



Trinity College Dublin

Coláiste na Tríonóide, Baile Átha Cliath

The University of Dublin

School of Computer Science and Statistics

CS7CS4 Machine Learning Final Assignment

Boris Flesch

20300025

January 25, 2021

MSc Computer Science, Intelligent Systems

Contents

1	Part 1 – Reviews predictions	2
1.1	Introduction	2
1.2	Data pre-processing	2
1.2.1	Parsing	2
1.2.2	Transforming reviews text to features	2
1.2.3	Tuning TfidfVectorizer	2
1.2.4	Improving the tokenizer	4
1.3	Selected classification models	6
1.3.1	Baseline	6
1.3.2	Logistic regression	6
1.3.3	k-Nearest Neighbours (kNN)	7
1.3.4	Multi-Layer Perceptron (MLP)	8
1.4	Predicting the review polarity	8
1.4.1	Experiments and results	8
1.4.2	Discussion of the results	12
1.5	Predicting if the review is for an "early access"	14
1.5.1	Experiments and results	14
1.5.2	Discussion of the results	16
1.6	Conclusions and reflections	16
2	Part 2.	18
2.1	Question (i)	18
2.2	Question (ii)	19
2.3	Question (iii)	20
2.4	Question (iv)	21
2.5	Question (v)	21
A	Appendix	22
A.1	Python Code	22
A.1.1	Class MyPreprocessor	22
A.1.2	Class MyModel	25
A.1.3	Class MyBaseline	26
A.1.4	Class MyLogisticRegression	27
A.1.5	Class MyKNeighborsClassifier	29
A.1.6	Class MyMLPClassifier	31
A.1.7	Main script	32

1 Part 1 – Reviews predictions

1.1 Introduction

Based on the provided video games reviews from Steam platform, our aim is to evaluate whether the review text can be used to (i) predict the review polarity and (ii) predict whether the review is for an “early access” version of a game or not.

We are first going to discuss data pre-processing. Then, we will mention the selected classification models and, finally, we will use these pre-processing steps and models to evaluate which predictions are possible based on the review text.

1.2 Data pre-processing

1.2.1 Parsing

To parse the JSON collection provided, we can use *jsonlines* that allows to parse seamlessly every item of the collection and store it in a Python list as following:

```
1 for item in jsonlines.open(path, 'r'):
2     self.reviews.append(item['text'])
3     self.voted_up.append(item['voted_up'])
4     self.early_access.append(item['early_access'])
```

Reviews will systematically be used as our input feature. Depending on the prediction to make, either *voted_up* or *early_access* can be used as the output.

1.2.2 Transforming reviews text to features

To convert a review text to an input feature X , we could use a bag of words with *one-hot encoding*: every single word would act as a feature vector and a 1 would be added in the bag of words (i.e. matrix) every time the word is used in a review.

A quite similar but more sophisticated approach is to use TF-IDF. That approach considers the importance of words using Term Frequency (the frequency of a term in a document) and Inverse Document Frequency (the frequency of a term in the overall collection of documents). $Tfidf(t, d)$ of token t in document d "is large for a token that occurs a lot in document d but only rarely in overall collection of documents" (Text Features lecture PDF, slide 12). Sklearn's `TfidfVectorizer` implements that TF-IDF approach. It can be used as following:

```
1 self.tfidf = TfidfVectorizer()
2 self.X = self.tfidf.fit_transform(self.reviews).toarray()
```

In this process of vectorization, we are not considering any language difference: all tokens are generated independently of the review language.

1.2.3 Tuning `TfidfVectorizer`

`TfidfVectorizer` provides parameters that can be tuned to improve the relevance of the generated features (i.e. vectors, tokens). These parameters are *min_df* and *max_df*, respectively ignoring "terms that have a document frequency strictly lower than the given threshold" and "terms that have a document frequency strictly higher than the given threshold" (Sklearn

documentation for TfidfVectorizer). These parameters can be either integers representing an absolute threshold or floats representing a portion of the documents. By default, *min_df* is set to 1 (i.e. not any term is excluded) and *max_df* is set to 1.0 (i.e. again, not any term is excluded by default). The approach chosen to tune these parameters is the following:

1. Use TfidfVectorizer with its default parameters
2. Use a test classification model such as Logistic Regression without any particular tuning (just for reference, i.e. to evaluate the relevance of changes made in data pre-processing)
3. Test a range of values for parameter *min_df* with a fixed *max_df* value and inversely
4. Cross-validate results (e.g. with F1-score) to choose appropriate *min_df* and *max_df*

That approach has been implemented with the following code (and a similar approach has been used for *cross_validate_max_df* method):

```

1 def cross_validate_min_df(self, min_df_range, max_df):
2     mean_scores = []; std_scores = []
3     for min_df in min_df_range:
4         self.preprocess_data(predict="voted_up", min_df=min_df,
5                               ↪ max_df=max_df)
6         X, y = self.get_data()
7         tmp_model = LogisticRegression(C=10, max_iter=1000)
8         scores = cross_val_score(tmp_model, X, y, cv=5, scoring='f1')
9         mean_scores.append(np.array(scores).mean())
10        std_scores.append(np.array(scores).std())

```

A first test has been done with a wide range of floats for *min_df*. However, 0.1 did already leave no term in the features. After focusing on a range of values < 0.1 , it appeared that very small values for *min_df* maximises the F1-score (Figure 1a). It is therefore interesting to test a range of integers representing only a very small amount of documents (Figure 1b):

```

1 # min_df_range = [0.001, 0.005, 0.01, 0.05, 0.075]
2 min_df_range = [1, 5, 10, 15, 20, 25, 50]
3 preprocessor.cross_validate_min_df(min_df_range, max_df=1.0)

```

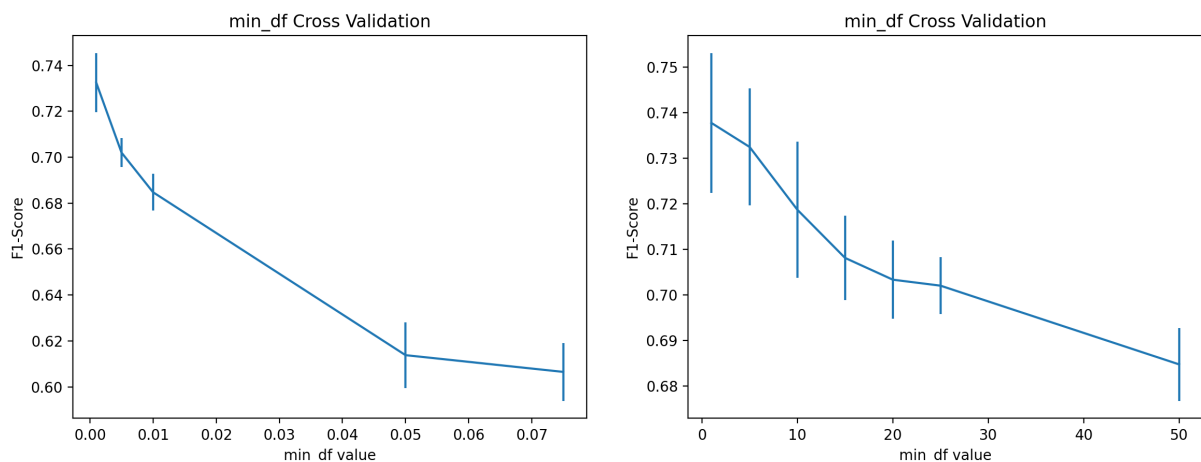


Figure 1: Cross-validation of *min_df*

It is worth noting that the default value of 1 which does not exclude any terms based on their low frequency actually maximises the F1-score. We can therefore argue that keeping $min_df = 1$ is a suitable choice.

A similar reasoning can be applied to max_df , for which choosing a value of 0.05 seems appropriate to maximise the F1-score while minimising its standard deviation.

```
1 # max_df_range = [0.1, 0.2, 0.5, 0.75, 0.9]
2 max_df_range = [0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.2]
3 preprocessor.cross_validate_max_df(min_df=1, max_df_range=max_df_range)
```

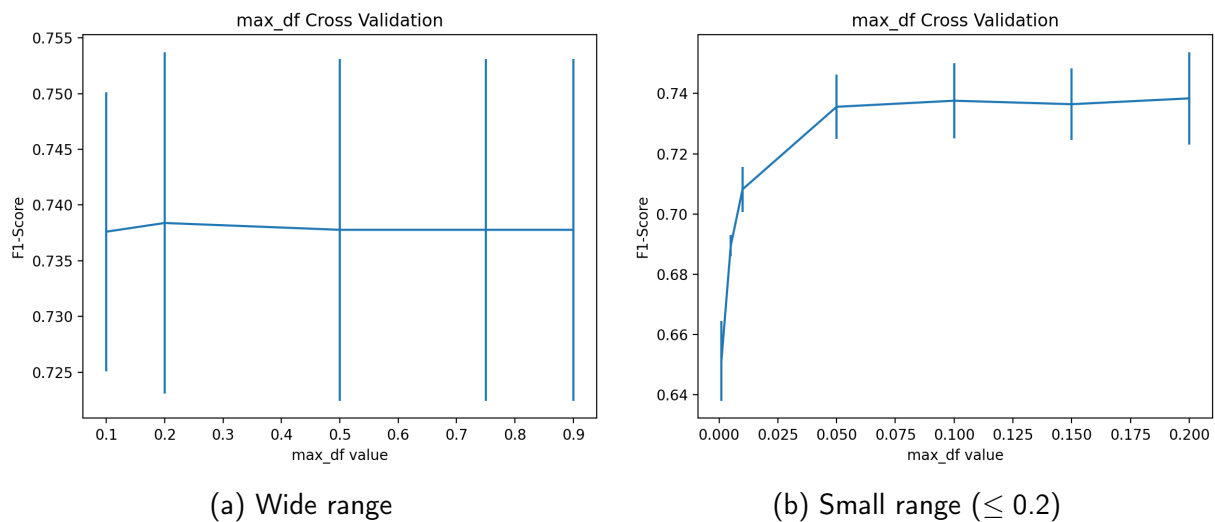


Figure 2: Cross-validation of max_df

1.2.4 Improving the tokenizer

Stopwords

NLTK library provides lists of common stop words in several languages that can be used as a parameter of Sklearn's `TfidfVectorizer`. However, our collection is composed of reviews in many languages and no flag indicates that language. Two approaches have been tested:

- (i) Focusing on stop words only in a single language, for example English
- (ii) Including stop words from several or even all languages available in NLTK

The following code has been used to test approach (i):

```
1 nltk.download('stopwords')
2 stop = list(stopwords.words('english'))
3 self.tfidf = TfidfVectorizer(stop_words=set(stop), min_df=min_df,
4                               ↪ max_df=max_df)
5 self.X = self.tfidf.fit_transform(self.reviews)
6 self.X = self.X.toarray()
```

Tested on a similar Logistic Regression classifier as for `TfidfVectorizer` tuning, approach (i) provided a F1-Score of 0.73, approach (ii) a score of 0.725 and not using any stop words list

provided a F1-Score of 0.73. As their impact is negligible, NLTK's stop words have not been used in this project.

