

Amélioration d'un cas de machine learning

**Challenge Data : Comment démasquer les fraudeurs ?
par BNP Paribas PF**

*Titre professionnel « Développeur en intelligence artificielle » de niveau 6
enregistré au RNCP sous le n°34757*

Passage par voie de formation

Parcours de 19 mois achevé le 10 mai 2023

Promotion SIMPLON – LANNION 2022 – 2023

Thomas JONCOURT

Table des matières

Contexte	3
Objectif	3
Base de données	3
Description du benchmark	4
<i>Métrique d'évaluation</i>	<i>4</i>
<i>Benchmark</i>	<i>5</i>
Pré-processing.....	5
<i>Nettoyage du texte</i>	<i>5</i>
<i>One hot encoding.....</i>	<i>6</i>
<i>Gestion de la taille des données à entraînés.....</i>	<i>6</i>
Choix de l'algorithme de machine learning	7
Gestion du déséquilibres des classes	8
Résultats	9
Conclusion.....	11

Contexte

La fraude est un problème majeur de plus en plus préoccupant pour les institutions financières du monde entier. Les criminels utilisent une grande variété de méthodes pour attaquer des organisations comme la BNP Paribas, quels que soient les systèmes, les canaux, les process ou les produits.

Le développement de méthodes de détection de la fraude est stratégique et essentiel pour les établissements bancaires. Les fraudeurs s'avèrent toujours très créatifs et ingénieux pour normaliser leurs comportements et les rendre difficilement identifiables. Une contrainte s'ajoute à cette problématique, la faible occurrence de la fraude dans la population.

Lien du challenge : <https://challengedata.ens.fr/challenges/104>

Objectif

L'objectif de ce challenge est de trouver la meilleure méthode pour transformer et agréger les données relatives au panier client d'un des partenaires de la BNP pour détecter les cas de fraude.

En utilisant ces données panier, les fraudeurs pourront être détectés, et ainsi refusés dans le futur.

Base de données

La base contient une liste d'achats effectués sur le site partenaire. Les informations décrivent exclusivement le contenu du panier.

Pour chaque observation de la base, il y a 147 colonnes dont 144 peuvent être regroupées en 6 catégories :

- Item : Catégorie du bien de l'item, exemple : Computer
- cash_price : Prix de l'item, exemple : 850
- make : Fabricant de l'Item, exemple : Apple
- model : Description du modèle de l'item, exemple : Apple Iphone XX
- goods_code : Code-barre de l'item, exemple : 2378284364
- Nbr_of_pro_purchas : Nombre de produits pour l'item, exemple : 2

Le panier se décompose au maximum en 24 items. Par exemple, si un panier contient 3 items alors toutes les informations relatives à ces 3 items seront renseignées dans les

colonnes item1, item2, item3, cash_price1, cash_price_2, cash_price3, make1, make 2, make3, model1, model2, model3, goods_code1, goods_code2, goods_code3, Nbr_of_prod_purchas1, Nbr_of_prod_purchas2 et Nbr_of_prod_purchas3. Les variables restantes (celles avec un indice > 3) seront vides .

Description du benchmark

Métrique d'évaluation

L'objectif est d'identifier une opération frauduleuse dans la population en prédisant un risque/probabilité de fraude. Par conséquent, la métrique à utiliser est l'aire sous la courbe Précision-Rappel, appelé également PR-AUC.

Cette métrique est appropriée pour évaluer correctement la performance d'un modèle sur la classe minoritaire dans le cadre d'une classification fortement déséquilibrée.

Le meilleur modèle correspondra à celui avec la valeur de PR-AUC la plus élevée.

Pour ce challenge, la PR-AUC sera estimée par la moyenne pondérée des précisions à chaque seuil avec le poids associé étant la variation du rappel entre le seuil précédent et le nouveau seuil :

Cette implémentation correspond à la métrique `average_precision_score` de scikit-learn.

