



# Implémentez un modèle de scoring



# Sommaire

1. Présentation de la problématique et du jeu de données
2. Présentation de la modélisation
3. Analyse du data drift
4. Présentation de la pipeline de déploiement
5. Démo du dashboard



## **Présentation de la problématique**

L'entreprise "Prêt à dépenser" propose des crédits à la consommation à des personnes ayant peu ou pas d'historique de prêt. Ici, la problématique va être de chercher à évaluer la probabilité de remboursement de chaque client et ainsi classer les demandes en crédits accordés ou refusés. C'est pourquoi, l'entreprise souhaite développer un algorithme de classification. Pour cela, elle se basera sur des sources de données variées, telles que les données comportementales et celles provenant d'autres institutions financières.

Comme les clients sont de plus en plus demandeurs de transparence en ce qui concerne les décisions d'octroi de crédit, l'entreprise a décidé de développer un dashboard interactif. Ce tableau de bord offrira une transparence maximale aux clients et permettra aux chargés de relation client d'expliquer facilement les décisions d'octroi de crédit.



## **Jeu de données mis à notre disposition:**

- 300 000 clients
- + de 300 variables
- Présente un gros déséquilibre

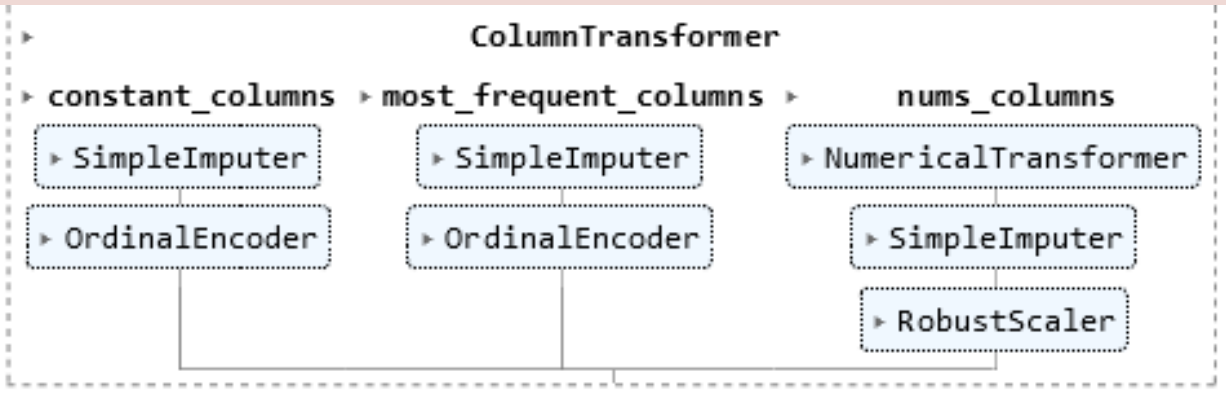


# Présentation de la modélisation

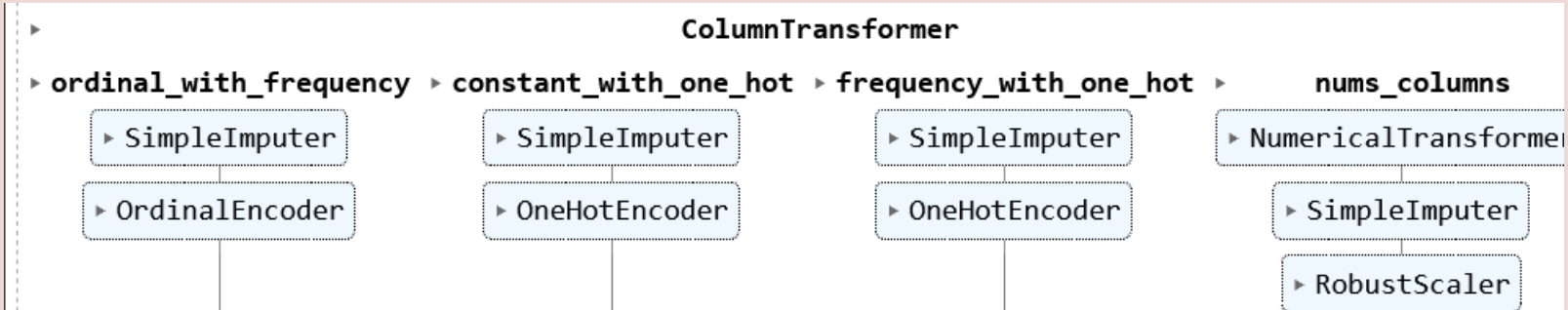


# Gestion des valeurs manquantes

## Modèle basée sur les arbres de décision

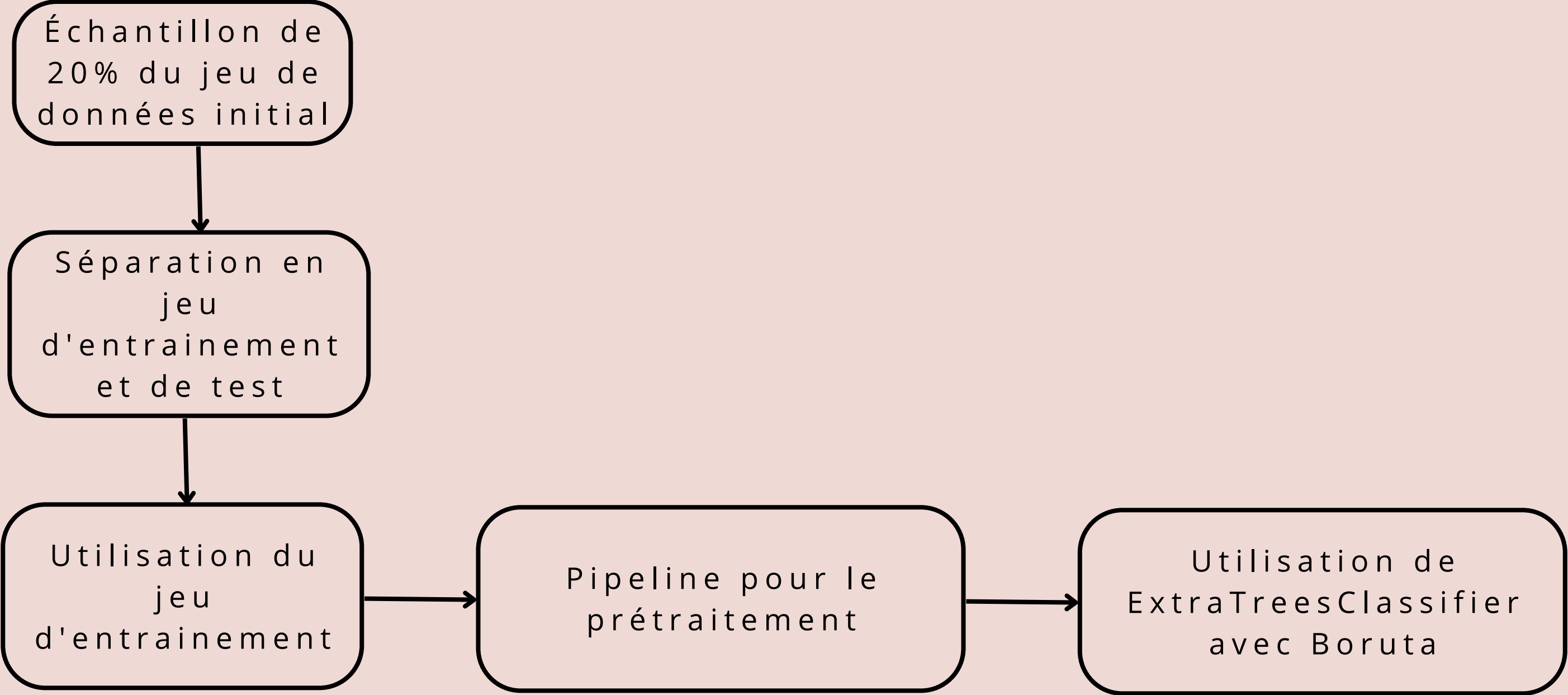


## Modèle linéaire





# Démarche pour la sélection de features





# Métrique choisis

	Pred Negative class	Pred Positive class
Negative class (Bon payeur)	TN (Bon payeur bien prédit)	FP (Bon payeur prédit mauvais payeur)
Positive class (Mauvais payeur)	FN (Mauvais payeur prédit bon payeur)	TP (Mauvais payeur bien prédit)

## Métrique technique choisis

- **L'aire sous la courbe (AUC)**: mesure la capacité du modèle à distinguer les mauvais payeurs (TP) des faux mauvais payeurs (FP)
- **Le fbeta score**: combine la précision et le rappel tout en accordant un poids à l'une ou l'autre de ces métriques. Un bêta de 2 sera défini.