Stemming

TfidfVectorizer accepts custom tokenizer. It is therefore possible to implement NLTK's stemmer which aim is to remove all affixes from words, keeping only their stem. NLTK provides two stemmers: PorterStemmer and SnowballStemmer. Two approaches have been tested to implement stemming:

- (i) Detecting the language of each review and applying only the appropriate stemmer
- (ii) Focusing on stemming only in a single language, for example English

Approach (i) seems to be the most appropriate approach. Some Python libraries such as *langdetect* are meant to detect the language of sentences:

```
1 def tokenize(self, text):
2     try:
3         lang = self.isoToLanguage(langdetect.detect(text))
4     except langdetect.lang_detect_exception.LangDetectException:
5         lang = None
6     tokens = nltk.word_tokenize(text)
7     stems = []
8     for item in tokens:
9         if lang:
10             stems.append(SnowballStemmer(lang).stem(item))
11         else:
12             stems.append(item)
13     return stems
```

After testing, *langdetect*'s approach often decreased F1-score (down to approximately 0.715) and precision, depending on the test sample. An hypothesis is that it could be due to classification errors in language prediction. That approach also increased processing time.

While SnowballStemmer is said to be "better" than the original PorterStemmer in NLTK's documentation, it appeared that applying a generic PorterStemmer to every token of the review text increased the F1-score in this case. Therefore, approach (ii) seemed to be more efficient, leading to a slightly increased F1-score of 0.745:

```
1 def tokenize(self, text):
2     tokens = nltk.word_tokenize(text)
3     stemmer = nltk.PorterStemmer()
4     stems = [stemmer.stem(item) for item in tokens]
5     return stems
```

In conclusion, approach (ii) seems to provide a more stable behaviour over the tests and increases our F1-score. We can therefore argue that keeping a generic PorterStemmer is an appropriate solution.

N-grams

The `TfidfVectorizer` approach provides feature vectors for each token which are, by default, 1-grams (i.e. words). However, that approach does not consider any relation between many words (i.e. sentences, paragraphs, 2 words, etc.).

It is possible to specify to `TfidfVectorizer` the range of n-gram to use. For instance, ranges of (1, 2), (2, 2) and (1, 3) have been tested. However, adding n-grams ranges drastically increases the number of feature vectors, so as the overall processing time.

Ranges of (1, 2) and (2, 2) did not improve either F1-score, accuracy or any other metric of the test Logistic Regression model. (1, 3) range slightly improved the F1-score but not enough to be relevant, especially considering the large increase of processing time (i.e. 477913 features for range (1, 3) against only 43625 for range (1, 1)). For instance, the default (1, 1) range is arguably sufficient.

1.3 Selected classification models

Different machine learning approaches have been tested to predict the review polarity and whether or not the review is for an "early access" version of a game based on the review text. For each model, we can first tune its hyperparameters using cross-validation and then evaluate its performance with different metrics (e.g. F1-score, confusion matrix, ROC/AUC, accuracy). As F1-score combines both precision and recall, it is a great metric to evaluate the performance of a classifier.

For convenience, a class *MyModel* has been defined with the following methods:

- *split(self, test_size=0.20)*: splits the data using Sklearn's *train_test_split* function (only training data is used in model's cross-validation and splitted in training/validation sets; so that test data remains unseen from the model during hyperparameters tuning)
- *train(self)*: trains the model using Sklearn's *model.fit()* function
- *save_model(self, path="")* and *load_model(self, path="")*

Abstract methods *confusion_matrix*, *print_report*, *plot_roc_curve* and *cross_validate_penalty* have also been defined. The idea of such an implementation is to provide a similar interface for every model to ease the usage of different models and the implementation of new ones.

1.3.1 Baseline

The "most frequent" strategy has been chosen as a baseline model. This baseline will act as a reference to evaluate the relevance and performance of actual models.

1.3.2 Logistic regression

Logistic regression classifier aim is to establish a decision boundary $\theta^T x = 0$ that allows it to predict the class of a given input with a certain confidence (i.e. probability). It relies on a cost function $J(\theta)$ that parameters θ need to minimise. It also takes advantage of a L2 penalty that encourages small (but non-zero) parameters values. That penalty can be tuned using C parameter.

The Logistic regression is implemented as following:

```

1 class MyLogisticRegression(MyModel):
2     def __init__(self, X, y):
3         # ...
4     def init_model(self, C=10, max_iterations=10000):
5         self.model = LogisticRegression(C=C, max_iter=max_iterations)
6
7     def cross_validate_penalty(self, C_range, k_fold_nb,
8         ↪ max_iterations=10000):
9         mean_scores = []; std_scores = []
10        for C in C_range:
11            tmp_model = LogisticRegression(C=C, max_iter=max_iterations)
12            scores = cross_val_score(tmp_model, self.X_train,
13            ↪ self.y_train, cv=k_fold_nb, scoring='f1')
14            mean_scores.append(np.array(scores).mean())
15            std_scores.append(np.array(scores).std())
16            plt.errorbar(C_range, mean_scores, yerr=std_scores)
17        # ...
18    def confusion_matrix(self):
19        disp = plot_confusion_matrix(self.model, self.X_test,
20        ↪ self.y_test)
21        # ...
22    def print_report(self):
23        print(classification_report(self.y_test, y_pred))
24        # ...
25    def plot_roc_curve(self, display_plot=True):
26        fpr, tpr, _ = roc_curve(self.y_test, decision_fct)
27        plt.plot(fpr, tpr)
28        # ...

```

Note: some methods displayed above have been shortened to keep them legible; full code is available in appendix. Also, as methods used to print models' metrics are very similar from one model another, they will not be displayed for each further model.

1.3.3 k-Nearest Neighbours (kNN)

kNN is an instance-based model: it bases its predictions directly on the training data. To predict the label of a point, kNN relies on the majority vote of its k closest training points. The distance is commonly measured as an Euclidian distance. Sklearn's KNeighborsClassifier provides a kNN classifier which is implemented in the following class:

```

1 class MyKNeighborsClassifier(MyModel):
2     def __init__(self, X, y):
3         # ...
4     def init_model(self, n=5):
5         self.model = KNeighborsClassifier(n_neighbors=n,
6         ↪ weights='uniform')
7
8     def cross_validate_n(self, n_range=[1, 10, 100, 1000], k_fold_nb=5):

```



```

8     mean_scores = []; std_scores = []
9     for n in n_range:
10         tmp_model = KNeighborsClassifier(n_neighbors=n,
11             ↪ weights='uniform')
12         scores = cross_val_score(tmp_model, self.X_train,
13             ↪ self.y_train, cv=k_fold_nb, scoring='f1')
14         mean_scores.append(np.array(scores).mean())
15         std_scores.append(np.array(scores).std())
16     plt.errorbar(n_range, mean_scores, yerr=std_scores)
17     # ...

```

1.3.4 Multi-Layer Perceptron (MLP)

A MLP is a neural network composed of an input layer, hidden layer(s) and an output layer. The number and size of hidden layers can be determined by cross-validation. In our case, we can use Sklearn's `MLPClassifier` with its standard solver *adam* implementing a "stochastic gradient-based optimizer" (Sklearn documentation for `MLPClassifier`) and the commonly used *ReLU* activation function:

```

1 class MyMLPClassifier(MyModel):
2     def __init__(self, X, y):
3         # ...
4     def init_model(self, n=(25,), max_iterations=10000):
5         self.model = MLPClassifier(hidden_layer_sizes=n,
6             ↪ max_iter=max_iterations)
7
8     def cross_validate_n(self, prev_hidden_layers=(), n_range=[1, 10,
9         ↪ 100, 1000], k_fold_nb=5, max_iterations=10000):
10         mean_scores = []; std_scores = []
11         for n in n_range:
12             hls = prev_hidden_layers + (n,)
13             tmp_model = MLPClassifier(hidden_layer_sizes=hls,
14                 ↪ max_iter=max_iterations)
15             scores = cross_val_score(tmp_model, self.X_train,
16                 ↪ self.y_train, cv=k_fold_nb, scoring='f1')
17             mean_scores.append(np.array(scores).mean())
18             std_scores.append(np.array(scores).std())
19         plt.errorbar(n_range, mean_scores, yerr=std_scores)
20         # ...

```

The cross-validation function takes both the range of new hidden layer sizes and the sizes of previous hidden layers (if there are any). The aim of that implementation is to be able to add and cross-validate easily new layers one after another.

1.4 Predicting the review polarity

1.4.1 Experiments and results

Pre-processing

By using parameter `predict="voted_up"`, `y` (i.e. our output variable) will contain the field `voted_up` of each review and `X` will contain the pre-processed reviews texts.

```
1 preprocessor = MyPreprocessor()
2 preprocessor.read_data('dataset/reviews_17.jl.json')
3 preprocessor.preprocess_data(predict="voted_up", min_df=1, max_df=0.05)
4 X, y = preprocessor.get_data(); preprocessor.print_report()
```

The short data report indicates that 2500 reviews have been `voted_up` over a total of 5000 reviews (i.e. 50%): the data is very well balanced, which is helpful to train classifiers.

Baseline

The baseline model can be used through `MyBaseline` class as following:

```
1 baseline = MyBaseline(X, y)
2 baseline.split(test_size=0.20)
3 baseline.init_model()
4 baseline.train()
5 baseline.confusion_matrix()
6 baseline.print_report()
```

This baseline provides the following metrics and the confusion matrix displayed in Figure 3: Precision = 0.25, Recall = 0.5, Accuracy = 0.5 and F1-score = 0.33.

It is interesting to note that the accuracy obtained reflects well our 50/50 balanced data.

