



Kaggle Titanic 入門 02

スコア向上策

庭師

庭子アバター v01.50
Avatar Museum12に展示中



0.73444 だった。gender_submission.csvよりもスコアが低いなぜ？

- スコアが 0.73444 で、gender_submission.csv(おそらく 0.76555 程度)より低い理由は、主に以下のような要因が考えられます：

✓ gender_submission.csv は「女性は生存、男性は死亡」のルールベース予測

- 非常にシンプルですが、Titanicデータでは 女性の生存率が高く、男性の生存率が低い ため、このルールだけでも そろそろ高い精度 を出せます。
- 実際、Kaggle公式のベースラインとしても提示されています。

🔍 あなたの RandomForest モデルのスコアが低い理由

1. 過学習 or 過剰な特徴量

- Title, FareBin, AgeBin, IsAlone などの特徴量を多数導入しましたが、適切に効いていない or ノイズになっている可能性があります。

2. 欠損値補完の方法が単純

- 年齢や運賃などを「中央値」や「最頻値」で一律補完していますが、より賢い方法（例: Sex や Pclass を使ったグループ中央値補完）だと精度が上がります。

3. 特徴量のスケールや調整がされていない

- RandomForestはスケールに強いですが、AgeBin や FareBin などカテゴリ化が精度を落としている可能性も。

4. パラメータチューニング未実施

- n_estimators, max_depth, min_samples_split などのチューニングが行われていません。
- デフォルト設定のままでは精度が出ないケースも多いです。

5. 学習データとテストデータで分布がズレている

- 特に FareBin や Title などがテスト側に未出現カテゴリを含んでいた可能性も。

✓ 改善案（提案）

1. シンプルな特徴量に絞る

- Pclass, Sex, Age, Fare などだけで

2. モデルをチューニング

- RandomForest のパラメータ調整
- または LightGBM に変更（高性能）

3. クロスバリデーションで検証

- train.csv 内で交差検証して精度を確認
- Public LB を基準にせず過学習を防ぐ

前回はこちらまで
でした！

この続きです！

← END 結論

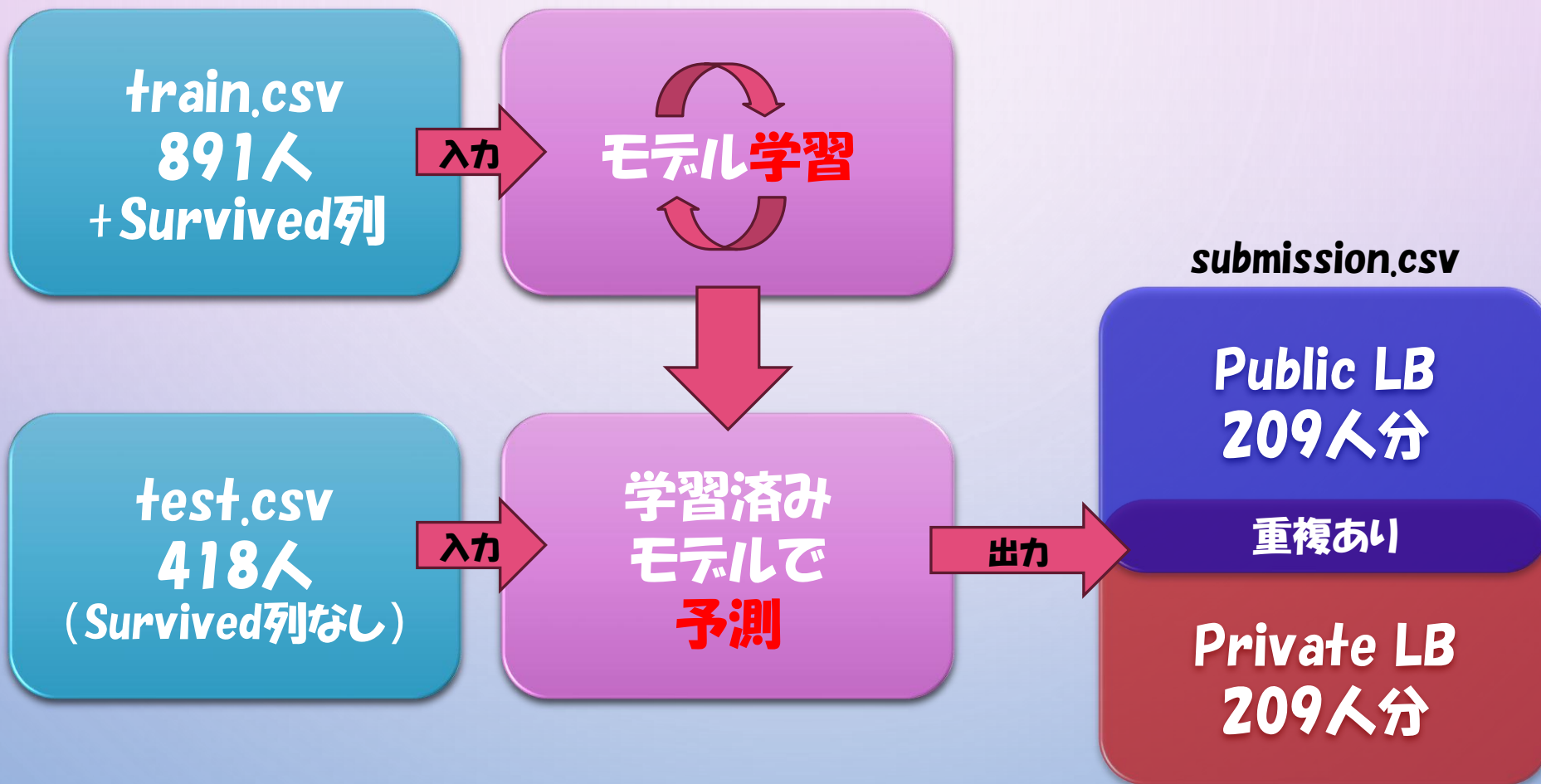
gender_submission.csv は実は「データに対して最適なバイアスのルール」に基づいていて、意外と強いベースラインです。

