

## Assignment 1

Analyze `weather.nominal.arff` using Weka 3.6.14

**Collaborators:** None.

### Exercise 1-1. Compute posterior distribution for the class

Use NaiveBayesSimple to find a Bayesian classifier (Laplace estimator = 1 is used in order to avoid frequency problems). Compute  $P(\text{"yes"})$  and  $P(\text{"no"})$  of the following instances.

(a) Windy = True

**Solution:** Given the evidence Windy = True, with independence assumption and Bayes rule, we have the probability of answer “yes” as:

$$P(\text{"yes"} | \text{Windy} = \text{True}) = \frac{P(\text{Windy} = \text{True} | \text{"yes"}) \times P(\text{"yes"})}{P(\text{Windy} = \text{True})}$$

Because of the independence assumption between attributes, we have the prior for Windy = True as follow:

$$\begin{aligned} P(\text{Windy} = \text{True}) &= P(\text{Windy} = \text{True} | \text{"yes"}) \times P(\text{"yes"}) + \\ &\quad P(\text{Windy} = \text{True} | \text{"no"}) \times P(\text{"no"}) \\ &= 0.3636 \times 0.625 + 0.5714 \times 0.375 \\ &= 0.4415 \end{aligned}$$

We have the posterior distribution for the class given Windy = True:

$$\begin{aligned} P(\text{"yes"} | \text{Windy} = \text{True}) &= \frac{P(\text{Windy} = \text{True} | \text{"yes"}) \times P(\text{"yes"})}{P(\text{Windy} = \text{True})} \\ &= \frac{0.3636 \times 0.625}{0.4415} \\ &= 0.5 \\ P(\text{"no"} | \text{Windy} = \text{True}) &= 1 - P(\text{"yes"} | \text{Windy} = \text{True}) \\ &= 0.5 \end{aligned}$$

(b) Humidity = High, Windy = True

**Solution:** Similar to question 1-(a), we have the join prior probability for the case Humidity = High (abbr. H = High), Windy = True (abbr. W = True) as follow:

$$\begin{aligned} P(H = \text{High}, W = \text{True}) &= P(H = \text{High}, W = \text{True} | \text{"yes"}) \times P(\text{"yes"}) + \\ &\quad P(H = \text{High}, W = \text{True} | \text{"no"}) \times P(\text{"no"}) \end{aligned}$$

Apply the attribute independence assumption, we have the following factorization of the conditional join probability:

$$\begin{aligned}
 P(H = \text{High}, W = \text{True}) &= P(H = \text{High} | \text{"yes"}) \times P(W = \text{True} | \text{"yes"}) \times P(\text{"yes"}) + \\
 &\quad P(H = \text{High} | \text{"no"}) \times P(W = \text{True} | \text{"no"}) \times P(\text{"no"}) \\
 &= 0.3636 \times 0.3636 \times 0.625 + 0.5714 \times 0.7142 \times 0.375 \\
 &= 0.2357
 \end{aligned}$$

We have the posterior distribution for the class given Windy = True and Humidity = High:

$$\begin{aligned}
 P(\text{"yes"} | H = \text{High}, W = \text{True}) &= \frac{P(W = \text{True}, H = \text{High} | \text{"yes"}) \times P(\text{"yes"})}{P(H = \text{High}, W = \text{True})} \\
 &= \frac{0.3636 \times 0.3636 \times 0.625}{0.2357} \\
 &= \mathbf{0.35} \\
 P(\text{"no"} | H = \text{High}, W = \text{True}) &= 1 - P(\text{"yes"} | \text{Windy} = \text{True}) \\
 &= \mathbf{0.65}
 \end{aligned}$$

(c) Temperature = Hot, Humidity = High, Windy = True

**Solution:** Similar to question 1-(a,b), we have the join prior probability for the case Temperature = Hot (abbr. T = Hot) Humidity = High (abbr. H = High), Windy = True (abbr. W = True) as follow:

$$\begin{aligned}
 P(H = \text{High}, W = \text{True}) &= P(H = \text{High}, W = \text{True} | \text{"yes"}) \times P(\text{"yes"}) + \\
 &\quad P(H = \text{High}, W = \text{True} | \text{"no"}) \times P(\text{"no"})
 \end{aligned}$$

