

EECS 391: Introduction to AI

Soumya Ray

Website: http://vorlon.case.edu/~sray/eecs391_sp12/index.html

Email: sray@case.edu

Office: Olin 516

Office hours: Tue 2:30-4:00 or by appointment

Binary Classification

	Has-fur?	Long-Teeth?	Scary?	<i>Lion?</i>
Animal ₁	Yes	No	No	No
Animal ₂	No	Yes	Yes	No
Animal ₃	Yes	Yes	Yes	Yes

Attribute-value representation
(**X**)

Binary Class Label (*y*)
(assigned by target concept/
teacher)

Naïve Bayes

- Simple probabilistic classifier for discrete data

$$p_{\mathbf{X},Y}(\mathbf{x}, y) = p(\mathbf{X} = \mathbf{x} | Y = y)p(Y = y)$$
$$= p(x_1, \dots, x_n | Y = y)p(Y = y)$$

Naïve Bayes assumption:
Attributes are conditionally independent given the class

$$\Rightarrow \prod_i p(X_i = x_i | Y = y)p(Y = y)$$

Naïve Bayes **parameters**: Instead of storing probabilities for each atomic event, we will only store these conditional probabilities and use this formula to calculate the probability for an atomic event (example).

Probability Estimation for Naïve Bayes

- Given a set of observations:

	Has-fur?	Long-Teeth?	Scary?	<i>Lion?</i>
Animal ₁	Yes	No	No	No
Animal ₂	No	Yes	Yes	No
Animal ₃	Yes	Yes	Yes	Yes

- Estimate** parameters $p(X_i=x_i/Y=y)$ and $p(Y=y)$
- We will use a method called “Maximum Likelihood Estimation”

Bayes Rule for Concept Learning

- Suppose we are given a set of examples D and we are considering a hypotheses space H
- Each hypothesis could have “caused” the data we observe
- The **posterior probability** of any hypothesis h in H is given by Bayes Rule:

$$\boxed{\text{Posterior}} \Pr(h | D) = \frac{\boxed{\text{Likelihood}} \Pr(D | h) \boxed{\text{Prior}} \Pr(h)}{\boxed{\text{Evidence}} \Pr(D)}$$

MAP Hypothesis

- Given: training sample D and hypothesis class H
- Do: Return the most probable hypothesis given the data---the **maximum a posteriori (MAP)** hypothesis

$$\begin{aligned}h_{MAP} &= \arg \max_{h \in H} \Pr(h | D) \\ &= \arg \max_{h \in H} \frac{\Pr(D | h) \Pr(h)}{\Pr(D)} \\ &= \arg \max_{h \in H} \Pr(D | h) \Pr(h)\end{aligned}$$

ML Hypothesis

- If every hypothesis in H has equal prior probability, only the first term matters
- This gives the **maximum likelihood (ML)** hypothesis

$$h_{ML} = \arg \max_{h \in H} \Pr(D | h)$$

Maximum Likelihood Estimation

- For naïve Bayes, a hypothesis is the vector of parameters, one for each of $p(X_i=x_i/Y=y)$ and $P(Y=y)$
- Assume X_i is binary (0/1)
 - Then $p(X_i=1/Y=1)$ is a parameter, call it θ_{i1}
 - There's another parameter for $p(X_i=1/Y=0)$, θ_{i0}
 - Finally there's a parameter for $p(Y=y)$, θ_y

Maximum Likelihood Estimation

$$h_{ML} = \arg \max_{h \in H} p(D | h)$$

$$p(D | h) = p(\{\mathbf{x}_d, y_d\}_{d=1 \dots m} | \{\theta_{i0}, \theta_{i1}\}_{i=1 \dots n}, \theta_y)$$

$$= \prod_{d=1}^m p(\mathbf{x}_d, y_d | \{\theta_{i0}, \theta_{i1}\}_{i=1 \dots n}, \theta_y)$$

$$= \prod_{d=1}^m \prod_{i=1}^n p(X_i = x_i | Y = y_d; \{\theta_{i0}, \theta_{i1}\}) \theta_{y_d}$$

	Has-fur? (f1)	Long-Teeth? (f2)	Scary? (f3)	Lion? (Y)
Animal ₁	Yes	No	No	No
Animal ₂	No	Yes	Yes	No
Animal ₃	Yes	Yes	Yes	Yes

$$\begin{aligned}
 p(D | h) &= \prod_{d=1}^m \prod_{i=1}^n p(X_i = x_i | Y = y_d; \{\theta_{i0}, \theta_{i1}\}) \theta_{y_d} \\
 &= \left[\theta_{10} (1 - \theta_{20}) (1 - \theta_{30}) \theta_{y_0} \right] \times \\
 &\quad \left[(1 - \theta_{10}) \theta_{20} \theta_{30} \theta_{y_0} \right] \times \\
 &\quad \left[\theta_{11} \theta_{21} \theta_{31} \theta_{y_1} \right]
 \end{aligned}$$

