

# Artificial Intelligence and Knowledge Engineering

## Assignment 4

Katarzyna Fojcik, Joanna Szolomicka, Piotr Miłkowski

May 15, 2024

### 1 Exercise goal

The purpose of the exercise is to provide a practical introduction to simple machine learning algorithms and the basic steps of implementing a machine learning project: data mining and problem definition, data preparation, selection of algorithms and hyperparameters, evaluation of algorithms, improvement of results, presentation of results.

### 2 Theoretical introduction

In this section, you will find information useful during solving the stated problems.

#### 2.1 Machine Learning

**Machine Learning** is one of the fields of artificial intelligence, a broad category of algorithms that aim to perform specific tasks automatically. Algorithms learn from data and then apply that knowledge to infer and solve the task.

**Machine learning types:**

- **Supervised learning** – a method of learning in which the training set, on which the algorithm learns, contains the solution to the problem in the form of labels or classes. The two main areas that use supervised learning are **regression** (value prediction) and **classification** (class prediction).
- **Unsupervised learning** – a method of learning in which the training set is not labeled. The input to the model is raw data, and the model has to find connections between the data (so-called learning without a teacher). Unsupervised algorithms are **clustering** and various algorithms for **data visualization** and **dimensionality reduction**.
- **Reinforcement learning** – its task is to interact with the environment based on the information collected. In contrast to methods mentioned earlier, in reinforcement learning, one does not prepare a set of training data, but an environment from which the model will collect data automatically. The aim of the model is to maximize the reward the environment returns. The environment may depend on the learning objective. In the case of teaching a program that plays games, the environment will be a game with its rules, and in the case of a program that learns to control a rover it will be the real world.

#### 2.2 Division of data into training and validation sets

**Training set** — a set of data used to train a model. Using this data, the model learns to classify properly and builds all dependencies. It predicts possible outcomes and makes decisions based on the data given to it.

**Validation set** – a dataset used to run an unbiased test of a model that was learning on training data. The test is performed during the selection of the model or hyperparameters. It is important

that the data in the validation set has not been used for the model training, as it will be unsuitable for objective, unbiased testing.

**Test set** – when the selected model is trained, it is time to test in on data from test set. It is very important that this data has not been previously used to train or validate the model. The purpose of testing the model is to see how well the chosen algorithm performs on data that it has never dealt with before.

**Cross-validation** is an alternative method of preparing training and testing sets for the model. It is **more efficient in assessing the actual accuracy of the model**, but it is less efficient than the naive method. This method consists of **dividing the dataset into subsets** and creating **validation and training datasets** out of them.

There are many types of cross-validation. Baseline K-fold cross-validation assumes that the data is independent and identically distributed (IID). It divides all samples into  $k$  groups of equal size if possible. The ML algorithm is trained using  $k - 1$  groups, and the skipped group is used for validation. Thus,  $k$  training sessions are conducted and the results from  $k$  validations are averaged.

The ratio of splitting the data into training, validation and test sets should depend on the sample size. There are different splitting practices, but it is generally accepted that for smaller sets we prioritize the need to train the algorithm on as many samples as possible. Then the validation (and test) sets are relatively smaller (in the case of cross-validation, we choose a larger  $k$ , e.g., 10). However, having several thousand samples at our disposal, we can decide to increase the number of samples in the validation set (alternatively, reduce the number of groups to  $k = 5$ ).

## 2.3 Naive Bayes Classifier

The Naive Bayes Classifier is a simple probabilistic classifier that naively assumes that features (predictors) for a fixed class label are independent and are equally significant. It is based on Bayes Theorem.

**Definition** (Bayes Theorem) *Let  $A$  i  $B$  be events,  $P(B) > 0$ ,  $P(A|B)$  is conditional probability, then:*

$$\begin{aligned} P(A \wedge B) &= P(A|B)P(B) = P(B|A)P(A) \\ P(B|A) &= \frac{P(A|B)P(B)}{P(A)} \end{aligned} \tag{1}$$

**Definition** (Naive Bayes Classifier) *Let  $X_1, \dots, X_n$  be features, then:*

$$\begin{aligned} P(X_1|X_2, Y) &= P(X_1|Y) \\ P(X_1, X_2, Y) &= P(X_1|X_2, Y)P(X_2, Y) = P(X_1|X_2, Y)P(X_2|Y)P(Y) \\ P(X_1, X_2|Y) &= P(X_1|Y)P(X_2|Y) \\ P(X_1, X_2, \dots, X_n|Y) &= \prod_{i=1}^n P(X_i|Y) \end{aligned} \tag{2}$$

## 2.4 Decision tree

A decision tree in machine learning is a representation of a classifier in the form of a tree that supports the decision process. Building such a tree is based on various algorithms and their parameters, which are calculated for each new node of the decision tree. The example parameters are entropy and information gain.

**Definition** (Entropy) *Entropy is a measure of the variability of the data given by formula 3.*

$$H(S) = - \sum_{x \in X} p(x) \log_2 p(x) \tag{3}$$

where  $S$  is a current dataset for which entropy is calculated (for each node of the tree it will be a different - suitably smaller dataset),  $X$  is a set of classes in the set  $S$ ,  $p(x)$  is the ratio of the number of elements in the class  $x$  to the number of elements in the set  $S$ .

**Definition** (Information Gain) *Information Gain  $G(A)$  is a measure of the difference in entropy before and after splitting dataset using a given attribute (feature).*

$$G(A) = H(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} H(S_v) \quad (4)$$

where  $H(S)$  is the entropy of the set  $S$ ,  $\text{Values}(A)$  is the set of all values of attribute  $A$ ,  $S_v$  is a subset of  $S$  such that  $S_v = \{s \in S : A(s) = v\}$ ,  $H(S_v)$  is the entropy of  $S_v$ .

