of mistakes

		Predicted label	
		+	-
True label	+	True Positive	False Negative (FN)
	-	False Positive (FP)	True Negative

Annotations:

- False Negative (FN): "swab negative for covid-19, but really have it" (blue arrow pointing to the bottom-left cell)
- False Positive (FP): "swab comes back pos., but don't really have covid" (blue arrow pointing to the bottom-left cell)

©2024 Emily Fox

CS 229: Machine Learning

88

## Cost of different types of mistakes can be different (& high) in some applications

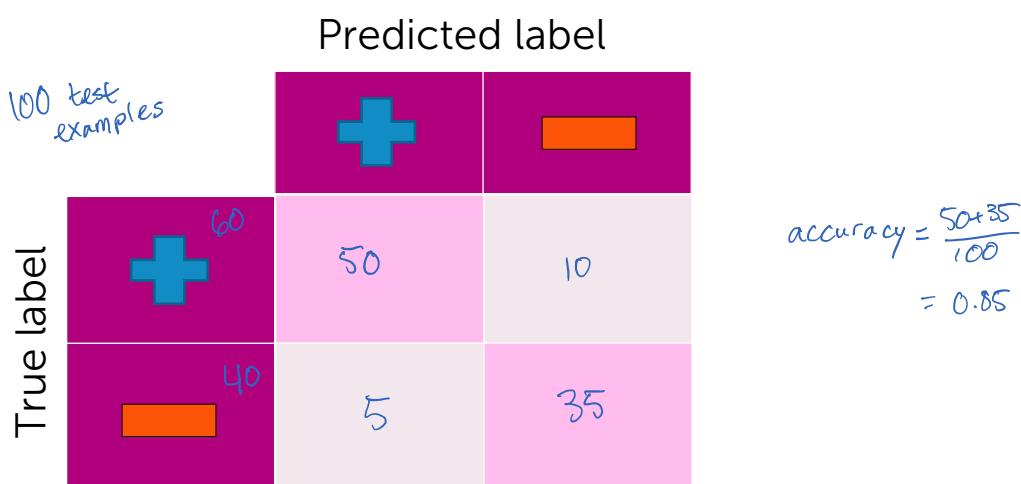
	Spam filtering	Medical diagnosis
False negative	Annoying	Disease not treated ↑ depends Wasteful treatment
False positive	Email lost <i>higher cost</i>	Anxiety...

©2024 Emily Fox

CS 229: Machine Learning

89

## Confusion matrix – binary classification



©2024 Emily Fox

CS 229: Machine Learning

90

## Encoding categorical inputs

**OPTIONAL**

©2024 Emily Fox

CS 229: Machine Learning

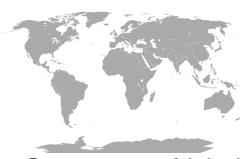
91

## Categorical inputs

- Numeric inputs:
  - #awesome, age, salary,...
  - Intuitive when multiplied by coefficient
    - e.g., 1.5 #awesome
- Categorical inputs:



Gender  
(Male, Female,...)



Country of birth  
(Argentina, Brazil, USA,...)



Zipcode  
(10005, 94305,...)

Numeric value, but should be interpreted as category  
(94305 not about 9x larger than 10005)

