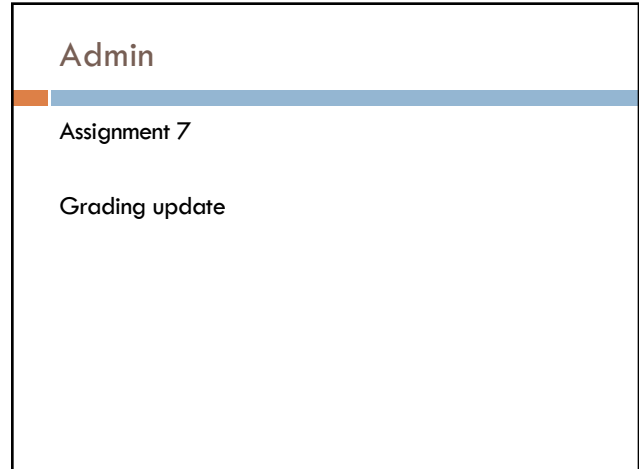


LOGISTIC REGRESSION

David Kauchak
CS158 – Fall 2019

1

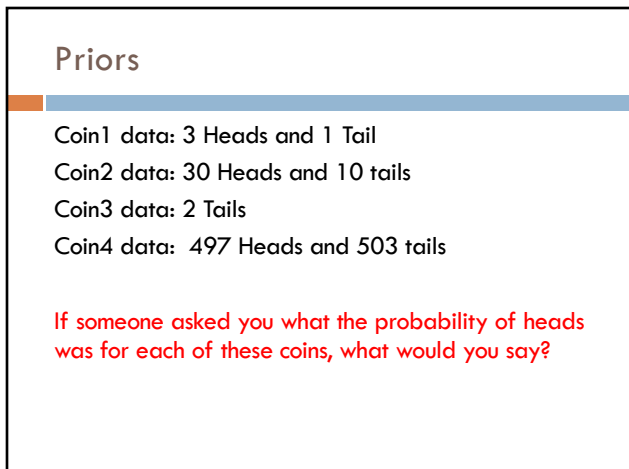


Admin

Assignment 7

Grading update

2

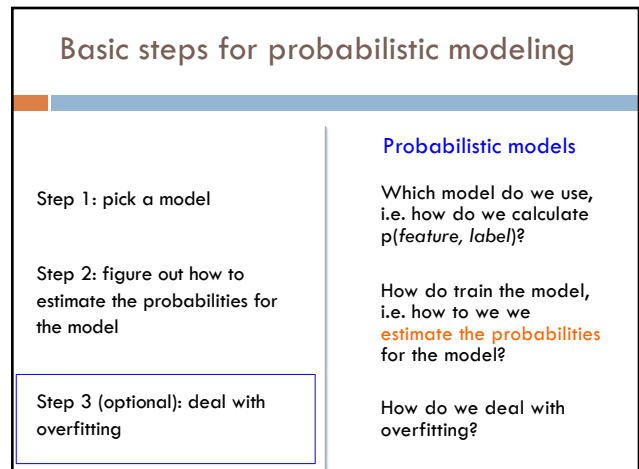


Priors

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?

3



Basic steps for probabilistic modeling

<p>Step 1: pick a model</p> <p>Step 2: figure out how to estimate the probabilities for the model</p> <p>Step 3 (optional): deal with overfitting</p>	<p>Probabilistic models</p> <p>Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?</p> <p>How do train the model, i.e. how to we we estimate the probabilities for the model?</p> <p>How do we deal with overfitting?</p>
---	--

4

Training revisited

What we're really doing during training is selecting the θ that maximizes:

$$p(\theta | data)$$

i.e.

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

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

5

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

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

(Hint: Bayes rule)

6

Estimating revisited

What are each of these probabilities?

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

7

Priors

likelihood of the data
under the model

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

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

probability of seeing the data
(regardless of model)

8

Priors

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

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

9

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?

10

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 Θ are equally likely!
- Relies solely on the likelihood

11

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 additional **bias**
into the model

12

Another view on the prior

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

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

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?

13

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 } prior regularizer } model bias

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

14

Prior for NB

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

Uniform prior

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

Dirichlet prior

$\lambda = 0$ → increasing

$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible_values_of_}x_i * \lambda}$

15

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?

Go to 0!

16

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!

17

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 λ additional examples

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

18

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

$$p(x_i | y) = \frac{\text{count}(x_i, y) + \lambda}{\text{count}(y) + \text{possible_values_of_}x_i * \lambda}$$

Does this do the right thing in these cases?