Let pos be the number of examples with $Y=1$ and d_i of those have $f_i=Yes$

$$p(D | h) = \prod_{d=1}^m \prod_{i=1}^n p(X_i = x_i | Y = y_d; \{\theta_{i0}, \theta_{i1}\}) \theta_{y_d}$$

$$\text{For } Y = 1, p(D | h) = \prod_{i=1}^n \theta_{i1}^{d_i} (1 - \theta_{i1})^{pos - d_i} \theta_{y_1}^{pos}$$

Number of examples with $Y=1$

$$\hat{\theta}_{i1} = \arg \max_{\theta_{i1}} \prod_{i=1}^n \theta_{i1}^{d_i} (1 - \theta_{i1})^{pos - d_i} \theta_{y_1}^{pos} = L(\theta_{i1})$$

$$LL(\theta_{i1}) = pos \log \theta_{y_1} + \sum_{i=1}^n d_i \log \theta_{i1} + (pos - d_i) \log(1 - \theta_{i1})$$

$$\frac{\partial LL}{\partial \theta_{i1}} = \frac{d_i}{\theta_{i1}} - \frac{(pos - d_i)}{(1 - \theta_{i1})} = 0$$

$$\text{or } d_i - d_i \theta_{i1} = pos \cdot \theta_{i1} - d_i \theta_{i1}$$

$$\text{or } d_i = pos \cdot \theta_{i1}$$

$$\text{or } \theta_{i1} = \frac{d_i}{pos}$$

Fraction of observed $Y=1$ examples where $X_i=1$!

Let pos be the number of examples with $Y=1$ and d_i of those have $f_i=Yes$

Naïve Bayes Parameter MLEs

$$\hat{p}(X_i = 1 | Y = 1) = \frac{\# \text{ observed examples with } X_i = 1 \text{ and } Y = 1}{\# \text{ observed examples with } Y = 1}$$

$$p(X_i = 1 | Y = 1) = \frac{p(X_i = 1, Y = 1)}{p(Y = 1)}$$

$$\hat{p}(Y = 1) = \frac{\# \text{ observed examples with } Y = 1}{\# \text{ observed examples}}$$

Example

	Has-fur?	Long-Teeth?	Scary?	<i>Lion?</i>
Animal ₁	Yes	No	No	No
Animal ₂	No	Yes	Yes	No
Animal ₃	Yes	Yes	Yes	Yes

$p(\text{Has-fur}=\text{Yes} \mid \text{Lion})=?$,

$p(\text{Has-fur}=\text{Yes} \mid \text{Not-Lion})=?$

$p(\text{Long-Teeth}=\text{Yes} \mid \text{Lion})=?$, $p(\text{Long-Teeth}=\text{Yes} \mid \text{Not-Lion})=?$

$p(\text{Scary}=\text{Yes} \mid \text{Lion})=?$,

$p(\text{Scary}=\text{Yes} \mid \text{Not-Lion})=?$

$p(\text{Lion})=?$

Smoothing probability estimates

- What happens if a certain value for a variable is not in our set of examples, for a certain class?
 - Suppose we're trying to classify lions and we've never seen a lion cub, so $\hat{p}(Scary = false | Lion) = 0$
 - When we see a cub, its probability of being a lion will be zero by our Naïve Bayes formula, even if it has long teeth and fur
 - It's a good idea to “smooth” our probability estimates to avoid this

m -Estimates

$$p(X_i = x_i | Y = y) = \frac{(\text{\# examples with } X_i = x_i \text{ and } Y = y) + mp}{(\text{\# examples with } Y = y) + m}$$

- p is our prior estimate of the probability
- m is called “Equivalent Sample Size” which determines the importance of p relative to the observations
- If variable has v values, the specific case of $m=v$, $p=1/v$ is called **Laplace smoothing**

Email Spam Filtering with Naïve Bayes

- Used very successfully to categorize documents
 - Is this document about “sports” or “finance”?
 - Is this email “spam” or “ham”?
- Given a vocabulary, each attribute X_i is the presence/absence of word i in the document
 - Ignores word order
 - “Bag-of-words” approach

Email Spam Filtering with Naïve Bayes

- Smoothed parameter estimates

$$p(\text{word}_k \text{ present} \mid Y = \text{spam}) = \boxed{\text{Bernoulli distribution}}$$
$$\frac{(\#\text{emails with } \text{word}_k \text{ present and } Y = \text{spam}) + 1}{(\#\text{emails with } Y = \text{spam}) + 2}$$

- Variations used by most major email clients/commercial spam filters e.g. DSPAM, SpamAssassin, SpamBayes, Bogofilter, ASSP
 - Read article on “Bayesian spam filtering” in wikipedia

Summary

- We learned about:
 - Maximum likelihood estimation
 - Smoothing parameters
- Next: Artificial Neural Networks