Benchmark

Le benchmark intègre plusieurs étapes de pré-processing et utilise un modèle de Machine Learning optimisé pour prédire le risque de fraude. Les détails de ces étapes ne sont pas précisés par la BNP.

PR-AUC2=0,14PR-AUC2=0,14

Pré-processing

Nettoyage du texte

Il y a dans le jeu quelques coquille dans certains noms de marques ou de catégories. Ils sont corrigés au cas par cas avec une fonction créé manuellement.

```
1 def clean(x):
2     if type(x)==str:
3         x = x.replace("& ", "")
4         x = x.replace(", ", "")
5         x = x.replace(" S ", "S ")
6         x = x.replace('TELEPHONESFAX', 'TELEPHONES FAX')
7         x = x.replace('CARPETS RUGS', 'CARPETS RUGS')
8         x = x.replace('BAGSWALLETS', 'BAGS WALLETS')
9         x = x.replace('!', '')
10        x = x.replace('BABYBJ RN', 'BABYBJÖRN')
11        x = x.replace('DR BROWN S', 'DR BROWNS')
12        x = x.replace('KIEHL S', 'KIEHLS')
13        x = x.replace('@', '')
14        x = x.replace('.', '')
15        x = x.replace('2', '')
16        x = x.replace(" + ", "")
17    return x

[ ] 1 for i in range(1,25):
2     train[f"item{i}"] = train[f"item{i}"].apply(clean)
3     train[f"make{i}"] = train[f"make{i}"].apply(clean)
4     train[f"model{i}"] = train[f"model{i}"].apply(clean)
5     test[f"item{i}"] = test[f"item{i}"].apply(clean)
6     test[f"make{i}"] = test[f"make{i}"].apply(clean)
7     train[f"model{i}"] = train[f"model{i}"].apply(clean)
```

One hot encoding

Avec ce jeu de données, le plus complexe est de transformer les informations contenues dans les colonnes `item{i}`, `make{i}` et `model{i}` (i de 1 à 24) qui contiennent des données textuelles. Ces données doivent être regroupées par catégories et transformées en valeurs numériques. On ne peut pas dans ce cas de figure utiliser une librairie clé en main comme scikit-learn car il faut agréger les valeurs des 24 colonnes de chaque type de colonne (exemple : item) et les « additionner » ensemble.

J'ai donc préféré écrire un script manuellement pour mieux contrôler le résultat du one hot encoding.

Exemple avec les colonnes `item{i}` :

```
1 train_item_price = pd.DataFrame(columns=items,index=train.index).fillna(0)
2 train_item_count = pd.DataFrame(columns=items,index=train.index).fillna(0)
3 test_item_price = pd.DataFrame(columns=items,index=test.index).fillna(0)
4 test_item_count = pd.DataFrame(columns=items,index=test.index).fillna(0)

1 for row in train.index:
2     for i in range(1,25):
3         item = train[f"item{i}"][row]
4         if item in items:
5             train_item_price[item][row] += train[f"cash_price{i}"][row] * train[f"Nbr_of_prod_purchas{i}"][row]
6             train_item_count[item][row] += train[f"Nbr_of_prod_purchas{i}"][row]
7
8 for row in test.index:
9     for i in range(1,25):
10        item = test[f"item{i}"][row]
11        if item in items:
12            test_item_price[item][row] += test[f"cash_price{i}"][row] * test[f"Nbr_of_prod_purchas{i}"][row]
13            test_item_count[item][row] += test[f"Nbr_of_prod_purchas{i}"][row]
```

2 types de matrices sont créés : une avec le nombre d'occurrence de chaque item/marque dans le panier et une avec le prix cumulé des produits de chaque item/marque dans le panier

