

RAPPORT DE KEVIN NGUETCHE POUR LE MINI-PROJET

A. Résumé du travail réalisé

A l'aide de GridSearchCV, nous avons sélectionné les meilleurs paramètres de 5 différents modèles de Machine Learning sur le jeu de donnée Fashion MNIST.

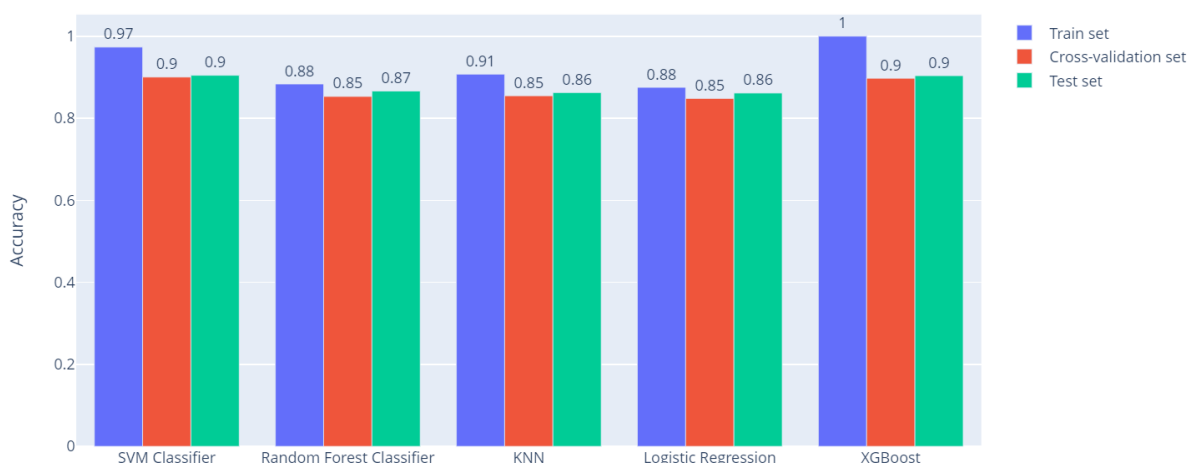
Trois évaluations différentes ont été effectuées pour chaque modèle à savoir : une **évaluation sur les données d'entraînement**, une **cross validation à 5 blocs** et enfin une **évaluation sur les données de tests**.

Le modèle ayant fourni de meilleures performances sur les données de test est **SVM** avec **0.9046** d'accuracy, suivi respectivement de **XGBOOST** avec **0.9037**, **Random Forest** avec **0.8899**, **KNN** avec **0.8625** et enfin **Regression logistic** avec **0.8614**.

Deux opérations de prétraitement ont été effectuées sur Fashion MNIST à savoir : **le mélange des données** (afin que toutes les classes soient représentées dans chaque portion du jeu de données) et **la normalisation Z-score** (afin de représenter les pixels sur un intervalle plus réduit que [0-255]).

B. Présentation des résultats

Accuracy Comparison By Kevin NGUETCHE



C. Résultats détaillés

Modèles	Accuracy Train set	Accuracy Cross-val	Accuracy Test set
SVM	0.9734	0.9003	0.9046
Random Forest	1.0000	0.8799	0.8899
KNN	0.9074	0.8547	0.8625
Logistic Regression	0.8752	0.8482	0.8614
XGBoost	1.0000	0.8974	0.9037

Dans les pages suivantes, nous présenterons plus en détail les étapes suivies.

I. JEU DE DONNEE FASHION MNIST

L'ensemble de données Fashion MNIST se compose de 70 000 images et chaque image a 784 caractéristiques (c'est-à-dire 28×28 pixels). Chaque pixel est une valeur de 0 à 255, décrivant l'intensité du pixel. 0 pour le blanc et 255 pour le noir.

Notre jeu de donnée a 10 classes d'étiquettes, numérotés de 0 à 9. Nous sommes donc en face d'un problème de classification multi-classes.

Suivant la consigne du projet, nous avons divisé notre jeu de données en 6/7 pour l'apprentissage soit 60 000 images, et 1/7 pour les tests soit 10 000 images.

