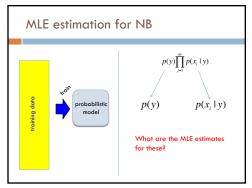
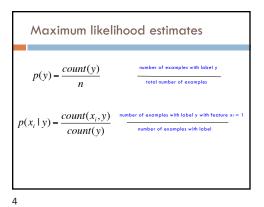
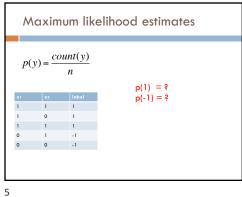


Admin
Assignment 7

1







Maximum likelihood estimates  $p(y) = \frac{count(y)}{n}$ p(1) = 3/5p(-1) = 2/51 1 1 1 1 1 0 1 -1 0 0 -1

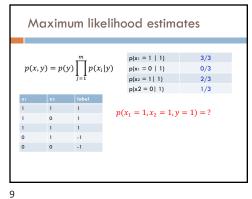
6

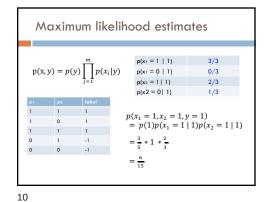
8

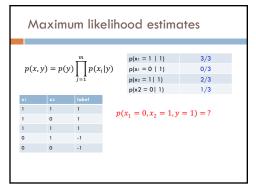
Maximum likelihood estimates  $p(x_i \mid y) = \frac{count(x_i, y)}{}$  $p(x_1 = 1 \mid 1)$  $p(x_1 = 0 | 1)$ ş count(y)  $p(x_2 = 1 \mid 1)$  ? p(x2 = 0 | 1)1 1 1 1 1 1 0 1 -1 0 -1

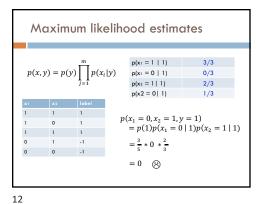
7

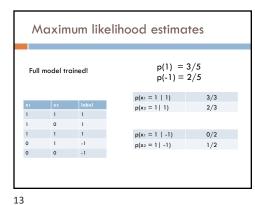
Maximum likelihood estimates  $p(x_i \mid y) = \frac{count(x_i, y)}{count(y)}$  $p(x_1 = 1 \mid 1)$  3/3  $p(x_1 = 0 \mid 1)$ 0/3  $p(x_2 = 1 \mid 1)$  2/3 p(x2 = 0 | 1) 1/3 1 1 1 1 0 1 1 1 1 0 1 -1 0 -1

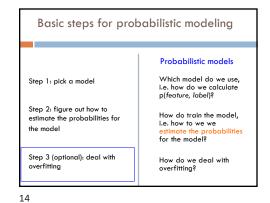


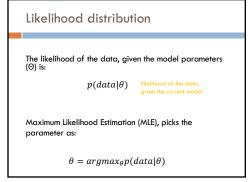


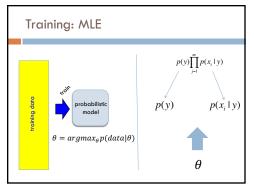


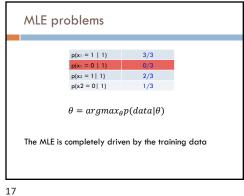


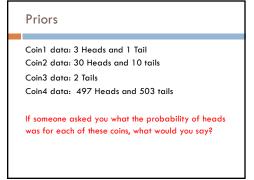


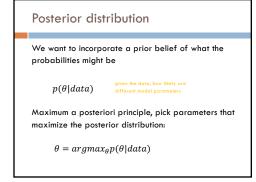


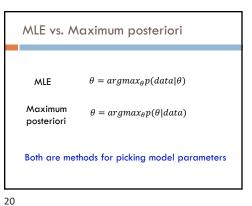










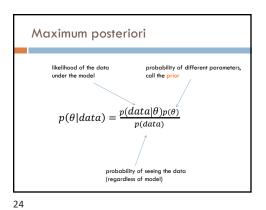


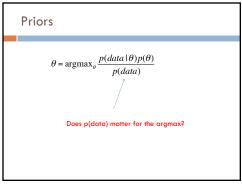
Maximum posteriori  $p(\theta|data) =$ ? (Hint: Bayes' rule) Maximum posteriori  $p(\theta|data) = \frac{p(data|\theta)p(\theta)}{p(data)}$ 

