

# LOGISTIC REGRESSION

David Kauchak  
CS158 – Fall 2023

1

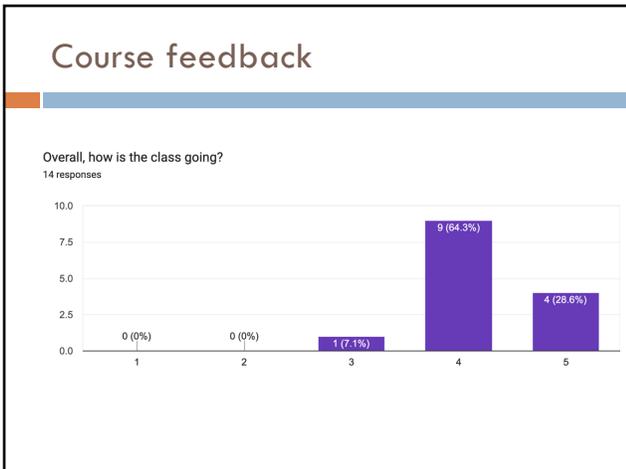
## Admin

Assignment 7

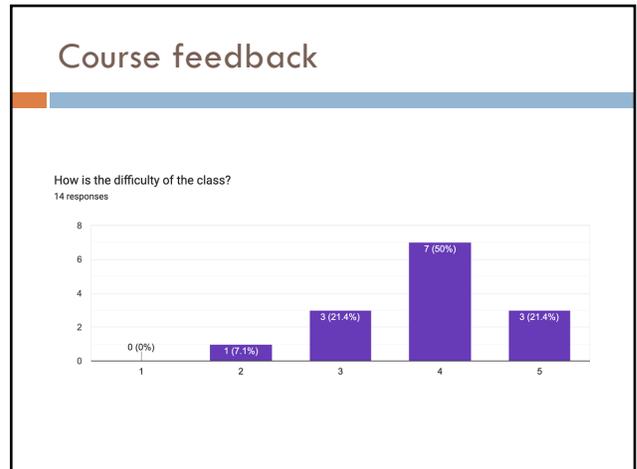
Grading update

Friday mentor hours: 6-8pm

2



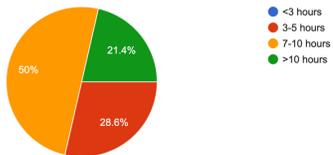
3



4

## Course feedback

About how many hours a week do you spend on this class (ignoring the DT assignment :)?  
14 responses



5

## Favorite thing

I feel like my coding skills in general are improving significantly. I am practicing concise documentation/efficient coding.

I like the coding part of the class!

getting better at java

6

## Improvements

Posting slides at the start of class.

I like to start the assignments as early as Monday night. I would love it if we had mentor sessions on Saturday as well. Or definitely both on Thursdays and Fridays the least, because Sunday is not enough if we are too far away from being done.

More mentor sessions :/ Not likely though. Maybe just three total would be awesome.

Releasing the assignments at the same time every week would help.

7

## Improvements

Give the big picture - how everything we are leading connects? the life cycle of ML project maybe?

I wished we went deeper in the math side of things. I think implementing ML algorithms is fun and cool, so don't change that!

8

## Improvements

Post the autograder score before we are done with the assignment. For several of these, we see our results, and they look good, but we don't know for sure if it's correct or not, so in a lot of cases we lose points on edge cases that we could've solved had we known they were problematic.

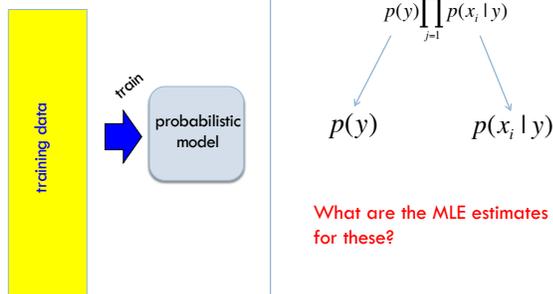
9

## Other comments

We lost points on the first three assignments to JavaDocs for other stylistic reasons, but we hadn't gotten our scores for the first assignment until after we had turned in the third assignment, so we didn't know we were supposed to do the JavaDocs and got dinged three times for it.