II. PRETRAITEMENT

Les opérations de prétraitement effectuées sont :

- [Le Mélange des données](#) : afin que toutes les classes soient représentées dans chaque portion du jeu de données ; cette opération améliore les performances de la validation.

```
1 # Prétraitement

1 # Permutation du jeu de donnée
2 import numpy as np
3 np.random.seed(42) # if you want reproducible results set the random seed value.
4 shuffle_index = np.random.permutation(60000)
5 x_train, y_train = x_train[shuffle_index], y_train[shuffle_index]
```

- [La normalisation Z-score](#) : afin de représenter les pixels sur un intervalle plus réduit que [0-255] (la plage [0-255] est très grande). Cette normalisation est donnée par la formule : **$x_{\text{normalisé}} = (x - \text{moyenne}) / \text{écart-type}$** .

```
# Normalisation Z-score du Jeu de donnée
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
x_train = scaler.fit_transform(x_train.astype(np.float64))
x_test = scaler.fit_transform(x_test.astype(np.float64))
```

III. SELECTION DES MEILLEURS PARAMETRES DE CHAQUE MODELE

Afin d'effectuer une comparaison équitable de nos différents modèles, nous les avons comparés en utilisant leurs meilleurs paramètres associés. GridSearchCV a été utilisé.

- Meilleurs paramètres de SVM sur Fashion MNIST

```
# Trouvons les meilleurs hyper-paramètres pour nos différents models sur nos données

# SVM

classifier_svm = svm.SVC()
print("Listes des paramètres de l'SVM", classifier_svm.get_params(), "\n")

parameters = [{'C': [1, 10, 100, 1000], 'kernel': ['linear', 'rbf'], 'gamma': [0.001, 0.0001], 'max_iter': [1000, 10000],
               'decision_function_shape': ['ovo', 'ovr']}]

grid_search = GridSearchCV(estimator=classifier_svm, param_grid=parameters, scoring='accuracy', cv=5)

grid_search.fit(x_train, y_train)

best_parameters = grid_search.best_params_
print("Meilleurs paramètres trouvés sont : %s" % best_parameters)

Meilleurs paramètres trouvés sont : {'C': 10, 'decision_function_shape': 'ovo', 'gamma': 0.001, 'kernel': 'rbf', 'max_iter': 10000}
```

- Meilleurs paramètres de Random Forest sur Fashion MNIST

```
1 # Random Forest
2
3 classifier_rf = RandomForestClassifier()
4
5 parameters = [{'random_state': [0, 10, 30, 42, 100], 'max_depth': [10, 50, 100], 'n_estimators': [10, 100, 1000]}]
6
7 grid_search = GridSearchCV(estimator=classifier_rf, param_grid=parameters, scoring='accuracy', cv=5)
8
9 grid_search.fit(x_train, y_train)
10
11 best_parameters = grid_search.best_params_
12 print("Meilleurs paramètres de Random forest trouvés sont : %s" % best_parameters)

Meilleurs paramètres de Random forest trouvés sont : {'max_depth': 50, 'n_estimators': 100, 'random_state': 0}
```

- Meilleurs paramètres de KNN sur Fashion MNIST

```
1 # KNN
2
3 classifier_knn = KNeighborsClassifier()
4
5 parameters = [{'n_neighbors': [3, 4, 5, 8, 12]}]
6
7 grid_search = GridSearchCV(estimator=classifier_knn, param_grid=parameters, scoring='accuracy', cv=5)
8
9 grid_search.fit(x_train, y_train)
10
11 best_parameters = grid_search.best_params_
12 print("Meilleurs paramètres du KNN trouvés sont : %s" % best_parameters)

Meilleurs paramètres du KNN trouvés sont : {'n_neighbors': 4}
```

– Meilleurs paramètres de Logistic Regression sur Fashion MNIST

```
1 # Logistic Regression
2
3 classifier_lr = LogisticRegression()
4
5 parameters = [{'random_state': [0,42], 'C': [1, 10, 100, 1000], 'multi_class' : ['multinomial'],
6               'solver' : ['lbfgs', 'liblinear']} ]
7
8 grid_search = GridSearchCV(estimator=classifier_lr, param_grid=parameters, scoring='accuracy', cv=5)
9
10 grid_search.fit(x_train, y_train)
11
12 best_parameters = grid_search.best_params_
13 print("Meilleurs paramètres de la regression logistique trouvés sont : %s" % best_parameters)
14
Meilleurs paramètres de la regression logistique trouvés sont : {'C': 100, 'multi_class': 'multinomial', 'random_state': 0,
'solver': 'lbfgs'}
```

– Meilleurs paramètres de XGBOOST sur Fashion MNIST

```
1 # XGBOOST
2 classifier_xgb = XGBClassifier()
3
4 parameters = [{'random_state': [0, 42], 'max_depth': [10, 50, 100], 'n_estimators': [10, 100, 1000]}]
5
6 grid_search = GridSearchCV(estimator=classifier_xgb, param_grid=parameters, scoring='accuracy', cv=5)
7
8 grid_search.fit(x_train, y_train)
9
10 best_parameters = grid_search.best_params_
11 print("Meilleurs paramètres de XGBOOST trouvés sont : %s" % best_parameters)
12
```

IV. Evaluation des 5 modèles

Maintenant que les meilleurs paramètres de chaque modèle ont été sélectionnés, évaluons ces derniers.