Gestion de la taille des données à entraînées

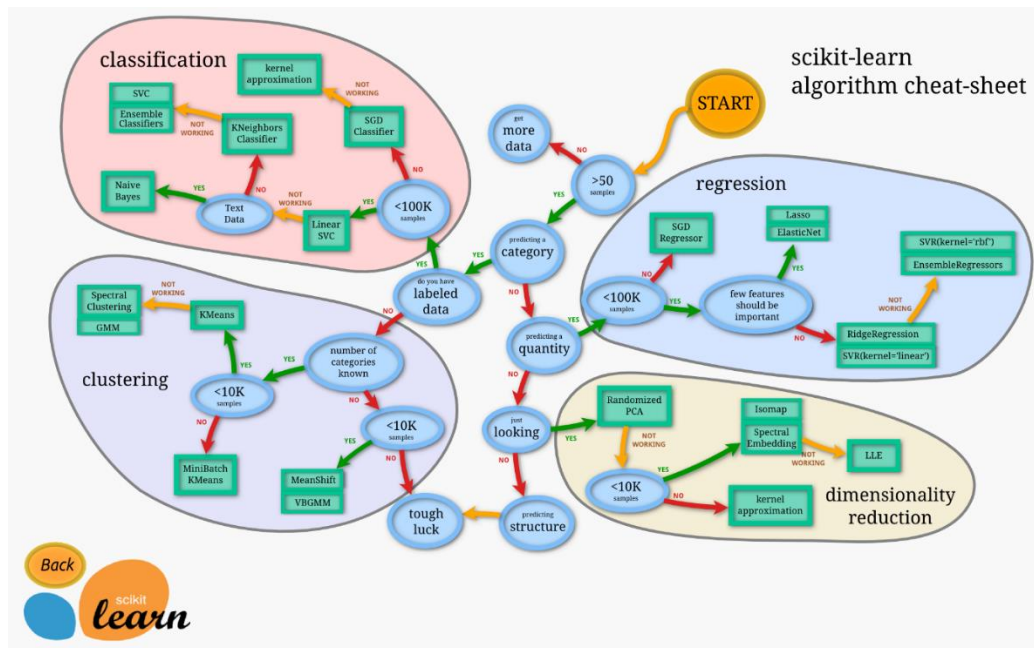
Il y a 184 items, 811 marques, 9678 modèles et 14880 code-barres différents dans ce jeu de données et 92790 échantillons. Afin de réduire la taille finale du jeu d'entraînement, les colonnes contenant les code-barres seront écartés. Ensuite les résultats des phases de one hot encoding pour les autres catégories seront

sauvegardés sous formes de matrice creuse afin d'optimiser la mémoire RAM et la puissance de calcul nécessaire à l'entraînement à l'aide de la librairie scipy.

```
1 sps.save_npz(f'drive/MyDrive/Challenges/bnp_dataset/new/matrice/train_item_price.npz', csr_matrix(train_item_price.fillna(0)))
2 sps.save_npz(f'drive/MyDrive/Challenges/bnp_dataset/new/matrice/test_item_price.npz', csr_matrix(test_item_price.fillna(0)))
3 sps.save_npz(f'drive/MyDrive/Challenges/bnp_dataset/new/matrice/train_item_count.npz', csr_matrix(train_item_count.fillna(0)))
4 sps.save_npz(f'drive/MyDrive/Challenges/bnp_dataset/new/matrice/test_item_count.npz', csr_matrix(test_item_count.fillna(0)))
```

Choix de l'algorithme de machine learning

Pour m'orienter dans le choix du modèle de machine learning, je me suis aidé de la carte proposée par scikit-learn suivante :



A noter que pour pouvoir évaluer le modèle à l'aide de la métrique `average_precision_score`, il est nécessaire de choisir un algorithme qui permette de prédire une probabilité pour chaque classe. J'utiliserais donc en priorité les modèles « Ensemble Classifiers » de scikit-learn.

