

COMP3308/COMP3608 Introduction to Artificial Intelligence

Weeks 6 Tutorial exercises Naïve Bayes. Evaluating Classifiers.

Exercise 1. Naïve Bayes (Homework)

Suppose you want to recognize good and bad items produced by your company. You are able to measure two properties of each item (P1 and P2) and express them with Boolean values. You randomly grab several items and test if they are good or bad, obtaining the following results:

P1	P2	result
Y	Y	good
Y	N	bad
N	N	good
N	Y	bad
Y	N	good
N	N	good

Use Naïve Bayes to predict the class, *good* or *bad*, of the following new item: P1=N, P2=Y. If there are ties, make a random choice.

Solution:

E is P1=N and P2=Y; E₁ is P1=N, E₂ is P2=Y

We need to compute P(good|E) and P(bad|E) and compare them.

$$P(\text{good} | E) = \frac{P(E_1 | \text{good}) P(E_2 | \text{good}) P(\text{good})}{P(E)}$$

$$P(\text{good}) = 4/6 = 2/3$$

$$P(E_1 | \text{good}) = P(P1=N | \text{good}) = 2/4 = 1/2$$

$$P(E_2 | \text{good}) = P(P2=Y | \text{good}) = 1/4$$

$$P(\text{bad}) = 2/6 = 1/3$$

$$P(E_1 | \text{bad}) = P(P1=N | \text{bad}) = 1/2$$

$$P(E_2 | \text{bad}) = P(P2=Y | \text{bad}) = 1/2$$

$$P(\text{good} | E) = \frac{\frac{1}{2} \frac{1}{4} \frac{2}{3}}{\frac{1}{12}} = \frac{1}{12}$$

$$P(\text{bad} | E) = \frac{\frac{1}{2} \frac{1}{2} \frac{1}{3}}{\frac{1}{12}} = \frac{1}{12}$$

The two probabilities are the same. To resolve the tie we randomly choose between the 2 classes => e.g. class *good*.

Exercise 2. Naïve Bayes

Why is the Naïve Bayesian classification called “naïve”?

Answer: Naïve Bayes assumes that the values of the attributes are independent of each other and that all attributes are equally important. These assumptions are unrealistic, that’s why it is called “Naïve”.

Exercise 3. Applying Naïve Bayes to data with both numerical and nominal attributes

Given is the training data in the table below (the *weather* data with some numerical attributes, *play* is the class). Predict the class of the following new example using the Naïve Bayes classification: outlook=overcast, temperature=60, humidity=62, windy=false.

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
overcast	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

Solution:

First, we need to calculate the mean μ and standard deviation σ values for the numerical attributes, for each class (yes and no) separately.

$$\mu = \frac{\sum_{i=1}^n X_i}{n} \quad \sigma^2 = \frac{\sum_{i=1}^n (X_i - \mu)^2}{n-1}$$

$X_i, i=1..n$ – the i -th measurement, n -number of measurements

You can use Excel, Python, Matlab etc. to calculate these values.

$\mu_{\text{temp_yes}}=73, \sigma_{\text{temp_yes}}=6.2; \quad \mu_{\text{temp_no}}=74.6, \sigma_{\text{temp_no}}=8.0$
 $\mu_{\text{hum_yes}}=79.1, \sigma_{\text{hum_yes}}=10.2; \quad \mu_{\text{hum_no}}=86.2, \sigma_{\text{hum_no}}=9.7$

Second, to calculate $f(\text{temperature}=60|\text{yes})$, $f(\text{temperature}=60|\text{no})$, $f(\text{humidity}=62|\text{yes})$ and $f(\text{humidity}=62|\text{no})$ using the probability density function for normal distribution:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$f(\text{temperature} = 60|\text{yes}) = \frac{1}{6.2\sqrt{2\pi}} e^{-\frac{(60-73)^2}{2*6.2^2}} = 0.0071$$

$$f(\text{temperature} = 60|\text{no}) = \frac{1}{8\sqrt{2\pi}} e^{-\frac{(60-74.6)^2}{2*8^2}} = 0.0094$$

$$f(\text{humidity} = 62|\text{yes}) = \frac{1}{10.2\sqrt{2\pi}} e^{-\frac{(62-79.1)^2}{2*10.2^2}} = 0.0096$$

$$f(\text{humidity} = 62|\text{no}) = \frac{1}{9.7\sqrt{2\pi}} e^{-\frac{(62-86.2)^2}{2*9.7^2}} = 0.0018$$

