

Chapter 9 Evidence and Probabilities

A broadly applicable framework both for **evaluating the evidence**, and for combining it to **estimate the resulting likelihood**

Combining Evidence Probabilistically

Bayes Rule for classification

$$p(C = c \mid \mathbf{E}) = \frac{p(\mathbf{E} \mid C = c) \cdot p(C = c)}{p(\mathbf{E})}$$

- Left hand side: probability of target variable C taking class of interest c, after taking the evidence E into account
- P(C=c) is the 'prior' probability of the class, could come from several places:
 - "subjective prior", meaning it is a belief of particular decision maker based on her knowledge and experience
 - 'prior' belief based on previous applications of Bayes Rule with other evidence
 - Unconditional probability inferred from data, usually from class prior / base rate of c, i.e. percentage of all examples in data that are of class c
- p(E | C = c) is the likelihood of seeing evidence E, given C=c.
 - Can be calculated from the percentage of examples of class c, which have feature vector of E
- P(E) is likelihood of evidence E in dataset
 - Can be calculated from percentage occurrence of E from among all examples

Conditional Independence and Naïve Bayes

Assume attributes are conditional independent, given the class
(So each evidence class can be decomposed)

Bayes Equation:

$$p(c \mid \mathbf{E}) = \frac{p(e_1 \mid c) \cdot p(e_2 \mid c) \cdots p(e_k \mid c) \cdot p(c)}{p(\mathbf{E})}$$

Naïve Bayes Classifier: estimate the probability that the example belongs to each class and reports the class with highest probability.

Advantages and Disadvantages of Naïve Bayes

- Efficient in terms of storage space and computation time
 - Training consists of only storage of counts of classes and feature occurrences
- Violation of independence assumption tend not to hurt classification performance.
 - E.g. double-counting two strongly correlated evidence would result in a more extreme classifier, but as long as the evidence shows right direction, the classifier would still work well
 - But need to be careful when we care the actual value of probability
E.g. when combining with cost-benefits
- Incremental learner
 - Update model one training example at a time. Doesn't need to re-process all previous training examples when new training data becomes available

A Model of Evidence Lift

Lift

- measures how much more prevalent the positive class is in the selected subpopulation over the prevalence in the population as a whole

$$\text{lift}_c(x) = \frac{p(x \mid c)}{p(x)}$$

If we assume full feature independence (Naïve-Naïve Bayes)

$$p(c \mid \mathbf{E}) = \frac{p(e_1 \mid c) \cdot p(e_2 \mid c) \cdots p(e_k \mid c) \cdot p(c)}{p(e_1) \cdot p(e_2) \cdots p(e_k)}$$

$$p(C = c \mid \mathbf{E}) = p(C = c) \cdot \text{lift}_c(e_1) \cdot \text{lift}_c(e_2) \cdots$$

An Example from Facebook "Likes"

Like	Lift	Like	Lift
<i>Lord Of The Rings</i>	1.69	Wikileaks	1.59
One Manga	1.57	Beethoven	1.52
Science	1.49	NPR	1.48
Psychology	1.46	<i>Spirited Away</i>	1.45
<i>The Big Bang Theory</i>	1.43	Running	1.41
Paulo Coelho	1.41	Roger Federer	1.40
<i>The Daily Show</i>	1.40	<i>Star Trek</i>	1.39
<i>Lost</i>	1.39	Philosophy	1.38
<i>Lie to Me</i>	1.37	<i>The Onion</i>	1.37
<i>How I Met Your Mother</i>	1.35	<i>The Colbert Report</i>	1.35
<i>Doctor Who</i>	1.34	<i>Star Trek</i>	1.32
<i>Howl's Moving Castle</i>	1.31	Sheldon Cooper	1.30
<i>Tron</i>	1.28	<i>Fight Club</i>	1.26
Angry Birds	1.25	<i>Inception</i>	1.25
<i>The Godfather</i>	1.23	<i>Weeds</i>	1.22

- The lift for High-IQ class
- (Proportion of "likes" among High-IQ class) / (Proportion of "likes" among all users)
- The base rate of High-IQ class is 0.14. If I give "like" to nothing, then probability of being high-IQ is just 0.14
- If I give like to Sheldon Cooper, then probability of High-IQ increases to 0.14 * 1.3 = 0.182
- If I have 3 likes, Sheldon Cooper, Star Trek and Lord of Rings, then my probability of being High-IQ is 0.14 * 1.3 * 1.39 * 1.69 = 0.4275