Republic of Yemen
IBB University
Faculty of Science
Department of IT
Data Mining



الجمهورية اليمنية جامعة إب كلية العلوم التطبيقية قسم : تقنية معلومات مقرر : تنقيب البيانات

الطلاب التالية أسمائهم

- ايمن محمد ناجي قمحان
  - ضياء فضل الحضرمي
    - حازم هزام العمري

#### **DEMAND FORECASTING**

Mission is to create a 3-month demand forecasting model for the relevant store chain using the following time series and machine learning techniques:

**Random Noise** 

**Lag Shifted Features** 

**Rolling Mean Features** 

**Exponentially Weighted Mean Features** 

**Custom Cost Function** 

**Model Validation with LightGBM** 

Solving the problem of multiple unclear products in some markets and some products required by customers are not available. This model predicts the impact of products on purchasing in a shopping period and a decision is made to purchase the most demanded products and reduce the demand for those with scarce demand.

Dataset is here: https://www.kaggle.com/c/demand-forecasting

#### **Dataset Description**

The objective of this competition is to predict 3 months of item-level sales data at different store locations.

File descriptions

train.csv - Training data

test.csv - Test data (Note: the Public/Private split is time based)

sample submission.csv - a sample submission file in the correct format

**Data fields** 

date - Date of the sale data. There are no holiday effects or store closures.

store - Store ID

item - Item ID

sales - Number of items sold at a particular store on a particular dat

In Arabic language

الهدف من المشروع كالتالي :

التنبؤ بالطلب للمتاجر

المهمة هي إنشاء نموذج تنبؤ بالطلب لمدة 3 أشهر لسلسلة المتاجر المعنية باستخدام السلاسل الزمنية وتقنيات التعلم الآلي التالية:

الضوضاء العشوائية

خصائص التأخير المُحوّل

خصائص المتوسط المنحدر

خصائص المتوسط المرجح أسيًا

دالة التكلفة المُخصصة

التحقق من صحة النموذج باستخدام LightGBM

حل مشكلة وهي تعدد المنتجات غير المرغوبة في بعض المتاجر وبعض المنتجات المطلوبة لدى العملاء لا يتم توفيرها فهذا المودل يتنبئ في المنتجات الأثر شراءً في المتجر في مده معينة ومنة يتم اتخاذ قرار لشراء المنتجات الأكثر طلباً وتقليل الطلب للمنتجات التي لديها طلب شحيح

مجموعة البيانات متوفرة هنا: https://www.kaggle.com/c/demand-forecasting

وصف مجموعة البيانات

الهدف من هذه المسابقة هو التنبؤ ببيانات مبيعات المنتجات على مدار ثلاثة أشهر في مواقع متاجر مختلفة.

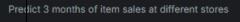
وصف الملفات

train.csv - بيانات التدريب

test.csv - بيانات الاختبار (ملاحظة: التقسيم بين عام وخاص يعتمد على الوقت)
sample\_submission.csv - نموذج لملف الإرسال بالتنسيق الصحيح
حقول البيانات
date - تاريخ بيانات البيع. لا توجد عطلات أو إغلاقات للمتاجر.
store - معرف المتجر
item - معرف المنتجات المباعة في متجر معين في تاريخ محدد.