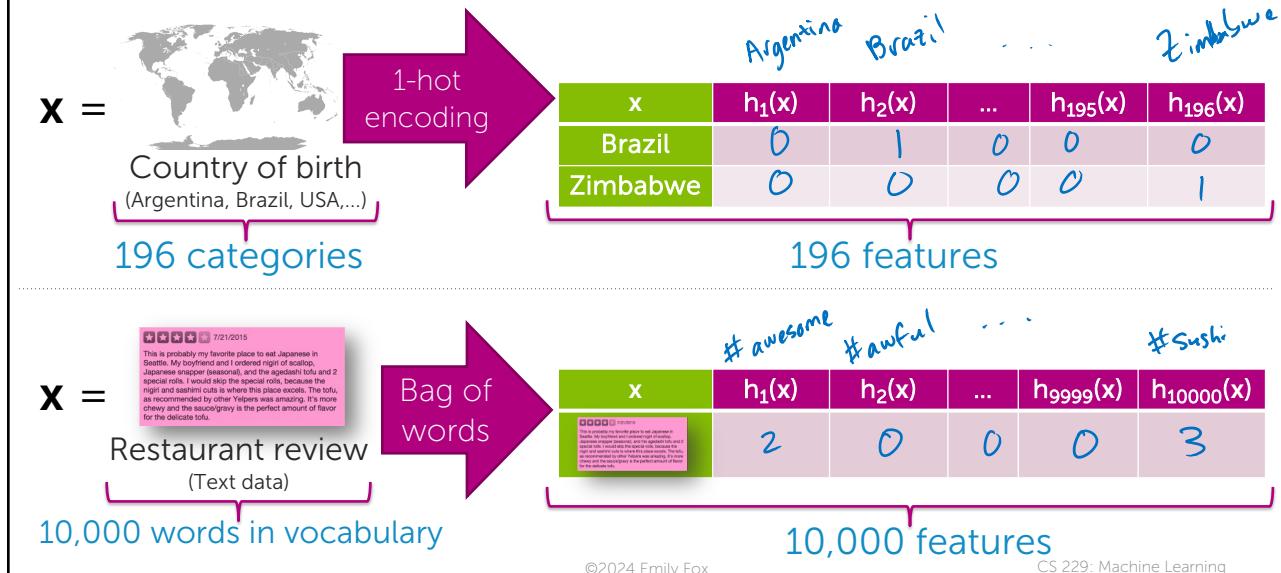
How do we multiply category by coefficient???  
Must convert categorical inputs into numeric features

©2024 Emily Fox

CS 229: Machine Learning

92

## Encoding categories as numeric features



©2024 Emily Fox

CS 229: Machine Learning

93

Multiclass classification  
using 1 versus all

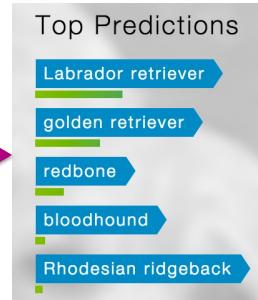
**OPTIONAL**

©2024 Emily Fox

CS 229: Machine Learning

94

## Multiclass classification



Input:  $x$   
Image pixels

Output:  $y$   
Object in image

©2024 Emily Fox

CS 229: Machine Learning

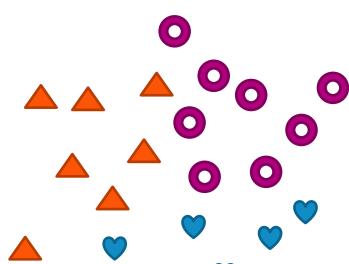
95

## Multiclass classification formulation

- C possible classes:  
–  $y$  can be 1, 2, ..., C      *before  $y \in \{-1, 1\}$*
- N datapoints:

Data point	$x[1]$	$x[2]$	$y$
$x_1, y_1$	2	1	▲
$x_2, y_2$	0	2	♥
$x_3, y_3$	3	3	○
$x_4, y_4$	4	1	○

Learn:  
 $\hat{P}(y=\blacktriangle | x)$   
 $\hat{P}(y=\heartsuit | x)$   
 $\hat{P}(y=\circlearrowleft | x)$



©2024 Emily Fox

CS 229: Machine Learning

96

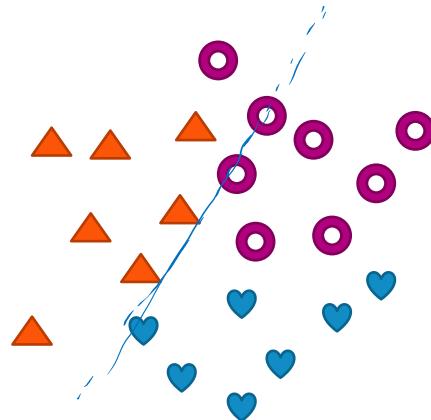
## 1 versus all:

Estimate  $\hat{P}(y=\triangle | x)$  using 2-class model

+1 class: points with  $y_i = \triangle$   
 -1 class: points with  $y_i = \heartsuit$  OR  $\circlearrowleft$

Train classifier:  $\hat{P}_{\triangle}(y=+1|x)$

Predict:  $\hat{P}(y=\triangle | x_i) = \hat{P}_{\triangle}(y=+1|x_i)$



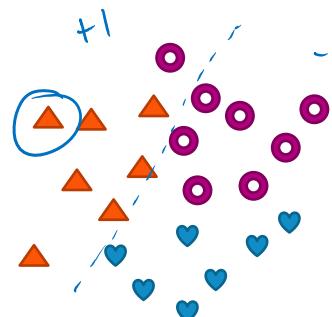
©2024 Emily Fox

CS 229: Machine Learning

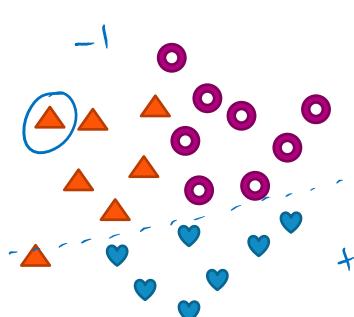
97

## 1 versus all: simple multiclass classification using C 2-class models

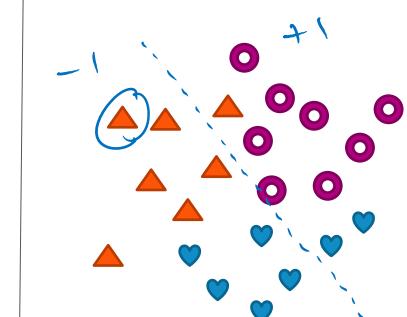
$$\hat{P}(y=\triangle | x_i) = \hat{P}_{\triangle}(y=+1 | x_i, w_{\triangle})$$



$$\hat{P}(y=\heartsuit | x_i) = \hat{P}_{\heartsuit}(y=+1 | x_i, w_{\heartsuit})$$



$$\hat{P}(y=\circlearrowleft | x_i) = \hat{P}_{\circlearrowleft}(y=+1 | x_i, w_{\circlearrowleft})$$



©2024 Emily Fox

CS 229: Machine Learning

98

**Multiclass training**

$\hat{P}_c(y=+1|x)$  = estimate of 1 vs all model for each class

**Predict most likely class**

```
max_prob = 0;  $\hat{y} = 0$ 
For c = 1,...,C:
  If  $\hat{P}_c(y=+1|x_i) > \text{max\_prob}$ :
     $\hat{y} = c$ 
    max_prob =  $\hat{P}_c(y=+1|x_i)$ 
```

©2024 Emily Fox

CS 229: Machine Learning

99

# Logistic Regression: Gradient descent details

👍
CS 229: Machine Learning
👎

Emily Fox
Stanford University

January 24, 2024
OPTIONAL

100

## Step 1: Rewrite log-likelihood

For simpler math, we'll rewrite likelihood with indicators:

$$\begin{aligned}\ell(\mathbf{w}) &= \sum_{i=1}^N \ln P(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \sum_{i=1}^N [\mathbb{1}[y_i = +1] \ln P(y = +1 | \mathbf{x}_i, \mathbf{w}) + \mathbb{1}[y_i = -1] \ln P(y = -1 | \mathbf{x}_i, \mathbf{w})]\end{aligned}$$

Indicator function:

$$\mathbb{1}[y_i = +1] = \begin{cases} 1 & \text{if } y_i = +1 \\ 0 & \text{if } y_i = -1 \end{cases}$$

©2024 Emily Fox

CS 229: Machine Learning

101

## Step 2: Express probabilities in terms of $\mathbf{w}$ and $h(\mathbf{x})$

Probability model predicts  $y=+1$ :

$$P(y=+1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\underline{\mathbf{w}^\top h(\mathbf{x})}}}$$