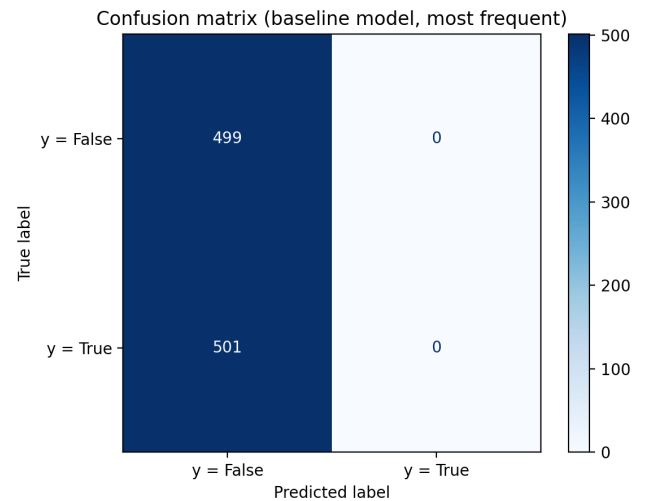


Figure 3: Baseline's confusion matrix

Logistic regression

We can use the following code to cross-validate hyperparameter `C` (see Figure 4 below):

```
1 logreg = MyLogisticRegression(X, y); logreg.split(test_size=0.20)
2 logreg.cross_validate_penalty(C_range=[0.1, 1, 10, 100, 1000],
  ↪ k_fold_nb=5)
```

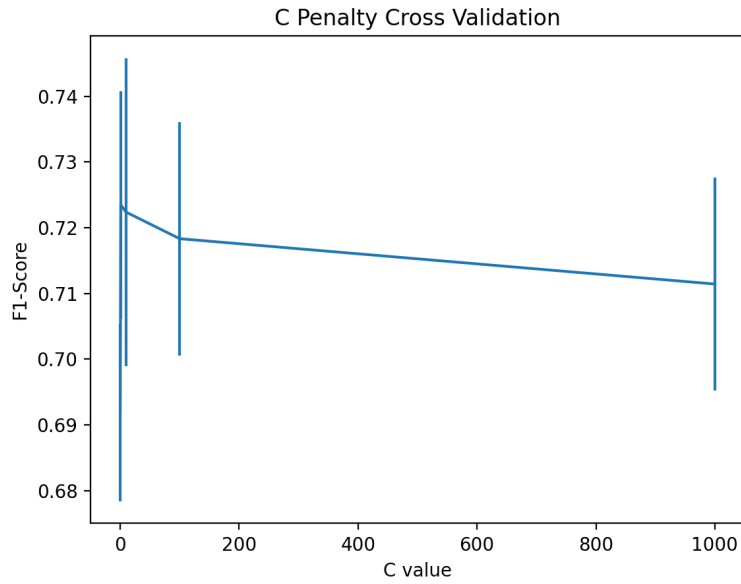
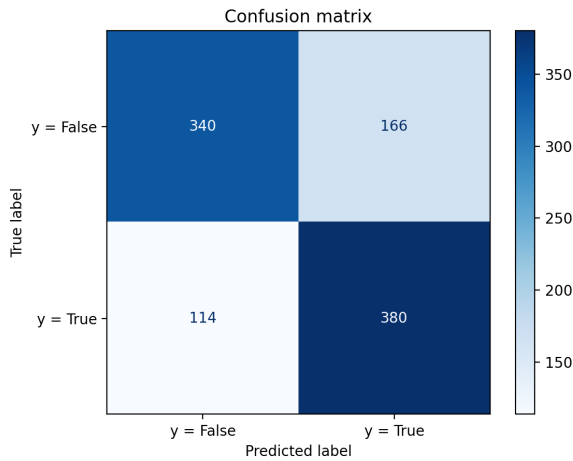


Figure 4: Logistic regression - C cross-validation

As $C = 10$ seems to maximise F1-score while preventing an important spread of the error (i.e. standard deviation), we can train and test the model with that hyperparameter value:

```
1 logreg.init_model(C=10, max_iterations=10000)
2 logreg.train(); logreg.confusion_matrix(); logreg.print_report()
```



(a) Confusion matrix

```
> Classification report:
```

	precision	recall	f1-score	support
Not voted up (y=-1)	0.78	0.64	0.70	508
Voted up (y=1)	0.68	0.82	0.74	492
accuracy			0.72	1000
macro avg	0.73	0.73	0.72	1000
weighted avg	0.73	0.72	0.72	1000

(b) Classification report

Figure 5: Logistic regression results

kNN

When cross-validating the kNN, it appears that a number of neighbours between 1 and 10 provides the highest F1-score as shown in Figure 6 below:

```
1 knn = MyKNeighborsClassifier(X, y); knn.split(test_size=0.20)
```

```

2 # knn.cross_validate_n(n_range=[1, 5, 10, 50, 100, 500, 1000],
  ↪ k_fold_nb=5)
3 knn.cross_validate_n(n_range=[1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
  ↪ k_fold_nb=5)

```

Therefore $n = 3$ neighbours is the most appropriate value as it maximises the F1-score while minimising the standard deviation. Our kNN can be trained using the code below (results Figure 7):

```

1 knn.init_model(n=3)
2 knn.train()
3 knn.confusion_matrix()
4 knn.print_report()

```

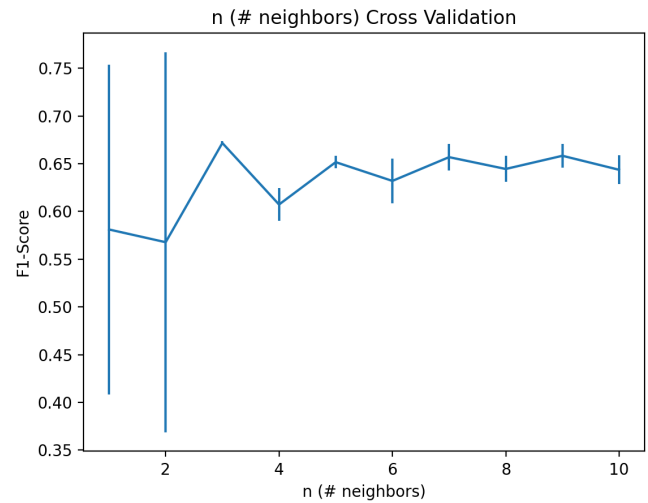
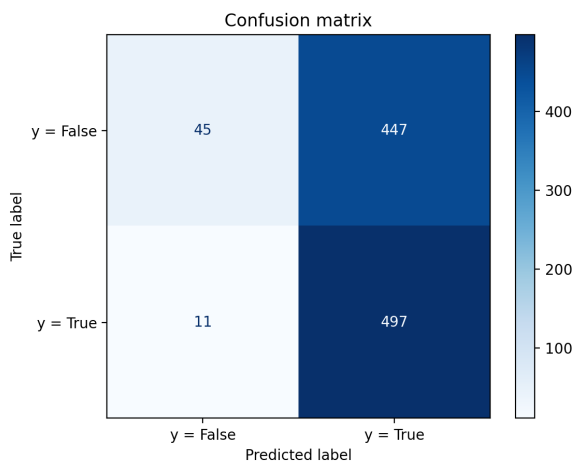


Figure 6: kNN Cross-validation



(a) Confusion matrix

	precision	recall	f1-score	support
Not voted up (y=-1)	0.81	0.11	0.20	496
Voted up (y=1)	0.53	0.97	0.68	504
accuracy			0.55	1000
macro avg	0.67	0.54	0.44	1000
weighted avg	0.67	0.55	0.44	1000

(b) Classification report

Figure 7: kNN results

MLP

MLP layers can be cross-validated using the following code. The idea is to begin by cross-validating the size of a single hidden layer. Then, we can iterate to cross-validate the size of a potential second layer using that first cross-validated layer size.

```

1 mlp = MyMLPClassifier(X, y); mlp.split(test_size=0.20)
2 # mlp.cross_validate_n(prev_hidden_layers=(), n_range=[1, 5, 10, 25, 50,
  ↪ 100, 500, 1000], max_iterations=10000)
3 mlp.cross_validate_n(prev_hidden_layers=(50,), n_range=[1, 5, 10, 25, 50,
  ↪ 100], max_iterations=10000)

```

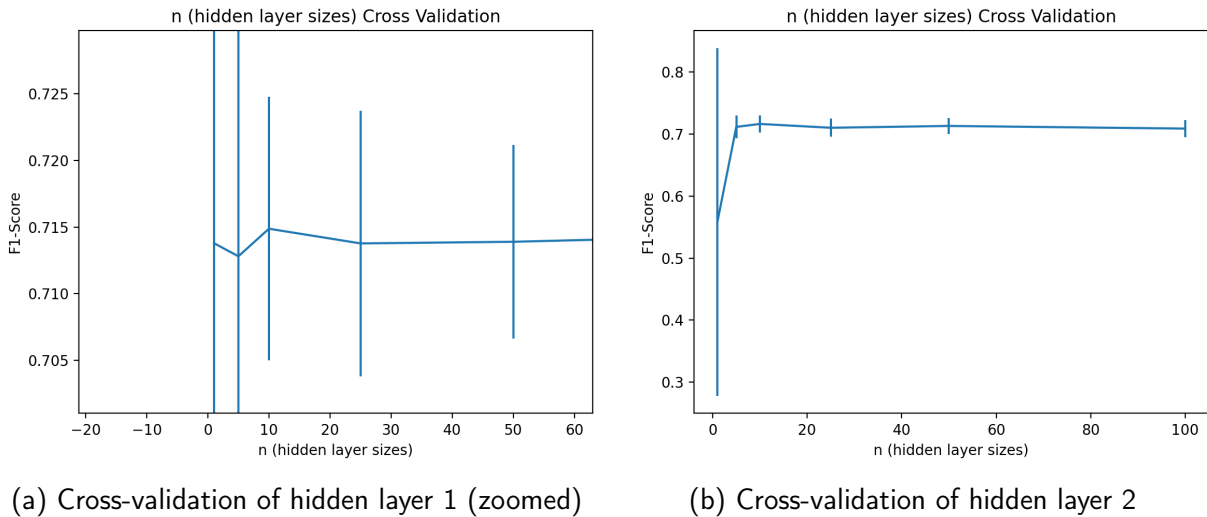


Figure 8: MLP Cross-validation

We can see from the cross-validation of the first layer (Figure 8a) that size ≥ 10 gives the highest F1-scores. To minimise standard deviation while keeping the model simple (i.e. which helps preventing over-fitting), we could choose $n = 50$ for the first layer. Adding new layers and/or nodes can allow our MLP to capture more behaviours/patterns in our training data (although too much nodes/layers can lead to over-fitting). As in this case adding a second hidden layer (Figure 8b) does not provide any improvement of the F1-score, we can keep a single hidden layer of size 50.

```
1 mlp.init_model(n=(50,))
2 mlp.train(); mlp.confusion_matrix(); mlp.print_report()
```

