Classification – Naïve Bayes

UROŠ KRČADINAC EMAIL: <u>uros@krcadinac.com</u> URL: <u>http://krcadinac.com</u>

Bayes rule

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

- H hypothesis
- E evidence related to the hypothesis H, i.e., the data to be used for validating (accepting/rejecting) the hypothesis H
- P(H) probability of the hypothesis (prior probability)
- P(E) probability of the evidence i.e., the state of the world described by the gathered data
- P(E|H) (conditional) probability of evidence E given that the hypothesis H holds
- P(H|E) (conditional) probability of the hypothesis H given the evidence E

Naive Bayes classifier

• Lets make an assumption that all attributes are mutually independent:

$$P(H|E) = \frac{P(E_1|H) * P(E_2|H) * ... * P(E_n|H) * P(H)}{P(E)}$$

Naive Bayes

- Makes two "naïve" assumptions over attributes:
 - all attributes are a priori equally important
 - all attributes are statistically independent (value of one attribute is not related to a value of another attribute)
- This assumptions mostly are not true, but in practice the algorithm gives good results

Example – Predicting whether a theater play will be performed

ToPlayOtNotToPlay.arff

| Outlook | Temp. | Humidity | Windy | Play |
|----------|-------|----------|-------|------|
| sunny | hot | high | false | no |
| sunny | hot | high | true | no |
| overcast | hot | high | false | yes |
| rainy | mild | high | false | yes |
| rainy | cool | normal | false | yes |
| rainy | cool | normal | true | no |
| overcast | cool | normal | true | yes |
| sunny | mild | high | false | no |
| sunny | cool | normal | false | yes |
| rainy | mild | normal | false | yes |
| sunny | mild | normal | true | yes |
| overcast | mild | high | true | yes |
| overcast | hot | normal | false | yes |
| rainy | mild | high | true | no |

Sunny weather

Suppose you know that it is sunny outside

Then 60% chance that Play = no

| Outlook | Temp. | Humidity | Windy | Play |
|----------|-------|----------|-------|------|
| sunny | hot | high | false | no |
| sunny | hot | high | true | no |
| overcast | hot | high | false | yes |
| rainy | mild | high | false | yes |
| rainy | cool | normal | false | yes |
| rainy | cool | normal | true | no |
| overcast | cool | normal | true | yes |
| sunny | mild | high | false | no |
| sunny | cool | normal | false | yes |
| rainy | mild | normal | false | yes |
| sunny | mild | normal | true | yes |
| overcast | mild | high | true | yes |
| overcast | hot | normal | false | yes |
| rainy | mild | high | true | no |

How well does outlook predict play?

| Outlook | Temp. | Humidity | V | Vindy | P | lay | |
|----------|-------|----------|------|-------|------|------|----|
| sunny | hot | high | 1 | false | r | າວ | |
| sunny | hot | high | true | | no | | |
| overcast | hot | high | 1 | false | у | es | |
| rainy | mild | high | 1 | false | у | es | |
| rainy | cool | normal | 1 | f-' | | | |
| rainy | cool | normal | | t | | Play | |
| overcast | cool | normal | | Outle | nok | Ves | no |
| sunny | mila | high | | | | yes | |
| sunny | cool | normal | | sunn | y 🛛 | 2 | 3 |
| rainy | mild | normal | 1 | ovor | | Л | 0 |
| sunny | mild | normal | | oven | Jasi | 4 | 0 |
| overcast | mild | high | | rainv | | 3 | 2 |
| overcast | hot | normal | | TOTA | | 0 | |
| rainy | mild | high | | | L | 9 | 5 |

How well does outlook predict play?

| | 2 | | | |
|----------|------|----|----------|---|
| | Play | |] | |
| Outlook | yes | no | 1 | |
| sunny | 2 | 3 | | |
| overcast | 4 | 0 | - | |
| rainy | 3 | 2 | 1 | |
| TOTAL | 9 | 5 | 1 | - |
| | | | For each | Ľ |

| Outlook | Temp. | Humidity | Windy | Play |
|----------|-------|----------|-------|------|
| sunny | hot | high | false | no |
| sunny | hot | high | true | no |
| overcast | hot | high | false | yes |
| rainy | mild | high | false | yes |
| rainy | cool | normal | false | yes |
| rainy | cool | normal | true | no |
| overcast | cool | normal | true | yes |
| sunny | mild | high | false | no |
| sunny | cool | normal | false | yes |
| rainy | mild | normal | false | yes |
| sunny | mild | normal | true | yes |
| overcast | mild | high | true | yes |
| overcast | hot | normal | false | yes |
| rainy | mild | high | true | no |