Probability model predicts  $y=-1$ :

$$\begin{aligned}P(y=-1|\mathbf{x}, \mathbf{w}) &= 1 - P(y=+1|\mathbf{x}, \mathbf{w}) \\ &= \frac{e^{-\underline{\mathbf{w}^\top h(\mathbf{x})}}}{1 + e^{-\underline{\mathbf{w}^\top h(\mathbf{x})}}}\end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

102

## Step 3: Plugging in for 1 data point

$$P(y = +1 | \mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}} \quad P(y = -1 | \mathbf{x}, \mathbf{w}) = \frac{e^{-\mathbf{w}^\top h(\mathbf{x})}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x})}}$$

$$\begin{aligned} \ell(\mathbf{w}) &= \mathbb{1}[y_i = +1] \ln P(y = +1 | \mathbf{x}_i, \mathbf{w}) + \mathbb{1}[y_i = -1] \ln P(y = -1 | \mathbf{x}_i, \mathbf{w}) \\ \ell(\mathbf{w}) &= -\mathbb{1}[y_i = +1] \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) + (1 - \mathbb{1}[y_i = +1]) (-\mathbf{w}^\top h(\mathbf{x}_i) - \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)})) \\ &= (1 - \mathbb{1}[y_i = +1]) \mathbf{w}^\top h(\mathbf{x}_i) - \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) \end{aligned}$$

↑ In num  
- ln den

©2024 Emily Fox

CS 229: Machine Learning

103

## Step 4: Gradient for 1 data point

$$\ell(\mathbf{w}) = -(1 - \mathbb{1}[y_i = +1]) \mathbf{w}^\top h(\mathbf{x}_i) - \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)})$$

$$\begin{aligned} \frac{\partial \ell(\mathbf{w})}{\partial w_j} &= -(1 - \mathbb{1}[y_i = +1]) h_j(\mathbf{x}_i) \\ &\quad + h_j(\mathbf{x}_i) (1 - P(y = +1 | \mathbf{x}_i, \mathbf{w})) \\ &= h_j(\mathbf{x}_i) [\mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w})] \end{aligned}$$

$$\textcircled{1} \quad \frac{\partial}{\partial w_j} w^\top h(\mathbf{x}_i) = h_j(\mathbf{x}_i)$$

$$\begin{aligned} \textcircled{2} \quad \frac{\partial}{\partial w_j} \ln(1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}) &= \frac{1}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} e^{-\mathbf{w}^\top h(\mathbf{x}_i)} (-h_j(\mathbf{x}_i)) \\ &= -h_j(\mathbf{x}_i) \frac{e^{-\mathbf{w}^\top h(\mathbf{x}_i)}}{1 + e^{-\mathbf{w}^\top h(\mathbf{x}_i)}} \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

104

## Step 5: Gradient over all data points

$$\frac{\partial \ell(\mathbf{w})}{\partial w_j} = \sum_{i=1}^N h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}) \right)$$

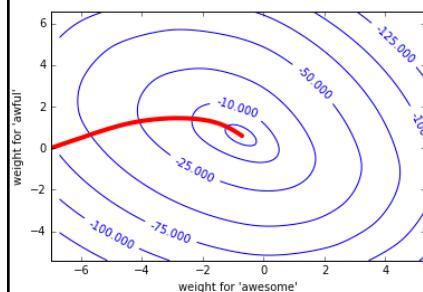
Sum over data points      Feature value      Difference between truth and prediction

©2024 Emily Fox

CS 229: Machine Learning

105

## Gradient ascent for logistic regression



```

init  $\mathbf{w}^{(1)} = 0$  (or randomly, or smartly),  $t=1$ 
while  $\|\nabla \ell(\mathbf{w}^{(t)})\| > \epsilon$ 
  for  $j=0, \dots, D$ 
    partial[j] =  $\sum_{i=1}^N h_j(\mathbf{x}_i) \left( \mathbb{1}[y_i = +1] - P(y = +1 | \mathbf{x}_i, \mathbf{w}^{(t)}) \right)$ 
     $\mathbf{w}_j^{(t+1)} \leftarrow \mathbf{w}_j^{(t)} + \eta \text{ partial}[j]$ 
   $t \leftarrow t + 1$ 
  Difference between truth and prediction