一方でモデルを使う場合は適切な特徴量選択・欠損処理・チューニングが必要です。

希望があれば、「精度を gender_submission.csv 以上に上げるための改良案」をコード付きで出せます。どうしますか？

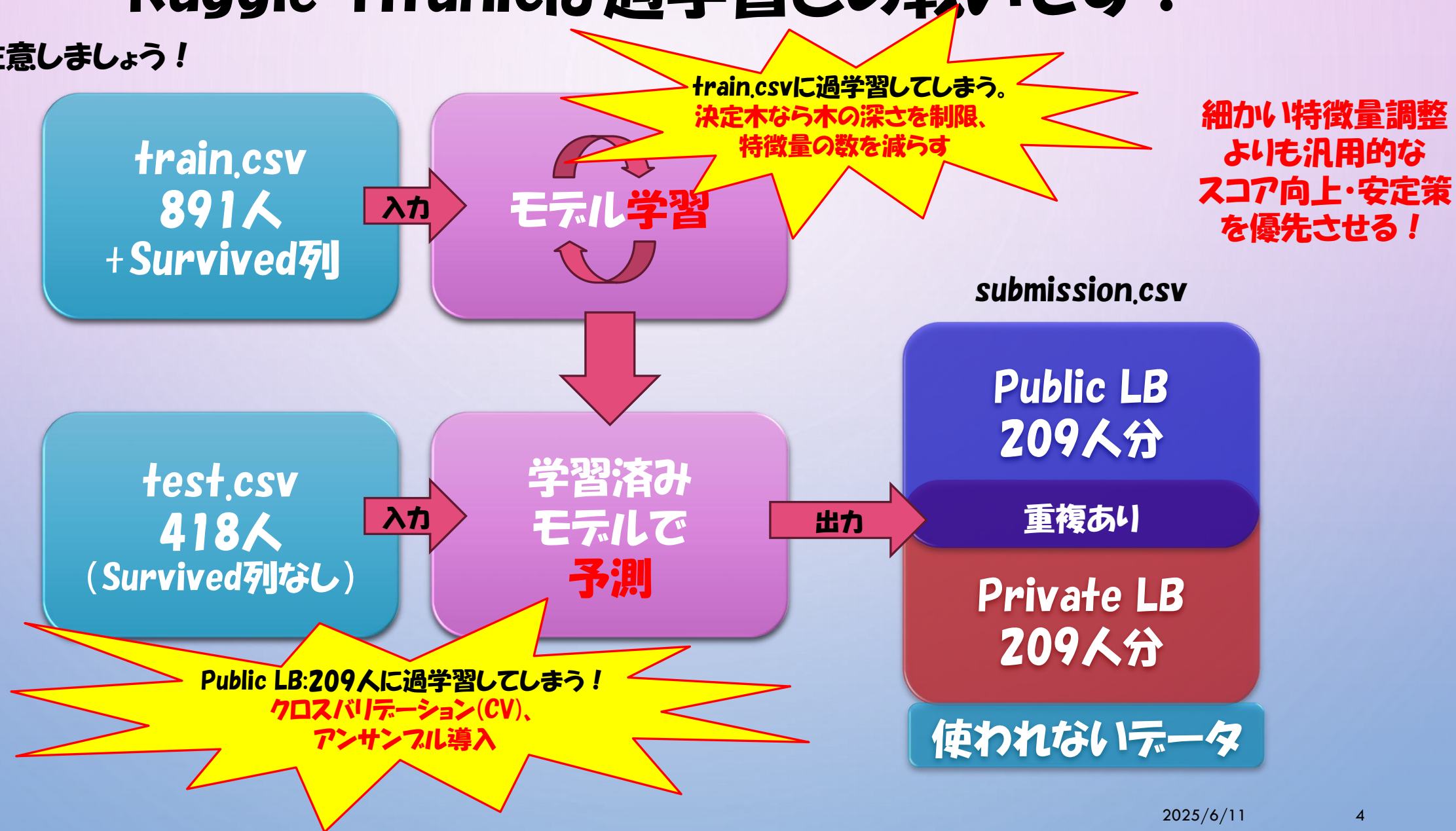
続く！！

train, test, Public LB, Private LB の関係



Kaggle Titanicは過学習との戦いです！

2カ所、注意しましょう！



戦略：まずスコアを安定させる！汎化対策！

- ChatGPTでtitanic予測させると、0.73~0.76くらいで大きくばらつきます。

【今回のシナリオ】

- さらに過学習(over fitting)のリスクがあるので、細かい特徴量調整やるまえに安定化・再現性をあげます。
 - `random_state` 固定 (`random_state=42` はよく使います。他の数字でもいい)
- その後汎用的なスコア向上策をいれます。
 - クロスバリデーション(CV)は早めに導入する。汎化対策。
 - 複数モデルの結果を集約するときは、平均より多数決(mode)が良さそう
- 段階的に特徴量エンジニアリングしていく

プログラムコード部は文字が小さいので
雰囲気みてやってください！
プログラムコードは配布します。

参考にさせてもらいました！

- KaggleチュートリアルTitanicで上位2%以内に入るノウハウ
- <https://qiita.com/jun40vn/items/d8a1f71fae680589e05c>



@jun40vn (Yuichi Sugita)

KaggleチュートリアルTitanicで上位2%以内に入るノウハウ

Python 機械学習 Kaggle titanic

投稿日 2019年11月21日

1.はじめに

機械学習を勉強し始めたので、なにはともあれ **Kaggle**の初心者向けチュートリアル、**タイタニック**にトライしてみました。

最初は、日本語のWeb情報を参考に、自分なりに特徴量を色々いじってみたわけですが、学習データでは良い精度が出ても、テストデータをSubmitすると、思った様に精度は上がらず、なかなか80%の壁を破れず悶々としていました。

そうした中で、英文なので取っ付き難いですが、Kaggle/Taitanicの**Notebook**にある先人の知恵を拝借して、やっと上位2%に入れたので、特に参考になった点を中心に備忘録として残します。



kaggle titanic をランダムフォレストで予測して

01_kaggle_rf_titanic_0.75358.py

gender_submission.csv :

02_kaggle_lgbm_titanic_0.74641.py

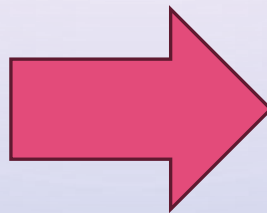
```
1 import pandas as pd
2 from sklearn.ensemble import RandomForestClassifier
3
4 # CSVファイルの読み込み
5 train_df = pd.read_csv("train.csv")
6 test_df = pd.read_csv("test.csv")
7
8 # 前処理
9 train_df["Sex"] = train_df["Sex"].map({"male": 0, "female": 1})
10 test_df["Sex"] = test_df["Sex"].map({"male": 0, "female": 1})
11
12 train_df["Age"].fillna(train_df["Age"].median(), inplace=True)
13 test_df["Age"].fillna(test_df["Age"].median(), inplace=True)
14 test_df["Fare"].fillna(test_df["Fare"].median(), inplace=True)
15
16 train_df["Embarked"].fillna("S", inplace=True)
17 test_df["Embarked"].fillna("S", inplace=True)
18 train_df["Embarked"] = train_df["Embarked"].map({"S": 0, "C": 1, "Q": 2})
19 test_df["Embarked"] = test_df["Embarked"].map({"S": 0, "C": 1, "Q": 2})
20
21 # 特徴量の指定
22 features = ["Pclass", "Sex", "Age", "SibSp", "Parch", "Fare", "Embarked"]
23 X = train_df[features]
24 y = train_df["Survived"]
25 X_test = test_df[features]
26
27 # モデル学習と予測
28 model = RandomForestClassifier(random_state=42)
29 model.fit(X, y)
30 predictions = model.predict(X_test)
31
32 # 提出ファイルの作成
33 submission = pd.DataFrame({
34     "PassengerId": test_df["PassengerId"],
35     "Survived": predictions
36 })
37 submission.to_csv("01_kaggle_rf_titanic.csv", index=False)
```