| | Play | | | Play | | Play | | Play | | | | Play | |
|----------|------|----|-------|------|----|--------|-----|------|-------|-----|----|-------|----|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 2 | 3 | hot | 2 | 2 | high | 3 | 4 | false | 6 | 2 | yes | 9 |
| overcast | 4 | 0 | mild | 4 | 2 | normal | 6 | 1 | true | 3 | 3 | no | 5 |
| rainy | 3 | 2 | cool | 3 | 1 | | | | | | | | |
| TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 14 |

Values to ratios

| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|----|-------|------|----|--------|------|----|-------|------|----|-------|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 2 | 3 | hot | 2 | 2 | high | 3 | 4 | false | 6 | 2 | yes | 9 |
| overcast | 4 | 0 | mild | 4 | 2 | normal | 6 | 1 | true | 3 | 3 | no | 5 |
| rainy | 3 | 2 | cool | 3 | 1 | | | | | | | | |
| TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 14 |

Covert values to ratios

| | Play | | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|------|------|--|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | | | |
| sunny | 0.22 | 0.60 | hot | 0.22 | 0.40 | high | 0.33 | 0.80 | false | 0.67 | 0.40 | yes | 0.64 | | |
| overcast | 0.44 | 0.00 | mild | 0.44 | 0.40 | normal | 0.67 | 0.20 | true | 0.33 | 0.60 | no | 0.36 | | |
| rainy | 0.33 | 0.40 | ool | 0.33 | 0.20 | | | | | | | | | | |

2 occurences of Play = no, where Outlook = rainy 5 occurences Play = no

Likelihood of playing under these weather conditions

Calculate the likelihood that: Outlook = sunny (0.22) Temperature = cool (0.33) Humidity = high (0.33) Windy = true (0.33) Play = yes (0.64)

Likelihood of playing under these weather conditions

| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|-----|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 0.22 | 0.60 | hot | 0.22 | 0.40 | high | 0.33 | 0.80 | false | 0.67 | 0.40 | yes | 0.64 |
| overcast | 0.44 | 0.00 | mild | 0.44 | 0.40 | normal | 0.67 | 0.20 | true | 0.33 | 0.60 | no | 0.36 |
| rainy | 0.33 | 0.40 | cool | 0.33 | 0.20 | | | | | | | | |

$0.22 \ge 0.33 \ge 0.33 \ge 0.33 \ge 0.64 = 0.0053$

Likelihood of NOT playing under these weather conditions

Calculate the likelihood that: Outlook = sunny (0.60) Temperature = cool (0.20) Humidity = high (0.80) Windy = true (0.60) Play = no (0.36)

Likelihood of **NOT** playing under these weather conditions

| | Play | | Play | | lay Play | | | | Play | | Play | | | | Play |
|----------|------|------|-------|------|----------|--------|------|------|-------|------|------|-----|------|--|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | | | |
| sunny | 0.22 | 0.60 | hot | 0.22 | 0.40 | high | 0.33 | 0.80 | false | 0.67 | 0.40 | yes | 0.64 | | |
| overcast | 0.44 | 0.00 | mild | 0.44 | 0.40 | normal | 0.67 | 0.20 | true | 0.33 | 0.60 | no | 0.36 | | |
| rainy | 0.33 | 0.40 | cool | 0.33 | 0.20 | | | | | | | | | | |

$0.60 \ge 0.20 \ge 0.80 \ge 0.60 \ge 0.36 = 0.0206$

The Bayes Theorem

```
Given these weather conditions:
Outlook = sunny
Temperature = cool
Humidity = high
Windy = true
```

Probability of Play = yes: $0.0053 \\ 0.0053 + 0.0206 \end{bmatrix} = 20.5\%$ Probability of Play = no: $0.0206 \\ 0.0053 + 0.0206 \end{bmatrix} = 79.5\%$

$$P(H|E) = \frac{P(E_1|H) * P(E_2|H) * \dots * P(E_n|H) * P(H)}{P(E)}$$

Likelihood of NOT playing under these weather conditions

Calculate the likelihood that:

Outlook = ovecast (0.00)Temperature = cool (0.20)Humidity = high (0.80)Windy = true (0.60)**Play = no** (0.36)



| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|-----|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 0.22 | 0.60 | hot | 0.22 | 0.40 | high | 0.33 | 0.80 | false | 0.67 | 0.40 | yes | 0.64 |
| overcast | 0.44 | 0.00 | mild | 0.44 | 0.40 | normal | 0.67 | 0.20 | true | 0.33 | 0.60 | no | 0.36 |
| rainy | 0.33 | 0.40 | cool | 0.33 | 0.20 | | | | | | | | |

The original dataset

| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|----|-------|------|----|--------|------|----|-------|------|----|-------|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 2 | 2 | hot | 2 | 2 | high | 3 | 4 | false | 6 | 2 | yes | 9 |
| overcast | 4 | 0 | mild | 4 | 2 | normal | 6 | 1 | true | 3 | 3 | no | 5 |
| rainy | 3 | | cool | 3 | 1 | | | | | | | | |
| TOTAL | 9 | | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 9 | 5 | TOTAL | 14 |

Laplace estimator: Add 1 to each count

After the Laplace estimator

| | Play | | | | Play | | | Play | | | Play | | | Play |
|----------|------|---|---|-------|------|----|--------|------|----|-------|------|----|-------|------|
| Outlook | yes | r | C | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 3 | Y | | hot | 3 | 3 | high | 4 | 5 | false | 7 | 3 | yes | 12 |
| overcast | 5 | 1 | | nild | 5 | 3 | normal | 7 | 2 | true | 4 | 4 | no | 8 |
| rainy | 4 | 3 | | cool | 4 | 2 | | | | | | | | |
| TOTAL | 12 | 8 | | TOTAL | 12 | 8 | TOTAL | 11 | 7 | TOTAL | 11 | 7 | TOTAL | 20 |

| | Play | | | Play | | Play | | Play | | | | Play | |
|----------|------|----|-------|------|----|--------|-----|------|-------|-----|----|-------|----|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 3 | 4 | hot | 3 | 3 | high | 4 | 5 | false | 7 | 3 | yes | 9 |
| overcast | 5 | 1 | mild | 5 | 3 | normal | 7 | 2 | true | 4 | 4 | no | 5 |
| rainy | 4 | 3 | cool | 4 | 2 | | | | | | | | |
| TOTAL | 12 | 8 | TOTAL | 12 | 8 | TOTAL | 11 | 7 | TOTAL | 11 | 7 | TOTAL | 14 |

Convert incremented counts to ratios after implementing the Laplace estimator

| Play | | | Play | | | Play | | | Play | | | Play | |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|------|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 0.25 | 0.50 | hot | 0.25 | 0.38 | high | 0.36 | 0.71 | false | 0.64 | 0.43 | yes | 0.64 |
| overcast | 0.42 | 0.13 | mild | 0.42 | 0.38 | normal | 0.64 | 0.29 | true | 0.36 | 0.57 | no | 0.36 |
| rainy | 0.33 | 0.38 | cool | 0.33 | 0.25 | | | | | | | | |

| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|-----|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 0.25 | 0.50 | hot | 0.25 | 0.38 | high | 0.36 | 0.71 | false | 0.64 | 0.43 | yes | 0.64 |
| overcast | 0.42 | 0.13 | mild | 0.42 | 0 38 | normal | 0.64 | 0.29 | true | 0.36 | 0.57 | no | 0.36 |
| rainy | 0.33 | 0.38 | cool | 0 33 | 0 25 | | | | | | | | |

Outlook = ovecast, Temperature = cool, Humidity = high, Windy = true

Play = no: $0.13 \ge 0.25 \ge 0.71 \ge 0.57 \ge 0.36 = 0.046$ **Play = yes**: $0.42 \ge 0.33 \ge 0.36 \ge 0.36 \ge 0.0118$

Probability of **Play = no**:

 $\frac{0.0046}{0.0046 + 0.0118} = 28\%$

Probability of **Play = yes**:

 $\frac{0.0118}{0.0046 + 0.0118} = 72\%$

Under these weather conditions: Temperature = cool Humidity = high Windy = true

NOT using Laplace estimator: Play = no: 79.5% Play = yes: 20.5% Using Laplace estimator: Play = no: 72.0% Play = yes: 28.0%

The effect of Laplace estimator has little effect as sample size grows.

