

# Theory

1. What is the difference between a rule-based system and a machine learning system?

A rule-based system is a method of setting rules artificially to get results for several specific application scenarios while a machine learning system is the method of learning rules to predict results from training data by algorithm models automatically.

2. What is the difference between unsupervised and supervised learning?

Supervised learning needs pairs of features and targets of the training dataset to train the model and predict results for the testing dataset.

Unsupervised learning does not have a training dataset but only one set of data, looking for rules., whose goal is not to generate a classification system, but to make decisions. The agent can motivate the correct behaviour and punish the wrong behaviour.

3. What do we mean when we say that a machine learning system is overfitting?

Overfitting means almost all the features of the training data were learned, including local features and false features due to noise.

Overfitting results in low "generalization" and the accuracy of the testing data.

## Practice

1. Your algorithm gets the following results in a classification experiment. Please compute the precision, recall, f-measure and accuracy \*manually\* (without the help of your computer/Python, please provide all steps and formulas). Include the process to get to the final result.

Handwritten calculations on grid paper:

TP (true positive) IF - 6  
FP (false positive) ~~T~~ F 3  
FN (false negative) F 3  
TN (true negative) IF F 8

precision  $P = \frac{TP}{TP+FP} = \frac{6}{6+3} = \frac{2}{3}$

recall  $R = \frac{TP}{TP+FN} = \frac{6}{6+3} = \frac{2}{3}$

$F = \frac{2PR}{P+R} = \frac{2 \times \frac{2}{3} \times \frac{2}{3}}{\frac{2}{3} + \frac{2}{3}} = \frac{2}{3}$

Accuracy  $= \frac{TP+TN}{TP+TN+FP+FN} = \frac{6+8}{20} = \frac{14}{20} = \frac{7}{10}$

2. You are given a dataset (named Wine data set) with different measured properties of different wines (dataset available in Learning Central). Your goal is to develop a machine learning model to predict the quality of an unseen wine given these properties. Train

two machine learning regression models and check their performance. Write, for each of the models, the main Python instructions to train and predict the labels (one line each, no need to include any data preprocessing) and the performance in the test set in terms of Root Mean Squared Error (RMSE)

Models:

```
01. lin_reg1 = SVR(kernel="linear", C=3, epsilon=1.5, gamma="auto")
02. lin_reg2=LinearRegression()
03. lin_reg3= Pipeline([
04.     ("poly_features", PolynomialFeatures(degree=3, interaction_only=True)),
05.     ("std_scaler", StandardScaler()),
06.     ("regul_reg", lin_reg2),
07. ])
08. plot_learning_curves(lin_reg1)
09. plot_learning_curves(lin_reg2)
10. plot_learning_curves(lin_reg3)
```

Train regression models and check performance using in terms of RMSE:

```
01. model.fit(X_train, y_train) # Train using different models
02. y_test_predict = model.predict(X_val) # predict results
03. Testscore=model.score(X_val,y_val) # compute score of model
04. TestMSE=mean_squared_error(y_val, y_test_predict) #compute MSE
05. print("test score:"+str(Testscore)) #print score of model
06. print("test RMSE:"+str(cmath.sqrt(TestMSE))) #print RMSE
```

	Score	RMSE
SVR with linear kernel	0.21	0.789
Linear Regression	0.24	0.774
Linear Regression After using StandardScaler() and PolynomialFeatures()	0.32	0.73

3. Train an SVM binary classifier using the Hateval dataset (available in Learning Central). The task consists of predicting whether a tweet represents hate speech or not. You can preprocess and choose the features freely. Evaluate the performance of your classifier in terms of accuracy using ten fold-cross validation. Write a table with the results of the classifier (accuracy, precision, recall and F-measure) in each of the folds and write a small summary (up

to 500 words) of how you preprocessed the data, chose the feature/s, and trained and evaluated your model.