```

©2024 Emily Fox

CS 229: Machine Learning

106

## Choosing the step size $\eta$

**OPTIONAL**

©2024 Emily Fox

CS 229: Machine Learning

107

### Simple rule of thumb for picking step size $\eta$

- Unfortunately, picking step size requires a lot of trial & error 😞
- Try a several values, exponentially spaced
  - **Goal:** plot learning curves to
    - find one  $\eta$  that is too small (smooth but moving too slowly)
    - find one  $\eta$  that is too large (oscillation or divergence)
- Try values in between to find “best”  $\eta$
- *Advanced tip:* can also try step size that decreases with iterations, e.g.,

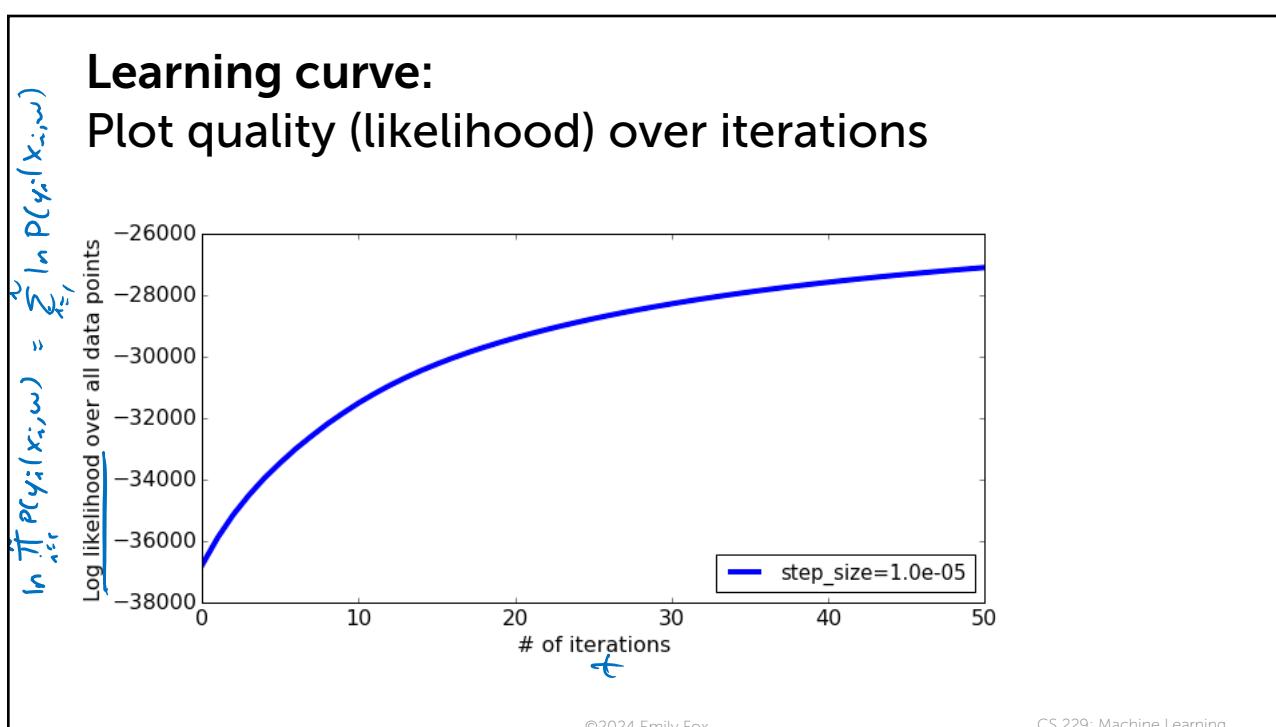
$$\eta_t = \frac{\eta_0}{t}$$



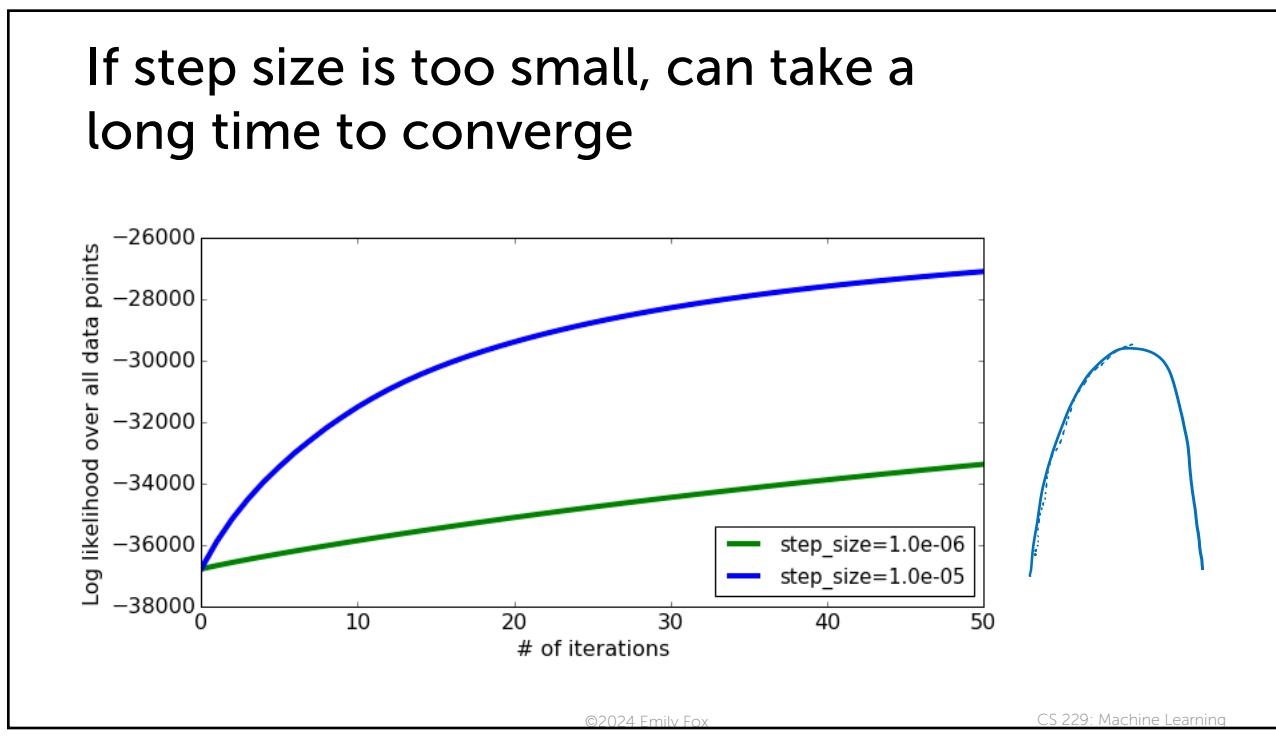
©2024 Emily Fox

CS 229: Machine Learning

108

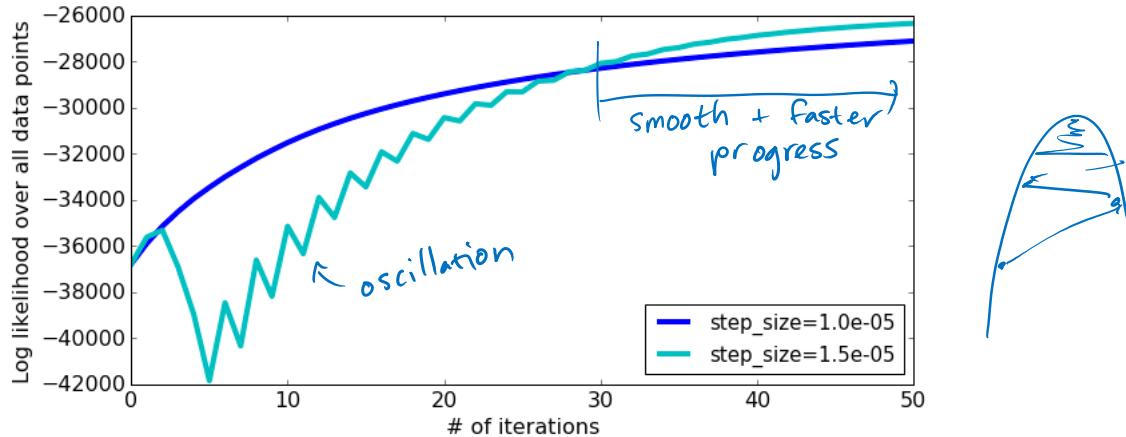