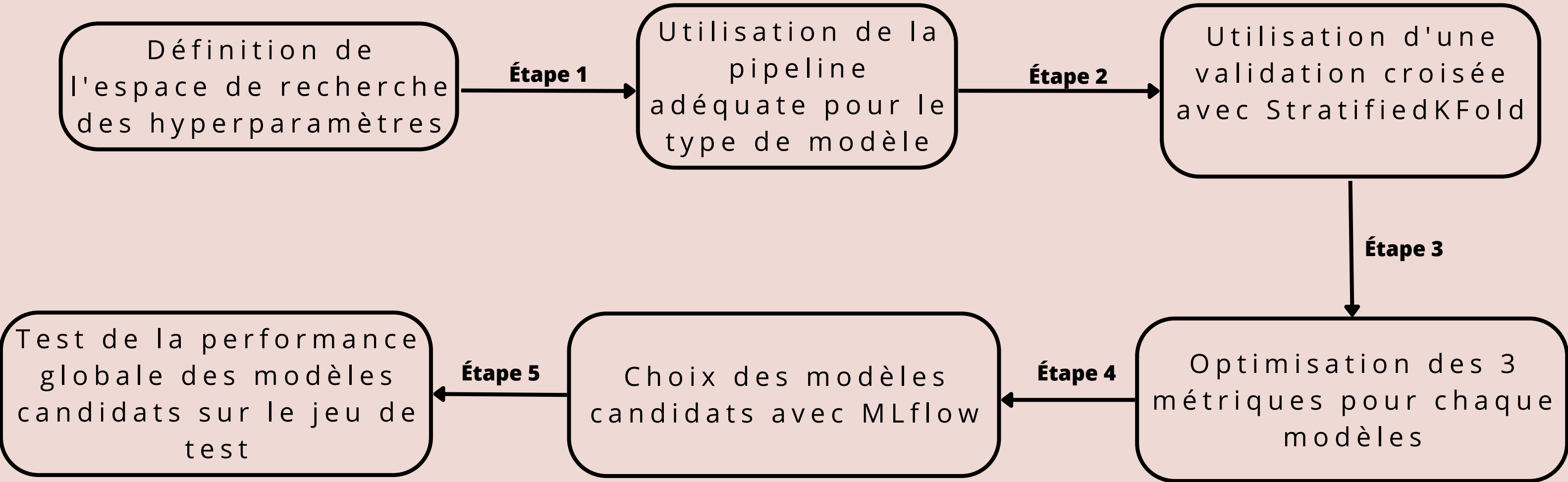
## Métrique métier

```
def cost_function(y, y_pred, cost):  
  
    Conf= confusion_matrix(y, y_pred)  
  
    tn, fp, fn, tp = Conf[0,0], Conf[0,1], Conf[1,0], Conf[1,1]  
  
    return (cost * fn + fp) / (tn + cost * fn + fp + tp)
```