Decision trees are very susceptible to **overfitting**. As the depth increases, the trees are able to adapt very well to the training data, but at the same time their ability to generalize deteriorates (the error on the test set increases).

## 2.5 PCA (Principal Component Analysis)

PCA is a statistical method of factor analysis: a dataset consisting of  $N$  observations, each containing  $K$  variables, can be interpreted as a cloud of  $N$  points in  $K$ -dimensional space. The aim of PCA is to rotate the coordinate system in such a way that the variance of the first coordinate is maximized first, then the second and so on. The transformed coordinate values are the principal components. In machine learning, it is used to reduce the dimension of the data.

## 2.6 Metrics

To evaluate the classification, the results of various metrics are used (usually only for the validation data):

### 2.6.1 Confusion matrix

Depending on the performance of the classifier, four cases are distinguished (Table 1):

- **TP (True Positive)** – number of truly recognized positive cases
- **FP (False Positive)** – number of falsely recognized positive cases
- **FN (False Negative)** – number of falsely recognized negative cases
- **TN (True Negative)** – number of truly recognized negative cases

		Real class	
		Positive class	Negative class
Prediction	Positive class	TP	FP
	Negative class	FN	TN

Table 1: Confusion matrix.

### 2.6.2 Accuracy

Accuracy ACC is the most common indicator that allows us to assess the quality of data classification. It specifies which part of all classified items has been correctly classified. That is, we divide the sum of the correct classifications from a given category (TP) and the correct ones from other categories (TN) by the number of all classified cases.

$$\text{ACC} = \frac{\text{TP} + \text{TN}}{P + N} \quad (5)$$

### 2.6.3 Recall/sensitivity/true positive rate

Recall TPR tells us what is the proportion of correctly predicted positive cases (TP) among all positive cases (including those misclassified as negative - FN).

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (6)$$

### 2.6.4 Precision/positive predictive rate

We divide the number of correctly predicted positive values (TP) by the sum of all positively predicted values (including those incorrectly predicted), i.e. TP+FP. As a result, we find how many positively predicted examples are actually positive.

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (7)$$

### 2.6.5 F1-score

F1-score is the harmonic mean between precision and recall. The closer it is to one, the better the classification model. In the best case, F1-score is equal to 1, when we are dealing with perfect sensitivity and precision.

$$\text{F1 - score} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}} \quad (8)$$

## 3 Task

Face the problem of identifying demand from the world of demand planning for T-shirts with given parameters. For this purpose, use the attached **t-shirts.csv**.

Inside the file there are 20,000 items labeled with the following parameters:

- size (XS, S, M, L, XL, XXL, 3XL)
- material (nylon, polyester, silk, cotton, linen)
- color (white, cream, blue, black, orange, green, yellow, red, violet, navy)
- sleeves (short, long)
- demand (low, medium, high)

The tasks should be done using Python and/or WEKA. **Scoring:**

1. data mining – provide basic statistics and comments on the features and labels in the dataset. (10 points)
2. data preparation – split the data into training and validation sets (alternatively use cross-validation), examine the impact of various types of data processing on the classification results (suggested: normalization, standardization, discretization, PCA) - compare the results without data processing with the results after data processing using at least 2 different methods (separately). (30 points)  
Bonus - remove 5% of the feature values and prepare the data using methods that deal with missing data. (5 points)
3. classification – test the classifiers and study the impact on the results: naive Bayes classifier and decision tree using at least 3 different sets of hyperparameters. (40 points)  
Bonus - Test more advanced algorithms (with understanding!), such as Random Forest or Support Vector Machine (SVM) classifier. (5 points)
4. classification evaluation – use the metrics for classification evaluation to compare the results of different types of data preparation and the selected classifier. Interpret the results. (20 points)

Prepare a report for the task, including a brief description of all the steps performed and the task results with an interpretation (preferably collected in a table). In your report, indicate the source materials used and briefly describe the libraries used in the implementation. Reports must be sent to the tutor at least 24h before the tasks' submission.

## 4 Appendix - filters and classifiers in Weka Explorer

### 4.1 Filters

- **attribute discretization:** filters → supervised → attribute → Discretize lub filters → unsupervised → attribute → Discretize
- **attribute normalization:** filters → unsupervised → attribute → Normalize
- **PCA:** filters → unsupervised → attribute → PrincipalComponents
- **deleting an attribute:** filters → unsupervised → attribute → Remove
- **completing the attribute value:** filters → unsupervised → attribute → ReplaceMissingValues
- **attribute standardization:** filters → unsupervised → attribute → Standardize
- **deleting instances with given attribute values:** filters → unsupervised → instance → RemoveWithValues

### 4.2 Attribute selection

Card *Select attributes* - selection by the button *Attribute Evaluator*:

- GainRatioAttributeEval
- InfoGainAttributeEval
- PrincipalComponents

### 4.3 Classifiers

Card *Classify*:

- **naive Bayes classifier:** classifiers → bayes → NaiveBayes or NaiveBayesMultinomial (depending on the needs)
- **decision tree C4.5:** classifiers → trees → J48
- **group of classifiers:**
  - **boosting:** classifiers → meta → AdaBoostM1
  - **bagging:** classifiers → meta → Bagging
  - **different types of classifiers:** classifiers → meta → Stacking

## 5 Appendix - useful libraries for Python

- from sklearn.preprocessing import Normalizer, StandardScaler
- from sklearn.decomposition import PCA
- from sklearn.model\_selection import train\_test\_split
- from sklearn.metrics import accuracy\_score, precision\_score, recall\_score, confusion\_matrix
- from sklearn.tree import DecisionTreeClassifier
- from sklearn.naive\_bayes import GaussianNB