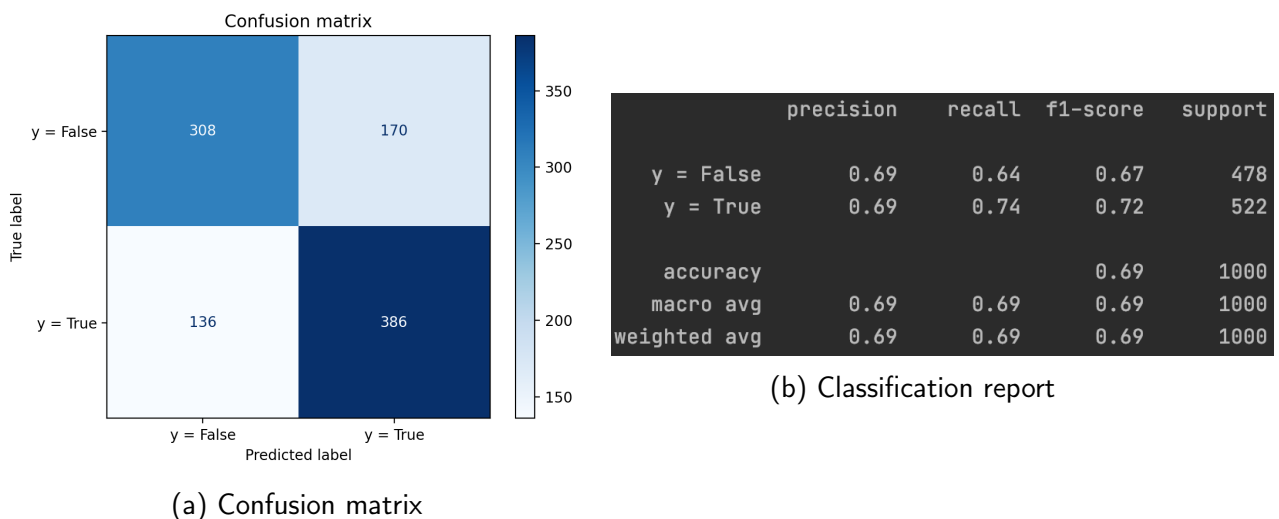


Figure 9: MLP results

1.4.2 Discussion of the results

ROC curves and AUC shown in Figure 10 beside are great metrics to evaluate the performance of our different models. The ROC curve is given by True Positive Rate vs False Positive Rate. As an ideal classifier would give 100% of true positive rate and 0% of false positive rate, we want a classifier which ROC curve comes as close as possible to the upper-left corner of the plot.

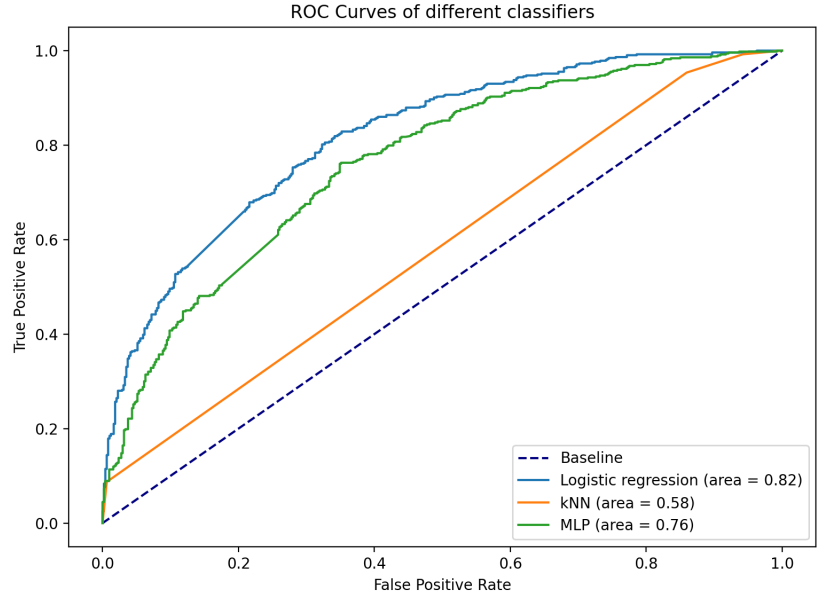


Figure 10: ROC Curves of different classifiers

We can see that the kNN provides only slightly better classification performances than our baseline model. Therefore, we can argue that kNN is not suitable to predict the polarity of a review.

MLP and Logistic Regression are much closer to each other while providing a largely better ROC curve with AUCs of respectively 0.76 and 0.82. We can establish a small comparison table from the previous results (Figures 5 and 9) and our baseline model:

	Baseline	MLP	Logistic Regression
Precision	0.25	0.69	0.73
Recall	0.5	0.69	0.73
F1-score	0.33	0.69	0.72
Accuracy	0.5	0.69	0.72

Table 1: Comparison of Baseline, MLP and Logistic Regression (review polarity prediction)

We can see that both MLP and Logistic Regression provide great improvements over the baseline model: they both offer an increase of at least 175% of the precision, 38% of the recall, 109% of the F1-score and 38% of the accuracy.

Logistic Regression seems to provide overall better classification performance than MLP. In addition to that, Logistic Regression provides a few advantages over MLP:

1. Logistic Regression is easier and faster to train
2. Logistic Regression is sort of interpretable given the values of θ (i.e. to evaluate which features/tokens of the text have more importance to predict the polarity of the review), whereas MLP acts more as a "black box" which cannot be interpreted

In conclusion, we can argue that the Logistic Regression model is a suitable classifier to predict the polarity of a review based on its text.

1.5 Predicting if the review is for an "early access"

The same approach as for previous section 1.4 will be used to determine whether or not a review is for an "early access" version of a game. Therefore, code snippets will not be displayed again for readability. The results of tested models will still be given and discussed in the last part of this section.

1.5.1 Experiments and results

Pre-processing

After pre-processing the data, it is interesting to note that only 518 reviews have `early_access = True` on a total of 5000 reviews (i.e. only 10%): the data is highly imbalanced. This bias may make the process of learning more difficult to classification models as they could easily tend toward the majority class.

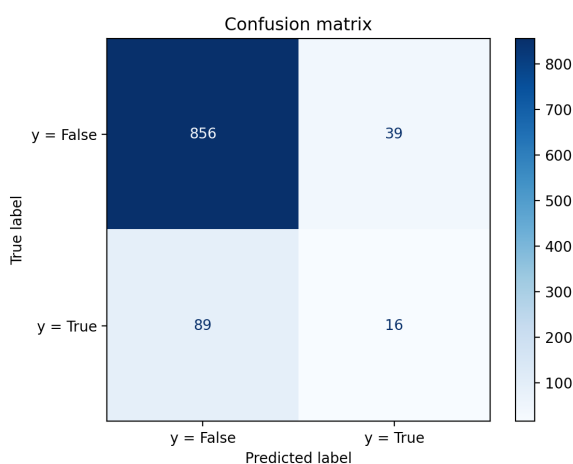
Baseline

The "most frequent" baseline model gives the following metrics for `early_access` predictions: Precision = 0.43, Recall = 0.5, Accuracy = 0.86 and F1-score = 0.46.

The accuracy of that baseline model also reflects well how imbalanced is our data. Indeed, we have only 318/5000 "early access" reviews. In other terms, around 90% of our reviews are not "early access" and can be easily classified by that baseline's "most frequent" strategy.

Logistic regression

Logistic regression classifier's penalty has been cross-validated on a wide range of values from 0 to 5000. $C \leq 100$ provide a F1-score close to 0, while $C \geq 1000$ is required to "maximise" F1-score which peak (averaged) is around 0.12 (i.e. we can expect poor predictions from such a low score). Once trained with $C = 1000$, we obtain the following results:



(a) Confusion matrix

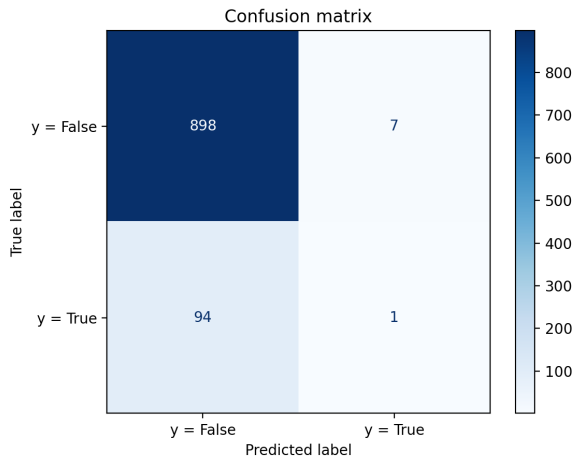
	precision	recall	f1-score	support
y = False	0.91	0.96	0.93	895
y = True	0.29	0.15	0.20	105
accuracy			0.87	1000
macro avg	0.60	0.55	0.57	1000
weighted avg	0.84	0.87	0.85	1000

(b) Classification report

Figure 11: Logistic Regression results

kNN

When cross-validating number of neighbours to consider in kNN classifier on a wide range of values (i.e. from 1 to 1000), the F1-score quickly drops to 0 when $C > 5$. When focusing on a shorter range of values, it seems that $n = 1$ provides the highest F1-score and gives the following classification results:



(a) Confusion matrix

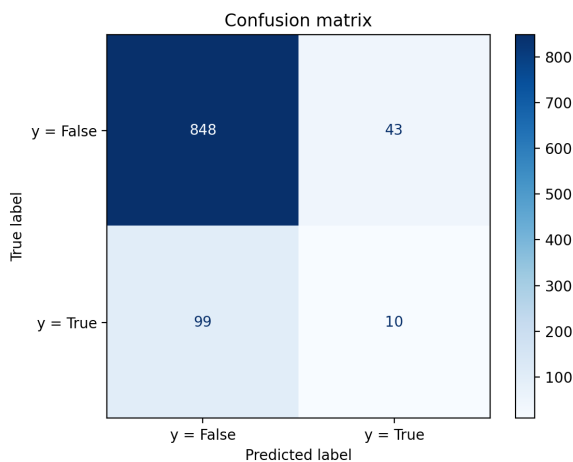
	precision	recall	f1-score	support
y = False	0.91	0.99	0.95	905
y = True	0.12	0.01	0.02	95
accuracy			0.90	1000
macro avg	0.52	0.50	0.48	1000
weighted avg	0.83	0.90	0.86	1000

(b) Classification report

Figure 12: kNN results

MLP

When determining the size of the first hidden layer of our MLP, using a size of $n = 50$ seems to maximise F1-score (approximately 0.125) while minimising the standard deviation. However, adding a new hidden layer does not improve the F1-score of our classifier: it is not able to capture further pattern/behaviour in our training data. Therefore, we can keep a single hidden layer of size $n = 50$ for this MLP classifier:



(a) Confusion matrix

	precision	recall	f1-score	support
y = False	0.90	0.95	0.92	891
y = True	0.19	0.09	0.12	109
accuracy			0.86	1000
macro avg	0.54	0.52	0.52	1000
weighted avg	0.82	0.86	0.84	1000

(b) Classification report

Figure 13: MLP results

1.5.2 Discussion of the results

Looking at the results of the tested models and the ROC curves displayed in Figure 14 beside, we can see that no model performs much better than the baseline.