22

21

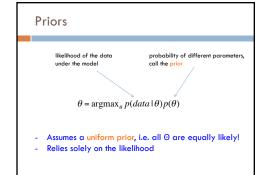
Maximum posteriori What are each of these probabilities?  $p(\theta|data) = \frac{p(data|\theta)p(\theta)}{p(data)}$ 23

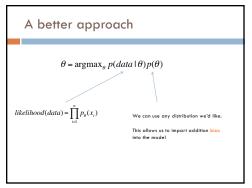


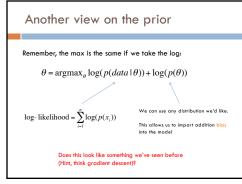


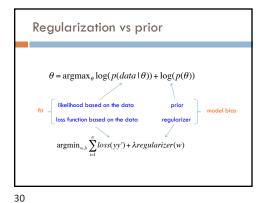
25

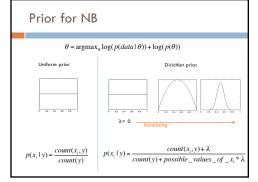
26

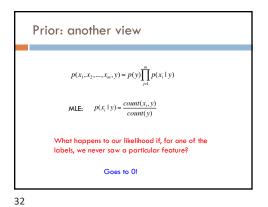


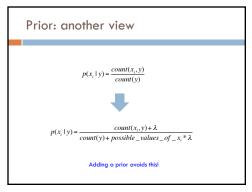


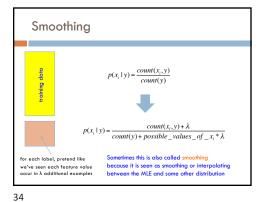












35

Priors

Coin1 data: 3 Heads and 1 Tail

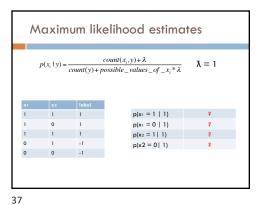
Coin2 data: 30 Heads and 10 tails

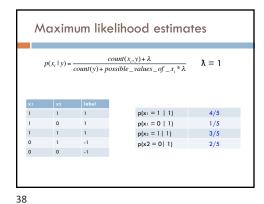
Coin3 data: 2 Tails

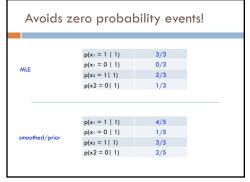
Coin4 data: 497 Heads and 503 tails  $p(heads) = \frac{count(heads) + \lambda}{totalflips + 2\lambda}$ Does this do the right thing in these cases?

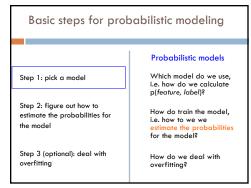
Priors

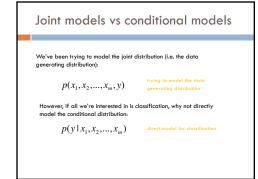
Coin1 data: 3 Heads and 1 Tail:  $p(heads) = \frac{4}{6} = 0.667$ Coin2 data: 30 Heads and 10 tails:  $p(heads) = \frac{31}{42} = 0.738$ Coin3 data: 2 Tails  $p(heads) = \frac{3}{4} = 0.75$ Coin4 data: 497 Heads and 503 tails  $p(heads) = \frac{498}{1002} = 0.497$   $p(heads) = \frac{count(heads) + \lambda}{totalflips + 2\lambda}$ 





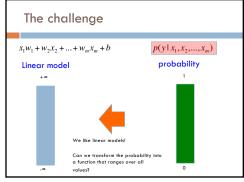






A first try: linear  $p(y \mid x_1, x_2, ..., x_m) = x_1 w_1 + w_2 x_2 + ... + w_m x_m + b$ Any problems with this?

