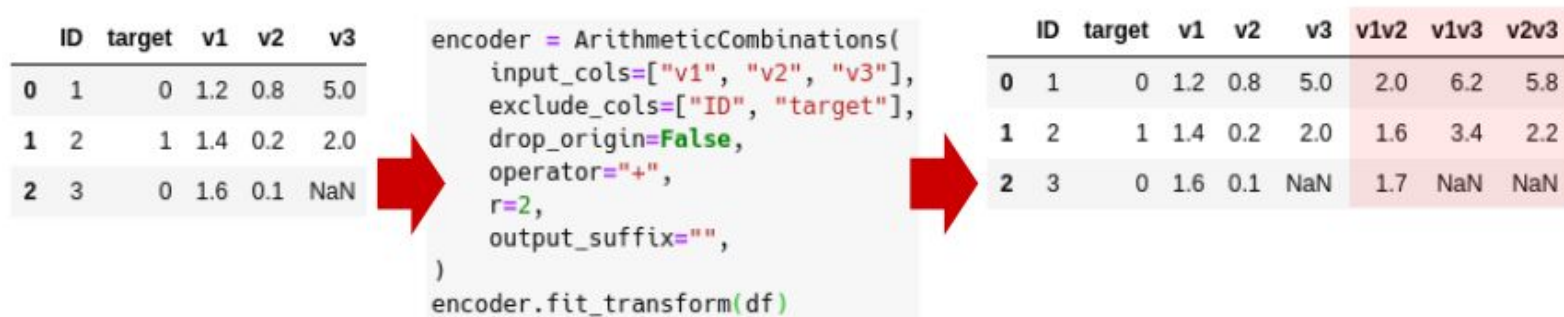# xfeat: Feature Engineering and Exploration Library

Kohei Ozaki at Preferred Networks

# Main Idea: DataFrame-IN, DataFrame-OUT

- xfeat provides feature encoders and selectors.
- DataFrame as Input. DataFrame as Output.
- Support both **pandas** and **cuDF** dataframe.

|   | ID | target | v1 | v2 | v3 |
|---|----|--------|-----|-----|-----|
| 0 | 1  | 0      | 1.2 | 0.8 | 5.0 |
| 1 | 2  | 1      | 1.4 | 0.2 | 2.0 |
| 2 | 3  | 0      | 1.6 | 0.1 | NaN |

```python
encoder = ArithmeticCombinations(
    input_cols=["v1", "v2", "v3"],
    exclude_cols=["ID", "target"],
    drop_origin=False,
    operator="+",
    r=2,
    output_suffix="",
)
encoder.fit_transform(df)
```

|   | ID | target | v1 | v2 | v3 | v1v2 | v1v3 | v2v3 |
|---|----|--------|-----|-----|-----|------|------|------|
| 0 | 1  | 0      | 1.2 | 0.8 | 5.0 | 2.0  | 6.2  | 5.8  |
| 1 | 2  | 1      | 1.4 | 0.2 | 2.0 | 1.6  | 3.4  | 2.2  |
| 2 | 3  | 0      | 1.6 | 0.1 | NaN | 1.7  | NaN  | NaN  |

# Main features: Encoders and Selectors

- Numerical:
  - SelectNumerical
  - ArithmeticCombinations
- Categorical:
  - SelectCategorical
  - LabelEncoder
  - Target Encoder
  - ConcatCombination
  - CountEncoder
  - UserDefinedLabelEncoder
- Feature Selection (Selectors): GBDTFeatureSelector
- Feature Exploration (Selectors): GBDTFeatureExploration

Generate features by arithmetic combinations of numerical columns.

v1 + v2

| | ID | target | v1 | v2 | v3 | v1v2 | v1v3 | v2v3 |
|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1.2 | 0.8 | 5.0 | 2.0 | 6.2 | 5.8 |
| 1 | 2 | 1 | 1.4 | 0.2 | 2.0 | 1.6 | 3.4 | 2.2 |
| 2 | 3 | 0 | 1.6 | 0.1 | NaN | 1.7 | NaN | NaN |

Generate features by concatenating string values of categoricals.