Prediction rules

| Outlook | Temp. | Humid. | Windy | Play |
|----------|-------|--------|-------|------|
| overcast | cool | high | false | no |
| overcast | cool | high | false | yes |
| overcast | cool | high | true | no |
| overcast | cool | high | true | yes |
| overcast | cool | normal | false | no |
| overcast | cool | normal | false | yes |
| overcast | cool | normal | true | no |
| overcast | cool | normal | true | yes |
| overcast | hot | high | false | no |
| overcast | hot | high | false | yes |
| overcast | hot | high | true | no |
| overcast | hot | high | true | yes |
| overcast | hot | normal | false | no |
| overcast | hot | normal | false | yes |
| overcast | hot | normal | true | no |
| overcast | hot | normal | true | yes |

Repeat previous calculation for all other combinations of weather conditions.

Calculate the rules for each pair.

Then throw out the rules with p < 0.5

Prediction rules

| | Play | | | Play | | | Play | | | Play | | | Play |
|----------|------|------|-------|------|------|--------|------|------|-------|------|------|-----|------|
| Outlook | yes | no | Temp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 0.25 | 0.50 | hot | 0.25 | 0.38 | high | 0.36 | 0.71 | false | 0.64 | 0.43 | yes | 0.64 |
| overcast | 0.42 | 0.13 | mild | 0.42 | 0.38 | normal | 0.64 | 0.29 | true | 0.36 | 0.57 | no | 0.36 |
| rainy | 0.33 | 0.38 | cool | 0.33 | 0.25 | | | | | | | | |



| Inst | Outlook | Temp. | Humid. | Windy | Play | Outlook | Temp. | Humid. | Windy | Play | Like. | Prob. |
|------|----------|-------|--------|-------|-------|---------|----------|--------|-------|------|--------|-------|
| | overcast | cool | high | false | no | 0.13 | 0.25 | 0.71 | 0.43 | 0.36 | 0.0034 | 14.2% |
| | overcast | cool | high | false | yes | 0.42 | 0.33 | 0.36 | 0.64 | 0.64 | 0.0207 | 85.8% |
| | overcast | cool | high | | | | | 0.71 | 0.57 | 0.36 | 0.0046 | 27.8% |
| | overcast | cool | high | | | probat | ollities | 0.36 | 0.36 | 0.64 | 0.0118 | 72.2% |
| | overcast | cool | normal | Iora | 11 30 | Combin | ations | 0.29 | 0.43 | 0.36 | 0.0014 | 3.6% |
| | overcast | cool | normal | false | yes | 0.42 | 0.33 | 0.64 | 0.64 | 0.64 | 0.0362 | 96.4% |
| | overcast | cool | normal | true | no | 0.13 | 0.25 | 0.29 | 0.57 | 0.36 | 0.0018 | 8.1% |
| 7 | overcast | cool | normal | true | yes | 0.42 | 0.33 | 0.64 | 0.36 | 0.64 | 0.0207 | 91.9% |
| | overcast | hot | high | false | no | 0.13 | 0.38 | 0.71 | 0.43 | 0.36 | 0.0051 | 24.9% |
| 3 | overcast | hot | high | false | yes | 0.42 | 0.25 | 0.36 | 0.64 | 0.64 | 0.0155 | 75.1% |

Prediction rules

| Inst | Outlook | Temp. | Humid. | Windy | Play | Prob. |
|------|----------|-------|--------|-------|------|-------|
| | overcast | cool | normal | false | yes | 96.4% |
| | overcast | mild | normal | false | yes | 95.7% |
| 13 | overcast | hot | normal | false | yes | 93.0% |
| 7 | overcast | cool | normal | true | yes | 91.9% |
| | overcast | mild | normal | true | yes | 90.4% |
| 5 | rainy | cool | normal | false | yes | 87.6% |
| | overcast | cool | high | false | yes | 85.8% |
| 10 | rainy | mild | normal | false | yes | 85.5% |
| | overcast | hot | normal | true | yes | 85.0% |
| 2 | sunny | hot | high | true | no | 83.7% |
| | overcast | mild | high | false | yes | 83.4% |
| 9 | sunny | cool | normal | false | yes | 79.9% |
| | rainy | hot | normal | false | yes | 77.9% |
| | sunny | mild | normal | false | yes | 76.8% |
| | sunny | mild | high | true | no | 75.5% |
| 3 | overcast | hot | high | false | yes | 75.1% |
| | rainy | cool | normal | true | yes | 75.1% |
| | rainy | hot | high | true | no | 74.3% |