10

## MLE estimation for NB



11

## Maximum likelihood estimates

$$p(y) = \frac{\text{count}(y)}{n} \quad \frac{\text{number of examples with label } y}{\text{total number of examples}}$$

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)} \quad \frac{\text{number of examples with label } y \text{ with feature } x_i = 1}{\text{number of examples with label } y}$$

12

### Maximum likelihood estimates

$$p(y) = \frac{\text{count}(y)}{n}$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(1) = ?$   
 $p(-1) = ?$

13

### Maximum likelihood estimates

$$p(y) = \frac{\text{count}(y)}{n}$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(1) = 3/5$   
 $p(-1) = 2/5$

14

### Maximum likelihood estimates

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(x_1 = 1   1)$	?
$p(x_1 = 0   1)$	?
$p(x_2 = 1   1)$	?
$p(x_2 = 0   1)$	?

15

### Maximum likelihood estimates

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

16

## Maximum likelihood estimates

$$p(x, y) = p(y) \prod_{j=1}^m p(x_j | y)$$

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$$p(x_1 = 1, x_2 = 1, y = 1) = ?$$

17

## Maximum likelihood estimates

$$p(x|y) = p(y) \prod_{j=1}^m p(x_j | y)$$

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$$\begin{aligned} p(x_1 = 1, x_2 = 1, y = 1) &= p(1)p(x_1 = 1 | 1)p(x_2 = 1 | 1) \\ &= \frac{3}{5} * 1 * \frac{2}{3} \\ &= \frac{6}{15} \end{aligned}$$

18

## Maximum likelihood estimates

$$p(x, y) = p(y) \prod_{j=1}^m p(x_j | y)$$

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$$p(x_1 = 0, x_2 = 1, y = 1) = ?$$

19

## Maximum likelihood estimates

$$p(x, y) = p(y) \prod_{j=1}^m p(x_j | y)$$

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$$\begin{aligned} p(x_1 = 0, x_2 = 1, y = 1) &= p(1)p(x_1 = 0 | 1)p(x_2 = 1 | 1) \\ &= \frac{3}{5} * 0 * \frac{2}{3} \\ &= 0 \quad \text{☹} \end{aligned}$$

20

## Maximum likelihood estimates

Full model trained!

$$p(1) = 3/5$$

$$p(-1) = 2/5$$

$x_1$	$x_2$	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(x_1 = 1   1)$	3/3
$p(x_2 = 1   1)$	2/3

$p(x_1 = 1   -1)$	0/2
$p(x_2 = 1   -1)$	1/2

21

## Priors

Coin1 data: 3 Heads and 1 Tail

Coin2 data: 30 Heads and 10 tails

Coin3 data: 2 Tails

Coin4 data: 497 Heads and 503 tails

If someone asked you what the probability of heads was for each of these coins, what would you say?

22

## Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

### Probabilistic models

Which model do we use, i.e. how do we calculate  $p(\text{feature}, \text{label})$ ?

How do train the model, i.e. how to we **estimate the probabilities** for the model?

How do we deal with overfitting?

23

## Training revisited

What we're really doing during training is selecting the  $\Theta$  that maximizes:

$$p(\theta | \text{data})$$

i.e.

$$\theta = \operatorname{argmax}_{\theta} p(\theta | \text{data})$$

That is, we pick the most likely model parameters given the data

24

## Estimating revisited

We want to incorporate a prior belief of what the probabilities might be

To do this, we need to break down our probability

$$p(\theta \mid data) = ?$$

(Hint: Bayes rule)

25

## Estimating revisited

What are each of these probabilities?

$$p(\theta \mid data) = \frac{p(data \mid \theta)p(\theta)}{p(data)}$$

26

## Priors

likelihood of the data  
under the model

probability of different parameters,  
call the **prior**

$$p(\theta \mid data) = \frac{p(data \mid \theta)p(\theta)}{p(data)}$$

probability of seeing the data  
(regardless of model)

27

## Priors

$$\theta = \operatorname{argmax}_{\theta} \frac{p(data \mid \theta)p(\theta)}{p(data)}$$

Does  $p(data)$  matter for the  $\operatorname{argmax}$ ?

