

10-701 Recitation (9/23/2010)

To cover:

MLE vs. MAP

NB example

Priors: Beta ( $\beta_H, \beta_T$ ) vs. Uniform

Relation between prior and # of data points

NB decision boundary

- Linear vs. quadratic

- connection to Bayes optimal classifier

Multinomial & Dirichlet

- Connection to Bernoulli and beta

LR

- connection to naive Bayes; be careful about conditions.

add- $\alpha$  smoothing

↓

## MLE vs. MAP

MLE and MAP are ways of estimating the parameters of a model

MLE - maximum likelihood estimate (frequentist)

MAP - maximum a posteriori estimate (Bayesian)

Example: Flipping a coin:

$X^{(1)}, \dots, X^{(n)} \sim \text{Bernoulli}(\theta)$   $\theta$ : probability coin comes up heads. Heads  $\Rightarrow X=1$

We are given a sample of  $n$  flips,  $X^1, \dots, X^n$ .

How do we guess the value of  $\theta$ ?

MLE: choose  $\hat{\theta}_{MLE}$  to maximize the likelihood of  $X^1, \dots, X^n$ .

$$\hat{\theta}_{MLE} = \arg \max_{\theta} p(X^1, \dots, X^n | \theta) \leftarrow \text{data likelihood}$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \prod_{i=1}^n \theta^{x^i} (1-\theta)^{(1-x^i)} \leftarrow \text{Maximize using Calculus.}$$

$$\text{In class, we saw } \hat{\theta}_{MLE} = \frac{\sum_{i=1}^n x^i}{n} = \frac{\# \text{ of heads}}{\# \text{ of flips}}$$

A difference

between frequentists

and Bayesians

is that

Bayesians treat

parameters

(like  $\theta$ )

as random

variables.

$\rightarrow$  MAP: Use Bayes' Rule to define a distribution over parameters  $\theta$ :

$$p(\theta | X^1, \dots, X^n) = \frac{\overset{\text{data}}{\downarrow \text{likelihood}} p(X^1, \dots, X^n | \theta) \overset{\text{prior}}{\downarrow} p(\theta)}{p(X^1, \dots, X^n)}$$

$\leftarrow$  posterior distribution

Note that  $p(\theta | X^1, \dots, X^n)$  is a distribution over possible parameters  $\theta$ .

MAP estimate chooses the  $\theta$  which maximizes

$p(\theta | X^1, \dots, X^n)$ :

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | X^1, \dots, X^n)$$

$$= \arg \max_{\theta} p(X^1, \dots, X^n | \theta) p(\theta)$$

Issue: we must choose the prior  $p(\theta)$

Usually we choose conjugate priors which make MAP estimation easy. Conjugate

of Bernoulli is Beta

$$\text{Beta}(\beta_H, \beta_T) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)}$$

So:  $\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \theta^{\beta_H-1 + \sum x_i} (1-\theta)^{\beta_T-1 + \sum (1-x_i)}$

And we saw that  $n$

$$\hat{\theta}_{\text{MAP}} = \frac{(\beta_H-1) + \sum x_i}{n + \beta_H + \beta_T - 2}$$

}  $\beta_H$  and  $\beta_T$  act like "virtual" flips!

Bayesian Inference in Pictures:



In this picture, the dark line is the final posterior distribution. The light lines show what the posterior looks like as we observe more data points. Note that these distributions converge to the final posterior.

A common choice for  $\beta_H, \beta_T$  is to set  $\beta_H = \beta_T = \alpha$ , in which case:

$$\hat{\theta}_{\text{MAP}} = \frac{(\alpha-1) + \sum x_i}{n + 2(\alpha-1)}$$

} This looks like something on your problem set... (Hint: set  $\alpha=2$ )

## Multinomial & Dirichlet Distributions

The multinomial distribution is like the Bernoulli distribution for more than 2 values.

Bernoulli - flipping a coin (2 outcomes)  
Multinomial - rolling a die (6 outcomes)

Let's say we're rolling an unfair die. How many parameters do we need to describe the distribution? Answer: 5.

In general:  $n$  values means  $n-1$  parameters.