هنا تم جلب البيانات

# **Store Item Demand Forecasting Challenge**



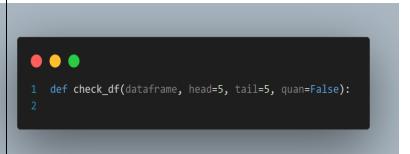


```
. .
    print(dataframe.shape)
print("###############################")
print(dataframe.dtypes)
         print("##################################")
         print(dataframe.tail(tail))
print("###################")
            def create_date_features(df):
        f create_date_features(df):
    df['month'] = df.date.dt.month
    df['day_of_month'] = df.date.dt.day
    df['day_of_wear'] = df.date.dt.dayofyear
    df['week_of_year'] = df.date.dt.dayofyear
    df['day_of_week'] = df.date.dt.dayofweek
    df['year'] = df.date.dt.year
    df['is_wfard'] = df.date.dt.year
    df['is_month_start'] = df.date.dt.is_month_start.astype(int)
    df['is_month_end'] = df.date.dt.is_month_end.astype(int)
    return df
    return df
    def random_noise(dataframe):
    return np.random.normal(scale=1.6, size=(len(dataframe),))
    # Lag/Shifted Features
def lag_features(dataframe, lags):
         for window in wind
              transform(
lambda x: x.shift(1).rolling(window=window, min_periods=10, win_type="triang").mean()) + random_noise(
dataframe)
return dataframe
54 def ewm_features(dataframe, alphas, lags):
55    for alpha in alphas:
              for lag in lags
                   def smape(preds, target):
    n = len(preds)
         n = Ien(Preds)
masked_arr = ~((preds == 0) & (target == 0))
preds, target = preds[masked_arr], target[masked_arr]
num = np.abs(preds - target)
         denom = np.abs(preds) + np.abs(target)
smape_val = (200 * np.sum(num / denom)) / n
          return smape val
    def lgbm_smape(preds, train_data):
    labels = train_data.get_label()
    smape_val = smape(np.expm1(preds), np.expm1(labels))
    return 'SMAPE', smape_val, False
         sns.set(font_scale=1)
sns.barplot(x="gain", y="feature", data=feat_imp[0:25])
              plt.title('feature')
plt.tight_layout()
              plt.show()
    for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

# ات من منصة Kaggel

يقوم هذا الكود بتهيئة بيئة العمل الخاصة بتحليل البيانات والتعلم الآلي عبر استيراد المكتبات التالية:

- numpy , pandas لتحليل البيانات ومعالجتها.
  - matplotlib, seaborn للرسم البياني.
  - Lightgbm لبناء نموذج تعلم آلي فعال.
- LabelEncoder لتحويل البيانات الفئوية إلى عددية.
- إعدادات عرض pandas وتحكم في رسائل التحذير لتحسين تجربة المستخدم.



import time

import numpy as np
import pandas as pd

import warnings

import seaborn as sns
import lightgbm as lgb

from matplotlib import pyplot as plt

pd.set\_option('display.max\_columns', None)

pd.set\_option('display.width', 500)
warnings.filterwarnings('ignore')

from sklearn.preprocessing import LabelEncoder

\*دالة لفحص وتحليل أولي لبيانات DataFrame تعرض التالى:

- الشكل العام للبيانات عدد الصفوف، الأعمدة
  - أنواع البيانات لكل عمود
- · أول وأخر عدد معين من الصفوف (افتراضي 5)
  - عدد القيم المفقودة في كل عمود
- إذا تم تفعيل quan=True، تُعرض أيضًا القيم المئوية(Quantiles)

```
def check_df(dataframe, head=5, tail=5, quan=False):
2
```

\* تُستخدم لاستخراج **ميزات مشتقة من التاريخ** من عمود date، مثل:

- الشهر، اليوم من الشهر، اليوم من السنة
  - الأسبوع، اليوم من الأسبوع، السنة
  - هل التاريخ في بداية أو نهاية الشهر؟
- هل اليوم عطلة (نهاية أسبوع)?
   هذه الميزات تساعد النماذج على فهم الأنماط
   الزمنية.

```
1 def random_noise(dataframe):
2
```

\* تُضيف ضجيج عشوائي (Noise) إلى البيانات، وهو أمر
 مفيد لمنع النماذج من التعلّم الزائد(Overfitting)

```
def lag_features(dataframe, lags):
2
```

تُنشئ ميزات الإزاحة الزمنية (Lag) للمبيعات السابقة، مثل مبيعات اليوم السابق، أو قبل يومين، إلخ.

- تقوم بذلك لكل مجموعة (store, item)
- تُستخدم هذه الميزات لتعليم النموذج العلاقة بين مبيعات اليوم والأيام السابقة

```
1 def roll_mean_features(dataframe, windows):
2
```

- تُنشئ ميزات المتوسط المتحرك (Rolling Mean)
   باستخدام نافذة window معينة.
- مفيدة لتحليل الاتجاهات على مدى زمني معين.
- تُستخدم دالة ()rolling مع نوع نافذة مثل"tiring

```
def ewm_features(dataframe, alphas, lags):
2
```

- \* تُنشئ ميزات ا**لمتوسط المرجح أسيًا(EWM)** ، الذي يعطي وزنًا أكبر للبيانات الأحدث.
  - ايتحكم بمقدار الوزن للبيانات الحديثة.
  - تُطبق لكل مجموعة (store, item) على بيانات المبيعات.