前処理を入れてくれている。

再現性確保のため random_state 固定
は必ず入れるように指示しましょう。



Kaggle code
(notebook)
向けLGBMにして



LGBMの方がスコア
が落ちる。チューニ
ングすれば、LGBM
の方が良いスコア
が出る可能性ある
が、ChatGPT内で
実行出来ることも
あり、今回は主にな
ランダムフォレスト
を採用。

```
1 import pandas as pd
2 import lightgbm as lgb
3
4 # CSVファイルの読み込み
5 train_df = pd.read_csv("/kaggle/input/titanic/train.csv")
6 test_df = pd.read_csv("/kaggle/input/titanic/test.csv")
7
8 # 前処理
9 train_df["Sex"] = train_df["Sex"].map({"male": 0, "female": 1})
10 test_df["Sex"] = test_df["Sex"].map({"male": 0, "female": 1})
11
12 train_df["Age"].fillna(train_df["Age"].median(), inplace=True)
13 test_df["Age"].fillna(test_df["Age"].median(), inplace=True)
14 test_df["Fare"].fillna(test_df["Fare"].median(), inplace=True)
15
16 train_df["Embarked"].fillna("S", inplace=True)
17 test_df["Embarked"].fillna("S", inplace=True)
18 train_df["Embarked"] = train_df["Embarked"].map({"S": 0, "C": 1, "Q": 2})
19 test_df["Embarked"] = test_df["Embarked"].map({"S": 0, "C": 1, "Q": 2})
20
21 # 特徴量の指定
22 features = ["Pclass", "Sex", "Age", "SibSp", "Parch", "Fare", "Embarked"]
23 X = train_df[features]
24 y = train_df["Survived"]
25 X_test = test_df[features]
26
27 # LightGBM モデルで学習と予測
28 model = lgb.LGBMClassifier(random_state=42)
29 model.fit(X, y)
30 preds = model.predict(X_test)
31
32 # 提出ファイルの作成
33 submission = pd.DataFrame({
34     "PassengerId": test_df["PassengerId"],
35     "Survived": preds
36 })
37 submission.to_csv("02_kaggle_lgbm_titanic.csv", index=False)
```

Kaggle code (notebook)
にコピーして実行

特徴量全部(LE:ラベルエンコーディング版)

03_rf_LabelE_nodrop_0.77272.py



特徴量全部入れてラベルエンコーディングしてからランダムフォレスト予測

このあたりから完全に同じコードは出てこないと思います。
必要に応じて、提供されているコードをChatGPTへアップロードして使ってください。

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.preprocessing import LabelEncoder
5
6 # データ読み込みと結合
7 train_df = pd.read_csv("train.csv")
8 test_df = pd.read_csv("test.csv")
9 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
10 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
11
12 # 欠損補完
13 df['Age'] = df['Age'].fillna(df['Age'].median())
14 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
15
16 # ラベルエンコーディング
17 label_cols = ['Sex', 'Embarked', 'Name', 'Ticket', 'Cabin']
18 for col in label_cols:
19     le = LabelEncoder()
20     df[col] = le.fit_transform(df[col].astype(str))
21
22 # X_train, y_train, X_test 分離
23 dfm = df # 列を削除しない
24 train = dfm[dfm['Survived'].notnull()]
25 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
26 X = train.drop('Survived', axis=1).values
27 y = train['Survived'].astype(int).values
28 X_test = test.values
29
30 # モデル学習・予測
31 model = RandomForestClassifier(random_state=42)
32 model.fit(X, y)
33 final_preds = model.predict(X_test)
34
35 # submissionファイル出力
36 submission = pd.DataFrame({
37     "PassengerId": test_df["PassengerId"],
38     "Survived": final_preds
39 })
40 submission.to_csv("03_rf_LabelE_nodrop.csv", index=False)
```

LE

特徴量数(Survivedを除く):11個

勝手に順番変えられる時がある。安定化のために
sort=False
を指定する。
順番変わるとLabel Encoding, Kaggle Score判定など挙動が変わる可能性あり。

gender_submission.csv :
0.76555 は超えた !

列名	データ型	説明
PassengerId	int	乗客ID。ユニークな識別番号。モデルには不要（識別用）。
Survived	int (0/1)	生存フラグ（0 = 死亡, 1 = 生存）。 目的変数 （ <code>train.csv</code> のみに存在）。
Pclass	int (1-3)	チケットクラス（1 = 上級, 2 = 中級, 3 = 下級）。社会階級の proxy。
Name	str	名前。敬称（ <code>Mr</code> , <code>Mrs</code> , <code>Miss</code> など）を抽出可能。
Sex	str	性別（ <code>male</code> / <code>female</code> ）。非常に重要な特徴。
Age	float	年齢。欠損あり（約20%）。子供・大人の判別などに有用。
SibSp	int	タイタニックに同乗していた 兄弟または配偶者 の数。家族構成の一部を表す。
Parch	int	タイタニックに同乗していた 親または子供 の数。これも家族構成を示す。
Ticket	str	チケット番号。情報抽出可能だが、ノイズも多い。
Fare	float	チケット料金。社会的地位や等級と相関がある。欠損値あり（ <code>test.csv</code> のみ）。
Cabin	str	キャビン番号。欠損多数。文字列の先頭（デッキレベル）を抽出して使う。
Embarked	str	乗船港（ <code>C</code> = Cherbourg, <code>Q</code> = Queenstown, <code>S</code> = Southampton）。欠損あり。

文字列はLE必要

Titanic 初期特徴量 おさらい

‘Braund, Mr. Owen Harris’ など

‘11774’, ‘STON/O2. 3101282’ など

‘B28’, ‘E46’, ‘F G73’ など

特徴量削減版(LE版)

特徴量数(*Survived*を除く):
7個

【削除した特徴量と理由】

• *PassengerId*: ただの固有ID。
いらない。

• *Name*, *Ticket*, *Cabin*:
単なるラベルエンコーディング
では意味が無い。

n_estimators=200はこの後の
CV安定させる大きめの木を採用。
デフォルトは100。

Titanicのような小さいテーブル
では、特徴量削減し、
*max_depth*制限すると、過学習
防げて、スコアが上がる。

n_estimators=200, *max_depth*=6

04_rf_LabelE_0.76315.py

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.preprocessing import LabelEncoder
5
6 # データ読み込みと結合
7 train_df = pd.read_csv("train.csv")
8 test_df = pd.read_csv("test.csv")
9 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
10 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
11
12 # 欠損補完
13 df['Age'] = df['Age'].fillna(df['Age'].median())
14 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
15
16 # ラベルエンコーディング
17 label_cols = ['Sex', 'Embarked']
18 for col in label_cols:
19     le = LabelEncoder()
20     df[col] = le.fit_transform(df[col].astype(str))
21
22 # X_train, y_train, X_test 分離
23 #['Pclass', 'Sex', 'Age', 'Fare', 'Embarked', 'SibSp', 'Parch']
24 dfm = df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], errors='ignore')
25 train = dfm[dfm['Survived'].notnull()]
26 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
27 X = train.drop('Survived', axis=1).values
28 y = train['Survived'].astype(int).values
29 X_test = test.values
30
31 # モデル学習・予測
32 model = RandomForestClassifier(random_state=42)
33 model.fit(X, y)
34 final_preds = model.predict(X_test)
35
36 # submissionファイル出力
37 submission = pd.DataFrame({
38     "PassengerId": test_df["PassengerId"],
39     "Survived": final_preds
40 })
41 submission.to_csv("04_rf_LabelE.csv", index=False)
```