Gestion du déséquilibres des classes

Dans ce jeu de données, il y a une problématique que je n'avais encore jamais rencontré ; le jeu de données est déséquilibré (seulement 1.4% de fraude). J'essayerais donc d'utiliser trois méthodes pour ré-équilibrer ce dataset :

- L'undersampling : Il s'agit ici simplement de retirer aléatoirement des instances de la classe majoritaire afin de ré-équilibrer les proportions. On perd toutefois de l'information, et il y a donc un risque d'underfitting. Cette méthode est préconisée lorsque l'on dispose d'un très grand nombre d'observations (**à minima > 10K**) comme c'est le cas pour notre dataset.
- L'oversampling : Il s'agit ici de dupliquer aléatoirement certaines instances des classes minoritaires, rendant ainsi leur signal plus puissant. Il y a toutefois ici un risque d'overfitting. Cette méthode est préconisée lorsque l'on dispose d'un nombre limité d'observations (< 10K), ou bien si le temps de calcul n'est pas un problème.
- La méthode des class weights permet de prendre en compte le caractère biaisé de la distribution du dataset et de créer un modèle pénalisé. Il s'agit ici de simplement attribuer des poids différents aux différentes classes de notre dataset, en donnant un poids plus important aux classes minoritaires, afin d'influencer le modèle lors de son entraînement. Nous pénalisons ainsi plus fortement une erreur de classification d'une classe minoritaire par rapport à une erreur de classification d'une classe majoritaire.

Il est à noter que la plupart des modèles disponibles dans la librairie [scikit-learn](https://scikit-learn.org/) ont un paramètre "[class_weight](#)" qui est par défaut initialiser à None, donnant alors des poids égaux pour toutes les classes. En mettant ce paramètre à "balanced", nous pouvons indiquer à scikit-learn de calculer automatiquement les poids pour chaque classe. Il est aussi possible de donner à scikit-learn un dictionnaire où les poids ont été entrés manuellement.

Résultats

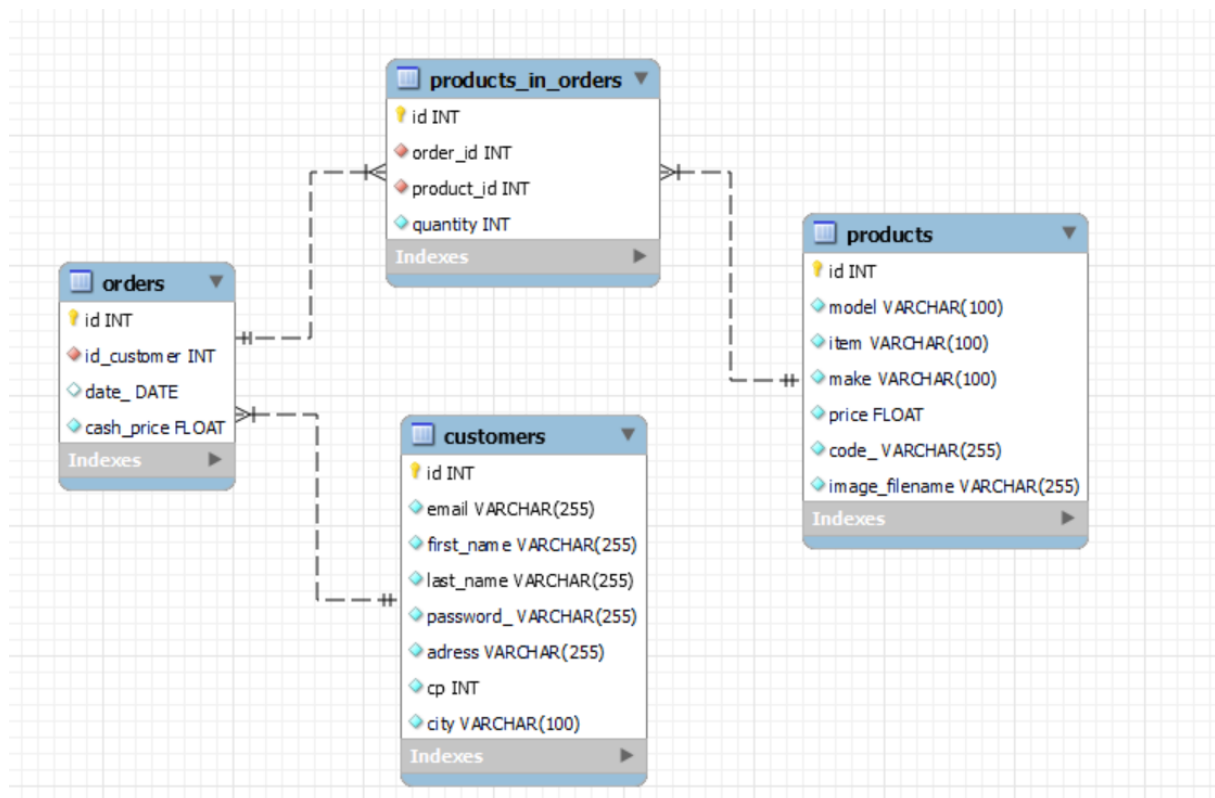
	class_weight = None	class_weight = « balanced »
RandomForestClassifier	0.184	0.120
GradientBoostingClassifier	0.126	
HistGradientBoostingClassifier	0.171	0.141
AdaBoostClassifier	0.090	
BaggingClassifier	0.159	
ExtraTreesClassifier	0.165	0.142