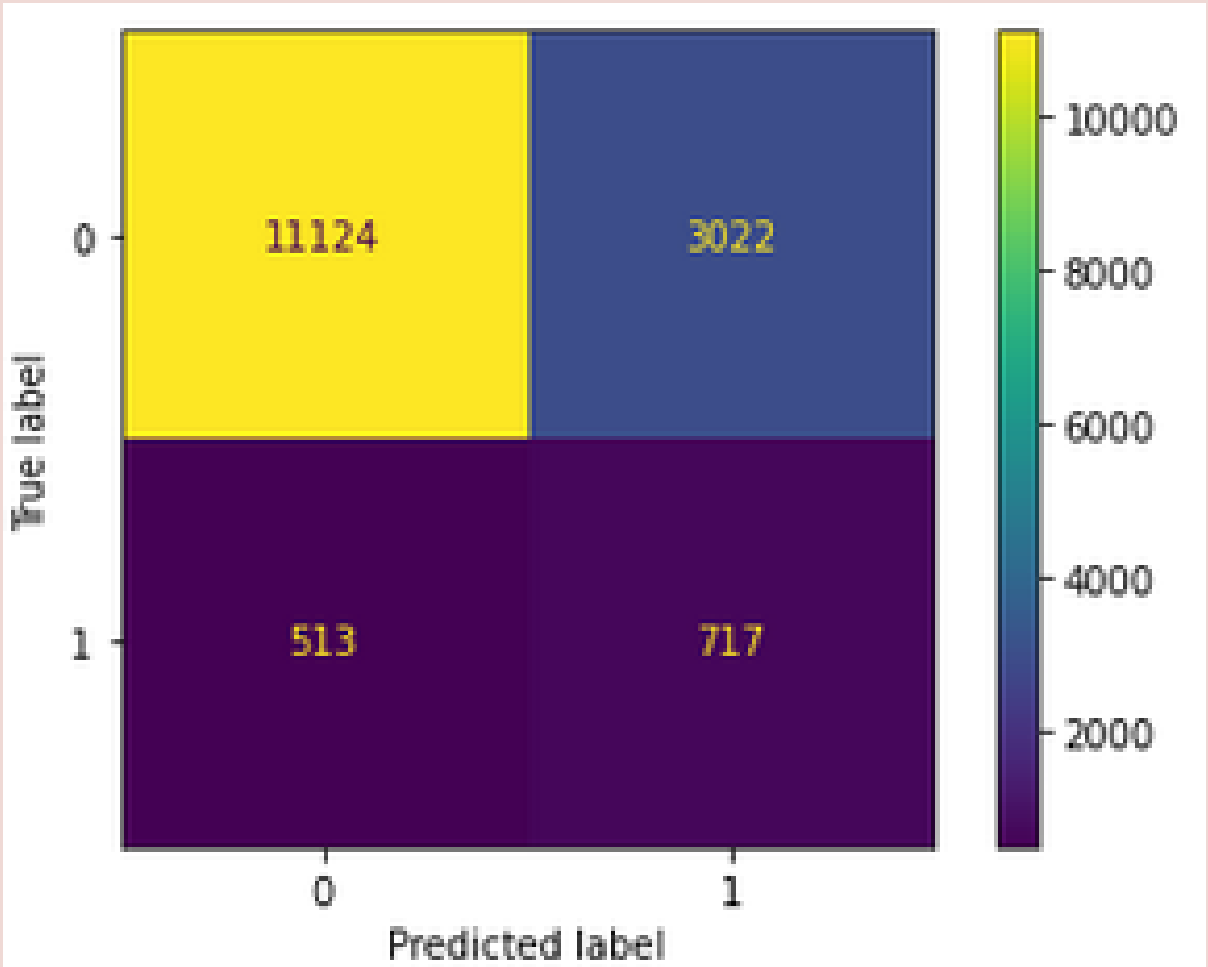
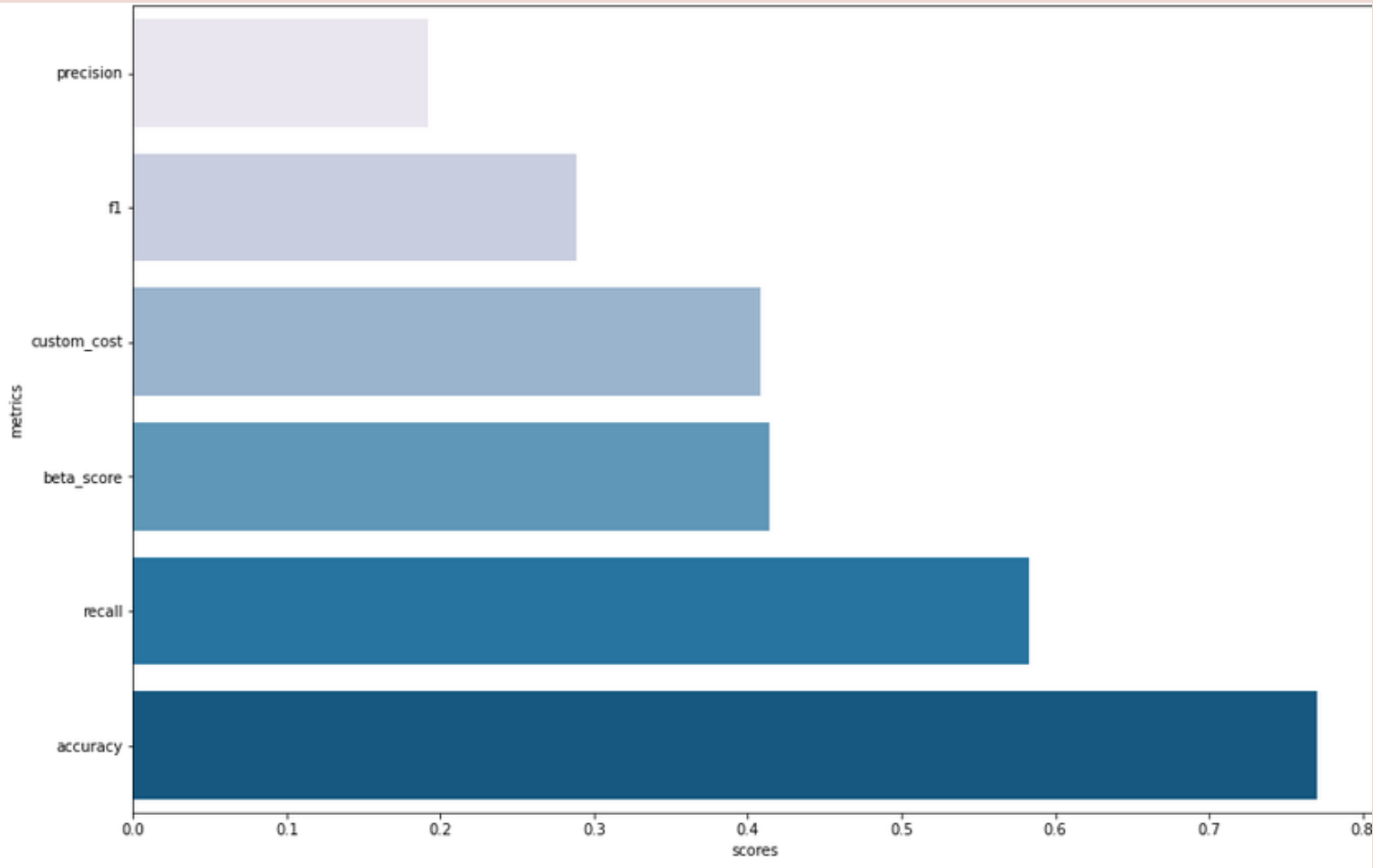