05_rf_LabelE_200-6_0.78229.py

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.preprocessing import LabelEncoder
5
6 # データ読み込みと結合
7 train_df = pd.read_csv("train.csv")
8 test_df = pd.read_csv("test.csv")
9 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
10 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
11
12 # 欠損補完
13 df['Age'] = df['Age'].fillna(df['Age'].median())
14 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
15
16 # ラベルエンコーディング
17 label_cols = ['Sex', 'Embarked']
18 for col in label_cols:
19     le = LabelEncoder()
20     df[col] = le.fit_transform(df[col].astype(str))
21
22 # X_train, y_train, X_test 分離
23 #['Pclass', 'Sex', 'Age', 'Fare', 'Embarked', 'SibSp', 'Parch']
24 dfm = df.drop(columns=['PassengerId', 'Name', 'Ticket', 'Cabin'], errors='ignore')
25 train = dfm[dfm['Survived'].notnull()]
26 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
27 X = train.drop('Survived', axis=1).values
28 y = train['Survived'].astype(int).values
29 X_test = test.values
30
31 # モデル学習・予測
32 model = RandomForestClassifier(n_estimators=200, max_depth=6, random_state=42)
33 model.fit(X, y)
34 final_preds = model.predict(X_test)
35
36 # submissionファイル出力
37 submission = pd.DataFrame({
38     "PassengerId": test_df["PassengerId"],
39     "Survived": final_preds
40 })
41 submission.to_csv("05_rf_LabelE_200-6.csv", index=False)
```

CV StratifiedKFold:5 説明

全データで学習/検証するため、trainに過学習しすぎる可能性あり。未知のデータに弱くなる→汎化性が低い

train.csv (Survivedあり): 891人

test.csv: 418人

04_rf_LabelE_0.76315.py
のケース

学習 / 検証 データ

学習済み
モデル

コンペtestデータ

学習データ

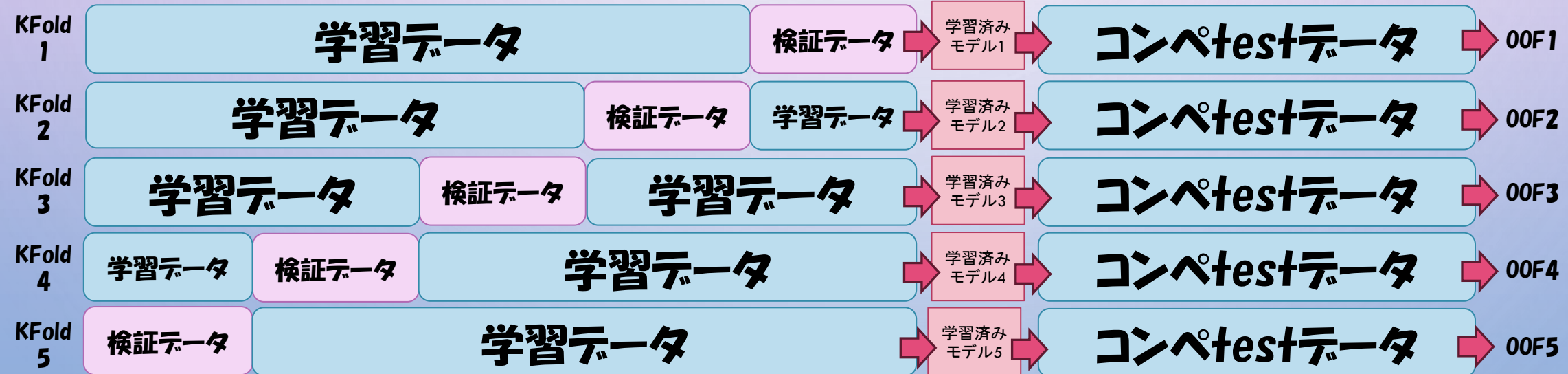
検証データ

学習済み
モデル

コンペtestデータ

汎化性は増すが、
学習に使われない
データが出てくる。

CV(Cross Validation) KFold:5



多数決
平均もあり

最終予測値

通常 StratifiedKFold() を使う。各特徴量を均等に配分してくれる。

00F: Out of Fold

CV:SKFold:5導入

CVスコアは高めに表示されます

CV SKFold:5 を導入して、少しスコアが下がったが・・・、Private LBでの安定性(スコアが下がらないよう)を考慮して、
CV SKFold:5、n_estimators:200で進める。

Fold 1 CV Accuracy: 0.84916
Fold 2 CV Accuracy: 0.80899
Fold 3 CV Accuracy: 0.81461
Fold 4 CV Accuracy: 0.84270
Fold 5 CV Accuracy: 0.85955
Mean CV Accuracy: 0.83500

06_rf_LabelE_200-6_skfold5_0.77990.py

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.metrics import accuracy_score
7 from scipy.stats import mode
8
9 # データ読み込みと結合
10 train_df = pd.read_csv("train.csv")
11 test_df = pd.read_csv("test.csv")
12 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
13 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
14
15 # Age, Fare補完 : base
16 df['Age'] = df['Age'].fillna(df['Age'].median())
17 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
18
19 # 特徴量とラベルエンコーディング, Drop 'PassengerId', 'Name', 'Ticket', 'Cabin'
20 features = ['Pclass', 'Sex', 'Age', 'Fare', 'Embarked', 'SibSp', 'Parch']
21 dfm = df[['Survived'] + features].copy()
22 for col in ['Sex', 'Embarked']:
23     le = LabelEncoder()
24     dfm[col] = le.fit_transform(dfm[col].astype(str))
25
26 # train/test 分割
27 train = dfm[dfm['Survived'].notnull()]
28 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
29 X = train.drop('Survived', axis=1).values
30 y = train['Survived'].astype(int).values
31 X_test = test.values
```

```
32
33 # StratifiedKFold + Mode voting(多数決)、CVスコア表示
34 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
35 test_preds = []
36 val_scores = []
37 for fold, (train_idx, val_idx) in enumerate(skf.split(X, y), 1):
38     model = RandomForestClassifier(n_estimators=200, max_depth=6, random_state=42)
39     model.fit(X[train_idx], y[train_idx])
40     val_pred = model.predict(X[val_idx])
41     val_acc = accuracy_score(y[val_idx], val_pred)
42     val_scores.append(val_acc)
43     print(f"Fold {fold} CV Accuracy: {val_acc:.5f}")
44     test_preds.append(model.predict(X_test))
45 print(f"Mean CV Accuracy: {np.mean(val_scores):.5f}")
46
47 # 多数決で最終予測を決定
48 test_preds = np.array(test_preds)
49 final_preds = mode(test_preds, axis=0, keepdims=False).mode
50
51 # サブミッション作成
52 submission = pd.DataFrame({
53     'PassengerId': test_df['PassengerId'],
54     'Survived': final_preds.astype(int)
55 })
56 submission.to_csv("06_rf_LabelE_200-6_skfold5.csv", index=False)
```