$\theta$  is our

parameter vector for the multinomial  $\rightarrow \theta_i$  is the probability of observing the  $i$ th outcome. (E.g.,  $\theta_1$  = probability of rolling 1 on the die)  
Constraint:  $\sum_{i=1}^n \theta_i = 1$

Conjugate of multinomial distribution is the Dirichlet distribution  
(Like conjugate of Bernoulli is Beta)

Note:

$\alpha$  is a vector  $\rightarrow \text{Dirichlet}(\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^n x_i^{\alpha_i - 1}$

As in Beta case, the  $\alpha_i$  values act like virtual observations!

$$\hat{\theta}_{\text{MLE}} = \frac{\text{Count}(X=i)}{(\text{Total \# of observations})}$$

$$\hat{\theta}_{\text{MAP}} = \frac{(\alpha_i - 1) + \text{Count}(X=i)}{(\text{total \# of observations}) + \sum_{i=1}^n (\alpha_i - 1)}$$

## Naïve Bayes Classifier

Let  $X = (X_1, \dots, X_d)^T$  be a feature vector,  $Y \in \{0, 1\}$  be a label.

Recall: classification means predicting  $Y$  from  $X$ , given training data

$$D = \{(X^1, Y^1), (X^2, Y^2), \dots, (X^n, Y^n)\}$$

Naïve Bayes Assumption: Probability of each feature  $X_i$  is conditionally independent given  $Y$ . Therefore, we can factor distribution  $P(Y, X)$ :

$$P(X, Y) = \underset{\uparrow}{P(Y)} \prod_{i=1}^d \underset{\uparrow}{P(X_i | Y)}$$

class prior      class conditional distributions

Class-conditional distributions  $P(X_i | Y)$  may have many different forms. E.g.,

$X_i$  continuous - Gaussian

$X_i$  boolean - Bernoulli

$X_i$  categorical - multinomial

(MLE) Parameter Estimation:  $\hat{\pi} = \frac{1}{n} \sum_{j=1}^n Y^j$  (Let  $P(Y=1) = \pi$ )

Estimate  $P(X_i | Y=y)$  based on the distribution.

Example:  $P(X_i | Y)$  is Gaussian with unit variance ( $\sigma^2=1$ )

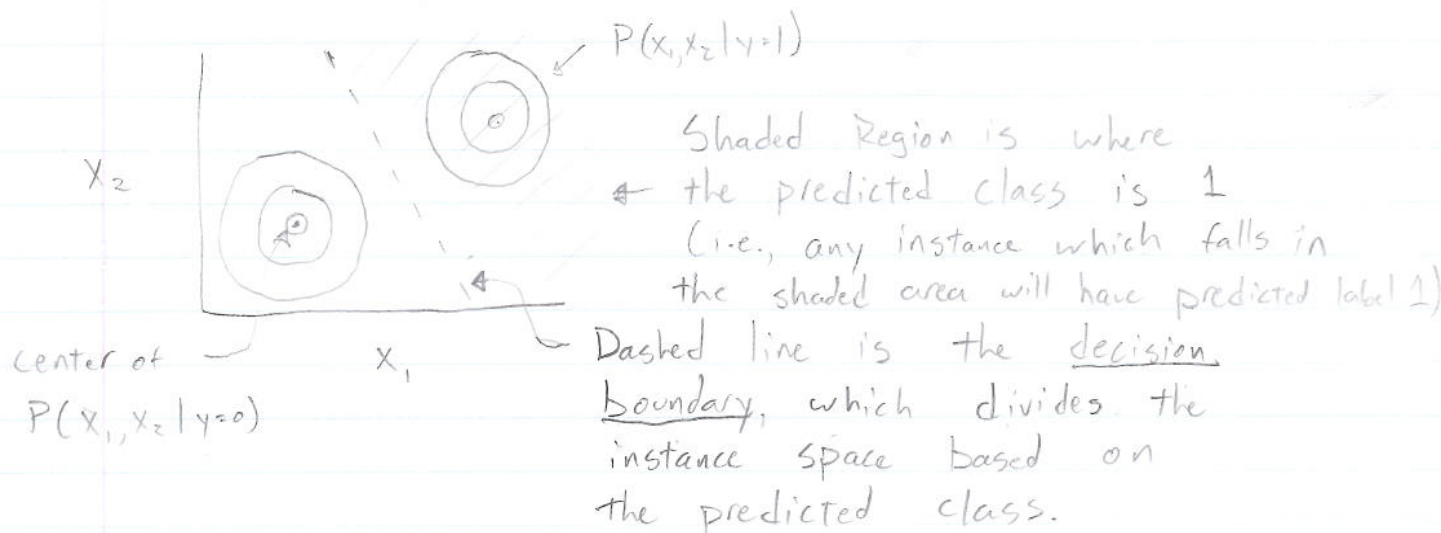
$$P(X_i | Y=y) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(X_i - \mu_{iy})^2}{2} \right\}$$