Apply the attribute independence assumption, we have the following factorization of the conditional join probability:

$$\begin{aligned}
 P(H = \text{High}, W = \text{True}) &= P(H = \text{High} | \text{"yes"}) \times P(W = \text{True} | \text{"yes"}) \times P(\text{"yes"}) + \\
 &\quad P(H = \text{High} | \text{"no"}) \times P(W = \text{True} | \text{"no"}) \times P(\text{"no"}) \\
 &= 0.3636 \times 0.3636 \times 0.625 + 0.5714 \times 0.7142 \times 0.375 \\
 &= 0.2357
 \end{aligned}$$

We have the posterior distribution for the class given Windy = True and Humidity = High:

$$\begin{aligned}
 P(\text{"yes"} | H = \text{High}, W = \text{True}) &= \frac{P(W = \text{True}, H = \text{High} | \text{"yes"}) \times P(\text{"yes"})}{P(H = \text{High}, W = \text{True})} \\
 &= \frac{0.3636 \times 0.3636 \times 0.625}{0.2357} \\
 &= \mathbf{0.35} \\
 P(\text{"no"} | H = \text{High}, W = \text{True}) &= 1 - P(\text{"yes"} | \text{Windy} = \text{True}) \\
 &= \mathbf{0.65}
 \end{aligned}$$

**Solution:** By equivalent translation to first-order logic, we have the following representation:

**Exercise 1-2.** [15 points] **Recurrence Relation Resolution**

For each of the following recurrence relations, pick the correct asymptotic runtime:

- (a) [5 points] Select the correct asymptotic complexity of an algorithm with runtime  $T(n, n)$  where

$$\begin{aligned} T(x, c) &= \Theta(x) && \text{for } c \leq 2, \\ T(c, y) &= \Theta(y) && \text{for } c \leq 2, \text{ and} \\ T(x, y) &= \Theta(x + y) + T(x/2, y/2). \end{aligned}$$

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

- (b) [5 points] Select the correct asymptotic complexity of an algorithm with runtime  $T(n, n)$  where

$$\begin{aligned} T(x, c) &= \Theta(x) && \text{for } c \leq 2, \\ T(c, y) &= \Theta(y) && \text{for } c \leq 2, \text{ and} \\ T(x, y) &= \Theta(x) + T(x, y/2). \end{aligned}$$

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

- (b) [5 points] Select the correct asymptotic complexity of an algorithm with runtime  $T(n, n)$  where

$$\begin{aligned} T(x, c) &= \Theta(x) && \text{for } c \leq 2, \\ T(x, y) &= \Theta(x) + S(x, y/2), \\ S(c, y) &= \Theta(y) && \text{for } c \leq 2, \text{ and} \\ S(x, y) &= \Theta(y) + T(x/2, y). \end{aligned}$$

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

## Peak-Finding

In Lecture 1, you saw the peak-finding problem. As a reminder, a *peak* in a matrix is a location with the property that its four neighbors (north, south, east, and west) have value less than or equal to the value of the peak. We have posted Python code for solving this problem to the website in a file called `ps1.zip`. In the file `algorithms.py`, there are four different algorithms which have been written to solve the peak-finding problem, only some of which are correct. Your goal is to figure out which of these algorithms are correct and which are efficient.

### Exercise 1-3. [16 points] Peak-Finding Correctness

(b) [4 points] Is `algorithm1` correct?

1. Yes.
2. No.

**Solution:** 1

(c) [4 points] Is `algorithm2` correct?

1. Yes.
2. No.

**Solution:** 1

(d) [4 points] Is `algorithm3` correct?

1. Yes.
2. No.

**Solution:** 1

(e) [4 points] Is `algorithm4` correct?

1. Yes.
2. No.

**Solution:** 1

**Exercise 1-4.** [16 points] **Peak-Finding Efficiency**

(f) [4 points] What is the worst-case runtime of `algorithm1` on a problem of size  $n \times n$ ?

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

(g) [4 points] What is the worst-case runtime of `algorithm2` on a problem of size  $n \times n$ ?

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

(h) [4 points] What is the worst-case runtime of `algorithm3` on a problem of size  $n \times n$ ?

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .
6.  $\Theta(2^n)$ .