SKFold:5
CVスコア表示なしだと
もっとシンプル。

多数決アンサンブル
meanもよく見る。

前処理導入

07_rf_LabelE_200-6_skfold5_Fe_0.78947.py

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.metrics import accuracy_score
7 from scipy.stats import mode
8
9 # データ読み込みと結合
10 train_df = pd.read_csv("train.csv")
11 test_df = pd.read_csv("test.csv")
12 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
13 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
14
15 # Age, Fare補完 : base
16 df['Age'] = df['Age'].fillna(df['Age'].median())
17 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
18
19 ##### 特徴量エンジニアリング : 全てまとめて 0.01 の精度向上
20 # Title作成 : 0.007 の精度向上, ['Master', 'Miss', 'Mr', 'Mrs', 'Officer', 'Royalty']に集約
21 df['Title'] = df['Name'].map(lambda x: x.split(',')[1].split('.')[0])
22 df['Title'].replace(['Capt', 'Col', 'Major', 'Dr', 'Rev'], 'Officer', inplace=True)
23 df['Title'].replace(['Don', 'Sir', 'the Countess', 'Lady', 'Dona'], 'Royalty', inplace=True)
24 df['Title'].replace(['Mme', 'Ms'], 'Mrs', inplace=True)
25 df['Title'].replace(['Mlle'], 'Miss', inplace=True)
26 df['Title'].replace(['Jonkheer'], 'Master', inplace=True)
27 # Cabin補完 : 0.005 の精度向上
28 df['Cabin'] = df['Cabin'].fillna('Unknown')
29 df['Cabin_L'] = df['Cabin'].str.get(0) # 最初の文字をCabinラベルとして使用
30 # Embarked補完 : 0.005 の精度向上
31 df['Embarked'] = df['Embarked'].fillna('S') # 'S'は最貧度
32 # Familyラベル : シンプルFamilySizeより0.015、なしより0.017 の精度向上
33 df['Family'] = df['SibSp'] + df['Parch'] + 1 # 家族人数を計算 (本人を含める)
34 df['Family_L'] = np.select([
35     (df['Family'] >= 2) & (df['Family'] <= 4), # 2:中規模
36     (df['Family'] == 1) | ((df['Family'] >= 5) & (df['Family'] <= 7)), # 1:単独 or やや多め
37     df['Family'] >= 8 # 0:超大規模
38 ], [2, 1, 0])
```

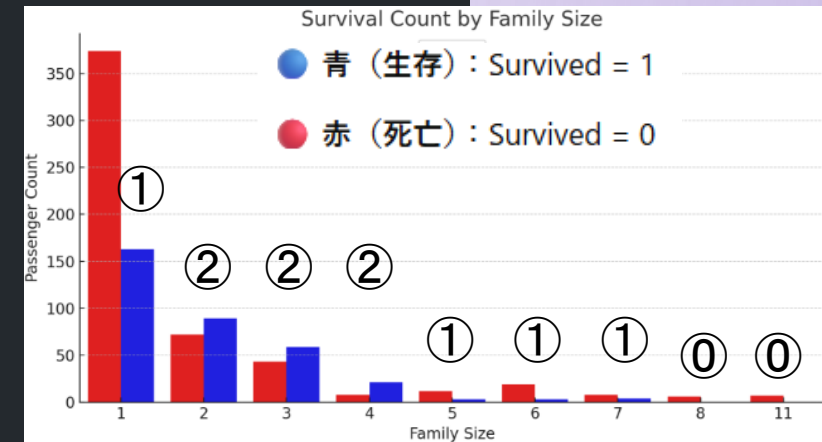
前処理追加

```
40 # 特徴量とラベルエンコーディング
41 features = ['Pclass', 'Sex', 'Age', 'Fare', 'Embarked', 'Title', 'Family_L', 'Cabin_L']
42 dfm = df[['Survived'] + features].copy()
43 for col in ['Sex', 'Embarked', 'Title', 'Cabin_L']:
44     le = LabelEncoder()
45     dfm[col] = le.fit_transform(dfm[col].astype(str))
46
47 # train/test 分割
48 train = dfm[dfm['Survived'].notnull()]
49 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
50 X = train.drop('Survived', axis=1).values
51 y = train['Survived'].astype(int).values
52 X_test = test.values
53
54 # StratifiedKFold + Mode voting(多数決)、CVスコア表示
55 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
56 test_preds = []
57 val_scores = []
58 for fold, (train_idx, val_idx) in enumerate(skf.split(X, y), 1):
59     model = RandomForestClassifier(n_estimators=200, max_depth=6, random_state=42)
60     model.fit(X[train_idx], y[train_idx])
61     val_pred = model.predict(X[val_idx])
62     val_acc = accuracy_score(y[val_idx], val_pred)
63     val_scores.append(val_acc)
64     print(f"Fold {fold} CV Accuracy: {val_acc:.5f}")
65     test_preds.append(model.predict(X_test))
66 print(f"Mean CV Accuracy: {np.mean(val_scores):.5f}")
67
68 # 多数決で最終予測を決定
69 test_preds = np.array(test_preds)
70 final_preds = mode(test_preds, axis=0, keepdims=False).mode
71
72 # サブミッション作成
73 submission = pd.DataFrame({
74     'PassengerId': test_df['PassengerId'],
75     'Survived': final_preds.astype(int)
76 })
77 submission.to_csv("07_rf_LabelE_200-6_skfold5_Fe.csv", index=False)
```

前処理説明

全体で0.01ほどのスコア向上。各特徴量を個別にon/offすると、下記コメント数値(非CV)でよい影響が大きいようだ。