– Evaluation de SVM

```
1 # SVM, Evaluation
2 classifier_svm = svm.SVC(C = 10.0, decision_function_shape = 'ovo', gamma = 0.001, kernel = 'rbf', max_iter = 10000)
3 classifier_svm.fit(x_train, y_train)
4 y_train_prd = classifier_svm.predict(x_train)
5
6 scores = cross_val_score(classifier_svm, x_train, y_train, cv=5)
7
8 y_test_prd = classifier_svm.predict(x_test)
9 acc_train_svm = accuracy_score(y_train, y_train_prd)
10 acc_val_svm = scores.mean()
11 acc_test_svm = accuracy_score(y_test, y_test_prd)
12
13 print("accuracy on train set:{:.4f}\naccuracy on cross validation set:{:.4f}\naccuracy on test set :{:.4f}".format(
14 acc_train_svm, acc_val_svm, acc_test_svm))
15
accuracy on train set:0.9734
accuracy on cross validation set:0.9003
accuracy on test set :0.9046
```

- Evaluation de Random Forest

```
1 # Random Forest, Evaluation
2 classifier_rf = RandomForestClassifier(random_state=0, n_estimators=100, max_depth=50)
3 classifier_rf.fit(x_train, y_train)
4 y_train_prd = classifier_rf.predict(x_train)
5
6 scores = cross_val_score(classifier_rf, x_train, y_train, cv=5)
7
8 y_test_prd = classifier_rf.predict(x_test)
9 acc_train_rf = accuracy_score(y_train, y_train_prd)
10 acc_val_rf = scores.mean()
11 acc_test_rf = accuracy_score(y_test, y_test_prd)
12
13 print("accuracy on train set:{:.4f}\naccuracy on cross validation set:{:.4f}\naccuracy on test set :{:.4f}".format(
14 acc_train_rf, acc_val_rf, acc_test_rf))
```

```
accuracy on train set:1.0000
accuracy on cross validation set:0.8799
accuracy on test set :0.8899
```

- Evaluation de KNN

```
1 # KNN, Evaluation
2 classifier_knn = KNeighborsClassifier(n_neighbors=4)
3 classifier_knn.fit(x_train, y_train)
4 y_train_prd = classifier_knn.predict(x_train)
5
6 scores = cross_val_score(classifier_knn, x_train, y_train, cv=5)
7
8 y_test_prd = classifier_knn.predict(x_test)
9 acc_train_knn = accuracy_score(y_train, y_train_prd)
10 acc_val_knn = scores.mean()
11 acc_test_knn = accuracy_score(y_test, y_test_prd)
12
13 print("accuracy on train set:{:.4f}\naccuracy on cross validation set:{:.4f}\naccuracy on test set :{:.4f}".format(
14 acc_train_knn, acc_val_knn, acc_test_knn))
```

```
accuracy on train set:0.9074
accuracy on cross validation set:0.8547
accuracy on test set :0.8625
```

- Evaluation de Logistic Regression

```
1 # Logistic Regression, Evaluation
2 classifier_lr = LogisticRegression(multi_class="multinomial", solver="lbfgs", C=100, random_state=0)
3 classifier_lr.fit(x_train, y_train)
4 y_train_prd = classifier_lr.predict(x_train)
5
6 scores = cross_val_score(classifier_lr, x_train, y_train, cv=5)
7
8 y_test_prd = classifier_lr.predict(x_test)
9 acc_train_lr = accuracy_score(y_train, y_train_prd)
10 acc_val_lr = scores.mean()
11 acc_test_lr = accuracy_score(y_test, y_test_prd)
12
13 print("accuracy on train set:{:.4f}\naccuracy on cross validation set:{:.4f}\naccuracy on test set :{:.4f}".format(
14 acc_train_lr, acc_val_lr, acc_test_lr))
```

```
accuracy on train set:0.8752
accuracy on cross validation set:0.8482
accuracy on test set :0.8614
```

