

# Lecture 8: Classification I - Naive Bayes

## Modeling Social Data, Spring 2017

### Columbia University

Vedant Dharnidharka

March 10, 2017

## 1 Classification

Classification problem deals with categorization of given examples into a set of categories.

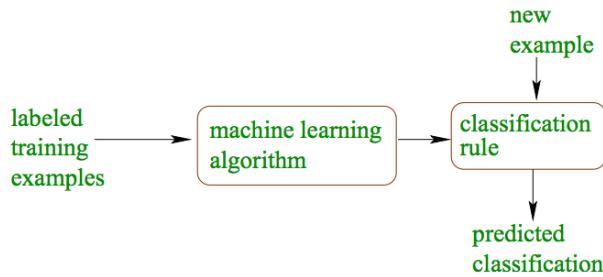


Figure 1: Classification

*Input:* As with regression, in a classification problem we start with measurements  $x_1, x_2, \dots, x_n$  in an input space  $X$ .

*Output:*

1.  $Y = \{-1, +1\}$  or  $\{0, 1\}$  is called binary classification.
2.  $Y = \{1, \dots, K\}$  is called multiclass classification.

Instead of a real-valued response, classification assigns  $x$  to a category. For pair  $(x, y)$ ,  $y$  is the class of  $x$ .

## 2 Naive Bayes Classifier

The Naive Bayes classifier is based on the Bayes rule of conditional probability. It makes use of all the attributes contained in the data, and analyses them individually as though they are equally important and independent of each other. In some cases it is also seen that Naive Bayes outperforms many other comparatively complex algorithms.

### 2.1 Bayes Theorem

A theorem describing how the conditional probability of each of a set of possible causes for a given observed outcome can be computed from knowledge of the probability of each cause and the conditional probability of the outcome of each cause.

Let  $(\Omega, P)$  be a probability space. ( $\Omega$  is the sample space;  $P$  is the probability distribution.)

For any event  $A, B$  in  $\Omega$

Formula:

$$P(A|B) = P(A) \frac{P(B|A)}{P(B)}, \quad (1)$$

Let  $E, H_0, H_1$  in  $\Omega$ .

Conditioned on  $E$ , of  $H_0$  and  $H_1$ , we find which one is more probable.

Compare  $P(H_0)P(E|H_0)$  to  $P(H_1)P(E|H_1)$

### 2.1.1 Example

Suppose result of test for genetic disease is correct with probability 95%, and suppose the disease is rare: any given person has disease with probability 1%.

**Question:** If test comes back positive for disease, is it more likely that you have disease or do not?

$E$  = test comes back positive for disease

$H_0$  = do not have disease

$H_1$  = have disease

$P(E | H_0) = 0.05$

$P(E | H_1) = 0.95$

$P(H_0) = 0.99$

$P(H_1) = 0.01$

Want to compare  $P(H_0 | E)$  to  $P(H_1 | E)$ , so compare

$P(H_0) \cdot P(E | H_0) = 0.99 \cdot 0.05$  and  $P(H_1) \cdot P(E | H_1) = 0.01 \cdot 0.95$ .

Figure 2: Bayes Theorem Example(Source: <http://www.cs.columbia.edu/~djhsu/coms4771-f16/lectures/slides-generative.4up.pdf>)

## 2.2 Maximum Likelihood Estimate

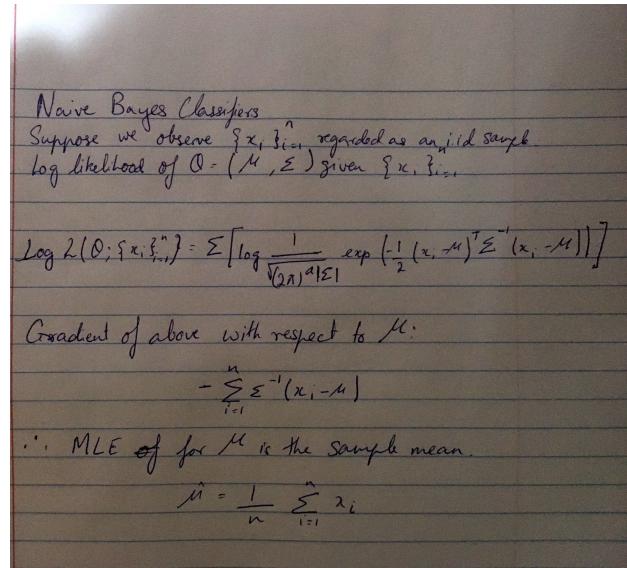


Figure 3: MLE for Naive Bayes

### 2.2.1 Advantages

- Simple, many variations
- Can leverage domain knowledge of class conditionals
- Can be very efficient when K is large.

### 2.2.2 Disadvantages

- Classifier relies on formula (via Bayes rule) that assumes data comes from estimated distribution, which is generally not true.
- Modeling P away from decision boundary between classes is a wasted effort.

## 3 Logistic Regression

In logistic regression probability of the response taking a particular value is modeled based on combination of values taken by the predictors.

Logistic Regression

Probability Distribution over  $\mathcal{X} \times \{0, 1\}$ ; let  $(x, y) \sim P$

Think of  $P$  comprised of 2 parts -

- (1) Marginal distribution of  $X$
- (2) Conditional distribution of  $Y$  given  $X = x$

$$\eta(x) = P(Y=1|X=x)$$

Bayes Classifier is

$$f^*(x) = \begin{cases} 0 & \text{if } \eta(x) \leq 1/2 \\ 1 & \text{if } \eta(x) > 1/2 \end{cases}$$

The log odds function at  $x$

$$x \rightarrow \log \frac{\eta(x)}{1-\eta(x)} \in [-\infty, +\infty]$$

Logistic Regression -

Feature space is  $\mathcal{X} \subseteq \mathbb{R}^d$

Statistical Model for  $Y|X=x$  for each  $x \in \mathcal{X}$

$$P = \{P_{(\beta_0, \beta)} : \beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^d\},$$

where.

$$\eta(\beta_0, \beta)(x) = P_{(\beta_0, \beta)}(Y=1|X=x) = \text{logistic}(\beta_0 + \langle \beta, x \rangle)$$

We know  $\text{logistic}(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$

Log Odds function of  $P_{(\beta_0, \beta)}$  is

$$x \rightarrow \log \frac{\eta(\beta_0, \beta)(x)}{1-\eta(\beta_0, \beta)(x)}$$

$$= \beta_0 + \langle \beta, x \rangle \text{ which is affine.}$$

Figure 4: Logistic Regression-Linear Classifier

### 3.1 Boosting

Boosting is nothing but using a learning algorithm that provides rough rules-of-thumb to construct a very accurate predictor.

Motivation: Easy to construct classification rules that are correct more-often-than-not (e.g., If over 5 percent of the

e-mail characters are dollar signs, then its spam.), but seems hard to find a single rule that is almost always correct.

Assumption: Availability of a base or weak learning algorithm which produces a weak classifier. Boosting improves the performance of the weak learning algorithm while treating it as a black box . Weak classifiers are not entirely trivial which means that the error rates are at least a better than a classifier whose every prediction is a random guess. The weak classifiers can be moderately inaccurate, but not as bad as random guessing.

When they are added, they are typically weighted in some way that is usually related to the weak learners' accuracy. After a weak learner is added, the data are assigned weights again. Therefore when an entry is wrongly classified, the classifier which correctly classifies has its weight increased.