29	24 janvier 2023 13:00	ThomasJoncourt	0,1745
30	27 février 2023 19:40	totof	0,1734
31	7 avril 2023 10:39	Squirrel & AntoineDP	0,1731
32	7 avril 2023 17:18	tlecornec & cleroy	0,1727
33	13 avril 2023 10:40	Tom2	0,1700
34	20 février 2023 14:02	VictorHoffmann	0,1692
35	28 mars 2023 14:20	AdamEmz	0,1689
36	15 février 2023 16:40	raoulaur@gmail.com	0,1681
37	9 avril 2023 17:37	SimonLeR & hugo.leveque	0,1676
38	10 mars 2023 14:46	VictorHoffmann2	0,1675
39	28 mars 2023 14:45	Khobby	0,1668
40	27 mars 2023 18:04	Pierre_Joly	0,1652
41	14 février 2023 17:12	neodelphis	0,1649
42	19 février 2023 02:00	tontontondavid	0,1626
43	29 mars 2023 00:26	Lokhy	0,1599
44	9 avril 2023 20:53	h-anhluong & hamyluong	0,1578
45	-	benchmark	0,1574
46	22 janvier 2023 20:48	Quentin	0,1569

PR-AUC2=0,14PR-AUC2=0,1745

Création d'une base de données relationnelle

Pour adapter le jeu de données afin d'être utilisé par une application de boutique en ligne, je l'ai transformé en une base de données SQL.