The kNN's ROC curve and AUC are so close to the baseline that it is not relevant to consider this classifier; although it provides a F1-score and an accuracy which are slightly better than the baseline model (respectively 0.48 and 0.90).

Actually, the training data is so imbalanced that the "most frequent" strategy of our baseline model already provides very high F1-score (0.46) and accuracy (0.86) compared to the tested models.

MLP and Logistic Regression classifiers seem to perform better than the kNN. However, they do still not provide a big difference compared to the baseline. Again, we can compare MLP and Logistic Regression to our baseline model:

	Baseline	MLP	Logistic Regression
Precision	0.43	0.54	0.6
Recall	0.5	0.52	0.55
F1-score	0.46	0.52	0.57
Accuracy	0.86	0.86	0.87

Table 2: Comparison of Baseline, MLP and Logistic Regression (early access prediction)

As well as in section 1.4, Logistic Regression provides overall better classification performance than MLP. However, that model gives at most at increase of 24% of the baseline's F1-score and only 1% increase in accuracy.

In conclusion, none of the models seem to be suitable to predict whether a review has been given to an "early access" version of a game or not.

1.6 Conclusions and reflections

- (i) We have shown that machine learning approaches such as MLP and Logistic Regression give significant classification performance to predict the review polarity (whether a game

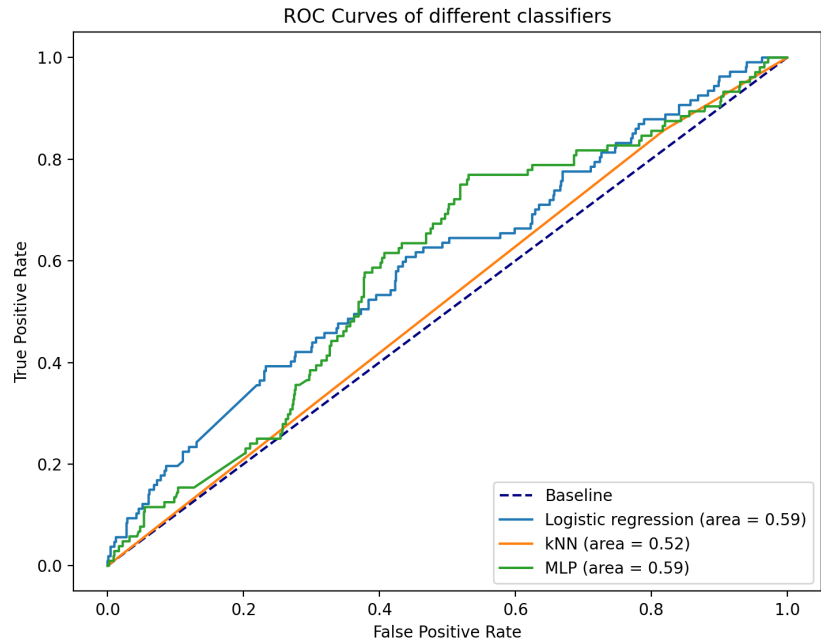


Figure 14: ROC Curves of different classifiers

has been “voted up” or not by the reviewer) based on the review text.

We can therefore conclude that the review text can be used to predict the review polarity.

However, these models could still be improved to offer even more robust classification performance. Further work would imply:

- Testing other approaches (e.g. SVM, Decision Tree Classifier, etc.).
 - Considering the reviews languages (e.g. using different models depending on the language and/or tuning the model dependently of the language).
 - Considering different pre-processing approaches such as Word2Vec or GloVe.
- (ii) None of the models tested have been able to accurately predict whether the review is for an “early access” version of a game or not based on the review text. It would be almost as accurate to say that a review is *always* for a non-early access version of a game than using one of the tested machine learning approaches.
- We can therefore conclude that the review text cannot be used to predict whether the review is for an “early access” version of a game or not.

However, it is important to note that the given data was highly imbalanced and that could have impacted our classifiers. Therefore, it would be possible to make further work on that question to potentially obtain better classification performance:

- Undersampling could be used to decrease the number of reviews where *early_access* = *False* so that the data becomes more balanced. However, undersampling needs to be carefully done to prevent from adding bias into the data (e.g. not removing all reviews in a specific language). Also, balancing the data by undersampling would lead to a total of 1036 reviews, which is quite small compared to the original 5000 reviews.
- A better approach would be oversampling: instead of decreasing the number of reviews where *early_access* = *False*, we could increase the number of reviews where *early_access* = *True* by including new reviews to the dataset.

2 Part 2.

2.1 Question (i)

We can speak of *under-fitting* when the model is too simple to catch the behaviour of our data (i.e. to fit to the shape of the data), thus leading to poor predictions. For example, trying to make prediction on quadratic-shaped data using a linear model would be a case of under-fitting. Using a more complex model, adding features (e.g. polynomial features) and tuning model's hyperparameters (e.g. penalty) can help preventing under-fitting.

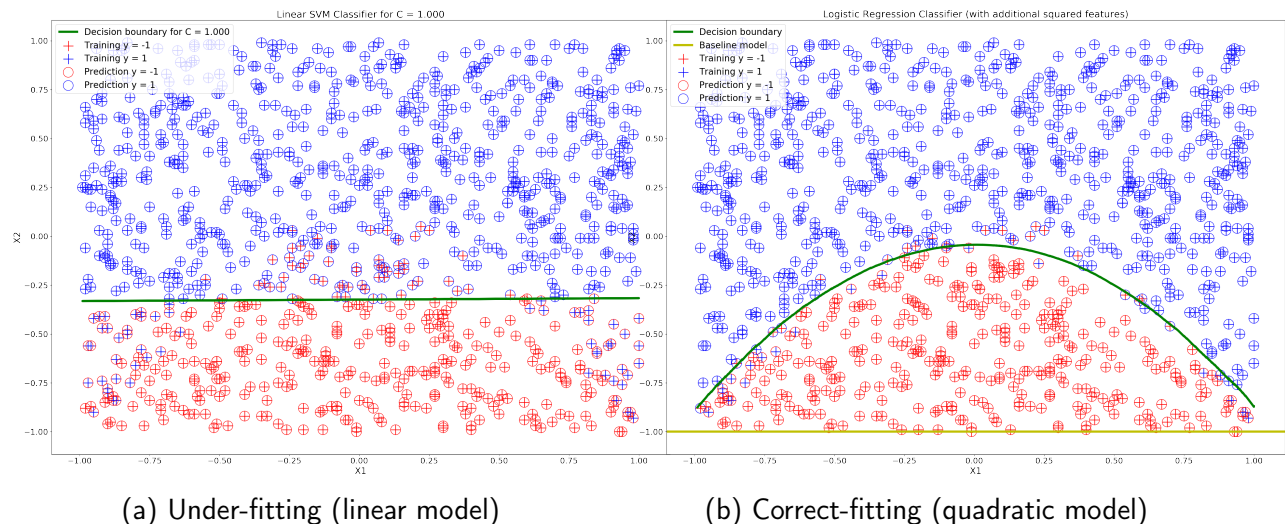
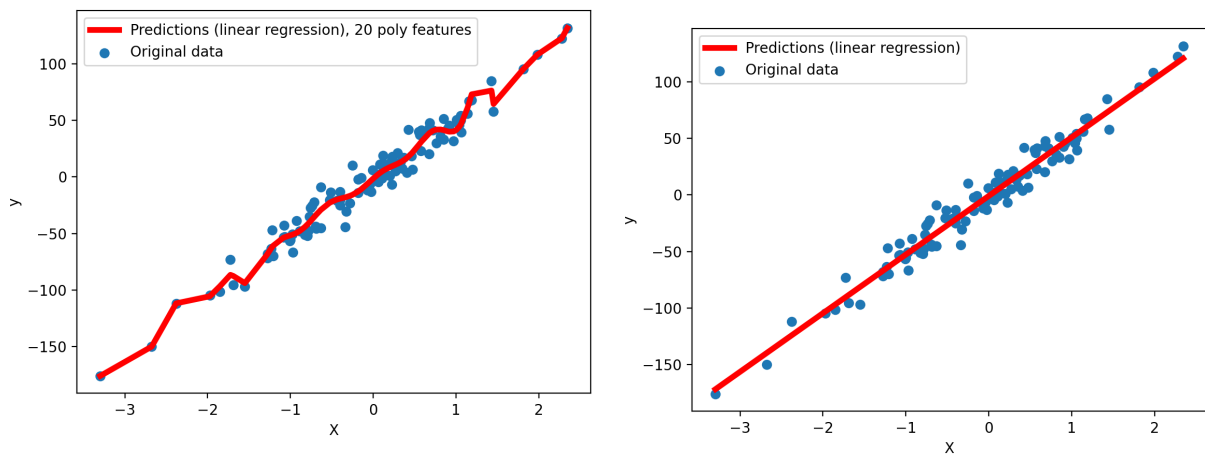


Figure 15: Example of under-fitting (*Source: Week 2 assignment*)

On the contrary, *over-fitting* occurs when the model is too complex and fits too much the noise of the training data. In other words, it means that the model is according too much importance to the noise of the training data and generalises poorly outside the scope of the training data. For example, using a model with too much polynomial features in a regression model can easily lead to over-fitting. Using a simpler model, using less features and tuning model's hyperparameters can help preventing over-fitting; a rule of thumb is to use the simplest model/parameters when many provide the same scores (e.g. F1-score, MSE, etc.).