```
1 def smape(preds, target):
2
```

\* دالة لحساب Percentage Error، وهي مقياس دقيق لتقييم النماذج التي تتنبأ بقيم مستمرة، خصوصًا في السلاسل الزمنية.

تعطي نتائج بين 0 و100 كلما اقتربت من 0 كان
 النموذج أفضل

```
def lgbm_smape(preds, train_data):
2
```

\* تُستخدم لربط SMAPE بنموذج LightGBM كتقييم مخصص (custom metric) أثناء التدريب.

تُحوّل القيم من Log إلى طبيعتها قبل التقييم
 باستخدام ()np.expm1

```
def plot_lgb_importances(model, plot=False, num=10):
```

\*تعرض أو تطبع أهمية الميزات (Feature Importance) في نموذج.LightGBM

- ❖ تعتمد على طريقتين: عدد الانقسامات (split) أو إجمالي
   التأثير.(gain)
- ❖ عند تفعیل plot-Trueترسم رسمًا بیانیًا لأهم 25 میزة.

```
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

#### جزء قراءة البيانات من كاجل

يُستخدم للحصول على قائمة بكل الملفات الموجودة في مجلد kaggle/input.

❖ مفيد عند استخدام منصة Kaggle للتأكد من الملفات المتاحة للقراءة.

```
1 df = pd.concat([train, test], sort=False)
2
```

```
    هذا القسم يقوم بتحميل البيانات وتحضيرها للمعالجة.
```

- تم دمج بيانات التدريب والاختبار لتسهيل معالجة موحدة لجميع البيانات .
- لاحقًا سيتم فصلهم مرة أخرى عند بناء النموذج والتنبؤ.

```
df["date"].min(), df["date"].max()
train["date"].min(), train["date"].max()
test["date"].min(), test["date"].max()
```

# التحليل والاستكشاف الأولى

\*نتحقق من نطاق التواريخ في كل من البيانات الكاملة (df) والتدريب (train) والاختبار.(test)

```
1 check_df(train)
2 check_df(test)
3 check_df(df)
4
```

نستخدم الدالة check\_df التي شرحناها سابقًا لفحص كل
 مجموعة بيانات: عدد الأعمدة، القيم المفقودة، أنواع
 البيانات، إلخ.

```
df["sales"].describe([0.10, 0.30, 0.50, 0.70, 0.80, 0.90, 0.95, 0.99])
```

نُحلل توزيع المبيعات باستخدام إحصائيات مئوية لتحديد
 القيم المتطرفة والوسيط.

```
1 df[["store"]].nunique()
2 df[["item"]].nunique()
3 df.groupby(["store"])["item"].nunique()
4
```

```
معلومات حول المتاجر والمنتجات
```

- ❖ نحصل على:
- عدد المتاجر الفريد
- ❖ عدد العناصر الفريدة
- عدد العناصر المتاحة في كل متجر

1 df.groupby(["store", "item"]).agg({"sales": ["sum", "mean", "median", "std"]})
2

```
    تحلیل المبیعات حسب کل متجر وغنصر: مجموع، متوسط، وسیط، وانحراف معیاري.
```

# df['date'] = pd.to\_datetime(df['date']) df = create\_date\_features(df) check\_df(df)

# استخراج ميزات تاريخية

- نحول العمود date إلى صيغة زمنية، ثم نستخدم دالة Vcreate\_date\_features
  - ❖ الشهر، الأسبوع، السنة، بداية/نهاية الشهر، إلخ.

```
df.groupby(["store", "item", "month"]).agg({"sales": ["sum", "mean", "median", "std"]})
2
```

💠 تحليل شهري للمبيعات حسب المتجر والعنصر.

```
df.sort_values(by=['store', 'item', 'date'], axis=0, inplace=True)
df = lag_features(df, [91, 98, 105, 112, 119, 126, 182, 364, 546, 728])
3
```

## ميزات الإزاحة الزمنية (Lag Features)

بعد فرز البيانات زمنيًا، ننشئ ميزات للمبيعات السابقة
 (Lagged Sales)، وهي مهمة جدًا للتنبؤ بناءً على القيم السابقة.

#### المتوسطات المتحركة

نُضيف ميزات للمتوسطات المتحركة طويلة المدى (سنة،
 سنة ونصف)، مما يساعد على التقاط، الاتجاهات.