# Démarche pour le choix du modèle finale





# Choix du seuil



- Seuil: 0.35



# Interprétation globale

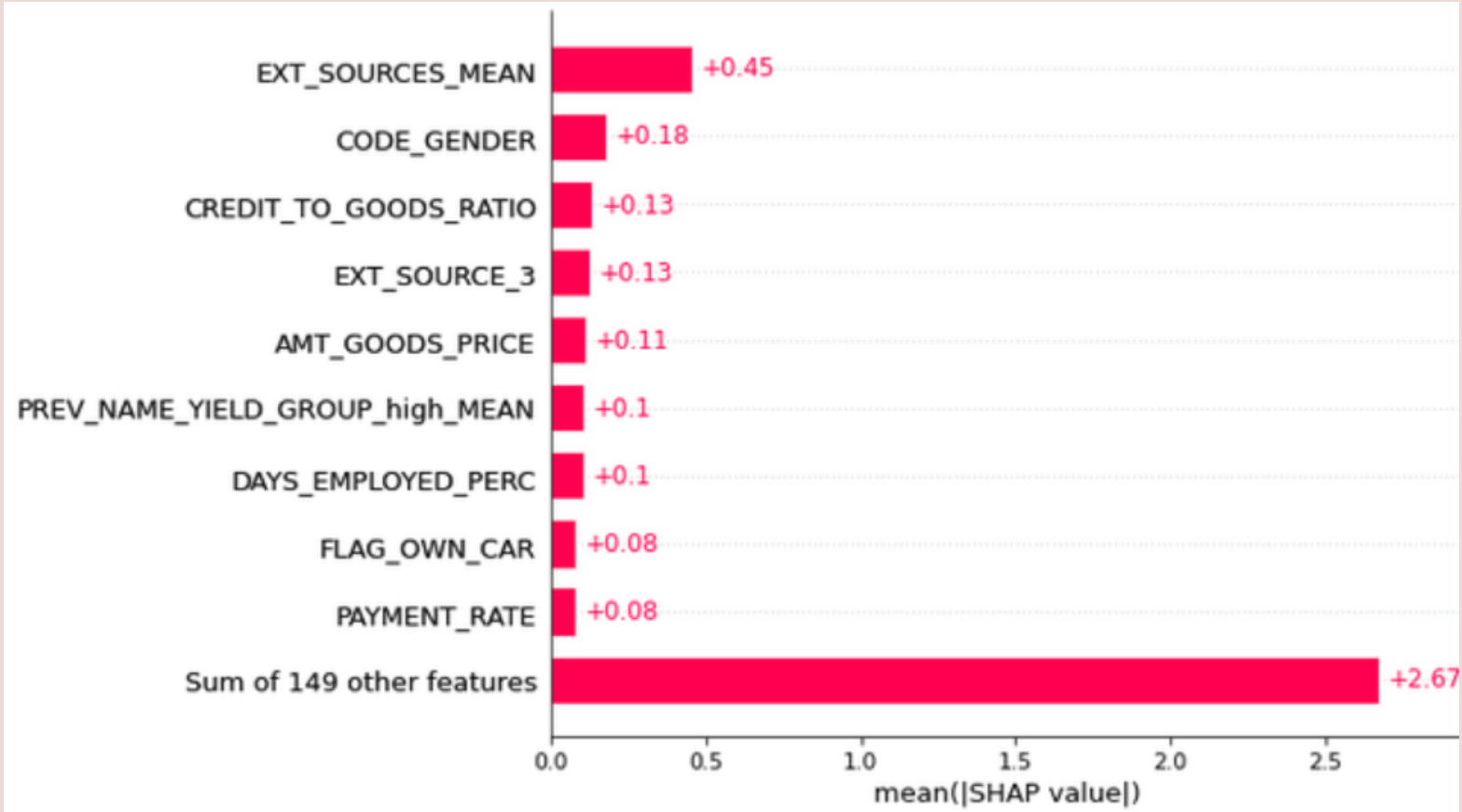


Figure 1

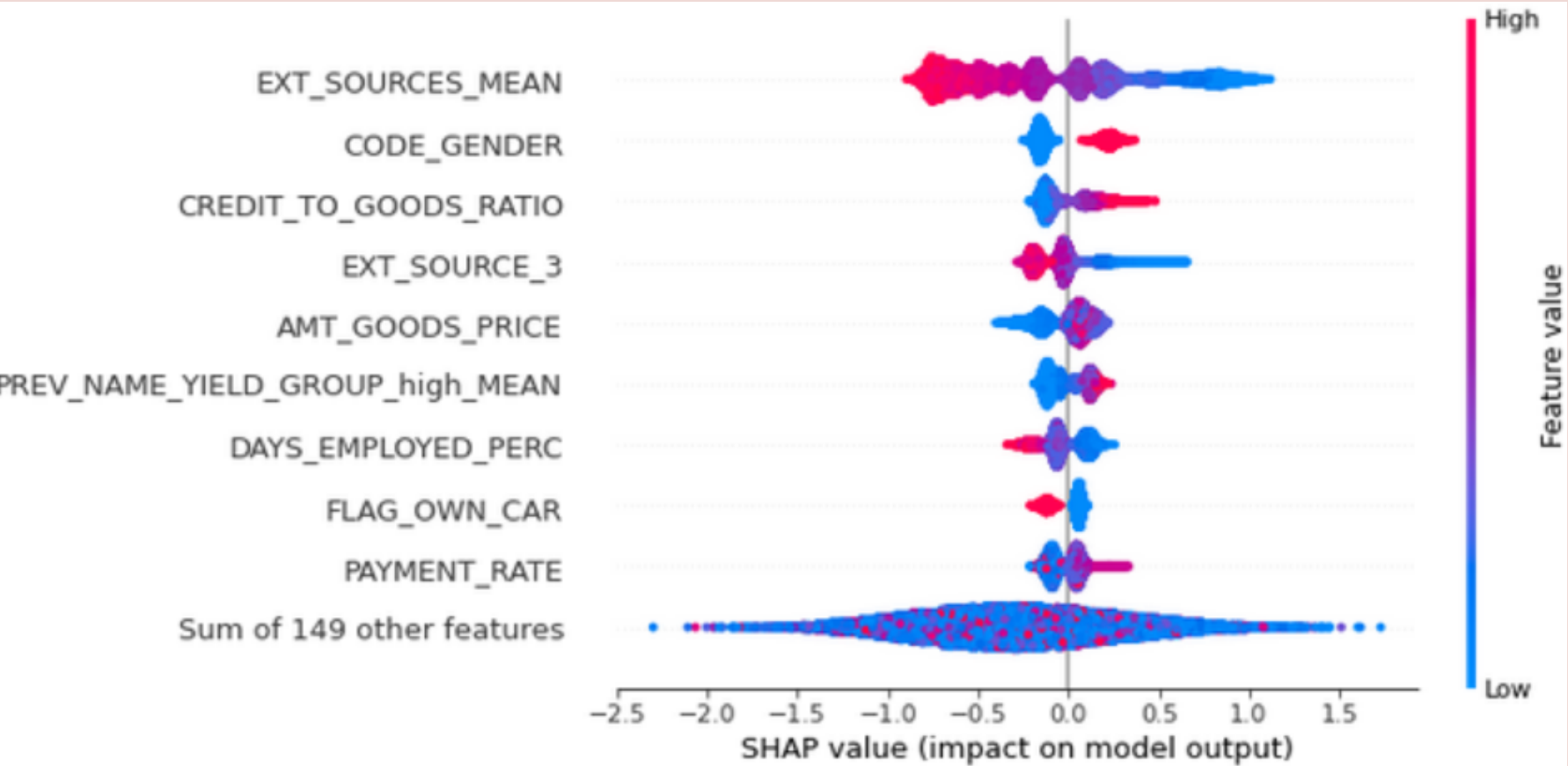
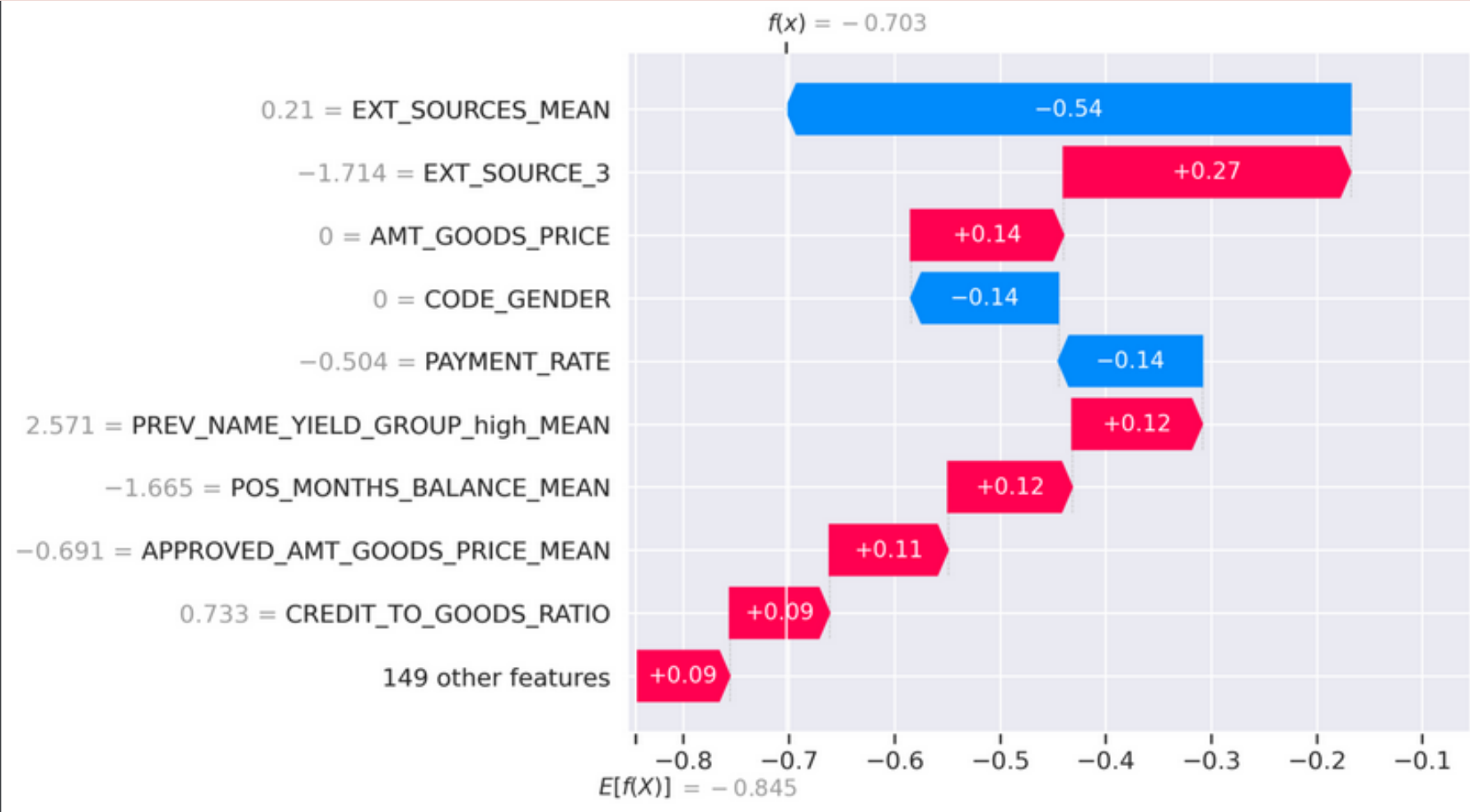