```
#### 特徴量エンジニアリング : 全てまとめて 0.01 の精度向上
# Title作成 : 0.007 の精度向上, ['Master', 'Miss', 'Mr', 'Mrs', 'Officer', 'Royalty']に集約
df['Title'] = df['Name'].map(lambda x: x.split(',')[1].split('.')[0])
df['Title'].replace(['Capt', 'Col', 'Major', 'Dr', 'Rev'], 'Officer', inplace=True)
df['Title'].replace(['Don', 'Sir', 'the Countess', 'Lady', 'Dona'], 'Royalty', inplace=True)
df['Title'].replace(['Mme', 'Ms'], 'Mrs', inplace=True)
df['Title'].replace(['Mlle'], 'Miss', inplace=True)
df['Title'].replace(['Jonkheer'], 'Master', inplace=True)
# Cabin補完 : 0.005 の精度向上
df['Cabin'] = df['Cabin'].fillna('Unknown')
df['Cabin_L'] = df['Cabin'].str.get(0) # 最初の文字をCabinラベルとして使用
# Embarked補完 : 0.005 の精度向上
df['Embarked'] = df['Embarked'].fillna('S') # 'S'は最貧度
# Familyラベル : シンプルFamilySizeより0.015、なしより0.017 の精度向上
df['Family'] = df['SibSp'] + df['Parch'] + 1 # 家族人数を計算 (本人を含める)
df['Family_L'] = np.select([
    (df['Family'] >= 2) & (df['Family'] <= 4), # 2:中規模
    (df['Family'] == 1) | ((df['Family'] >= 5) & (df['Family'] <= 7)), # 1:単独 or やや多め
    df['Family'] >= 8 # 0:超大規模
], [2, 1, 0])
```

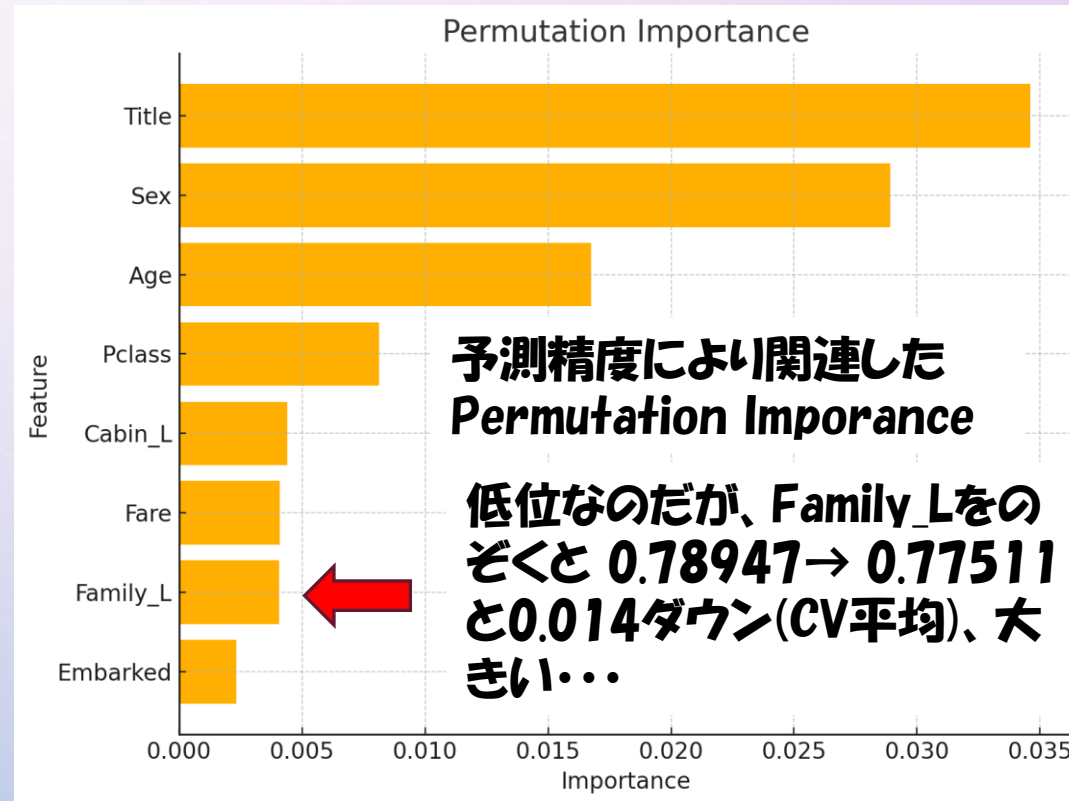
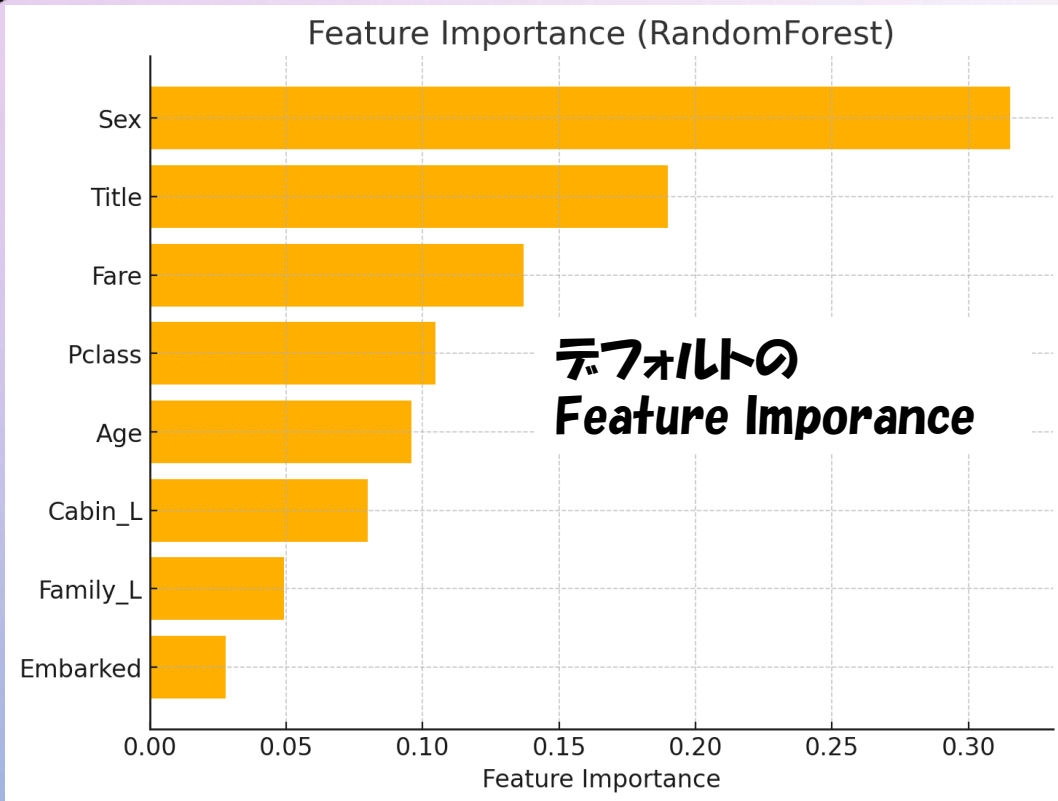




Feature Importance

feature importanceを横棒グラフで表示して

LGBMでは
左がsplit(デフォルト)、
右がgain。



① Mean Decrease in Impurity (MDI)

(決定木の分割で減少したジニ不純度やエントロピーの合計)

- scikit-learn の `feature_importances_` 属性がこれ
- 各特徴量が木の分割にどれだけ貢献したかを測る
- 計算が高速
- バイアスがある (カテゴリ特徴量や多くの分割候補がある特徴量に有利になる傾向)

② Permutation Importance

(特徴量の値をシャッフルして予測精度の低下を測定)

- 特徴量の値をシャッフル → 精度低下を測定 → 重要度を算出
- モデル訓練後に外部から計算可能
- 特徴量のスケールや型に依存しにくい
- ノイズに強く、実用上こちらを採用する論文やKaggleも多い
- 計算コストは高め

後処理(ルールベース)導入

08_rf_LabelE_200-6_skfold5_FeAll_0.80382.py

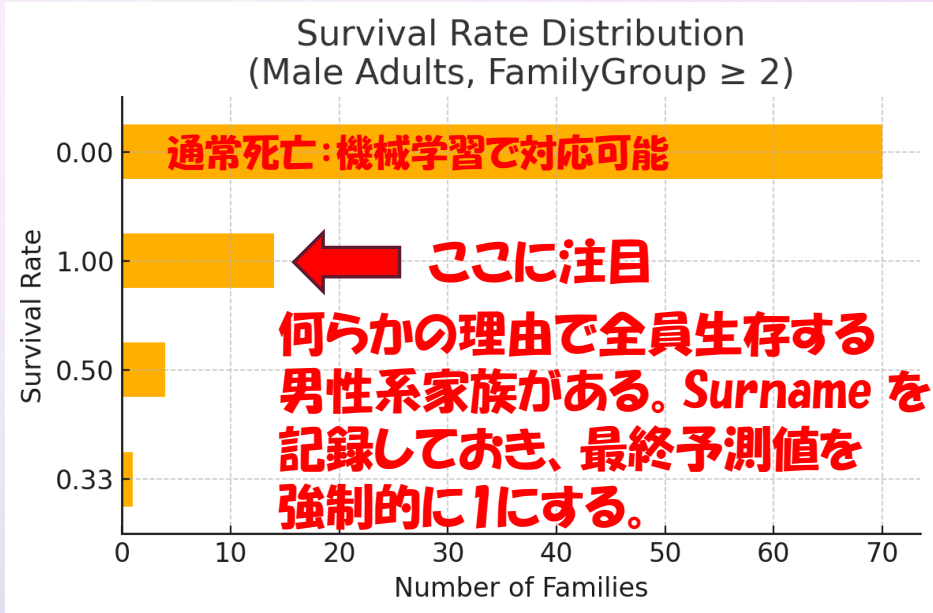
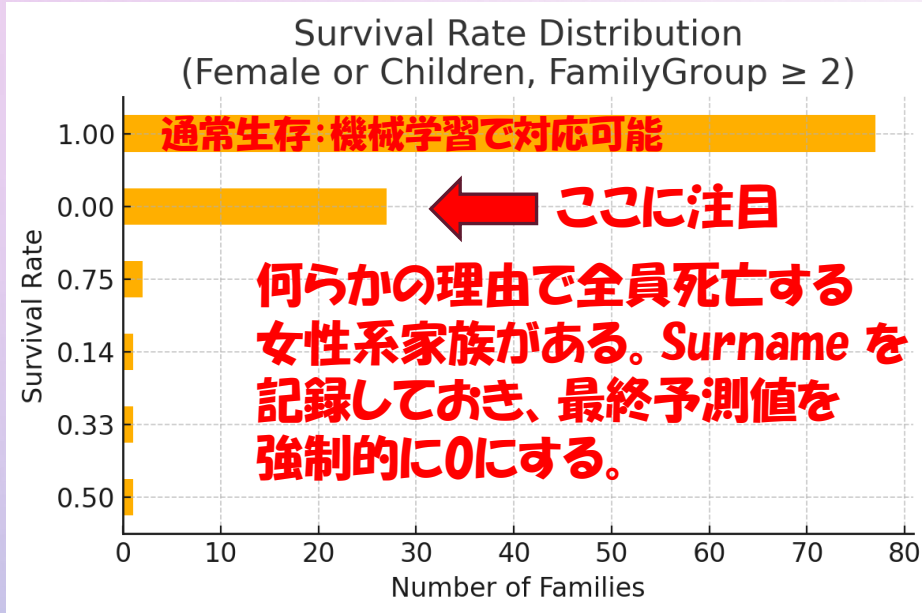
0.80越え!

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import StratifiedKFold
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.metrics import accuracy_score
7 from scipy.stats import mode
8
9 # データ読み込みと結合
10 train_df = pd.read_csv("train.csv")
11 test_df = pd.read_csv("test.csv")
12 test_df['Survived'] = np.nan # 後で分離できるように np.nan(null) を入れておく
13 df = pd.concat([train_df, test_df], ignore_index=True, sort=False)
14
15 # Age, Fare補完 : base
16 df['Age'] = df['Age'].fillna(df['Age'].median())
17 df['Fare'] = df['Fare'].fillna(df['Fare'].median())
18
19 ##### 特徴量エンジニアリング : 全てまとめて 0.01 の精度向上
20 # Title作成 : 0.007 の精度向上, ['Master', 'Miss', 'Mr', 'Mrs', 'Officer', 'Royalty']に集約
21 df['Title'] = df['Name'].map(lambda x: x.split(',')[1].split('.')[0])
22 df['Title'].replace(['Capt', 'Col', 'Major', 'Dr', 'Rev'], 'Officer', inplace=True)
23 df['Title'].replace(['Don', 'Sir', 'the Countess', 'Lady', 'Dona'], 'Royalty', inplace=True)
24 df['Title'].replace(['Mme', 'Ms'], 'Mrs', inplace=True)
25 df['Title'].replace(['Mlle'], 'Miss', inplace=True)
26 df['Title'].replace(['Jonkheer'], 'Master', inplace=True)