Table 1. The performance of the classifier in 10 fold cross-validation

Average Accuracy: 0.767				
	precision	recall	f1	accuracy
1	0.759032	0.760441	0.759699	0.768889
2	0.746116	0.750333	0.747062	0.750000
3	0.736743	0.735590	0.736142	0.747778
4	0.752314	0.756914	0.753437	0.756667
5	0.753298	0.758570	0.754744	0.758889
6	0.769067	0.768011	0.768503	0.773333
7	0.756119	0.751786	0.753339	0.758889
8	0.734555	0.728520	0.730232	0.736667
9	0.731099	0.734003	0.732331	0.741111
10	0.761527	0.762203	0.761850	0.766667

## Preprocessing

### Lemmatization

```
01. lemmatizer = nltk.stem.WordNetLemmatizer()
```

### Removing stop words

```
01. stopwords=set(nltk.corpus.stopwords.words('english'))
02. stopwords.add(".")
03. stopwords.add(",")
04. stopwords.add("!")
05. stopwords.add(":")
06. stopwords.add("?")
```

## Choice of features

### Bag of words

```
01. # Function to get the list of tokens from a string
02. def get_list_tokens(string):
03.     sentence_split=nlk.tokenize.sent_tokenize(string)
04.     list_tokens=[]
05.     for sentence in sentence_split:
06.         list_tokens_sentence=nlk.tokenize.word_tokenize(sentence)
07.         for token in list_tokens_sentence:
08.             list_tokens.append(lemmatizer.lemmatize(token).lower())
09.     return list_tokens
10.
11. # Function to get feature vectors
12. def get_vector_text(list_vocab,string):
13.     vector_text=np.zeros(len(list_vocab))
14.     list_tokens_string=get_list_tokens(string) ##Use function
15.     for i, word in enumerate(list_vocab):
16.         if word in list_tokens_string:
17.             #print (list_tokens_string.count(word))
18.             vector_text[i]=list_tokens_string.count(word) #Count words
19.     return vector_text
20.
21.
22. # Function to retrieve features based on bag of words
23. def get_vocabulary(training_set, num_features):
24.     vocabulary_inner=[]
25.     dict_word_frequency={}
26.     for instance in training_set:
27.         sentence_tokens=get_list_tokens(instance[0])##Use function
28.         for word in sentence_tokens:
29.             if word in stopwords: continue
30.             if word not in dict_word_frequency: dict_word_frequency[word]=1
31.             else: dict_word_frequency[word]+=1
32.     sorted_list = sorted(dict_word_frequency.items(), key=operator.itemgetter(1), reverse=True)[:num_features]
33.     for word,frequency in sorted_list:
34.         vocabulary_inner.append(word)
35.     return vocabulary_inner
36.
37. X_train=[]
38. Y_train=[]
39. for instance in training_set:
40.     vector_instance=get_vector_text(vocabulary,instance[0])
41.     X_train.append(vector_instance)
42.     Y_train.append(instance[1])
```

Each token occurrence frequency is treated as a feature, then a collection of text documents is transformed into numerical feature vectors as input data.

## Train the model

```
01. def train_svm_classifier(training_set, vocabulary):
02.     X_train=[]
03.     Y_train=[]
04.     for instance in training_set:
05.         vector_instance=get_vector_text(vocabulary,instance[0])
06.         X_train.append(vector_instance)
07.         Y_train.append(instance[1])
08.     #test_svm_function(X_train,Y_train) # Parameters Adjustment
09.     scaler = StandardScaler()
10.     #print(scaler)
11.     X_scaled = scaler.fit_transform(X_train) #StandScaler Training dataset
12.     #print(X_scaled)
13.     #test_svm_function(X_train,Y_train)
14.     #svm_clf=sklearn.svm.SVC(kernel="rbf",gamma=3.59) #train:99.7%, test:55.8%
15.     svm_clf=sklearn.svm.SVC(kernel="linear",C=1,gamma='auto') #train:85% test:76%
16.     #svm_clf=sklearn.svm.SVC(kernel="poly",degree=2,gamma='auto') #train:59% test 58%
17.     svm_clf.fit(X_scaled,Y_train)
18.     print("Train data score:"+str(svm_clf.score(X_scaled,Y_train)))
19.
20.     return svm_clf
```