| Inst | Outlook | Temp. | Humid. | Windy | Play | Prob. |
|------|----------|-------|--------|-------|------|-------|
| | overcast | cool | high | true | yes | 72.2% |
| | sunny | cool | high | true | no | 72.0% |
| | rainy | mild | normal | true | yes | 71.6% |
| 1 | sunny | hot | high | false | no | 68.8% |
| 12 | overcast | mild | high | true | yes | 68.4% |
| | sunny | hot | normal | false | yes | 66.5% |
| 14 | rainy | mild | high | true | no | 63.5% |
| | sunny | cool | normal | true | yes | 63.0% |
| | rainy | cool | high | false | yes | 61.7% |
| | rainy | hot | normal | true | yes | 60.2% |
| | rainy | cool | high | true | no | 59.1% |
| 11 | sunny | mild | normal | true | yes | 58.6% |
| 4 | rainy | mild | high | false | yes | 57.3% |
| 8 | sunny | mild | high | false | no | 57.0% |
| | overcast | hot | high | true | yes | 56.4% |
| | rainy | hot | high | false | no | 55.4% |
| | sunny | hot | normal | true | no | 54.0% |
| | sunny | cool | high | false | no | 52.4% |

The instance 6 is missing

Rules predicting class for all combinations of attributes

Comparing the prediction with the original data

| Inst | Outlook | Temp. | Humid. | Windy | Play | Prob. | Actual |
|------|----------|-------|--------|-------|------|-------|--------|
| 1 | sunny | hot | high | false | no | 72.6% | no |
| 2 | sunny | hot | high | true | no | 86.1% | no |
| 3 | overcast | hot | high | false | yes | 71.6% | yes |
| 4 | rainy | mild | high | false | yes | 52.8% | yes |
| 5 | rainy | cool | normal | false | yes | 85.5% | yes |
| 6 | rainy | cool | normal | true | yes | 75.1% | no |
| 7 | overcast | cool | normal | true | yes | 90.4% | yes |
| 8 | sunny | mild | high | false | no | 61.4% | no |
| 9 | sunny | cool | normal | false | yes | 76.8% | yes |
| 10 | rainy | mild | normal | false | yes | 83.0% | yes |
| 11 | sunny | mild | normal | true | yes | 54.2% | yes |
| 12 | overcast | mild | high | true | yes | 64.3% | yes |
| 13 | overcast | hot | normal | false | yes | 91.7% | yes |
| 14 | rainy | mild | high | true | no | 67.6% | no |

Naïve Bayes in Weka



Predictions over training dataset

ToPlayOtNotToPlay.arff dataset

| | Weka Explorer |
|---|-------------------------------------|
| Preprocess Classify | Classifier evaluation options |
| Classifier Choose Naive Payee | 🥑 Output model |
| Test options | ☑ Output per-class stats |
| • Use training set | Output entropy evaluation measures |
| O Supplied test set Set | Output confusion matrix |
| O Cross-validation Polds 10 O Percentage split % 66 | Store predictions for visualization |
| More options | Output predictions |

Classification results

| \square | Play | | | Play | | Play | | | Play | | | Play | |
|-----------|--------------------------|----|-------|------|----|--------|-----|----|------------|--------|----|------|------|
| Outlook | yes | no | ſemp. | yes | no | Humid. | yes | no | Windy | yes | no | | |
| sunny | 3 | 4 | not | 3 | 3 | high | 4 | 5 | false | 7 | 3 | yes | 12 |
| overcast | 5 | 1 | nild | 5 | 3 | normal | 7 | 2 | Classifier | output | | | |
| rainy | 4 | 3 | loo: | 4 | 2 | | | | Attrib | nte | | 20 | URA |
| TOTAL | 12 | 8 | TOTAL | 12 | 8 | TOTAL | 11 | 7 | ACCIL | (| | | .63) |
| | _ | | | | | | | | | | | | |
| | | | | | | | | | Outloo | k | | | |
| | | | | | | | | 1 | sunn | Y | | 4.0 | 3.0 |
| | | | | | | | | | over | cast | | 1.0 | 5.0 |
| | | | | | | | | | rain | Y | | 3.0 | 4.0 |
| | | | | | | | | | [tot | al] | | 8.0 | 12.0 |
| | | | | | | | | | Temp. | | | | |
| | | | | | | | | | hot | | | 3.0 | 3.0 |
| | | | | | | | | | mild | l I | | 3.0 | 5.0 |
| | | | | | | | | | cool | | | 2.0 | 4.0 |
| | | | | | | | | | [tot | al] | | 8.0 | 12.0 |
| | | | | | | | | | Humidi | tv | | | |
| . | | | | | | | | | high | | | 5.0 | 4.0 |
| The L | The Laplace estimator is | | | | | | | | norm | al | | 2.0 | 7.0 |
| auto | automatically applied | | | | | | | | [tot | al] | | 7.0 | 11.0 |