Third, we can calculate the probabilities for the nominal attributes:

$$P(\text{yes})=9/14=0.643$$

$$P(\text{no})=5/14=0.357$$

$$P(\text{outlook}=\text{overcast}|\text{yes})=4/9=0.444$$

$$P(\text{outlook}=\text{overcast}|\text{no})=1/5=0.2$$

$$P(\text{windy}=\text{false}|\text{yes})=6/9=0.667$$

$$P(\text{windy}=\text{false}|\text{no})=2/5=0.4$$

Fourth, we can calculate the final probabilities:

$$P(\text{yes} | E) = \frac{0.444 * 0.0071 * 0.0096 * 0.667 * 0.643}{P(E)} = \frac{12.97 * 10^{-6}}{P(E)}$$

$$P(\text{no}|E) = \frac{0.2 * 0.0094 * 0.0018 * 0.4 * 0.357}{P(E)} = \frac{4 * 10^{-7}}{P(E)}$$

Therefore, the Naïve Bayes classifier predicts $\text{play}=\text{yes}$ for the new example.

Exercise 4. Bayes Theorem (Advanced only)

Suppose that the fraction of undergraduate students who smoke is 15% and the fraction of graduate students who smoke is 23%. If 1/5 of the University students are graduate students and the rest are undergraduates, what's the probability that a student who smokes is a graduate student?

Hint: Use the Bayes Theorem; you will need to calculate the denominator using the *law of total probability*, see its Wikipedia description.

Solution:

Suppose that:

X represents if the student smokes: {smoker, non-smoker} or abbreviated {S, NS} and

Y represents the type of student: {undergraduate, graduate} or abbreviated {UG, G}

Given: $P(G)=1/5=0.2$, $P(UG)=4/5=0.8$, $P(S|UG)=0.15$, $P(S|G)=0.23$

$P(G|S)=?$

$$P(G | S) = \frac{P(S | G) P(G)}{P(S)}$$

To calculate $P(S)$ we will use the law of total probability. If $\{Y_1, Y_2, \dots, Y_k\}$ is the set of mutually exclusive and exhaustive outcomes of Y, then:

$$P(X) = \sum_{i=1}^k P(X, Y_i) = \sum_{i=1}^k P(X | Y_i) P(Y_i)$$

$$\Rightarrow P(S) = P(S|G)P(G) + P(S|UG)P(UG) = 0.23*(1/5) + 0.15*(4/5) = 0.166$$

$$\Rightarrow P(G|S) = (0.23*0.2)/0.166 = 0.277$$

Exercise 5. Using Weka – Comparing Classifiers

1. Load the iris dataset
 2. Choose “Percentage split” mode for evaluation: 66% training set, 33% testing set
 3. Run the Naïve Bayes and review Weka’s output
 4. For comparison, also run k-nearest neighbor with k=1 and 3 (IB1 and IBk), OneR and ZeroR. Which is the most accurate classifier?
 5. Change the test mode to “Cross validation”. Apply 10-fold cross validation instead of percentage split as evaluation mode and compare the classifiers.
 - Which classifier produced the most accurate classification?
 - Which evaluation strategy (percentage split or 10-fold cross validation) produced better results?
 - Which evaluation strategy, percentage split or cross validation, is more statistically reliable and why?
 6. Apply leave-one-out cross validation. Tip: You need to specify the number of folds in the WEKA’s cross validation box.
 7. Check the confusion matrix printed by WEKA for one of the classifiers, e.g. Naïve Bayes, and verify the accuracy, recall, precision and F1 measure. Note: Weka shows recall, precision and F1 for each class separately.
- Answer:** Leave-one-out is n-fold cross validation where n is the number of examples. Thus, the number of folds should be set to 150 for iris data.