SVC models with linear, "RBF" and "poly" kernel provided by SVM are used. Also, methods of the learning curve and grid search are used to find the best parameters for each model. Function, as follows, are carried out for parameter adjustment.

```
01. #learning curve
02. def test_svm_function(X_train,Y_train):
03.     gamma_range=np.logspace(0,1,10)
04.     score=[]
05.     for i in gamma_range:
06.         svm_clf=SVC(kernel="rbf",gamma=i,cache_size=5000).fit(X_train,Y_train)
07.         score.append(svm_clf.score(X_train,Y_train))
08.     print(max(score),gamma_range[score.index(max(score))])
09.     plt.plot(gamma_range,score)
10.     return gamma_range[score.index(max(score))]
11.
12. #grid search
13. def test_svm_function2(X_train,Y_train):
14.     gamma_range=np.logspace(-10,1,20)
15.     coef0_range=np.linspace(0,5,10)
16.
17.     param_grid=dict(gamma=gamma_range,coef0=coef0_range)
18.
19.     cv=StratifiedShuffleSplit(n_splits=10,test_size=0.3,random_state=420)
20.     grid=GridSearchCV(SVC(kernel="rbf"),param_grid=param_grid,cv=cv)
21.     grid.fit(X_train,Y_train)
22.     print("The best parameters are %s with a score of %0.5f"%(grid.best_params_,grid.best_score_))
```

Finally, SVC with a linear kernel is the best model among them and kept in codes.

## Evaluation using 10-fold cross-validation

```
01. def Kfold():
02.     table=pd.DataFrame(columns=['precision','recall','f1','accuracy'])
03.     kf = KFold(n_splits=10)
04.     print (kf)
05.     random.shuffle(dataset_full)
06.     kf.get_n_splits(dataset_full)
07.     index=0
08.     for train_index, test_index in kf.split(dataset_full):
09.         train_set_fold=[]
10.         test_set_fold=[]
11.         accuracy_total=0.0
12.         index+=1
13.         for i,instance in enumerate(dataset_full):
14.             if i in train_index:
15.                 train_set_fold.append(instance)
16.             else:
17.                 test_set_fold.append(instance)
18.
19.         vocabulary_fold=get_vocabulary(train_set_fold, 1000)
20.         svm_clf_fold=train_svm_classifier(train_set_fold, vocabulary_fold)
21.         X_test_fold=[]
22.         Y_test_fold=[]
23.         for instance in test_set_fold:
24.             vector_instance=get_vector_text(vocabulary_fold,instance[0])
25.             X_test_fold.append(vector_instance)
26.             Y_test_fold.append(instance[1])
27.         Y_test_fold_gold=np.asarray(Y_test_fold)
28.         X_test_fold=np.asarray(X_test_fold)
29.         Y_test_predictions_fold=svm_clf_fold.predict(X_test_fold)
30.         print ("Fold completed.")
31.         ## compute different value ##
32.         print("Test data score:"+str(svm_clf_fold.score(X_test_fold,Y_test_fold_gold)))
33.         accuracy_fold=accuracy_score(Y_test_fold_gold, Y_test_predictions_fold)
34.         accuracy_total+=accuracy_fold
35.         precision_fold=precision_score(Y_test_fold_gold, Y_test_predictions_fold, average='macro')
36.         recall_fold=recall_score(Y_test_fold_gold, Y_test_predictions_fold, average='macro')
37.         f1_fold=f1_score(Y_test_fold_gold, Y_test_predictions_fold, average='macro')
38.         ## draw the graph ##
39.         table=table.append(pd.DataFrame({'precision':precision_fold,'recall':recall_fold,'f1':f1_fold,'accuracy':accuracy_fold},index=[index]))
40.         average_accuracy=accuracy_total/10
41.         print ("Average Accuracy: "+str(round(average_accuracy,3)))
42.         print (table)
43.     Kfold()
```

Firstly, the dataset is split into ten pairs of training dataset and testing dataset. Then vectorize these datasets using the same features based on word frequency. Finally, the model is evaluated by accuracy, precision, recall and f1 scores. Then average accuracy is computed by ten folds.

Table 1 shows the performance of the model and average accuracy: 0.767