The following code can be used to generate an example of over-fitting with too many polynomial features (see result Figure 16):

```
1 X, y = make_regression(n_samples=100, n_features=1, noise=10)
2 X, y = zip(*sorted(zip(X, y))) # Sort data
3 reg = LinearRegression().fit(X, y)
4 plt.scatter(X, y)
5 plt.plot(X, reg.predict(X))
6 # ...
7 reg = make_pipeline(PolynomialFeatures(20), LinearRegression())
8 reg.fit(X, y)
9 plt.scatter(X, y)
10 plt.plot(X, reg.predict(X))
```



(a) Over-fitting (polynomial features up to degree 20)

(b) Correct-fitting (no polynomial features)

Figure 16: Example of over-fitting

2.2 Question (ii)

The following algorithm shows an implementation of k-fold cross-validation from scratch.

Notes:

- `array[x : y]` retrieves indexes x (included) to y (excluded) of the array.
- *hyperparameter* is the hyperparameter to cross-validate (e.g. a penalty, number of neighbours, etc.)
- *scoreFunction* represents a scoring function such as MSE, F1-score, etc.
- *meanScores* is used in this example but other metrics such as standard deviation of the scores or errors could also be considered

Algorithm 1 k-fold cross-validation

```

1:  $X, y \leftarrow$  Original data
2:  $nbDataPoints \leftarrow \text{length}(X)$ 
3:  $k \leftarrow$  number of folds
4:
5: // 1. Split data into k folds
6: for  $i = 0, \dots, k - 1$  do
7:    $offset \leftarrow i * nbDataPoints / k$ 
8:    $maxIndex \leftarrow offset + nbDataPoints / k$ 
9:    $folds[i][X] \leftarrow X[offset : maxIndex]$ 
10:   $folds[i][y] \leftarrow y[offset : maxIndex]$ 
11: end for
12:
13: // 2. Perform cross-validation
14:  $meanScores \leftarrow []$ 
15:  $range \leftarrow [1, 10, 100, 1000]$ 
16: for  $value$  in  $range$  do
17:    $tmpScores \leftarrow []$ 
18:    $model \leftarrow \text{TheModel}(\text{hyperparameter} = value),$ 
19:   for  $i = 0, \dots, k - 1$  do
20:      $train \leftarrow \text{concatenate}(folds[0 : i], folds[i + 1 : k])$ 
21:      $test \leftarrow folds[i]$ 
22:      $model.train(train[X], train[y])$ 
23:      $predictions \leftarrow model.predict(test[X])$ 
24:      $tmpScores.append(\text{scoreFunction}(test[y], predictions))$ 
25:   end for
26:    $meanScores.append(\text{mean}(tmpScores))$ 
27: end for
28:
29: // 3. Evaluate
30:  $\text{Plot}(meanScores)$ 

```

2.3 Question (iii)

Tuning model's hyperparameter(s) can help finding the right trade-off between over- and under-fitting, for example by adding penalty to some parameters of the model. A straightforward approach consists in splitting the data in a training set and a testing set (e.g. 80/20 split) to then train and evaluate the model over a range of hyperparameters values. However, if either the training data or the test data is noisy, hyperparameters' may be tuned inadequately, thus leading to poor predictions.

K-fold cross-validation provides a way to evaluate a model by equally considering the data provided. That is, training the model each time with a different slice (i.e. fold) of the original data. That process allows to smooth out the noise by averaging the scores (e.g. MSE, F1-score) of the trained model. Therefore, k-fold cross-validation is a great technique for tuning model's hyperparameters which are directly correlated with the balance between over- and under-fitting.

2.4 Question (iv)

An advantage of a kNN classifier over a logistic regression classifier is that kNNs are very easy to use, simple and effective. Indeed, a kNN only requires its number of neighbours k to be set properly and it then provides predictions based directly on the training data (i.e. kNNs are instance- or data-based models). Logistic regression classifiers are a bit more complex as they rely on the minimisation of a cost function to make predictions.

However, the fact that kNNs rely on the data itself is also a downside compared to logistic regressions. As a kNN requires to consider all data points available to make predictions, it is a good solution for small data but it can quickly become much slower than a logistic regression classifier because of its large computation cost at runtime if the amount of data is important.

Finally, another con of kNNs is that they are not good at making predictions when stepping away from training data (i.e. outliers). This is due to the fact that kNNs rely directly on the training data and "stepping away" from it implies that less data points are available in that area; that is, less relevant neighbours to consider. On the contrary, a logistic regression classifier can easily extrapolate on the training data to make predictions above that scope (i.e. by "extending" its decision boundary based on its parameters θ).

2.5 Question (v)

kNN classifiers are highly sensitive to imbalanced data. When making predictions on imbalanced data, a kNN would tend to predict the majority class.

That sensitivity to imbalanced data is directly correlated with the measurement of distances and the fact that kNNs are instance-based models: the more imbalanced is the original data, the more nearest neighbours of one class we can potentially find close to the data point to predict; thus leading to tend to predict the majority class.

A kNN would also give inaccurate predictions for high-dimensions data (i.e. data having an important number of features). For instance, if we consider two data points of dimension 2 $x_1 = (14, 13)$ and $x_2 = (19, 0)$, they are both at a distance of approximately $d = 19$ from the origin $(0, 0)$. However, we might want our kNN to find points which are closer on every axis instead of just a single one. The more features we add, the more difficult it becomes for our kNN to determine the closest (relevant) neighbours of a point. This challenge is also known as the "Curse of Dimensionality".

That may also explain why kNNs have not been able to provide accurate predictions for both parts (i) and (ii) of the first part of this assignment, considering the large number of features provided.

A Appendix

A.1 Python Code

A.1.1 Class MyPreprocessor

```

1  import jsonlines
2  import numpy as np
3  from sklearn.feature_extraction.text import TfidfVectorizer
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.model_selection import cross_val_score
6  import matplotlib.pyplot as plt
7  import nltk
8  from nltk.stem.snowball import SnowballStemmer
9  import langdetect
10
11
12  class MyPreprocessor:
13      def __init__(self):
14          self.reviews = []
15          self.voted_up = []
16          self.early_access = []
17          self.X = None
18          self.y = None
19          self.tfidf = None
20
21      def read_data(self, path):
22          print("> Reading data")
23          for item in jsonlines.open(path, 'r'):
24              self.reviews.append(item['text'])
25              self.voted_up.append(item['voted_up'])
26              self.early_access.append(item['early_access'])
27
28      def cross_validate_min_df(self, min_df_range, max_df):
29          mean_scores = []; std_scores = []
30          for min_df in min_df_range:
31              self.preprocess_data(predict="voted_up", min_df=min_df,
32                                  ↪ max_df=max_df)
33              X, y = self.get_data()
34              print("\t> For min_df = %.2f..." % min_df)
35              tmp_model = LogisticRegression(C=10, max_iter=1000)
36              scores = cross_val_score(tmp_model, X, y, cv=5, scoring='f1')
37              mean_scores.append(np.array(scores).mean())
38              std_scores.append(np.array(scores).std())
39
40          print("> min_df cross validation done")
41          plt.figure()
42          plt.errorbar(min_df_range, mean_scores, yerr=std_scores)

```

```
42     plt.title("min_df Cross Validation")
43     plt.xlabel("min_df value")
44     plt.ylabel("F1-Score")
45     plt.show()
46
47     def cross_validate_max_df(self, min_df, max_df_range):
48         mean_scores = []; std_scores = []
49         for max_df in max_df_range:
50             self.preprocess_data(predict="voted_up", min_df=min_df,
51                                   ↪ max_df=max_df)
52             X, y = self.get_data()
53             print("\t> For max_df = %.2f..." % max_df)
54             tmp_model = LogisticRegression(C=10, max_iter=1000)
55             scores = cross_val_score(tmp_model, X, y, cv=5, scoring='f1')
56             mean_scores.append(np.array(scores).mean())
57             std_scores.append(np.array(scores).std())
58
59         print("> max_df cross validation done")
60         plt.figure()
61         plt.errorbar(max_df_range, mean_scores, yerr=std_scores)
62         plt.title("max_df Cross Validation")
63         plt.xlabel("max_df value")
64         plt.ylabel("F1-Score")
65         plt.show()
66
67     def isoToLanguage(self, lang):
68         if lang == "da":
69             return "danish"
70         elif lang == "nl":
71             return "dutch"
72         elif lang == "en":
73             return "english"
74         elif lang == "fi":
75             return "finnish"
76         elif lang == "fr":
77             return "french"
78         elif lang == "de":
79             return "german"
80         elif lang == "hu":
81             return "hungarian"
82         elif lang == "it":
83             return "italian"
84         elif lang == "pt":
85             return "portuguese"
86         elif lang == "ro":
87             return "romanian"
88         elif lang == "ru":
89             return "russian"
```



```
89         elif lang == "es":
90             return "spanish"
91         elif lang == "sv":
92             return "swedish"
93         else:
94             return None
95
96     def tokenize(self, text, use_lang_detect=False):
97         tokens = nltk.word_tokenize(text)
98         stems = []
99
100        if use_lang_detect:
101            try:
102                lang = langdetect.detect(text)
103                lang = self.isoToLanguage(lang)
104            except langdetect.lang_detect_exception.LangDetectException:
105                lang = None
106
107            for item in tokens:
108                if lang:
109                    stems.append(SnowballStemmer(lang).stem(item))
110                else:
111                    stems.append(item)
112        else:
113            stemmer = nltk.PorterStemmer()
114            stems = [stemmer.stem(item) for item in tokens]
115
116        return stems
117
118    def preprocess_data(self, predict="voted_up", min_df=1, max_df=1.0):
119        self.tfidf = TfidfVectorizer(tokenizer=self.tokenize,
120        ↪ min_df=min_df, max_df=max_df, ngram_range=(1, 1))
121        self.X = self.tfidf.fit_transform(self.reviews)
122        self.X = self.X.toarray()
123
124        if predict == "voted_up":
125            self.y = self.voted_up
126        else:
127            self.y = self.early_access
128
129    def get_data(self):
130        return self.X, self.y
131
132    def print_report(self, print_features_names=False):
133        print("> Features:")
134        print("\t> Number:", len(self.tfidf.get_feature_names()))
135        if print_features_names:
136            print("\t> Names:", self.tfidf.get_feature_names())
```

```

136
137     print("\n> Total number of reviews:", len(self.reviews))
138
139     voted_up_reviews = 0
140     for voted_up in self.voted_up:
141         if voted_up:
142             voted_up_reviews += 1
143
144     print("> Total number of \"voted_up\" reviews:",
145           ↪ voted_up_reviews)
146
147     early_access_reviews = 0
148     for early_access in self.early_access:
149         if early_access:
150             early_access_reviews += 1
151
152     print("> Total number of \"early_access\" reviews:",
153           ↪ early_access_reviews)