$$\hat{\mu}_{iy} = \frac{\sum_{j=1}^n \mathbb{1}(Y^j=y) X_i^j}{\sum_{j=1}^n \mathbb{1}(Y^j=y)}$$

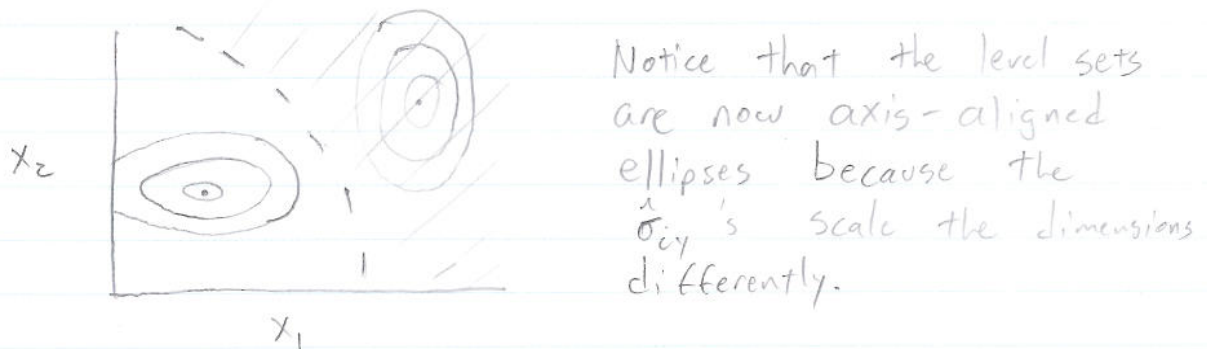


## Decision Boundaries

Consider the GNB classifier from the previous example. Assume we have two features  $x_1, x_2$ . The classifier will estimate one spherical Gaussian per class. Consider the level sets of these Gaussians:



As another example, let's say we now estimate a standard deviation  $\hat{\sigma}_{iy}$  for each class conditional distribution:



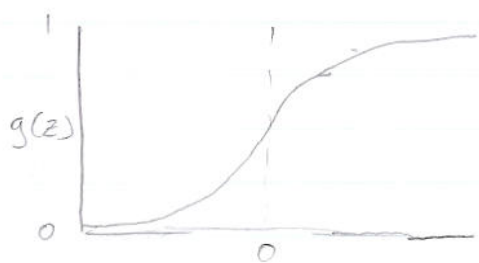
In our first example, our decision boundary is linear, while in the second example it is a curve (actually a quadratic curve). The shape of a classifier's decision boundary is an important characteristic of the classifier.

## Logistic Regression (LR):

LR models  $P(Y|X)$  instead of  $P(X, Y)$ .

$$P(Y=1|X) = \frac{\exp \{ w_0 + \sum_{i=1}^n w_i X_i \}}{1 + \exp \{ w_0 + \sum_{i=1}^n w_i X_i \}}$$

N.B.  $g(z) = \frac{\exp \{ z \}}{1 + \exp \{ z \}}$  is the logistic function.



Logistic function has a value between 0 and 1.  
 $g(z) \rightarrow 0$  as  $z \rightarrow -\infty$   
 $g(z) \rightarrow 1$  as  $z \rightarrow \infty$   
 $g(0) = 0.5$

LR's decision boundary is linear:

$$P(Y=1|X) = 0.5 = \frac{\exp \{ w_0 + \sum_{i=1}^n w_i X_i \}}{1 + \exp \{ w_0 + \sum_{i=1}^n w_i X_i \}}$$
$$\Rightarrow 0 = w_0 + \sum_{i=1}^n w_i X_i$$

linear equation for the decision boundary.