27 # Cabin補完 : 0.005 の精度向上
28 df['Cabin'] = df['Cabin'].fillna('Unknown')
29 df['Cabin_L'] = df['Cabin'].str.get(0) # 最初の文字をCabinラベルとして使用
30 # Embarked補完 : 0.005 の精度向上
31 df['Embarked'] = df['Embarked'].fillna('S') # 'S'は最貧度
32 # Familyラベル : シンプルFamilySizeより0.015、なしより0.017 の精度向上
33 df['Family'] = df['SibSp'] + df['Parch'] + 1 # 家族人数を計算 (本人を含める)
34 df['Family_L'] = np.select([
35     (df['Family'] >= 2) & (df['Family'] <= 4), # 2:中規模
36     (df['Family'] == 1) | ((df['Family'] >= 5) & (df['Family'] <= 7)), # 1:単独 or やや多め
37     df['Family'] >= 8 # 0:超大規模
38 ], [2, 1, 0])
39
40 # 特徴量とラベルエンコーディング
41 features = ['Pclass', 'Sex', 'Age', 'Fare', 'Embarked', 'Title', 'Family_L', 'Cabin_L']
42 dfm = df[['Survived'] + features].copy()
43 for col in ['Sex', 'Embarked', 'Title', 'Cabin_L']:
44     le = LabelEncoder()
45     dfm[col] = le.fit_transform(dfm[col].astype(str))
46
```

```
47 # train/test 分割
48 train = dfm[dfm['Survived'].notnull()]
49 test = dfm[dfm['Survived'].isnull()].drop('Survived', axis=1)
50 X = train.drop('Survived', axis=1).values
51 y = train['Survived'].astype(int).values
52 X_test = test.values
53
54 # StratifiedKFold + Mode voting(多数決)、CVスコア表示
55 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
56 test_preds = []
57 val_scores = []
58 for fold, (train_idx, val_idx) in enumerate(skf.split(X, y), 1):
59     model = RandomForestClassifier(n_estimators=200, max_depth=6, random_state=42)
60     model.fit(X[train_idx], y[train_idx])
61     val_pred = model.predict(X[val_idx])
62     val_acc = accuracy_score(y[val_idx], val_pred)
63     val_scores.append(val_acc)
64     print(f"Fold {fold} CV Accuracy: {val_acc:.5f}")
65     test_preds.append(model.predict(X_test))
66 print(f"Mean CV Accuracy: {np.mean(val_scores):.5f}")
67
68 # 多数決で最終予測を決定
69 test_preds = np.array(test_preds)
70 final_preds = mode(test_preds, axis=0, keepdims=False).mode
71
72 # 特殊家族補正
73 # 通常予測とはことなる特殊な家族の補正 (最終予測用) : 0.014 の精度向上
74 surnames = df['Name'].str.split(',').str[0].str.strip() # 名前から姓 (Surname) を取り出す
75 family_n = surnames.map(surnames.value_counts()) # 同じ姓の人数をカウント (家族単位)
76 # 「2人以上の家族」かつ「16歳以下または女性」で構成(通常生存)され「全員死亡」している特異な家族リスト
77 Dead_list = set(surnames[ (family_n>=2)&((df['Age']<=16)|(df['Sex']=='female'))&(df['Survived']==0) ])
78 # 「2人以上の家族」かつ「男性の大人」で構成(通常死亡)され「全員生存」している特異な家族リスト
79 Surv_list = set(surnames[ (family_n>=2)&((df['Age']> 16)&(df['Sex']=='male' ))&(df['Survived']==1) ])
80 # 予測データに同じ姓があれば、そのパターンに従って生死を補正する
81 for i, s in enumerate(df.loc[df['Survived'].isnull(), 'Name'].str.split(',').str[0].str.strip()):
82     final_preds[i] = 0 if s in Dead_list else 1 if s in Surv_list else final_preds[i]
83
84 # サブミッション作成
85 submission = pd.DataFrame({
86     'PassengerId': test_df['PassengerId'],
87     'Survived': final_preds.astype(int)
88 })
89 submission.to_csv("08_rf_LabelE_200-6_skfold5_FeAll.csv", index=False)
```