```

A.1.2 Class MyModel

```

1  import abc
2  from joblib import dump, load
3  from sklearn.model_selection import train_test_split
4
5
6  class MyModel:
7      def save_model(self, path=""):
8          if not path:
9              path = "saved-models/" + self.__class__.__name__ + ".joblib"
10
11          dump(self.model, path)
12
13      def load_model(self, path=""):
14          if not path:
15              path = "saved-models/" + self.__class__.__name__ + ".joblib"
16
17          self.model = load(path)
18
19      def split(self, test_size=0.20):
20          self.X_train, self.X_test, self.y_train, self.y_test =
21          ↪ train_test_split(self.X, self.y, test_size=test_size) #
22          ↪ random_state=42
23
24      def train(self):
25          self.model.fit(self.X_train, self.y_train)
26
27      @abc.abstractmethod

```

```
26     def init_model(self):
27         pass
28
29     @abc.abstractmethod
30     def confusion_matrix(self):
31         pass
32
33     @abc.abstractmethod
34     def print_report(self):
35         pass
36
37     @abc.abstractmethod
38     def plot_roc_curve(self):
39         pass
40
41     @abc.abstractmethod
42     def cross_validate_penalty(self):
43         pass
```

A.1.3 Class MyBaseline

```
1  from models.MyModel import MyModel
2  from sklearn.dummy import DummyClassifier
3  from sklearn.metrics import accuracy_score
4  from sklearn.metrics import precision_recall_fscore_support
5  import matplotlib.pyplot as plt
6  from sklearn.metrics import plot_confusion_matrix
7
8
9  class MyBaseline(MyModel):
10     def __init__(self, X, y):
11         self.X = X
12         self.y = y
13         self.X_train = None
14         self.X_test = None
15         self.y_train = None
16         self.y_test = None
17         self.model = None
18
19     def init_model(self):
20         self.model = DummyClassifier(strategy="most_frequent")
21
22     def cross_validate_penalty(self):
23         print("No cross-validation available for Baseline")
24
25     def confusion_matrix(self):
26         print("> Confusion matrix & scores")
```

```

27     disp = plot_confusion_matrix(self.model, self.X_test,
    ↪     self.y_test, display_labels=["y = False", "y = True"],
    ↪     cmap=plt.cm.Blues, values_format='d')
28     disp.ax_.set_title("Confusion matrix (baseline model, most
    ↪     frequent)")
29     plt.show()
30
31     def print_report(self):
32         y_pred = self.model.predict(self.X_test)
33         precision, recall, fscore, _ =
    ↪         precision_recall_fscore_support(self.y_test, y_pred,
    ↪         average="macro")
34         print("> Baseline scores (strategy: most frequent):")
35         print("Precision (avg):", precision, ", Recall (avg):", recall,
    ↪         ", F1-Score (avg):", fscore)
36         print("Accuracy:", accuracy_score(self.y_test, y_pred))
37
38     def plot_roc_curve(self, display_plot=True):
39         """Plot ROC curve and ROC area of the model"""
40         plt.figure(123)
41         plt.plot([0, 1], [0, 1], color='navy', linestyle='--',
    ↪         label="Baseline")
42         plt.xlabel('False Positive Rate')
43         plt.ylabel('True Positive Rate')
44         plt.title('ROC Curve - Baseline')
45         plt.legend(loc="lower right")
46         if display_plot:
47             plt.show()

```

A.1.4 Class MyLogisticRegression

```

1  from models.MyModel import MyModel
2  from sklearn.model_selection import cross_val_score
3  import matplotlib.pyplot as plt
4  from sklearn.linear_model import LogisticRegression
5  from sklearn.metrics import plot_confusion_matrix
6  import numpy as np
7  from sklearn.metrics import classification_report
8  from sklearn.metrics import roc_curve, auc
9
10
11  class MyLogisticRegression(MyModel):
12      def __init__(self, X, y):
13          self.X = X
14          self.y = y
15          self.X_train = None
16          self.X_test = None
17          self.y_train = None

```

```

18     self.y_test = None
19     self.model = None
20
21     def init_model(self, C=10, max_iterations=10000):
22         self.model = LogisticRegression(C=C, max_iter=max_iterations)
23
24     def cross_validate_penalty(self, C_range=[1, 10, 100, 1000],
25     ↪ k_fold_nb=5, max_iterations=10000):
26         print("> C Cross validation (range:", C_range, ")")
27         mean_scores = []; std_scores = []
28
29         for C in C_range:
30             print("\t> For C = %.2f..." % C)
31             tmp_model = LogisticRegression(C=C, max_iter=max_iterations)
32             scores = cross_val_score(tmp_model, self.X_train,
33             ↪ self.y_train, cv=k_fold_nb, scoring='f1')
34             mean_scores.append(np.array(scores).mean())
35             std_scores.append(np.array(scores).std())
36
37         print("> C Cross validation done")
38         plt.figure()
39         plt.errorbar(C_range, mean_scores, yerr=std_scores)
40         plt.title("C Penalty Cross Validation")
41         plt.xlabel("C value")
42         plt.ylabel("F1-Score")
43         plt.show()
44
45     def confusion_matrix(self):
46         print("> Confusion matrix:")
47         plt.figure()
48         disp = plot_confusion_matrix(self.model, self.X_test,
49         ↪ self.y_test, display_labels=["y = False", "y = True"],
50         ↪ cmap=plt.cm.Blues,
51         ↪ values_format='d')
52         disp.ax_.set_title("Confusion matrix")
53         plt.show()
54
55     def print_report(self):
56         y_pred = self.model.predict(self.X_test)
57         print("> Classification report:")
58         print(classification_report(self.y_test, y_pred, target_names=["y
59         ↪ = False", "y = True"]))
60
61         print("> Intercept:")
62         print(self.model.intercept_)
63         print("> Coefficients:")
64         print(self.model.coef_)

```

```

61     def plot_roc_curve(self, display_plot=True):
62         """Plot ROC curve and ROC area of the model"""
63         decision_fct = self.model.decision_function(self.X_test)
64         fpr, tpr, _ = roc_curve(self.y_test, decision_fct)
65         plt.figure(123)
66         plt.plot(fpr, tpr, label='Logistic regression (area = %0.2f)' %
67                 ↪ auc(fpr, tpr))
68         plt.xlabel('False Positive Rate')
69         plt.ylabel('True Positive Rate')
70         plt.title('ROC Curve - Logistic Regression')
71         plt.legend(loc="lower right")
72         if display_plot:
73             plt.show()

```

A.1.5 Class MyKNeighborsClassifier

```

1  from models.MyModel import MyModel
2  from sklearn.model_selection import cross_val_score
3  import matplotlib.pyplot as plt
4  from sklearn.neighbors import KNeighborsClassifier
5  from sklearn.metrics import plot_confusion_matrix
6  import numpy as np
7  from sklearn.metrics import classification_report
8  from sklearn.metrics import roc_curve, auc
9
10
11  class MyKNeighborsClassifier(MyModel):
12      def __init__(self, X, y):
13          self.X = X
14          self.y = y
15          self.X_train = None
16          self.X_test = None
17          self.y_train = None
18          self.y_test = None
19          self.model = None
20
21      def init_model(self, n=5):
22          self.model = KNeighborsClassifier(n_neighbors=n,
23          ↪ weights='uniform')
24
25      def cross_validate_penalty(self):
26          print("No cross-validation available for kNN")
27
28      def cross_validate_n(self, n_range=[1, 10, 100, 1000], k_fold_nb=5):
29          print("> n cross validation (range:", n_range, ")")
30          mean_scores = []; std_scores = []
31
32          for n in n_range:

```

```

32     print("\t> For n = %d..." % n)
33     tmp_model = KNeighborsClassifier(n_neighbors=n,
34     ↪ weights='uniform')
35     scores = cross_val_score(tmp_model, self.X_train,
36     ↪ self.y_train, cv=k_fold_nb, scoring='f1')
37     mean_scores.append(np.array(scores).mean())
38     std_scores.append(np.array(scores).std())
39
40     print("> n cross validation done")
41     plt.figure()
42     plt.errorbar(n_range, mean_scores, yerr=std_scores)
43     plt.title("n (# neighbors) Cross Validation")
44     plt.xlabel("n (# neighbors)")
45     plt.ylabel("F1-Score")
46     plt.show()
47
48     def confusion_matrix(self):
49         print("> Confusion matrix:")
50         plt.figure()
51         disp = plot_confusion_matrix(self.model, self.X_test,
52         ↪ self.y_test, display_labels=["y = False", "y = True"],
53         ↪ cmap=plt.cm.Blues,
54         ↪ values_format='%d')
55         disp.ax_.set_title("Confusion matrix")
56         plt.show()
57
58     def print_report(self):
59         y_pred = self.model.predict(self.X_test)
60         print("> Scores:")
61         print(classification_report(self.y_test, y_pred, target_names=["y
62         ↪ = False", "y = True"]))
63
64     def plot_roc_curve(self, display_plot=True):
65         """Plot ROC curve and ROC area of the model"""
66         y_pred_proba = self.model.predict_proba(self.X_test)
67         fpr, tpr, _ = roc_curve(self.y_test, y_pred_proba[:, 1])
68         plt.figure(123)
69         plt.plot(fpr, tpr, label='kNN (area = %0.2f)' % auc(fpr, tpr))
70         plt.xlabel('False Positive Rate')
71         plt.ylabel('True Positive Rate')
72         plt.title('ROC Curve - kNN')
73         plt.legend(loc="lower right")
74         if display_plot:
75             plt.show()