28

### Priors

likelihood of the data under the model

probability of different parameters, call the **prior**

$$\theta = \operatorname{argmax}_{\theta} p(\text{data} | \theta)p(\theta)$$

What does MLE assume for a prior on the model parameters?

29

### Priors

likelihood of the data under the model

probability of different parameters, call the **prior**

$$\theta = \operatorname{argmax}_{\theta} p(\text{data} | \theta)p(\theta)$$

- Assumes a **uniform prior**, i.e. all  $\Theta$  are equally likely!
- Relies solely on the **likelihood**

30

### A better approach

$$\theta = \operatorname{argmax}_{\theta} p(\text{data} | \theta)p(\theta)$$

$$\text{likelihood}(\text{data}) = \prod_{i=1}^n p_{\theta}(x_i)$$

We can use any distribution we'd like. This allows us to impart addition **bias** into the model

31

### Another view on the prior

Remember, the max is the same if we take the log:

$$\theta = \operatorname{argmax}_{\theta} \log(p(\text{data} | \theta)) + \log(p(\theta))$$

$$\text{log-likelihood} = \sum_{i=1}^n \log(p(x_i))$$

We can use any distribution we'd like. This allows us to impart addition **bias** into the model

Does this look like something we've seen before?

32

### Regularization vs prior

$$\theta = \operatorname{argmax}_{\theta} \log(p(\text{data} | \theta)) + \log(p(\theta))$$

fit { likelihood based on the data  
loss function based on the data } model bias { prior  
regularizer }

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \text{loss}(yy') + \lambda \text{regularizer}(w)$$

33

### Prior for NB

$$\theta = \operatorname{argmax}_{\theta} \log(p(\text{data} | \theta)) + \log(p(\theta))$$

Uniform prior | Dirichlet prior

$\lambda = 0$  increasing

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible\_values\_of\_}x_i * \lambda}$$

34

### Prior: another view

$$p(x_1, x_2, \dots, x_m, y) = p(y) \prod_{j=1}^m p(x_j | y)$$

MLE:  $p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$

What happens to our likelihood if, for one of the labels, we never saw a particular feature?

Goes to 0!

35

### Prior: another view

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

↓

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible\_values\_of\_}x_i * \lambda}$$

Adding a prior avoids this!

36

### Smoothing

training data

$$p(x_i | y) = \frac{\text{count}(x_i, y)}{\text{count}(y)}$$

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible\_values\_of\_}x_i * \lambda}$$

for each label, pretend like we've seen each feature value occur in  $\lambda$  additional examples

Sometimes this is also called **smoothing** because it is seen as smoothing or interpolating between the MLE and some other distribution

37

### Priors

Coin1 data: 3 Heads and 1 Tail  
 Coin2 data: 30 Heads and 10 tails  
 Coin3 data: 2 Tails  
 Coin4 data: 497 Heads and 503 tails

$$p(\text{heads}) = \frac{\text{count}(\text{heads}) + \lambda}{\text{totalflips} + 2\lambda}$$

Does this do the right thing in these cases?

38

### Maximum likelihood estimates

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible\_values\_of\_}x_i * \lambda} \quad \lambda = 1$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(x_1 = 1   1)$	?
$p(x_1 = 0   1)$	?
$p(x_2 = 1   1)$	?
$p(x_2 = 0   1)$	?

39

### Maximum likelihood estimates

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible\_values\_of\_}x_i * \lambda} \quad \lambda = 1$$

x1	x2	label
1	1	1
1	0	1
1	1	1
0	1	-1
0	0	-1

$p(x_1 = 1   1)$	4/5
$p(x_1 = 0   1)$	1/5
$p(x_2 = 1   1)$	3/5
$p(x_2 = 0   1)$	2/5

40

## Avoids zero probability events!

$p(x_1 = 1   1)$	3/3
$p(x_1 = 0   1)$	0/3
$p(x_2 = 1   1)$	2/3
$p(x_2 = 0   1)$	1/3

smoothed/prior

$p(x_1 = 1   1)$	4/5
$p(x_1 = 0   1)$	1/5
$p(x_2 = 1   1)$	3/5
$p(x_2 = 0   1)$	2/5

41

## Basic steps for probabilistic modeling

Step 1: pick a model

### Probabilistic models

Which model do we use, i.e. how do we calculate  $p(\text{feature}, \text{label})$ ?