後処理追加

後処理(ルールベース)説明



将来は、こういう特徴量を自動発見して学習する機械学習アルゴリズムが生まれるでしょう。

```
72 # 特殊家族補正
73 # 通常予測とはことなる特殊な家族の補正（最終予測用）： 0.014 の精度向上
74 surnames = df['Name'].str.split(',').str[0].str.strip() # 名前から姓（Surname）を取り出す
75 family_n = surnames.map(surnames.value_counts()) # 同じ姓の人数をカウント（家族単位）
76 # 「2人以上の家族」かつ「16歳以下または女性」で構成(通常生存)され「全員死亡」している特異な家族リスト
77 Dead_list = set(surnames[ (family_n>=2)&((df['Age']<=16)|(df['Sex']=='female'))&(df['Survived']==0) ])
78 # 「2人以上の家族」かつ「男性の大人」で構成(通常死亡)され「全員生存」している特異な家族リスト
79 Surv_list = set(surnames[ (family_n>=2)&((df['Age']> 16)&(df['Sex']=='male' ))&(df['Survived']==1) ])
80 # 予測データに同じ姓があれば、そのパターンに従って生死を補正する
81 for i, s in enumerate(df.loc[df['Survived'].isnull(), 'Name'].str.split(',').str[0].str.strip()):
82     final_preds[i] = 0 if s in Dead_list else 1 if s in Surv_list else final_preds[i]
83
```

さらなるスコア向上案

- 特徴量エンジニアリングのさらなる探求
- LGBMでtuning
- アンサンブル/スタッキング
 - たとえば下記アルゴリズムはscikit learnに入っているので、ChatGPT内でアンサンブル実行可能。
 - LogisticRegression
 - Random Forests
 - Gradient Boosting Trees
 - Kaggle code実行などでより強いアルゴリズムでアンサンブル
 - LGBM
 - CatBoost
 - XGBoost, etc

やってみたが...

アンサンブルでスコアは上がらなかった:0.79186 (汎化性は高いかも)
各アルゴリズム毎に、特徴量前処理、ハイパーパラメータの
チューニングが必要かも知れない。


まとめ

**Kaggle面白いですよ！
ChatGPT(生成AI)
駆使して、楽しみましょう！**


- まだスコア向上策はあいそう
- 最後のルールベース追加は、様々な議論あるだろう
 - **Kaggleとしては合法ルール**
 - 人間が特徴量エンジニアリングで気づくべき派
 - ひらめき・特徴量エンジニアリングで気づいた者が勝利
 - ルールベース追加や、この予測が入るように前処理するなど
 - 機械学習アルゴリズムが自動導入すべき派
 - 人間に頼らず、強い機械学習アルゴリズムが自動で導入すべきという意見もある。
 - よいより機械学習アルゴリズムの探索・開発へ
 - なのでTitanicのような古いコンペでも新しいアルゴリズムが試されている
- ChatGPTは万能ではない。しかしAIペアプログラミング相手として優秀！
 - 細かいSyntax Error・凡ミスは、発生しても自動修正してくれる。
 - **人間は本質的なデータサイエンス思考を止めなくて良い。細かいプログラミングに煩わされない。**
 - 今のChatGPTがAIエージェント的にどこまで使えるか、楽しみながら確認してみてください！

本資料、各プログラムは 菜園の柱 ブログに置きます

- <https://garden-pillars.blog/>

**菜園の柱** 趣味のCGとChiaのブログ

ホーム | 庭子アバター 配布所 | 庭子アバター制作まとめ | Chia関連まとめ | 登場人物紹介



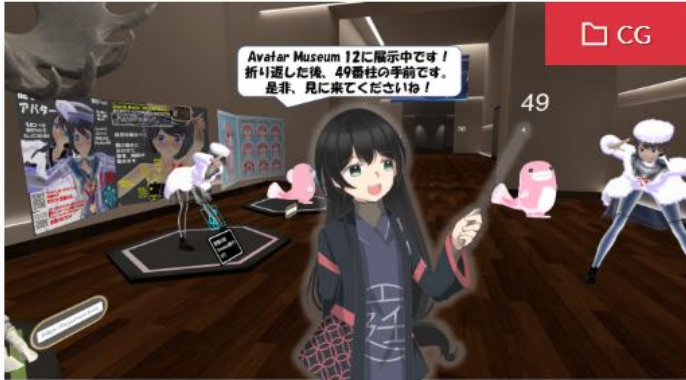
Kaggle Titanic 入門 01
Kaggle を始めてみよう！
庭子

庭子アバター v0150
Avatar Museum 12に展示中

📅 2025.05.30

Kaggle Titanic 入門01

データサイエンティスト集会 in VRC でChatGPT使いまくって、初心者がKaggle TitanicコンペやるためのLTを講演しました。その時の発表資料を置いておきます。



Avatar Museum 12に展示中です！
折り返した後、49番柱の手前です。
是非、見に来てくださいね！

49

CG

📅 2025.05.10

Avatar Museum 12 に展示中！

庭子 v0150（Booth）をAvatar Museum 12に展示中です。一番奥まで行って、折り返した後、49番柱の手前です。主なupdateは下記：是非、見に来てくださいね！