Figure 2



# Interprétation locale





# Analyse du data drift



# Analyse du data drift

Dataset Drift

Dataset Drift is NOT detected. Dataset drift detection threshold is 0.5

158

Columns

12

Drifted Columns

0.0759

Share of Drifted Columns

Data Drift Summary

Drift is detected for 7.595% of columns (12 out of 158).

Q Search

×

Column	Type	Reference Distribution	Current Distribution	Data Drift	Stat Test	Drift Score
> PAYMENT_RATE	num			Detected	Wasserstein distance (normed)	0.577878
> AMT_GOODS_PRICE	num			Detected	Wasserstein distance (normed)	0.212276
> NAME_CONTRACT_TYPE	cat			Detected	Jensen-Shannon distance	0.14782
> PREV_NAME_PORTFOLIO_XNA_MEAN	num			Detected	Wasserstein distance (normed)	0.136548
> PREV_NAME_YIELD_GROUP_XNA_MEAN	num			Detected	Wasserstein distance (normed)	0.133758
> FLAG_EMAIL	num			Detected	Jensen-Shannon distance	0.124014
> PHONE_TO_BIRTH_RATIO	num			Detected	Wasserstein distance (normed)	0.123733
> PREV_NAME_CONTRACT_STATUS_Approved_MEAN	num			Detected	Wasserstein distance (normed)	0.117387
> CLOSED_DAYS_CREDIT_UPDATE_MEAN	num			Detected	Wasserstein distance (normed)	0.112929
> PREV_NAME_PAYMENT_TYPE_XNA_MEAN	num			Detected	Wasserstein distance (normed)	0.111889