Step 2: figure out how to estimate the probabilities for the model

How do train the model, i.e. how to we we **estimate the probabilities** for the model?

Step 3 (optional): deal with overfitting

How do we deal with overfitting?

42

## Joint models vs conditional models

We've been trying to model the joint distribution (i.e. the data generating distribution):

$$p(x_1, x_2, \dots, x_m, y)$$

However, if all we're interested in is classification, why not directly model the conditional distribution:

$$p(y | x_1, x_2, \dots, x_m)$$

43

## A first try: linear

$$p(y | x_1, x_2, \dots, x_m) = x_1 w_1 + w_2 x_2 + \dots + w_m x_m + b$$

Any problems with this?

- Nothing constrains it to be a probability
- Could still have combination of features and weight that exceeds 1 or is below 0

44

### The challenge

$x_1w_1 + w_2x_2 + \dots + w_mx_m + b$

**Linear model**

$+\infty$

$-\infty$

**probability**

$p(y | x_1, x_2, \dots, x_m)$

1

0

We like linear models!

Can we transform the probability into a function that ranges over all values?

45

### Odds ratio

Rather than predict the probability, we can predict the ratio of 1/0 (positive/negative)

Predict the **odds** that it is 1 (true): **How much more likely is 1 than 0.**

Does this help us?

$$\frac{P(1|x_1, x_2, \dots, x_m)}{P(0|x_1, x_2, \dots, x_m)} = \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = x_1w_1 + w_2x_2 + \dots + w_mx_m + b$$

46

### Odds ratio

$x_1w_1 + w_2x_2 + \dots + w_mx_m + b$

**Linear model**

$+\infty$

$-\infty$

**odds ratio**

$\frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)}$

$+\infty$

0

Where is the dividing line between class 1 and class 0 being selected?

47

### Odds ratio

$\frac{P(1|x_1, x_2, \dots, x_m)}{P(0|x_1, x_2, \dots, x_m)} > \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)}$

$P(1|x_1, x_2, \dots, x_m) > 1 - P(1|x_1, x_2, \dots, x_m)$

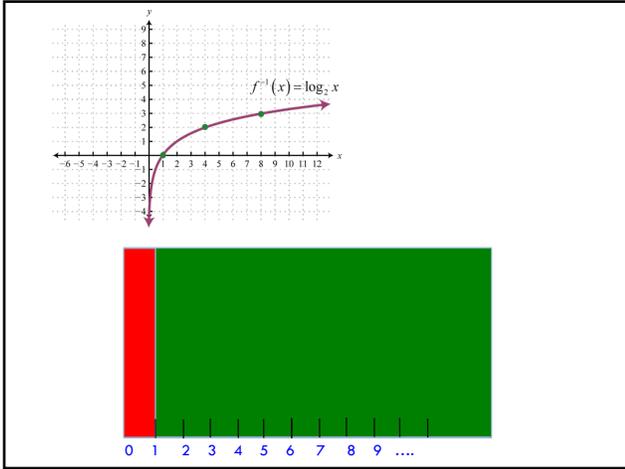
We're trying to find some transformation that transforms the odds ratio to a number that is  $-\infty$  to  $+\infty$

Does this suggest another transformation?

odds ratio

0 1 2 3 4 5 6 7 8 9 ...

48



49

### Log odds (logit function)

$$x_1 w_1 + w_2 x_2 + \dots + w_m x_m + b = \log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)}$$

Linear regression odds ratio

+∞

-∞

How do we get the probability of an example?