# (EWM) المتوسطات الأسية المرجحة

تُنشأ ميزات تأخذ في الاعتبار الأوزان الزمنية للأسابيع أو
 الأشهر الماضية، حيث تُعطى أهمية أكبر للقيم الأحدث.

# كود التحضير للنموذج

- نحول المتغيرات الفئوية (... store, item, ...) إلى تمثيل
   رقمى باستخدام .One-Hot Encoding
- نحول المبيعات إلى القيمة اللوغاريتمية لتقليل تأثير القيم المتطرفة.

# -\* نقسم البيانات زمنياً للتدريب :

- مجموعة التدريب حتى نهاية 2016
- مجموعة التحقق هي أول 3 شهور من 2017، لمحاكاة
   التنبؤ في 2018

```
1 df = roll_mean_features(df, [365, 546])
2
```

```
1 alphas = [0.95, 0.9, 0.8, 0.7, 0.5]
2 lags = [91, 98, 105, 112, 180, 270, 365, 546, 728]
3 df = ewm_features(df, alphas, lags)
4
```

```
1 df = pd.get_dummies(df, columns=['store', 'item', 'day_of_week', 'month'])
2 df['sales'] = np.logip(df["sales"].values)
3
```

```
1 train = df.loc[(df["date"] < "2017-01-01"), :]
2 val = df.loc[(df["date"] >= "2017-01-01") & (df["date"] < "2017-04-01"), :]
3</pre>
```

```
1 Y_train = train['sales']
2 X_train = train[cols]
3 Y_val = val['sales']
4 X_val = val[cols]
5
```

+- نحدد المتغير المستهدف (sales) والمتغيرات المدخلة cols) يجب أن تحتوي على أسماء الأعمدة التي تم استخراجها.(

```
1 Y_train.shape, X_train.shape, Y_val.shape, X_val.shape
```

\* نتأكد من صحة أبعاد البيانات قبل التدريب.

```
1 lgb_params = {
2    'metric': {'mae'},
3    'num_leaves': 10,
4    'learning_rate': 0.02,
5    'feature_fraction': 0.8,
6    'max_depth': 5,
7    'verbose': 0,
8    'num_boost_round': 1000,
9    'early_stopping_rounds': 200,
10    'nthread': -1
11 }
12
```

- إعدادات نموذج:LightGBM

```
❖ mae مقياس الخطأ
```

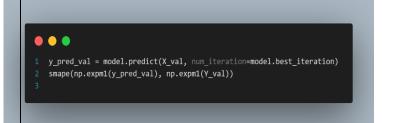
- num\_leaves, max\_depth خاتمونج النموذج
- 💸 early\_stopping\_rounds توقف التدريب عند عدم التحسن
  - learning\_rate معدل التعلم الصغير لتجنب التذبذب

```
1 lgbtrain = lgb.Dataset(data=X_train, label=Y_train, feature_name=cols)
2 lgbval = lgb.Dataset(data=X_val, label=Y_val, reference=lgbtrain, feature_name=cols)
3
```

♦ إعداد بيانات التدريب والتحقق لاستخدامها في نموذج LightGBM.

```
1 model = lgb.train(
2    lgb_params,
3    lgbtrain,
4    valid_sets=[lgbtrain, lgbval],
5    num_boost_round=lgb_params['num_boost_round'],
6    callbacks=[lgb.early_stopping(lgb_params['early_stopping_rounds']), lgb.log_evaluation(100)],
7    feval=lgbm_smape
8 )
9
```

- تدریب النموذج باستخدام:
  - نه التحقق المبكر(Early stopping)
    - ❖ تقییم مخصص عبرSMAPE
    - طباعة النتائج كل 100 جولة

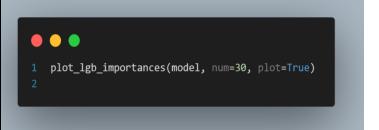


# تقييم النموذج على مجموعه التحقق

- نقوم بالتنبؤ على مجموعة التحقق (X\_val) باستخدام أفضل تكرار من التدريب(best\_iteration)
- \* ثم نحسب SMAPEعلى القيم الأصلية (بعد عكس log1p باستخدام (expm1 عليم دقة النموذج.