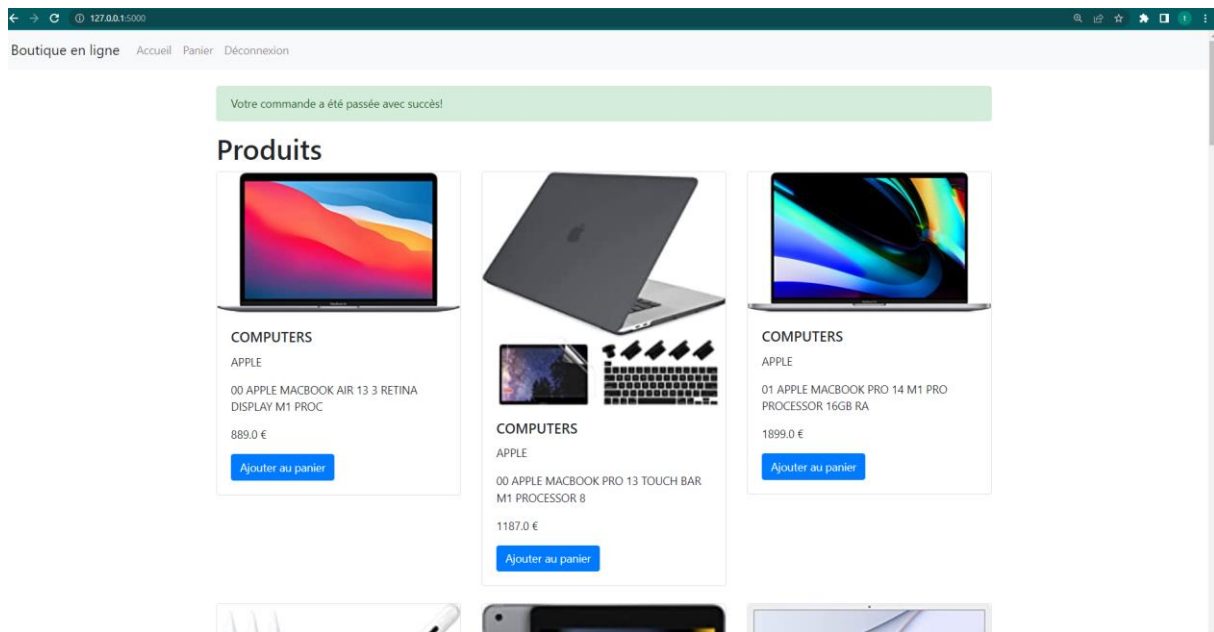
Le script SQL de création de la base de données est donné en Annexe 1.

Le script d'import des données dans la table «products» est donné en Annexe 2.

Application

Afin de simuler une utilisation pratique de notre modèle de détection de fraude, j'ai créé une application Flask qui sera relié à notre base de données SQL. Les informations de l'utilisateur et ses commandes seront stockées dans la base puis traité pour calculer la probabilité de fraude sur chaque commande.

Le script du fichier main.py de l'application est donné en annexe 3



Conclusion

Le score de précision moyenne obtenue est supérieur à celui du benchmark, j'ai donc réussi à proposer une solution plus performante sur la détection de fraude. Le classement de l'algorithme parmi tous ceux proposés dans le challenge est plutôt correct même s'il existe apparemment de bien meilleures solutions.

Les solutions pour essayer de rééquilibrer le dataset ne semblent pas fonctionner. Je ne sais pas si c'est une mauvaise utilisation de ma part ou si ce dataset n'est tout simplement pas adapté au rééquilibrage. On peut supposer que dans notre cas précis, vu que le taux de fraude est une observation réelle, il ne doit pas être modifié pour ne pas s'éloigner de la réalité. Le rééquilibrage serait donc une solution adaptée à d'autres cas de machine Learning de classification où la proportion de chaque classe n'est pas associée à un phénomène réel comme des datasets de classification d'image par exemple.

Il est certainement possible d'obtenir de meilleurs résultats avec des modèles plus complexes comme des modèles de Deep Learning notamment.