+∞

-∞

50

### Log odds (logit function)

$$\log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b$$

$$\frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = e^{w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b}$$

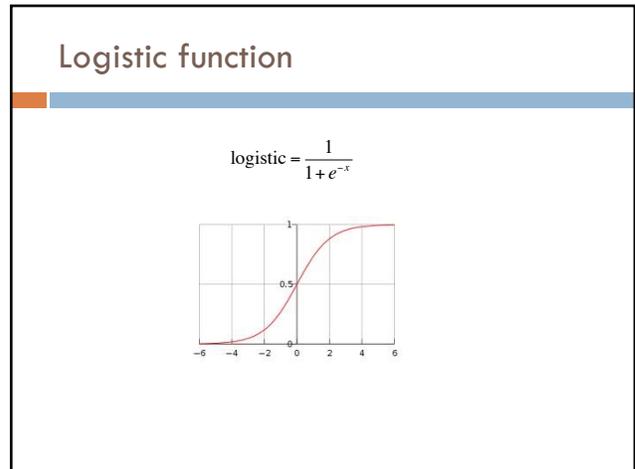
$$P(1|x_1, x_2, \dots, x_m) = (1 - P(1|x_1, x_2, \dots, x_m)) e^{w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b}$$

...

$$P(1|x_1, x_2, \dots, x_m) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b)}}$$

anyone recognize this?

51



52

## Logistic regression

How would we classify examples once we had a trained model?

$$\log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b$$

If the sum  $> 0$  then  $p(1)/p(0) > 1$ , so positive

if the sum  $< 0$  then  $p(1)/p(0) < 1$ , so negative

Still a *linear* classifier (decision boundary is a line)

53

## Training logistic regression models

How should we learn the parameters for logistic regression (i.e. the  $w$ 's and  $b$ )?

$$\log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b$$

parameters

$$P(1|x_1, x_2, \dots, x_m) = \frac{1}{1 + e^{-(w_1 x_1 + w_2 x_2 + \dots + w_m x_m + b)}}$$

54

## MLE logistic regression

Find the parameters that maximize the likelihood (or log-likelihood) of the data:

$$\begin{aligned} \text{log-likelihood} &= \sum_{i=1}^n \log(p(x_i)) \\ &= \sum_{i=1}^n \log\left(\frac{1}{1 + e^{-y_i(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)}}\right) \quad \text{assume labels } 1, -1 \\ &= \sum_{i=1}^n -\log(1 + e^{-y_i(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)}) \end{aligned}$$

55

## MLE logistic regression

$$\text{log-likelihood} = \sum_{i=1}^n -\log(1 + e^{-y_i(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)})$$

We want to maximize, i.e.

$$\begin{aligned} MLE(\text{data}) &= \operatorname{argmax}_{w, b} \text{log-likelihood}(\text{data}) \\ &= \operatorname{argmax}_{w, b} \sum_{i=1}^n -\log(1 + e^{-y_i(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)}) \\ &= \operatorname{argmin}_{w, b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1 x_{i1} + w_2 x_{i2} + \dots + w_m x_{im} + b)}) \end{aligned}$$

Look familiar? Hint: anybody reading the book?

56

### MLE logistic regression

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)})$$

Surrogate loss functions:

- Zero/one:  $\ell^{(0/1)}(y, \hat{y}) = \mathbf{1}[y\hat{y} \leq 0]$
- Hinge:  $\ell^{(\text{hinge})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$
- Logistic:  $\ell^{(\text{log})}(y, \hat{y}) = \frac{1}{\log 2} \log(1 + \exp[-y\hat{y}])$
- Exponential:  $\ell^{(\text{exp})}(y, \hat{y}) = \exp[-y\hat{y}]$
- Squared:  $\ell^{(\text{sq})}(y, \hat{y}) = (y - \hat{y})^2$

57

### logistic regression: three views

$$\log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m$$

linear classifier

$$P(1|x_1, x_2, \dots, x_m) = \frac{1}{1 + e^{-(w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m)}}$$

conditional model  
logistic

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)})$$

linear model  
minimizing logistic loss

58

### Overfitting

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)})$$

If we minimize this loss function, in practice, the results aren't great and we tend to overfit

Solution?

59

### Regularization/prior

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \text{regularizer}(w, b)$$

or

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) - \log(p(w, b))$$

What are some of the regularizers we know?