109



110

## Compare converge with different step sizes

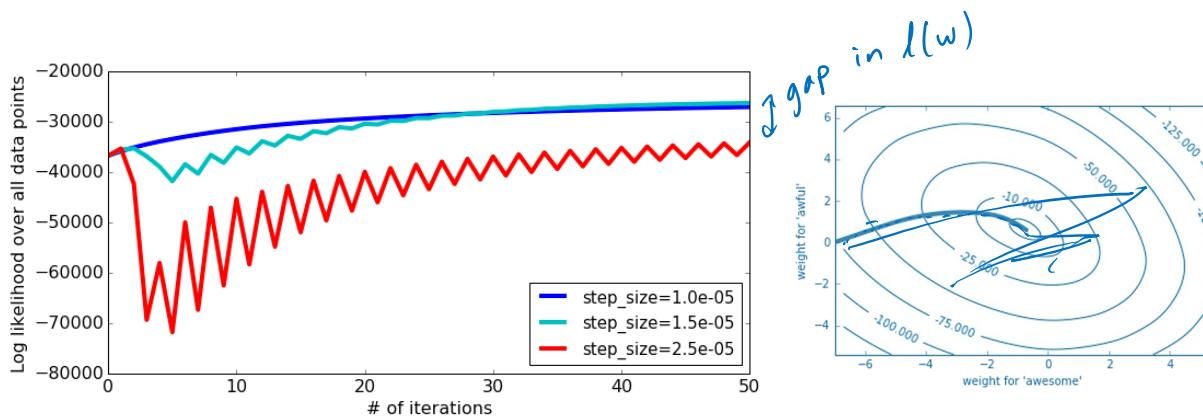


©2024 Emily Fox

CS 229: Machine Learning

111

## Careful with step sizes that are too large

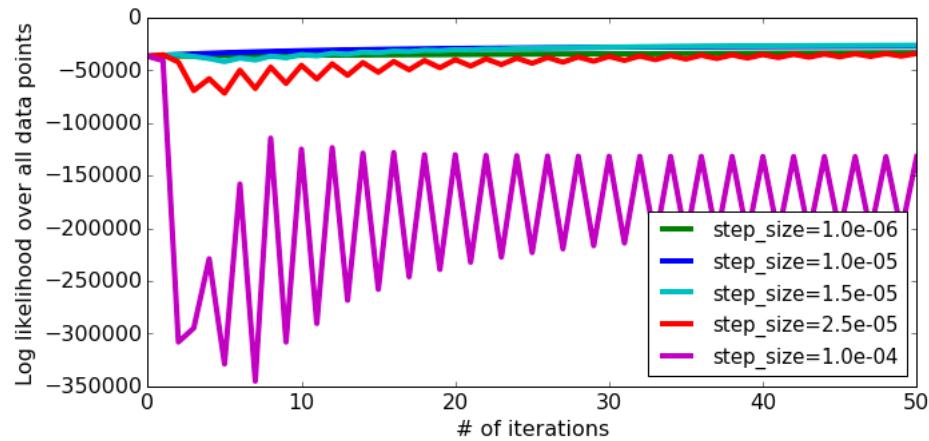


©2024 Emily Fox

CS 229: Machine Learning

112

Very large step sizes can even cause divergence or wild oscillations



©2024 Emily Fox

CS 229: Machine Learning