Annexe

Annexe 1 : Script SQL de création de la base de données

```
1 • drop database if exists bnp2;
2   -- Crée la base de données
3 • CREATE DATABASE bnp2;
4   -- Utilise la base de données
5 • USE bnp2;
6   -- Crée la table des produits
7 • CREATE TABLE products (
8     id INT AUTO_INCREMENT PRIMARY KEY,
9     model VARCHAR(100) NOT NULL,
10    item VARCHAR(100) NOT NULL,
11    make VARCHAR(100) NOT NULL,
12    price FLOAT NOT NULL,
13    code_ VARCHAR(255) NOT NULL,
14    image_filename VARCHAR(255) NOT NULL
15  );
16  -- Crée la table des produits
17 • CREATE TABLE customers (
18     id INT AUTO_INCREMENT PRIMARY KEY,
19     email VARCHAR(255) NOT NULL,
20     first_name VARCHAR(255) NOT NULL,
21     last_name VARCHAR(255) NOT NULL,
22     password_ VARCHAR(255) NOT NULL,
23     adress VARCHAR(255) NOT NULL,
24     cp INT NOT NULL,
25     city VARCHAR(100) NOT NULL
26  );
27  -- Crée la table des commandes
28 • CREATE TABLE orders (
29     id INT AUTO_INCREMENT PRIMARY KEY,
30     id_customer INT NOT NULL,
31     date_ DATE,
32     cash_price FLOAT NOT NULL,
33     FOREIGN KEY (id_customer) REFERENCES customers(id)
34  );
```

```

35     -- Crée la table des éléments de commande (order_items)
36 • CREATE TABLE products_in_orders (
37     id INT AUTO_INCREMENT PRIMARY KEY,
38     order_id INT NOT NULL,
39     product_id INT NOT NULL,
40     quantity INT NOT NULL,
41     FOREIGN KEY (order_id) REFERENCES orders(id),
42     FOREIGN KEY (product_id) REFERENCES products(id)
43 );

```

Annexe 2 : Script d'import des données dans la table « products »

```

5 train = pd.read_csv("/content/drive/MyDrive/Challenges/bnp_dataset/x_train.csv")
6
7 models = []
8 for i in range(1,25):
9     models += train[f"model{i}"].values.tolist()
10
11 df_models = pd.DataFrame({"model":pd.Series(models).value_counts().index,"count":pd.Series(models).value_counts().values})
12
13 df_models["make"] = None
14 df_models["item"] = None
15 df_models["price"] = None
16 df_models["good_code"] = None
17
18 for row in df_models.index:
19     model = df_models["model"][row]
20     mask = train[[f"model{i}" for i in range(1,25)].isin([model]).any(axis=1)
21     for i in range(1,25):
22         if train[mask].iloc[0][f"model{i}"] == model:
23             df_models["make"][row] = train[mask].iloc[0][f"make{i}"]
24             df_models["item"][row] = train[mask].iloc[0][f"item{i}"]
25             df_models["price"][row] = train[mask].iloc[0][f"cash_price{i}"]
26             df_models["good_code"][row] = train[mask].iloc[0][f"goods_code{i}"]
27             break
28
29 df_models["code_"] = df_models["good_code"]
30 df_models["image_filename"] = df_models["model"] + ".jpg"
31 df_models = df_models[["model", "item", "make", "price", "code_", "image_filename"]]
32
33 # Connect to the MySQL database
34 cnx = mysql.connector.connect(user='root', password='root', host='localhost', database='bnp2')
35
36
37 sql = "INSERT INTO products (model, item, make, price, code_, image_filename) VALUES (%s, %s, %s, %s, %s, %s)"
38 val = [[df_models["model"][row], df_models["item"][row], df_models["make"][row], str(df_models["price"][row]), df_models["code_"][row], df_models["image_filename"][row]] for row in df_models.iloc[0:20].index]
39
40 # Create a cursor object
41 cursor = cnx.cursor()
42
43 cursor.executemany(sql, val)
44
45 cnx.commit()