60

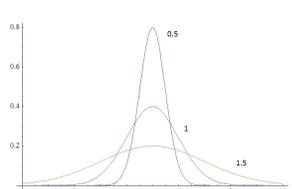
### Regularization/prior

**L2 regularization:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \|w\|^2$$

**Gaussian prior:**  
Gaussians are defined by a mean ( $\mu$ ) and a variance ( $\sigma^2$ )

$p(w,b) \sim$



61

### Regularization/prior

**L2 regularization:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \|w\|^2$$

**Gaussian prior:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \frac{1}{2\sigma^2} \|w\|^2$$

Does the  $\lambda$  make sense?  $\lambda = \frac{1}{2\sigma^2}$

62

### Regularization/prior

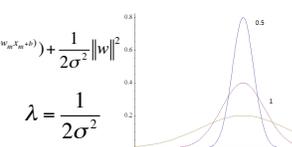
**L2 regularization:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \|w\|^2$$

**Gaussian prior:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \frac{1}{2\sigma^2} \|w\|^2$$

$\lambda = \frac{1}{2\sigma^2}$



63

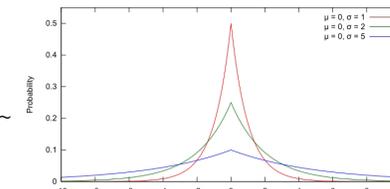
### Regularization/prior

**L1 regularization:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \|w\|$$

**Laplacian prior:**

$p(w,b) \sim$



64

## Regularization/prior

**L1 regularization:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \lambda \|w\|$$

**Laplacian prior:**

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) + \frac{1}{\sigma} \|w\|$$

$$\lambda = \frac{1}{\sigma}$$

65

## L1 vs. L2

L1 = Laplacian prior

L2 = Gaussian prior

66

## Logistic regression

Why is it called logistic regression?  
It is a classifier??

$$\log \frac{P(1|x_1, x_2, \dots, x_m)}{1 - P(1|x_1, x_2, \dots, x_m)} = w_1x_1 + w_2x_2 + \dots + w_mx_m + b$$

67

## A digression: regression vs. classification

<table border="0"> <tr><th>Raw data</th><th>Label</th></tr> <tr><td style="text-align: center;">■</td><td style="text-align: center;">0</td></tr> <tr><td style="text-align: center;">■</td><td style="text-align: center;">0</td></tr> <tr><td style="text-align: center;">■</td><td style="text-align: center;">1</td></tr> <tr><td style="text-align: center;">■</td><td style="text-align: center;">1</td></tr> <tr><td style="text-align: center;">■</td><td style="text-align: center;">0</td></tr> </table>	Raw data	Label	■	0	■	0	■	1	■	1	■	0	➔	<table border="0"> <tr><th>features</th><th>Label</th></tr> <tr><td style="text-align: center;">f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>n</sub></td><td style="text-align: center;">■</td></tr> <tr><td style="text-align: center;">f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>n</sub></td><td style="text-align: center;">■</td></tr> <tr><td style="text-align: center;">f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>n</sub></td><td style="text-align: center;">■</td></tr> <tr><td style="text-align: center;">f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>n</sub></td><td style="text-align: center;">■</td></tr> <tr><td style="text-align: center;">f<sub>1</sub>, f<sub>2</sub>, f<sub>3</sub>, ..., f<sub>n</sub></td><td style="text-align: center;">■</td></tr> </table>	features	Label	f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■	f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■	f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■	f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■	f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■
Raw data	Label																									
■	0																									
■	0																									
■	1																									
■	1																									
■	0																									
features	Label																									
f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■																									
f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■																									
f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■																									
f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■																									
f <sub>1</sub> , f <sub>2</sub> , f <sub>3</sub> , ..., f <sub>n</sub>	■																									

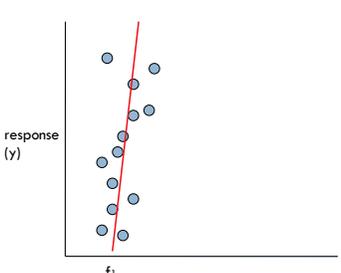
extract features

classification:  
discrete (some finite set of labels)

regression: real value

68

### linear regression



Given some points, find the **line** that best fits/explains the data

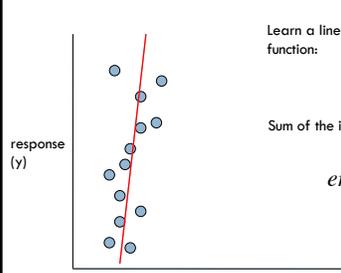
Our model is a line, i.e. we're assuming a linear relationship between the feature and the label value

$$h(y) = w_1 x_1 + b$$

How can we find this line?