# Conclusion

- Bonnes performances globales
- Mieux comprendre le jeu de données
- Tester d'autres méthodes de gestion des déséquilibres des classes
- Affiner la recherche d'hyperparamètres
- Affiner la métrique avec l'équipe métier



# Présentation de la pipeline de déploiement





# Github actions

## build

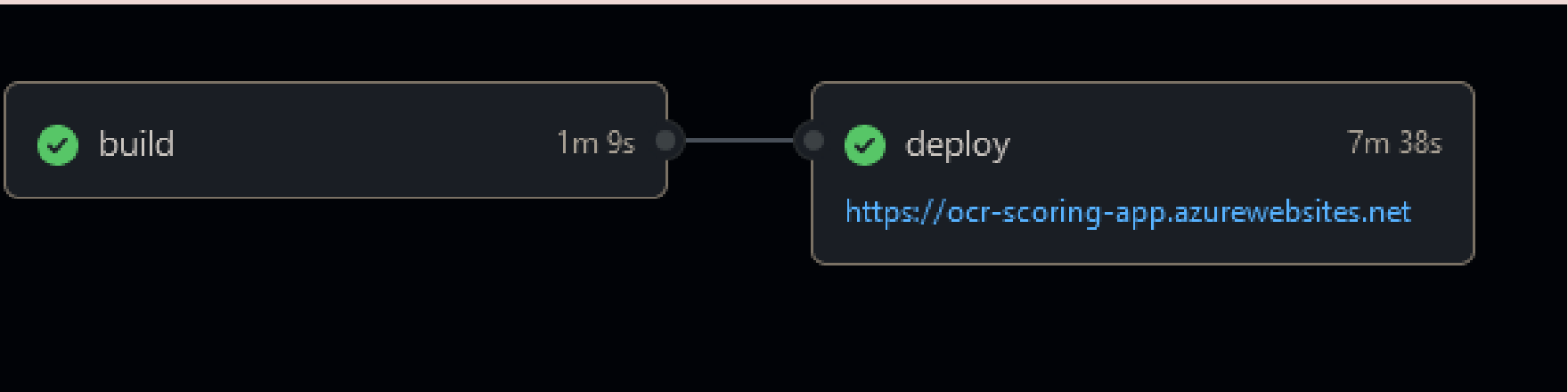
succeeded yesterday in 1m 9s

- > ✓ Set up job
- > ✓ Run actions/checkout@v2
- > ✓ Set up Python version
- > ✓ Create and start virtual environment
- > ✓ Install dependencies
- > ✓ Test with pytest
- > ✓ Upload artifact for deployment jobs
- > ✓ Post Run actions/checkout@v2
- > ✓ Complete job

## deploy

succeeded yesterday in 7m 38s

- > ✓ Set up job
- > ✓ Download artifact from build job
- > ✓ Deploy to Azure Web App
- > ✓ Complete job





# Tests unitaires

```
21 def test_check_columns():
22     df = pd.read_json(test_data)
23     columns = df.columns
24     assert all(item in columns for item in columns_names)
25
26 def test_check_pipeline_predict():
27     df = pd.read_json(test_data)
28     y_pred = pipeline.predict(df)
29     y_proba = pipeline.predict_proba(df)
30     assert y_pred[0] in [0, 1]
31
32     assert len([y_proba[0][1]]) == 1
33
34     assert (y_proba[0][1] >= 0) & (y_proba[0][1] <= 1)
35
36 def test_shape_features_prep():
37     df = pd.read_json(test_data)
38     actual_shape = pipeline[0].transform(df).shape
39     expected_shape = (1, 158)
40     assert actual_shape == expected_shape
41
42 def test_check_trans_nums():
43     df = pd.read_json(test_data)
44     data = pipeline_nums.transform(df[nums_columns_name])
45     assert data.shape == (1, 147)
46
```

```
✓ Test with pytest

1 ▶ Run pytest
6 ===== test session starts =====
7 platform linux -- Python 3.9.16, pytest-7.2.1, pluggy-1.0.0
8 rootdir: /home/runner/work/scoring-app/scoring-app
9 collected 4 items
10
11 tests/test_api.py .... [100%]
12
13 ===== 4 passed in 1.76s =====
```



# Démo du dashboard



MERCI DE VOTRE

ATTENTION

---