```

A.1.6 Class MyMLPClassifier

```

1  from models.MyModel import MyModel
2  from sklearn.model_selection import cross_val_score
3  import matplotlib.pyplot as plt
4  from sklearn.neural_network import MLPClassifier
5  from sklearn.metrics import plot_confusion_matrix
6  import numpy as np
7  from sklearn.metrics import classification_report
8  from sklearn.metrics import roc_curve, auc
9
10
11 class MyMLPClassifier(MyModel):
12     def __init__(self, X, y):
13         self.X = X
14         self.y = y
15         self.X_train = None
16         self.X_test = None
17         self.y_train = None
18         self.y_test = None
19         self.model = None
20
21     def init_model(self, n=(25,), max_iterations=10000):
22         self.model = MLPClassifier(hidden_layer_sizes=n,
23                                     ↪ max_iter=max_iterations)
24
25     def cross_validate_penalty(self):
26         print("No cross-validation available for Baseline")
27
28     def cross_validate_n(self, prev_hidden_layers=(), n_range=[1, 10,
29 ↪ 100, 1000], k_fold_nb=5, max_iterations=10000):
30         # Add parameters "prev_hidden_layers" (i.e. to cross-validate a
31 ↪ new layer's size)
32         print("> n cross validation (range:", n_range, ")")
33         mean_scores = []; std_scores = []
34
35         for n in n_range:
36             hls = prev_hidden_layers + (n,)
37             print("\t> For hidden_layer_sizes =", hls, "...")
38             tmp_model = MLPClassifier(hidden_layer_sizes=hls,
39 ↪ max_iter=max_iterations)
40             scores = cross_val_score(tmp_model, self.X_train,
41 ↪ self.y_train, cv=k_fold_nb, scoring='f1')
42             mean_scores.append(np.array(scores).mean())
43             std_scores.append(np.array(scores).std())
44
45         print("> n cross validation done")
46         plt.figure()

```



```

42     plt.errorbar(n_range, mean_scores, yerr=std_scores)
43     plt.title("n (hidden layer sizes) Cross Validation")
44     plt.xlabel("n (hidden layer sizes)")
45     plt.ylabel("F1-Score")
46     plt.show()
47
48     def confusion_matrix(self):
49         print("> Confusion matrix:")
50         plt.figure()
51         disp = plot_confusion_matrix(self.model, self.X_test,
52             ↪ self.y_test, display_labels=["y = False", "y = True"],
53             ↪ cmap=plt.cm.Blues,
54                                     values_format='d')
55         disp.ax_.set_title("Confusion matrix")
56         plt.show()
57
58     def print_report(self):
59         y_pred = self.model.predict(self.X_test)
60         print("> Scores:")
61         print(classification_report(self.y_test, y_pred, target_names=["y
62             ↪ = False", "y = True"]))
63
64     def plot_roc_curve(self, display_plot=True):
65         """Plot ROC curve and ROC area of the model"""
66         y_pred_proba = self.model.predict_proba(self.X_test)
67         fpr, tpr, _ = roc_curve(self.y_test, y_pred_proba[:, 1])
68         plt.figure(123)
69         plt.plot(fpr, tpr, label='MLP (area = %0.2f)' % auc(fpr, tpr))
70         plt.xlabel('False Positive Rate')
71         plt.ylabel('True Positive Rate')
72         plt.title('ROC Curve - MLP Classifier')
73         plt.legend(loc="lower right")
74         if display_plot:
75             plt.show()

```

A.1.7 Main script

```

1  from models.MyBaseline import MyBaseline
2  from models.MyLogisticRegression import MyLogisticRegression
3  from models.MyPreprocessor import MyPreprocessor
4  from models.MyKNeighborsClassifier import MyKNeighborsClassifier
5  from models.MyMLPClassifier import MyMLPClassifier
6  import matplotlib.pyplot as plt
7
8  part_i = True
9  part_ii = True
10
11  if part_i:

```

```

12  # (i) predict the review polarity (where a game has been "voted up"
    ↪ or not by the reviewer)
13
14  # (i)(i) Preprocessing data
15  print("> Preprocessing data")
16  preprocessor = MyPreprocessor()
17  preprocessor.read_data('dataset/reviews_17.jl.json')
18  preprocessor.preprocess_data(predict="voted_up", min_df=1,
    ↪ max_df=0.05)
19  X, y = preprocessor.get_data()
20
21  preprocessor.print_report()
22
23  # print("\n=== min_df cross-validation ===")
24  # min_df_range = [1, 5, 10, 15, 20, 25, 50]
25  # print("> min_df cross validation (range:", min_df_range, ")")
26  # preprocessor.cross_validate_min_df(min_df_range=min_df_range,
    ↪ max_df=1.0) # max_df=0.150
27
28  # print("\n=== max_df cross-validation ===")
29  # max_df_range = [0.001, 0.005, 0.01, 0.05, 0.1, 0.15, 0.2]
30  # print("> max_df cross validation (range:", max_df_range, ")")
31  # preprocessor.cross_validate_max_df(min_df=1,
    ↪ max_df_range=max_df_range)
32
33  # (i)(ii) Baseline model
34  print("\n=== BASELINE MODEL ===")
35  baseline = MyBaseline(X, y)
36  baseline.split(test_size=0.20)
37  baseline.init_model()
38  baseline.train()
39  # baseline.confusion_matrix()
40  baseline.print_report()
41  baseline.plot_roc_curve(display_plot=False)
42
43  # (i)(iii) Logistic regression
44  print("\n=== LOGISTIC REGRESSION ===")
45  logreg = MyLogisticRegression(X, y)
46  logreg.split(test_size=0.20)
47  # logreg.cross_validate_penalty(C_range=[0.1, 1, 10, 100, 1000],
    ↪ k_fold_nb=5)
48  logreg.init_model(C=10, max_iterations=10000)
49  logreg.train()
50  # logreg.confusion_matrix()
51  logreg.print_report()
52  logreg.plot_roc_curve(display_plot=False)
53
54  # (i)(iv) kNN

```

```

55 print("\n=== K NEAREST NEIGHBORS CLASSIFIER ===")
56 knn = MyKNeighborsClassifier(X, y)
57 knn.split(test_size=0.20)
58 # knn.cross_validate_n(n_range=[1, 5, 10, 50, 100, 500, 1000],
   ↪ k_fold_nb=5)
59 # knn.cross_validate_n(n_range=[1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
   ↪ k_fold_nb=5)
60 knn.init_model(n=3)
61 knn.train()
62 knn.confusion_matrix()
63 knn.print_report()
64 knn.plot_roc_curve(display_plot=False)
65
66 # (i)(v) MLP
67 print("\n=== MLP CLASSIFIER ===")
68 mlp = MyMLPClassifier(X, y)
69 mlp.split(test_size=0.20)
70 # mlp.cross_validate_n(prev_hidden_layers=(), n_range=[1, 5, 10, 25,
   ↪ 50, 100, 500, 1000], max_iterations=10000)
71 mlp.cross_validate_n(prev_hidden_layers=(50,), n_range=[1, 5, 10, 25,
   ↪ 50, 100], max_iterations=10000)
72 # mlp.cross_validate_n(prev_hidden_layers=(50, 5), n_range=[1, 5, 10,
   ↪ 25, 50, 100], max_iterations=10000)
73 mlp.init_model(n=(50,))
74 mlp.train()
75 # mlp.save_model("saved-models/MyMLPClassifier-50.joblib")
76 # mlp.load_model()
77 mlp.confusion_matrix()
78 mlp.print_report()
79 mlp.plot_roc_curve(display_plot=False)
80
81 # Print all ROC Curves in a single figure
82 plt.figure(123) # ROC Curves are added by each model on figure #123
83 plt.title("ROC Curves of different classifiers")
84 plt.show()
85
86
87 if part_ii:
88     # (ii)(i) Preprocessing data
89     print("> Preprocessing data")
90     preprocessor = MyPreprocessor()
91     preprocessor.read_data('dataset/reviews_17.jl.json')
92     preprocessor.preprocess_data(predict="early_access", min_df=1,
   ↪ max_df=0.05)
93     X, y = preprocessor.get_data()
94
95     preprocessor.print_report()
96

```

```

97     # (ii)(ii) Baseline model
98     print("\n=== BASELINE MODEL ===")
99     baseline = MyBaseline(X, y)
100    baseline.split(test_size=0.20)
101    baseline.init_model()
102    baseline.train()
103    # baseline.confusion_matrix()
104    baseline.print_report()
105    baseline.plot_roc_curve(display_plot=False)
106
107    # (ii)(iii) Logistic regression
108    print("\n=== LOGISTIC REGRESSION ===")
109    logreg = MyLogisticRegression(X, y)
110    logreg.split(test_size=0.20)
111    # logreg.cross_validate_penalty(C_range=[0.1, 1, 10, 100, 1000,
112        ↪ 5000], k_fold_nb=5)
113    logreg.init_model(C=1000, max_iterations=10000)
114    logreg.train()
115    # logreg.confusion_matrix()
116    logreg.print_report()
117    logreg.plot_roc_curve(display_plot=False)
118
119    # (ii)(iv) kNN
120    print("\n=== K NEAREST NEIGHBORS CLASSIFIER ===")
121    knn = MyKNeighborsClassifier(X, y)
122    knn.split(test_size=0.20)
123    # knn.cross_validate_n(n_range=[1, 5, 10, 50, 100, 500, 1000],
124        ↪ k_fold_nb=5)
125    # knn.cross_validate_n(n_range=[1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11,
126        ↪ 12], k_fold_nb=5)
127    knn.init_model(n=1)
128    knn.train()
129    # knn.confusion_matrix()
130    knn.print_report()
131    knn.plot_roc_curve(display_plot=False)
132
133    # (ii)(v) MLP
134    print("\n=== MLP CLASSIFIER ===")
135    mlp = MyMLPClassifier(X, y)
136    mlp.split(test_size=0.20)
137    # mlp.cross_validate_n(prev_hidden_layers=(), n_range=[1, 5, 10, 25,
138        ↪ 50, 100, 500, 1000], max_iterations=10000)
139    mlp.init_model(n=(50,))
140    mlp.train()
141    # mlp.save_model("saved-models/MyMLPClassifier-ii-50.joblib")
142    # mlp.load_model()
143    # mlp.confusion_matrix()
144    mlp.print_report()

```

```
141     mlp.plot_roc_curve(display_plot=False)
142
143     # Print all ROC Curves in a single figure
144     plt.figure(123) # ROC Curves are added by each model on figure #123
145     plt.title("ROC Curves of different classifiers")
146     plt.show()
```