19

Basic steps for probabilistic modeling

<p>Step 1: pick a model</p> <hr/> <p>Step 2: figure out how to estimate the probabilities for the model</p> <hr/> <p>Step 3 (optional): deal with overfitting</p>	<p>Probabilistic models</p> <p>Which model do we use, i.e. how do we calculate $p(\text{feature}, \text{label})$?</p> <p>How do train the model, i.e. how to we we estimate the probabilities for the model?</p> <p>How do we deal with overfitting?</p>
---	---

20

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:

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

21

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

22

The challenge

$$x_1 w_1 + w_2 x_2 + \dots + w_m x_m + b$$

Linear model

$+\infty$



$-\infty$

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

probability

1



0



We like linear models!

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

23

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?

$$\frac{P(1 | x_1, x_2, \dots, x_m)}{P(0 | x_1, x_2, \dots, x_m)} = \frac{P(1 | x_1, x_2, \dots, x_m)}{1 - P(1 | x_1, x_2, \dots, x_m)} = x_1 w_1 + w_2 x_2 + \dots + w_m x_m + b$$

24

Odds ratio

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

Linear model

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

odds ratio

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

25

Odds ratio

$P(1|x_1, x_2, \dots, x_m) > P(0|x_1, x_2, \dots, x_m)$

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

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

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

26

$f^{-1}(x) = \log_e x$

27

Log odds (logit function)

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

Linear regression

odds ratio

How do we get the probability of an example?

28

Log odds (logit function)

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

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

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

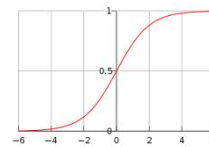
...

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

29

Logistic function

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



30

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)

31

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

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

32

MLE logistic regression

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

$$\begin{aligned} \text{log-likelihood} &= \sum_{i=1}^n \log(p(x_i)) \\ &= \sum_{i=1}^n \log\left(\frac{1}{1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}}\right) \quad \text{assume labels } 1, -1 \\ &= \sum_{i=1}^n -\log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) \end{aligned}$$

33

MLE logistic regression

$$\text{log-likelihood} = \sum_{i=1}^n -\log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)})$$

We want to maximize, i.e.

$$\begin{aligned} MLE(data) &= \operatorname{argmax}_{w,b} \text{log-likelihood}(data) \\ &= \operatorname{argmax}_{w,b} \sum_{i=1}^n -\log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) \\ &= \operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) \end{aligned}$$

Look familiar? Hint: anybody read the book?

34

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: $\ell^{(\text{hin})}(y, \hat{y}) = \max\{0, 1 - y\hat{y}\}$

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

35

logistic regression: three views

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

$$P(1|x_1, x_2, \dots, x_m) = \frac{1}{1 + e^{-(w_0 + w_1x_1 + w_2x_2 + \dots + w_mx_m)}} \quad \begin{array}{l} \text{conditional model} \\ \text{logistic} \end{array}$$

$$\operatorname{argmin}_{w,b} \sum_{i=1}^n \log(1 + e^{-y_i(w_1x_{i1} + w_2x_{i2} + \dots + w_mx_{im} + b)}) \quad \begin{array}{l} \text{linear model} \\ \text{minimizing logistic loss} \end{array}$$

36

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?

37

Regularization/prior

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

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?

38

Regularization/prior

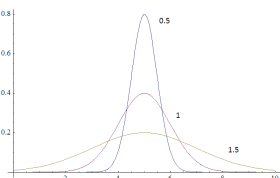
L2 regularization:

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

Gaussian prior:

Gaussians are defined by a mean (μ) and a variance (σ^2)

$p(w,b) \sim$



39

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 λ make sense?

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

40

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:

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

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

41

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:

$$p(w,b) \sim$$

42

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:

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

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

43

L1 vs. L2

L1 = Laplacian prior

L2 = Gaussian prior

44

Logistic regression

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

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

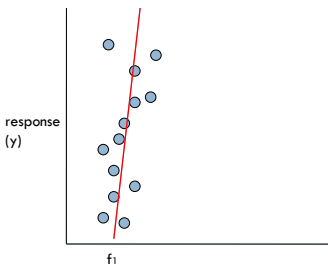
45

A digression: regression vs. classification

<table border="0"> <tr> <th style="text-align: left;">Raw data</th> <th style="text-align: left;">Label</th> </tr> <tr> <td style="text-align: center;">■</td> <td>0</td> </tr> <tr> <td style="text-align: center;">■</td> <td>0</td> </tr> <tr> <td style="text-align: center;">■</td> <td>1</td> </tr> <tr> <td style="text-align: center;">■</td> <td>1</td> </tr> <tr> <td style="text-align: center;">■</td> <td>0</td> </tr> </table>	Raw data	Label	■	0	■	0	■	1	■	1	■	0	<div style="color: blue; font-size: 2em;">➔</div> extract features	<table border="0"> <tr> <th style="text-align: left;">features</th> <th style="text-align: left;">Label</th> </tr> <tr> <td style="text-align: left;">f₁, f₂, f₃, ..., f_n</td> <td style="text-align: center;">■</td> </tr> <tr> <td style="text-align: left;">f₁, f₂, f₃, ..., f_n</td> <td style="text-align: center;">■</td> </tr> <tr> <td style="text-align: left;">f₁, f₂, f₃, ..., f_n</td> <td style="text-align: center;">■</td> </tr> <tr> <td style="text-align: left;">f₁, f₂, f₃, ..., f_n</td> <td style="text-align: center;">■</td> </tr> <tr> <td style="text-align: left;">f₁, f₂, f₃, ..., f_n</td> <td style="text-align: center;">■</td> </tr> </table>	features	Label	f ₁ , f ₂ , f ₃ , ..., f _n	■	f ₁ , f ₂ , f ₃ , ..., f _n	■	f ₁ , f ₂ , f ₃ , ..., f _n	■	f ₁ , f ₂ , f ₃ , ..., f _n	■	f ₁ , f ₂ , f ₃ , ..., f _n	■	<div style="color: blue;"> classification: discrete (some finite set of labels) </div> <div style="color: blue; margin-top: 10px;"> regression: real value </div>
Raw data	Label																										
■	0																										
■	0																										
■	1																										
■	1																										
■	0																										
features	Label																										
f ₁ , f ₂ , f ₃ , ..., f _n	■																										
f ₁ , f ₂ , f ₃ , ..., f _n	■																										
f ₁ , f ₂ , f ₃ , ..., f _n	■																										
f ₁ , f ₂ , f ₃ , ..., f _n	■																										
f ₁ , f ₂ , f ₃ , ..., f _n	■																										

46

linear regression



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

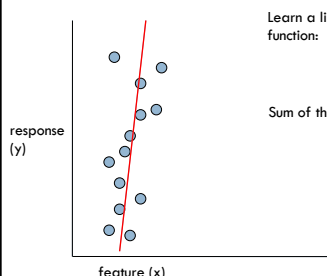
Our model is a line, i.e. we're assuming a linear relationship between the feature and the label value

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

How can we find this line?

47

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!

48

Error minimization

How do we find the minimum of an equation?

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

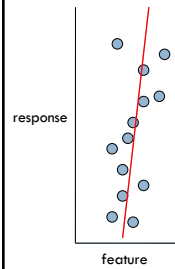
Take the derivative, set to 0 and solve (going to be a min or a max)

Any problems here?

Ideas?

49

Linear regression



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



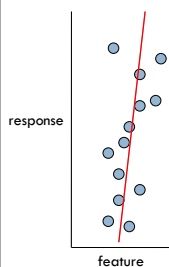
$$error(h) = \sum_{i=1}^n (y_i - h(f_i))^2$$

squared error is convex!

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

50

Linear regression



Learn a line h that minimizes an error function:

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

in the case of a 2d line:

$$error(h) = \sum_{i=1}^n (y_i - \underbrace{(w_1 x_i + w_0)}_{\text{function for a line}})^2$$

51

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

52

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

53

Multiple linear regression

We can still calculate the squared error like before

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

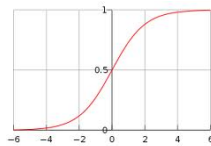
$$error(h) = \sum_{i=1}^n (y_i - (w_0 + w_1 f_{i1} + w_2 f_{i2} + \dots + w_m f_{im}))^2$$

Still can solve this exactly!

54

Logistic function

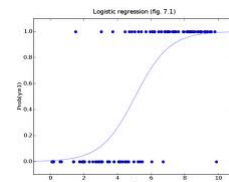
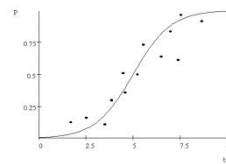
$$\text{logistic} = \frac{1}{1 + e^{-x}}$$



55

Logistic regression

Find the best fit of the data based on a logistic



56

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 we **estimate the probabilities** for the model?

How do we deal with overfitting?

Probabilistic models summarized

Two classification models:

- ▣ Naïve 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

57

58