Classification results

| Classifier output | | | | | | | | Instance 6 is marked as a wrong identified instance |
|-------------------------------------|--|---|---|------------------------|---|---|--------|---|
| === Predictions on training set === | | | | | | | | |
| = i | === Pro inst#, 1 2 3 4 5 6 7 8 9 10 11 12 13 | actual, 1:no 1:no 2:yes 2:yes 2:yes 1:no 2:yes 1:no 2:yes 2:yes 2:yes 2:yes 2:yes 2:yes | n training predicted, 1:no 1:no 2:yes 2:yes 2:yes 2:yes 1:no 2:yes 2:yes 2:yes 2:yes 2:yes 2:yes 2:yes | set === error, + | probabi *0.704 *0.847 0.263 0.446 0.133 0.263 0.087 *0.588 0.214 0.155 0.432 0.333 0.075 | lity dis 0.296 0.153 *0.737 *0.554 *0.867 *0.737 *0.913 0.412 *0.786 *0.845 *0.845 *0.667 *0.925 | tribut | Probability of each instance in the dataset |
| | 14 | 1:no | 1:no | | *0 . 652 | 0.348 | | |

Naïve Bayes features

- Intended primarily for the work with nominal attributes
- In case of numeric attributes:
 - Use the probability distribution of attributes (Normal distribution is default) for probability estimation for the each attribute
 - Discretize the attribute's values

Example 2 – Eatable Mushrooms dataset

EdibleMushrooms.arff

- Eatable Mushrooms dataset based on "National Audubon Society Field Guide to North American Mushrooms"
- Hypothetical samples with descriptions corresponding to 23 species of mushrooms
- There are 8124 instances with 22 nominal attributes which describe mushroom characteristics; one of which is whether a mushroom is eatable or not
- Our goal is to predict whether a mushroom is eatable or not

Data in this dataset are hypothetical and these results are not to be used in real life!

Baseline classifier

diabetes.arff

- There are total of 768 instances (500 negative, 268 positive)
- A priori probabilities for classes negative and positive are

$$Negativan = \frac{500}{768} \cdot 100\% = 65.1\%$$
$$Pozitivan = \frac{268}{768} \cdot 100\% = 34.9\%$$

- Baseline classifier classifies every instances to the dominant class, the class with the highest probability
- In Weka, the implementation of baseline classifier is: rules -> ZeroR

Baseline classifier in Weka: rules -> ZeroR



Baseline classifier

- Open dataset diabetes.arff
- Test option: Percentage split 66%
- Test classifiers:

| • | rules -> ZeroR | 65% |
|---|---------------------|-----|
| • | trees -> J48 | 76% |
| • | bayes -> NaiveBayes | 77% |
| • | lazy -> Ibk | 73% |

• For every classification problem test first whether the tested classifier performs better than the baseline classifier

Example 3 – Supermarket dataset

supermarket.arff

- Dataset describes data about the article sales in a local supermarket in New Zealand in one day.
- Attributes are nominal and describes different store departments and different article categories (e.g.. "bread and cake' refer to the group of baking products).
- Value "t" of an attributes means that the shopping cart contained at least one product for the specific department or at least one product from the product category.
- Class has values "low" and "high" determining whether a byer spent less or more than 100\$ for the shopping

Recommendations and credits

Weka Tutorials and Assignments @ The Technology Forge

Link: <u>http://www.technologyforge.net/WekaTutorials/</u>

"Data Mining with Weka" and "More Data Mining with Weka": MOOCs from the University of Waikato. A self-paced session of "Data Mining with Weka" runs until 23 October June 2015.

• Link: <u>https://www.youtube.com/user/WekaMOOC/</u>

(Anonymous) survey for your comments and suggestions: http://goo.gl/cqdp3l

ANY QUESTIONS?

UROŠ KRČADINAC EMAIL: <u>uros@krcadinac.com</u> URL: <u>http://krcadinac.com</u>