

## تصنيف الرسائل الإلكترونية للكشف عن البريد المزعج

### ➤ مقدمة:

مع الاستخدام الواسع النطاق للبريد الإلكتروني في كل من التواصل الشخصي والمهني، أصبحت مشكلة رسائل البريد المزعج مصدر قلق كبير. فالرسائل الغير مرغوب فيها لا تعبر فقط عن البريد العشوائي بل تشكل أيضاً مخاطر أمنية وتهديدات للخصوصية للمستخدمين. تعاني الطرق التقليدية لفرز البريد المزعج، من قيود في القابلية للتكيف والفعالية. لذلك، اكتساب تقنيات التعلم الآلي لتصنيف البريد الإلكتروني قد اكتسب أهمية بسبب قدرتها الكبيرة على تحسين الدقة والكفاءة في اكتشاف البريد المزعج.

في هذا التقرير، سنقوم بدراسة استخدام النماذج الإحصائية Naive Bayes و Logistic Regression لتصنيف البريد الإلكتروني إلى فئات "spam" و "ham" (غير مؤذي). سنقوم بتحليل أداء النماذج وفقاً لمعايير مثل الدقة ومصفوفة الارتباك.

### ➤ الطريقة:

في هذا المشروع، نهدف إلى تطوير برنامج لتصنيف الرسائل الإلكترونية إلى فئتين: البريد المزعج والشرعي (البريد الحقيقي). تتضمن الطريقة المقترحة الخطوات التالية:

1. استيراد المكتبات اللازمة.
2. يتم تعريف نموذجين للتصنيف:
  - الأول يستخدم Naive Bayes.
  - الثاني يستخدم Logistic Regression.
3. تعريف الكلاس الرئيسي EmailClassifier: يقوم بقراءة البيانات، تنظيفها، تدريب النماذج، وطباعة النتائج.
4. قراءة وتنظيف البيانات: يتم قراءة البيانات من ملف CSV ومعالجتها لإزالة القيم المفقودة.
5. تدريب النماذج: يتم تدريب النماذج باستخدام تقنية K-Fold Cross Validation.
6. تقييم الأداء: يتم قياس الدقة ومصفوفة الارتباك لكل نموذج.

### ➤ التجربة:

في نصف المشروع، قمنا بإكمال تجهيز البيانات وتدريب النموذج بنجاح. تم تنظيف المجموعة البيانات، واستخراج الميزات من النصوص البريدية. أظهرت التجارب الأولية مع خوارزمية Naive Bayes نتائج واعدة من حيث الدقة والكفاءة. ومع ذلك، يتطلب الأمر إجراء تجارب وتقييم إضافي لضبط النموذج وتحسين أدائه.

## ➤ الاستنتاج:

في الختام، فإن تطوير برنامج لتصنيف الرسائل الإلكترونية إلى فئتي البريد المزعج والشرعي باستخدام تقنيات التعلم الآلي له القدرة على تعزيز أمان البريد الإلكتروني وتجربة المستخدم. من خلال تجهيز البيانات بشكل منهجي، واختيار الميزات المناسبة، وتدريب النماذج، نهدف إلى إنشاء نظام تصنيف فعال ودقيق. ستسهم التجارب المستمرة والتعديلات في تحسين فعالية الطريقة المقترحة في مكافحة رسائل البريد المزعج.

## ➤ المراجع:

[1] مجموعة بيانات Email Spam Classification Dataset. متوفر على Kaggle:

<https://www.kaggle.com/datasets/purusinghvi/email-spam-classification-dataset>

## البرنامج:

```
1 import pandas as pd
2 from sklearn.model_selection import KFold
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
7
8 # Naive Bayes المتعلق النموذج تعريف
9 class NaiveBayesModel:
10     def __init__(self):
11         def __init__(self):
12             self.model = MultinomialNB() # Naive Bayes استخدام نموذج
13             self.vectorizer = CountVectorizer() # Bag-of-Words استخدام تقنية
14
15     def train(self, X_train, y_train):{}
16         # تحويل النصوص إلى تمثيل رقمي
17         # تدريب النموذج
18
19     def predict(self, X_test):{}
20         # تحويل النصوص إلى تمثيل رقمي
21         # التنبؤ بالتصنيف
22
23 # Logistic Regression المتعلق النموذج تعريف
24 class LogisticRegressionModel:
25     def __init__(self):
26         self.model = LogisticRegression(max_iter=1000) # Logistic Regression استخدام نموذج مع زيادة عدد الدورات
27         self.vectorizer = CountVectorizer() # Bag-of-Words استخدام تقنية
28
29     def train(self, X_train, y_train):{}
30         # تدريب
31
32
33     def predict(self, X_test):{}
34         # تبيؤ
35
36
37 # تعريف الكلاس الرئيسي لتصنيف البريد الإلكتروني
38 class EmailClassifier:
39     def __init__(self, data_path):
40         self.data_path = data_path
41         self.data = None
42         self.models = {'Naive Bayes': NaiveBayesModel(), 'Logistic Regression': LogisticRegressionModel()}
43
44     def read_data(self):
45         self.data = pd.read_csv(self.data_path) # قراءة البيانات من المسار المحدد
46
47     def clean_data(self):
48         self.data = self.data.fillna('') # ملء القيم المفقودة بقيم فارغة
49         self.data = self.data[self.data['text'].notna()] # إزالة الصفوف التي تحتوي على قيمة مفقودة في العمود 'text'
50
51     def train_models(self):{}
52         #class في ال تقوم بتدريب النماذج الموجودة
53     def print_results(self, results):{}
54         # طباعة متوسط مصفوفات الارتباك لكل نموذج
55
56 if __name__ == "__main__":
57     data_path = "C:\\Users\\sulim\\Downloads\\DataSets_for_Spam-Emails\\combined_data.csv" # مسار ملف البيانات
58     classifier = EmailClassifier(data_path)
59     classifier.read_data() # قراءة البيانات
60     classifier.clean_data() # تنظيف البيانات
61     results = classifier.train_models() # تدريب النماذج
62     classifier.print_results(results) # طباعة النتائج
```