# عرض اهم الميزات

نستخدم دالة plot\_lgb\_importances الميزات في النموذج النهائي، مما يساعدنا على فهم ما
 يؤثر أكثر على التنبؤ.



```
1 train = df.loc[~df.sales.isna()]
2 test = df.loc[df.sales.isna()]
3
```

## تدريب النموذج النهائي

يتم استخدام كامل بيانات التدريب لتدريب النموذج
 النهائي، واستبعاد الصفوف التي لا تحتوي على sales
 )وهى مجموعة الاختبار.(

```
1 lgb_params = {..., "num_boost_round": model.best_iteration}
2 model = lgb.train(lgb_params, lgbtrain_all, num_boost_round=model.best_iteration)
3
```

نستخدم أفضل عدد تكرارات تم الحصول عليه سابقًا
 لتدريب نموذج قوى على كامل البيانات.

```
submission_df = test.loc[:, ['id', 'sales']]
submission_df['sales'] = np.expm1(test_preds)
submission_df['id'] = submission_df.id.astype(int)
submission_df.to_csv('submission_demand.csv', index=False)
```

#### إنشاء ملف التقديم Submission

- نحضر الملف حسب تنسيق مسابقة كاجل:
  - نستخدم pid نستخدم
- نحول القيم من log إلى القيم الأصلية

```
with open('/content/drive/My Drive/lgbm_model.pkl', 'wb') as model_file:
pickle.dump(model, model_file)
submission_df.to_csv('/content/drive/My Drive/submission_demand.csv', index=False)
4
```

#### حفض النموذج وملف التقييم

#### نحفظ:

- النموذج المدرب بصيغة pkl.
  - ❖ ملف النتائج بصيغة csv.
- \* يتم حفظ كل ذلك في Google Drive لاستخدامه لاحقًا أو تحميله.

```
with open('/content/drive/MyDrive/lgbm_model.pkl', 'rb') as f:
model = pickle.load(f)
```

# تحميل النموذج نخليه لاحقأ

نستخدم مكتبة pickle لإعادة تحميل النموذج عند
 الحاجة دون إعادة تدريبه.

# التنبؤ اليدوي لمتجر معين

- تُجهز هذه الدالة إدخال يدوي لتنبؤ مبيعات منتج معين في
   متجر معين في تاريخ معين:
  - پتم استخراج خصائص التاریخ
  - one-Hot Encoding پتم تنفیذ
    - ❖ ملء الأعمدة الناقصة بـ 0
  - ترتیب الأعمدة تمامًا كما في التدریب

input\_df = prepare\_input(store=3, item=25, date\_str='2018-01-01', feature\_columns=cols)
prediction = model.predict(input\_df)[0]

1 def prepare\_input(store, item, date\_str, feature\_columns): ...

يتم تمرير البيانات إلى النموذج للحصول على توقع المبيعات لليوم المحدد.

```
plt.bar(['Predicted Sales'], [prediction], color='green')
2
```

```
    يتم رسم الرسم البيائي الذي يُظهر قيمة المبيعات المتوقعة بشكل مرئى أنيق.
```

# خلاصة المشروع: توقع الطلب باستخدام LightGBM

- قمنا ببناء نموذج دقيق لتوقع المبيعات اليومية لكل منتج في كل متجر باستخدام تقنيات تعلم الآلة.
  - اعتمدنا على تحويل البيانات الزمنية إلى ميزات ذكية مثل:
    - o الإزاحات الزمنية (Lag Features)
    - o المتوسطات المتحركة (Rolling Mean)
    - المتوسطات الأسية المرجحة (EWM)
  - مؤشرات التاريخ (الشهر، اليوم، بداية/نهاية الشهر...إلخ)
  - استخدمنا نموذج LightGBM نظراً لسرعته ودقته مع البيانات الكبيرة.
    - تم تقييم النموذج باستخدام مقياس SMAPE لضمان جودة التنبؤ.
  - بعد التحقق درّبنا نموذجاً نهائيًا على كامل البيانات وولّدنا ملف التقديم النهائي.
  - أنشأنا دالة للتنبؤ اليدوي بأي تاريخ/منتج/متجر، مع عرض النتائج بطريقة رسومية بسيطة.

النتيجة	
نموذج فعال قادر على دعم اتخاذ القرار في التخطيط للمخزون، التسعير، وإدارة الطلبات بكفاءة عالية.	
اشراف المهندس / فؤاد الــــدغار	