**Solution:** 1

(i) [4 points] What is the worst-case runtime of `algorithm4` on a problem of size  $n \times n$ ?

1.  $\Theta(\log n)$ .
2.  $\Theta(n)$ .
3.  $\Theta(n \log n)$ .
4.  $\Theta(n \log^2 n)$ .
5.  $\Theta(n^2)$ .

6.  $\Theta(2^n)$ .

**Solution:** 1

**Exercise 1-5.** [19 points] **Peak-Finding Proof**

Please modify the proof below to construct a proof of correctness for the *most efficient correct algorithm* among `algorithm2`, `algorithm3`, and `algorithm4`.

The following is the proof of correctness for `algorithm1`, which was sketched in Lecture 1.

We wish to show that `algorithm1` will always return a peak, as long as the problem is not empty. To that end, we wish to prove the following two statements:

**1. If the peak problem is not empty, then `algorithm1` will always return a location.** Say that we start with a problem of size  $m \times n$ . The recursive subproblem examined by `algorithm1` will have dimensions  $m \times \lfloor n/2 \rfloor$  or  $m \times (n - \lfloor n/2 \rfloor - 1)$ . Therefore, the number of columns in the problem strictly decreases with each recursive call as long as  $n > 0$ . So `algorithm1` either returns a location at some point, or eventually examines a subproblem with a non-positive number of columns. The only way for the number of columns to become strictly negative, according to the formulas that determine the size of the subproblem, is to have  $n = 0$  at some point. So if `algorithm1` doesn't return a location, it must eventually examine an empty subproblem.

We wish to show that there is no way that this can occur. Assume, to the contrary, that `algorithm1` does examine an empty subproblem. Just prior to this, it must examine a subproblem of size  $m \times 1$  or  $m \times 2$ . If the problem is of size  $m \times 1$ , then calculating the maximum of the central column is equivalent to calculating the maximum of the entire problem. Hence, the maximum that the algorithm finds must be a peak, and it will halt and return the location. If the problem has dimensions  $m \times 2$ , then there are two possibilities: either the maximum of the central column is a peak (in which case the algorithm will halt and return the location), or it has a strictly better neighbor in the other column (in which case the algorithm will recurse on the non-empty subproblem with dimensions  $m \times 1$ , thus reducing to the previous case). So `algorithm1` can never recurse into an empty subproblem, and therefore `algorithm1` must eventually return a location.

**2. If `algorithm1` returns a location, it will be a peak in the original problem.** If `algorithm1` returns a location  $(r_1, c_1)$ , then that location must have the best value in column  $c_1$ , and must have been a peak within some recursive subproblem. Assume, for the sake of contradiction, that  $(r_1, c_1)$  is not also a peak within the original problem. Then as the location  $(r_1, c_1)$  is passed up the chain of recursive calls, it must eventually reach a level where it stops being a peak. At that level, the location  $(r_1, c_1)$  must be adjacent to the dividing column  $c_2$  (where  $|c_1 - c_2| = 1$ ), and the values must satisfy the inequality  $val(r_1, c_1) < val(r_1, c_2)$ .

Let  $(r_2, c_2)$  be the location of the maximum value found by `algorithm1` in the dividing column. As a result, it must be that  $val(r_1, c_2) \leq val(r_2, c_2)$ . Because the algorithm

chose to recurse on the half containing  $(r_1, c_1)$ , we know that  $val(r_2, c_2) < val(r_2, c_1)$ . Hence, we have the following chain of inequalities:

$$val(r_1, c_1) < val(r_1, c_2) \leq val(r_2, c_2) < val(r_2, c_1)$$

But in order for `algorithm1` to return  $(r_1, c_1)$  as a peak, the value at  $(r_1, c_1)$  must have been the greatest in its column, making  $val(r_1, c_1) \geq val(r_2, c_1)$ . Hence, we have a contradiction.

**Solution:** Write your proof here.

### Exercise 1-6. [19 points] Peak-Finding Counterexamples

For each incorrect algorithm, upload a Python file giving a counterexample (i.e. a matrix for which the algorithm returns a location that is not a peak).

**Solution:**

```
problemMatrix = [
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0]
]
```

```
problemMatrix = [
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0]
]
```