```

Annexe 3 : Script de fichier main.py de l'application Flask

```
1  from flask import Flask, render_template, request, redirect, url_for, flash, session
2  from flask_sqlalchemy import SQLAlchemy
3  from flask_login import LoginManager, login_user, logout_user, login_required, current_user
4  from models import db, Customers, Products, Orders, ProductsInOrder
5  from werkzeug.security import check_password_hash
6  from werkzeug.security import generate_password_hash
7
8  app = Flask(__name__)
9  app.config['SQLALCHEMY_DATABASE_URI'] = 'mysql+pymysql://root:root@localhost/bnp2'
10 app.config['SQLALCHEMY_TRACK_MODIFICATIONS'] = False
11 app.secret_key = 'your_secret_key'
12
13 db.init_app(app)
14
15 login_manager = LoginManager()
16 login_manager.init_app(app)
17 login_manager.login_view = 'login'
18
19 @login_manager.user_loader
20 def load_user(user_id):
21     return Customers.query.get(int(user_id))
22
23 # Routes à créer ici (index, register, login, etc.)
24
25 @app.route('/')
26 def index():
27     products = Products.query.all()
28     return render_template('index.html', products=products)
29
```

```
30 @app.route('/register', methods=['GET', 'POST'])
31 def register():
32     if current_user.is_authenticated:
33         return redirect(url_for('index'))
34
35     if request.method == 'POST':
36         email = request.form['email']
37         first_name = request.form['first_name']
38         last_name = request.form['last_name']
39         password = request.form['password']
40         adress = request.form['adress']
41         cp = request.form['cp']
42         city = request.form['city']
43
44         # Vérifiez si l'utilisateur existe déjà
45         user = Customers.query.filter_by(email=email).first()
46         if user:
47             flash('Cet email est déjà utilisé. Veuillez en utiliser un autre.', 'danger')
48             return redirect(url_for('register'))
49
50         hashed_password = generate_password_hash(password)
51         new_customer = Customers(email=email, first_name=first_name, last_name=last_name, password=hashed_password, adress=adress, cp=cp, city=city)
52         db.session.add(new_customer)
53         db.session.commit()
54
55         flash('Inscription réussie! Vous pouvez maintenant vous connecter.', 'success')
56         return redirect(url_for('login'))
57
58     return render_template('register.html')
59
```

```

60 @app.route('/login', methods=['GET', 'POST'])
61 def login():
62     if current_user.is_authenticated:
63         return redirect(url_for('index'))
64
65     next_page = request.args.get('next', url_for('index'))
66
67     if request.method == 'POST':
68         email = request.form['email']
69         password = request.form['password']
70         remember_me = request.form.get('remember_me')
71
72         user = Customers.query.filter_by(email=email).first()
73
74         if user and check_password_hash(user.password_, password):
75             login_user(user, remember=remember_me)
76             session['cart'] = {} # Initialise le panier
77             return redirect(next_page)
78         else:
79             flash('Identifiants incorrects. Veuillez réessayer.', 'danger')
80
81     return render_template('login.html', next=next_page)
82
83 @app.route('/logout')
84 def logout():
85     logout_user()
86     return redirect(url_for('index'))

```

```

101 @app.route('/cart', methods=['GET', 'POST'])
102 #@login_required
103 def cart():
104     if request.method == 'POST':
105         # Créez la commande pour le client actuel
106         order = Orders(id_customer=current_user.id, cash_price=0)
107         db.session.add(order)
108         db.session.commit()
109
110         # Parcourez les éléments du panier et ajoutez-les à la commande
111         cart_items = request.form.getlist('product_id')
112         cart_quantities = request.form.getlist('quantity')
113         total_price = 0
114
115         for i, product_id in enumerate(cart_items):
116             product = Products.query.get(product_id)
117             quantity = int(cart_quantities[i])
118             total_price += product.price * quantity
119
120         # Ajoutez l'élément de commande à la base de données
121         order_item = ProductsInOrder(order_id=order.id, product_id=product_id, quantity=quantity)
122         db.session.add(order_item)
123
124         # Mettez à jour le prix total de la commande
125         order.cash_price = total_price
126         db.session.commit()
127
128         flash("Votre commande a été passée avec succès!", "success")
129         return redirect(url_for('index'))
130
131     cart = session.get('cart') or {}
132     cart_items = []
133     for product_id, quantity in cart.items():
134         product = Products.query.get(product_id)
135         cart_items.append({'product': product, 'quantity': quantity})
136
137     return render_template('cart.html', cart_items=cart_items)
138
139 if __name__ == '__main__':
140     app.run(debug=True)

```