69

### Linear regression



Learn a line  $h$  that minimizes some loss/error function:

$$error(h) = ?$$

Sum of the individual errors:

$$error(h) = \sum_{i=1}^n |y_i - h(f_i)|$$

0/1 loss!

70

### Error minimization

How do we find the minimum of an equation?

$$error(h) = \sum_{i=1}^n |y_i - h(f_i)|$$

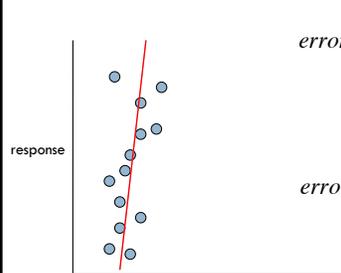
Take the derivative, set to 0 and solve (going to be a min or a max)

Any problems here?

Ideas?

71

### Linear regression



$$error(h) = \sum_{i=1}^n |y_i - h(f_i)|$$

↓

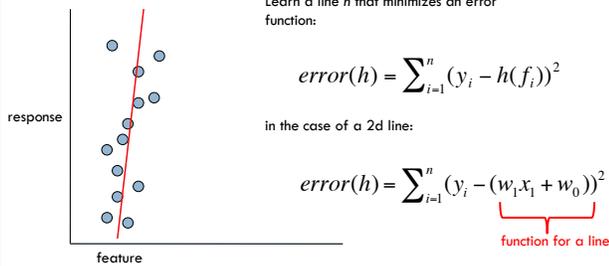
$$error(h) = \sum_{i=1}^n (y_i - h(f_i))^2$$

squared error is convex!

Squared:  $\ell^{(sq)}(y, \hat{y}) = (y - \hat{y})^2$

72

## Linear regression



73

## Linear regression

We'd like to *minimize* the error

Find  $w_1$  and  $w_0$  such that the error is minimized

$$error(h) = \sum_{i=1}^n (y_i - (w_1 f_i + w_0))^2$$

We can solve this in closed form

74

## Multiple linear regression

If we have  $m$  features, then we have a line in  $m$  dimensions

$$h(\vec{f}) = w_0 + w_1 f_1 + w_2 f_2 + \dots + w_m f_m$$

weights

75

## Multiple linear regression

We can still calculate the squared error like before

$$h(\vec{f}) = w_0 + w_1 f_1 + w_2 f_2 + \dots + w_m f_m$$

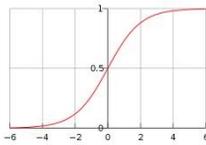
$$error(h) = \sum_{i=1}^n (y_i - (w_0 + w_1 f_{i1} + w_2 f_{i2} + \dots + w_m f_{im}))^2$$

Still can solve this exactly!

76

## Logistic function

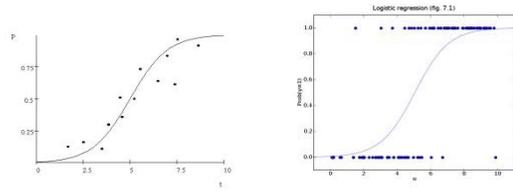
$$\text{logistic} = \frac{1}{1 + e^{-x}}$$



77

## Logistic regression

Find the best fit of the data based on a logistic



78

## Basic steps for probabilistic modeling

Step 1: pick a model

Step 2: figure out how to estimate the probabilities for the model

Step 3 (optional): deal with overfitting

### Probabilistic models

Which model do we use, i.e. how do we calculate  $p(\text{feature}, \text{label})$ ?

How do train the model, i.e. how to we **estimate the probabilities** for the model?

How do we deal with overfitting?

79

## Probabilistic models summarized

Two classification models:

- Naive Bayes (models **joint** distribution)
- Logistic Regression (models **conditional** distribution)
  - In practice this tends to work better if all you want to do is classify

Priors/smoothing/regularization

- Important for both models
- In theory: allow us to impart some prior knowledge
- In practice: avoids overfitting and often tune on development data

80