Where does the LR form come from?

Consider our Gaussian Naive Bayes classifier with unit variance for each class:

$$P(Y=0|X) = \frac{P(Y=0) P(X|Y=0)}{P(Y=0) P(X|Y=0) + P(Y=1) P(X|Y=1)}$$
$$= \frac{1}{1 + \frac{P(Y=1) P(X|Y=1)}{P(Y=0) P(X|Y=0)}}$$

Consider the second term in the denominator:

$$\frac{P(Y=1) P(X|Y=1)}{P(Y=0) P(X|Y=0)} = \exp \left\{ \ln \left( \frac{P(Y=1)}{P(Y=0)} \frac{P(X|Y=1)}{P(X|Y=0)} \right) \right\}$$

Our NB assumption says:  $P(X|Y) = \prod_{i=1}^n P(x_i|Y)$   
so:

$$\frac{P(Y=1) P(X|Y=1)}{P(Y=0) P(X|Y=0)} = \exp \left\{ \ln \left( \frac{P(Y=1)}{P(Y=0)} \right) + \sum_{i=1}^n \ln \left( \frac{P(x_i|Y=1)}{P(x_i|Y=0)} \right) \right\}$$

Consider the terms in the sum:

$$\ln \left( \frac{P(x_i|Y=1)}{P(x_i|Y=0)} \right) = \ln \left( \frac{\frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (x_i - \hat{\mu}_{i1})^2 \right\}}{\frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} (x_i - \hat{\mu}_{i0})^2 \right\}} \right)$$

$$= \ln \exp \left\{ -\frac{1}{2} (x_i - \hat{\mu}_{i0})^2 + \frac{1}{2} (x_i - \hat{\mu}_{i1})^2 \right\}$$

$$= \frac{1}{2} \left( (x_i^2 - 2x_i \hat{\mu}_{i0} + \hat{\mu}_{i0}^2) - (x_i^2 - 2x_i \hat{\mu}_{i1} + \hat{\mu}_{i1}^2) \right)$$

$$= (\hat{\mu}_{i1} - \hat{\mu}_{i0}) x_i + \frac{\hat{\mu}_{i0}^2 - \hat{\mu}_{i1}^2}{2}$$

Note that the above equation is linear  
in  $x_i$ ! So, if we set:

$$w_i = \hat{\mu}_{i1} - \hat{\mu}_{i0}$$

$$w_0 = \ln \left( \frac{P(Y=1)}{P(Y=0)} \right) + \sum_{i=1}^n \left( \frac{\hat{\mu}_{i0}^2 - \hat{\mu}_{i1}^2}{2} \right)$$

Then:

$$P(Y=0|X) = \frac{1}{1 + \exp \left\{ w_0 + \sum_{i=1}^n w_i x_i \right\}}$$



WARNING: This derivation does not work for all types of NB classifiers.  
 Example: our earlier GNB classifier with estimated per-class variances.  
 (Recall the decision boundary was a quadratic curve, while LR's decision boundary is always linear!)

### Multiclass Logistic Regression

What if we have  $K$  labels,  $y_1, \dots, y_K$  instead of just 2 labels?

$$P(Y=y_j | X) = \frac{\exp \left\{ w_{0j} + \sum_{i=1}^d w_{ij} X_i \right\}}{1 + \sum_{l=1}^K \exp \left\{ w_{0l} + \sum_{i=1}^d w_{il} X_i \right\}}$$

$$P(Y=y_K | X) = \frac{1}{1 + \sum_{l=1}^K \exp \left\{ w_{0l} + \sum_{i=1}^d w_{il} X_i \right\}}$$

Decision Rule: choose  $y_K$  with maximum probability.

$$\begin{aligned} \hat{y} &= \arg \max_y P(Y=y | X) \\ &= \arg \max_j \left( 0, w_{0j} + \sum_{i=1}^d w_{ij} X_i \right) \end{aligned}$$

Ex. 3 classes decision boundary