- Evaluation de XGBOOST

```

1 # XGBOOST, Evaluation
2 classifier_xgb = XGBClassifier(random_state=0, n_estimators=100, max_depth=50)
3 classifier_xgb.fit(x_train, y_train)
4 y_train_prd = classifier_xgb.predict(x_train)
5
6 scores = cross_val_score(classifier_xgb, x_train, y_train, cv=5)
7
8 y_test_prd = classifier_xgb.predict(x_test)
9 acc_train_xgb = accuracy_score(y_train, y_train_prd)
10 acc_val_xgb = scores.mean()
11 acc_test_xgb = accuracy_score(y_test, y_test_prd)
12
13 print("accuracy on train set:{:.4f}\naccuracy on cross validation set:{:.4f}\naccuracy on test set :{:.4f}".format(
14 acc_train_xgb, acc_val_xgb, acc_test_xgb))

```

accuracy on train set:1.0000
accuracy on cross validation set:0.8974
accuracy on test set :0.9037

V. COMPARAISON DES MODELES

Visualisons à présent nos modèles en fonction de leur évaluation sur les données d'entrainements, les données avec Cross validation à 5 blocs et les données de Test.

```

1 # Comparaison de nos modèles
2 acc_combine = {'Model': ['SVM Classifieur', 'Random Forest Classifieur', 'KNN', 'Logistic Regression', 'XGBoost'],
3 'Accuracy_Tra': [acc_train_svm, acc_train_rf, acc_train_knn, acc_train_lr, acc_train_xgb],
4 'Accuracy_Val': [acc_val_svm, acc_val_rf, acc_val_knn, acc_val_lr, acc_val_xgb],
5 'Accuracy_Test': [acc_test_svm, acc_test_rf, acc_test_knn, acc_test_lr, acc_test_xgb]
6 }
7
8

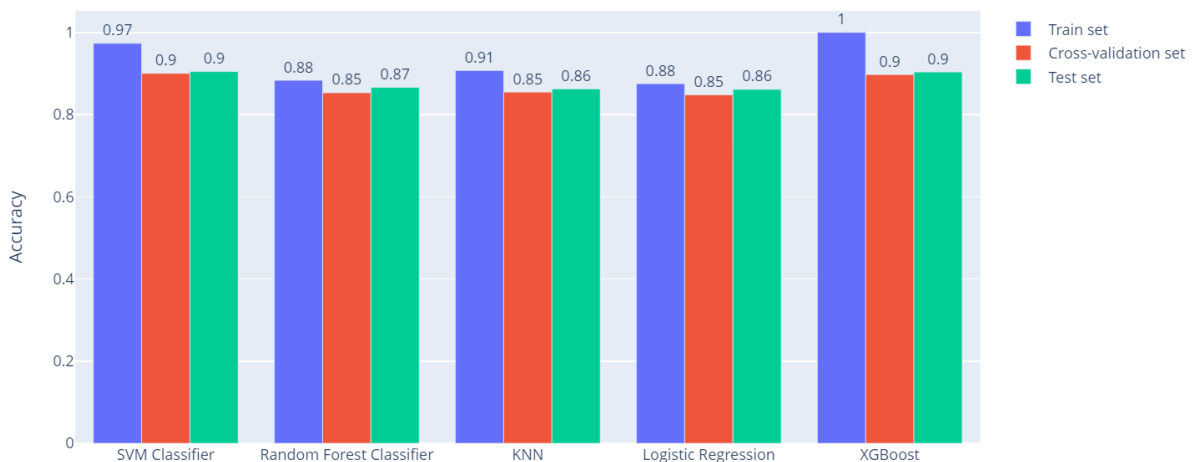
```

```

1 # Schéma de comparaison
2
3 fig = go.Figure(data=[
4 go.Bar(name='Train set', x=acc_combine['Model'], y=acc_combine['Accuracy_Tra'], text=np.round(acc_combine['Accuracy_Tra']
5 go.Bar(name='Cross-validation set', x=acc_combine['Model'], y=acc_combine['Accuracy_Val'], text=np.round(acc_combine['Acc
6 go.Bar(name='Test set', x=acc_combine['Model'], y=acc_combine['Accuracy_Test'], text=np.round(acc_combine['Accuracy_Test'
7 ]))
8
9 fig.update_layout(barmode='group', title_text='Accuracy Comparison By Kevin NGUETCHE', yaxis=dict(
10 title='Accuracy'))
11 fig.show()
12

```

Accuracy Comparison By Kevin NGUETCHE



On peut constater à travers la figure que la performance maximale sur les données Fashion MNIST est de **90%**. Elle est atteinte en utilisant soit l'approche One vs One de **SVM** ; soit en fixant le nombre d'estimateurs de **XGBOOST** à 100 pour un random_state à 0.