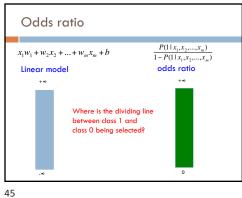
- Nothing constrains it to be a probability
- Could still have combination of features and weight that exceeds 1 or is below 0

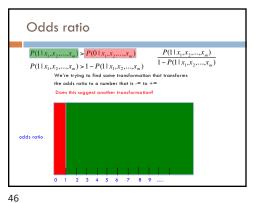


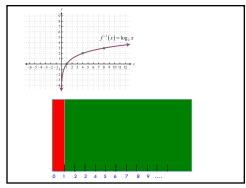
Rather than predict the probability, we can predict the ratio of 1/0 (positive/negative)

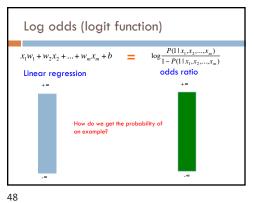
Predict the **odds** that it is 1 (true): How much more likely is 1 than 0.

Does this help us?  $\frac{P(1|x_1,x_2,...,x_m)}{P(0|x_1,x_2,...,x_m)} = \frac{P(1|x_1,x_2,...,x_m)}{1-P(1|x_1,x_2,...,x_m)} = x_1w_1 + w_2x_2 + ... + w_mx_m + b$ 









Log odds (logit function)

$$\begin{split} \log \frac{P(1 \mid x_1, x_2, ..., x_m)}{1 - P(1 \mid x_1, x_2, ..., x_m)} &= w_1 x_2 + w_2 x_2 + ... + w_m x_m + b \\ \\ \frac{P(1 \mid x_1, x_2, ..., x_m)}{1 - P(1 \mid x_1, x_2, ..., x_m)} &= e^{w_1 x_2 + w_2 x_2 + ... + w_m x_m + b} \\ P(1 \mid x_1, x_2, ..., x_m) &= (1 - P(1 \mid x_1, x_2, ..., x_m)) e^{w_1 x_2 + w_2 x_2 + ... + w_m x_m + b} \\ \\ P(1 \mid x_1, x_2, ..., x_m) &= \frac{1}{1 + e^{-(w_1 x_2 + w_2 x_2 + ... + w_m x_m + b)}} \end{split}$$

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

49

51

50

Logistic regression

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

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

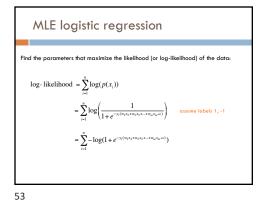
If the sum > 0 then p(1)/p(0) > 1, so positive

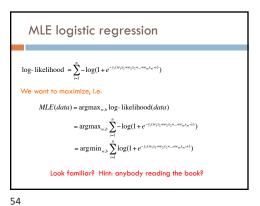
if the sum < 0 then p(1)/p(0) < 1, so negative

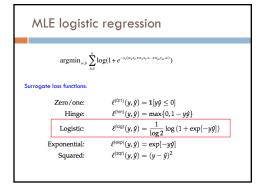
Still a linear classifier (decision boundary is a line)

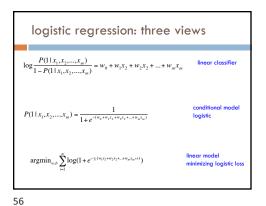
Training logistic regression models

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

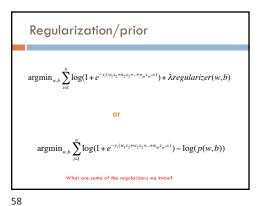










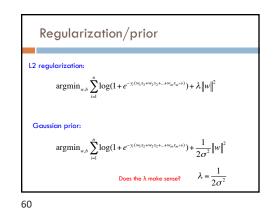


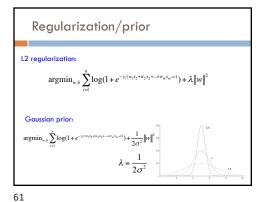
Regularization/prior

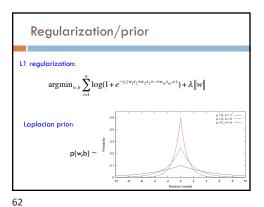
1.2 regularization:  $\arg\min_{w,b} \sum_{l=1}^{n} \log(1 + e^{-y_{l}(w_{l}x_{2} + w_{2}x_{2} + \dots + w_{w}x_{m} + b)}) + \lambda \|w\|^{2}$ Gaussian prior:

Gaussians ore defined by a mean (µ) and a variance (a?)  $p(w,b) \sim \sum_{l=1}^{n} \log(1 + e^{-y_{l}(w_{l}x_{2} + w_{2}x_{2} + \dots + w_{w}x_{m} + b)}) + \lambda \|w\|^{2}$ 

59

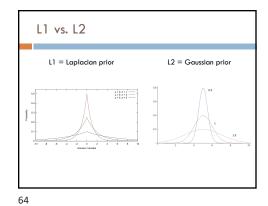


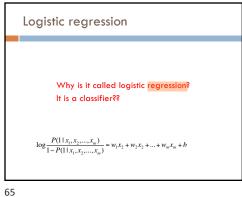


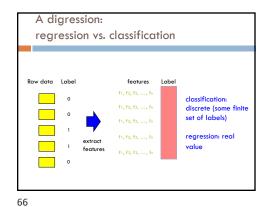


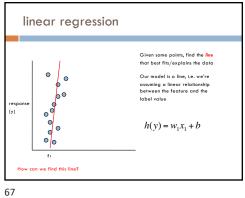
Regularization/prior

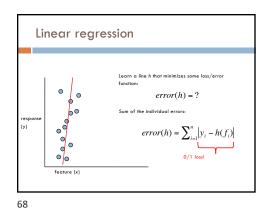
L1 regularization:  $\operatorname{argmin}_{w,b} \sum_{i=1}^{n} \log(1 + e^{-y_i(w_i x_2 + w_2 x_2 + \dots + w_m x_m + b)}) + \lambda \|w\|$ Laplacian prior:  $\operatorname{argmin}_{w,b} \sum_{i=1}^{n} \log(1 + e^{-y_i(w_i x_2 + w_2 x_2 + \dots + w_m x_m + b)}) + \frac{1}{\sigma} \|w\|$   $\lambda = \frac{1}{\sigma}$ 

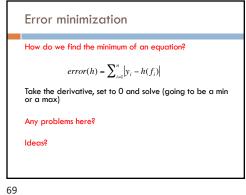


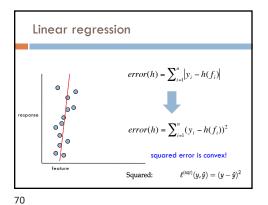




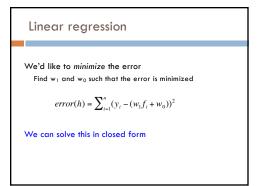








Linear regression Learn a line h that minimizes an error  $error(h) = \sum_{i=1}^{n} (y_i - h(f_i))^2$  $error(h) = \sum_{i=1}^{n} (y_i - (w_1 x_1 + w_0))^2$ function for a line



## Multiple linear regression

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

$$h(\bar{f}) = w_0 + w_1 f_1 + w_2 f_2 + \ldots + w_m f_m$$
 weights

## Multiple linear regression

We can still calculate the squared error like before

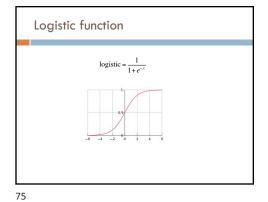
$$h(\bar{f}) = w_0 + w_1 f_1 + w_2 f_2 + ... + w_m f_m$$

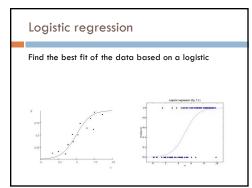
$$error(h) = \sum_{i=1}^{n} (y_i - (w_0 + w_1 f_1 + w_2 f_2 + ... + w_m f_m))^2$$

Still can solve this exactly!

73

74





## 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 Basic steps for probabilistic models Which model do we use, i.e. how do we calculate p(feature, label)? How do train the model, i.e. how to we we estimate the probabilities for the model? How do we deal with overfitting?

77

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