Filter-based method using LightGBM. Its hyperparameter is tuned by Optuna.

# Main features: Pipeline

Combine multiple encoders into a single encoder object.

```python
encoder = Pipeline([
    SelectCategorical(exclude_cols=["id", "user_id"]),

    # If there are many categorical columns,
    # users can specify the columns to be combined with `input_cols` kwargs.
    # `r=2` specifies the number of columns to combine the columns.
    ConcatCombination(drop_origin=True, output_suffix="", r=2),

    LabelEncoder(output_suffix=""),
])
encoder.fit_transform(df).head()
```

# Main features: Serialize/Deserialize

Serialize/Deserialize Encoders with pickle.

This makes it easier to separate training step and inference step.

1. Fit parameters on a train set.

```python
encoder = Pipeline([
    SelectCategorical(exclude_cols=["id", "user_id"]),
    LabelEncoder(output_suffix=""),
])
df_train_encoded = encoder.fit_transform(df_train)

with open("label_encoder.pkl", "wb") as f:
    pickle.dump(encoder, f)
```
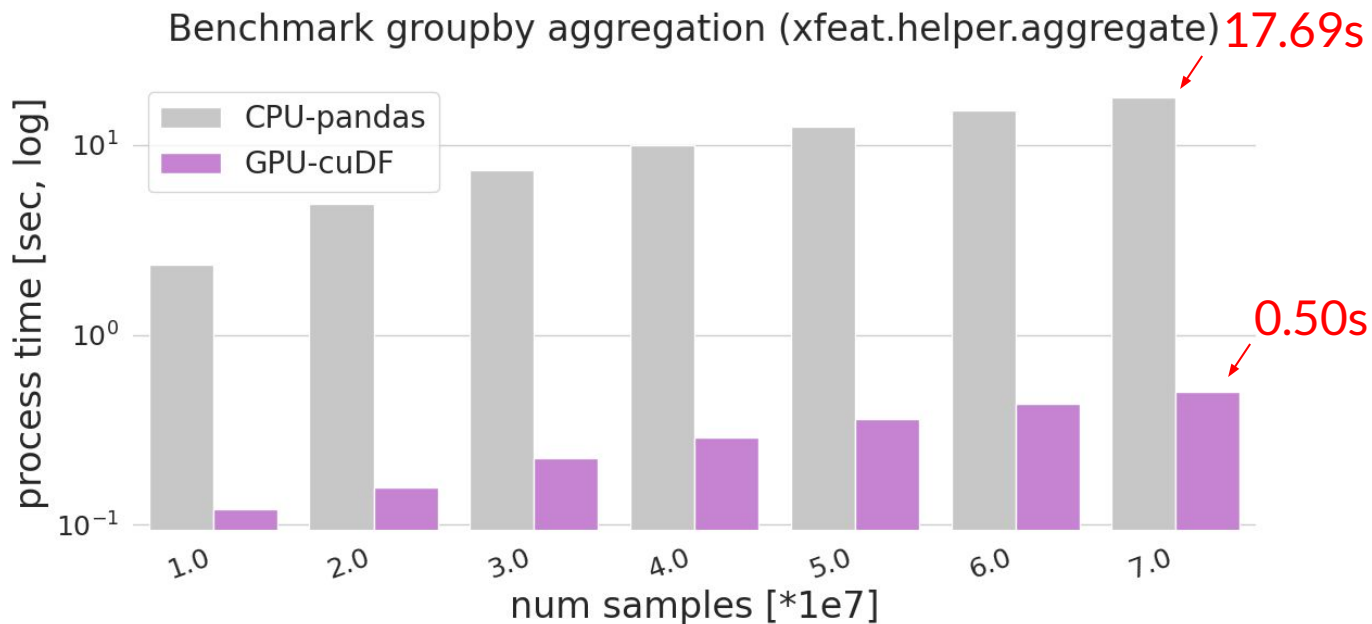
2. Transform on a test set later.

```python
with open("label_encoder.pkl", "rb") as f:
    encoder = pickle.load(f)

encoder.transform(df_test).head()
```
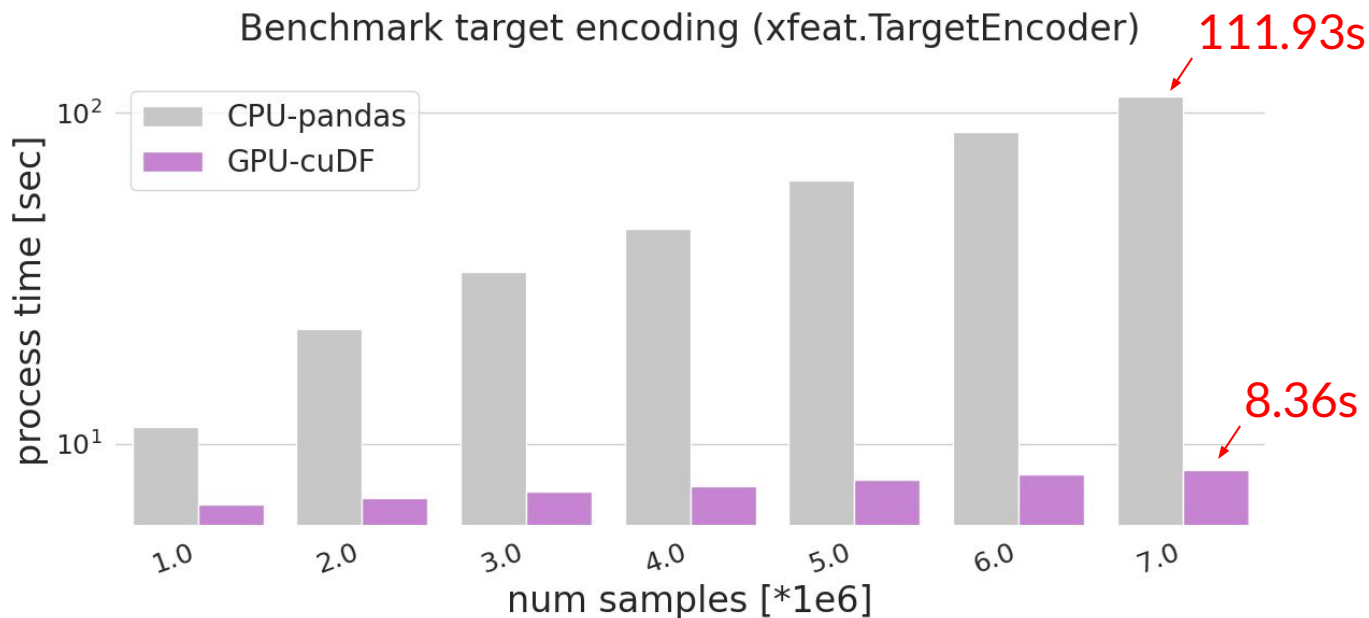
# Benchmark: Group-by aggregation

Encoders can be greatly accelerated with cuDF and CuPy.



Benchmark groupby aggregation (xfeat.helper.aggregate)

17.69s

0.50s

# Benchmark: Target Encoding

In Target Encoding, aggregation is computed on the GPUs using CuPy.



Benchmark target encoding (xfeat.TargetEncoder)

# Practical result on Kaggle competition

With using xfeat, I got 12th place LB score with less than 200 lines.
I plan to upload the solution code on kaggle.com after releasing xfeat.

The initial OSS release is planned on June.

# Appendix: solution code (1 of 3)

```python
52  def feature_engineering():
53      # (1) Save numerical features
54      SelectNumerical().fit_transform(pd.read_feather("train_test.ftr")).reset_index(
55          drop=True
56      ).to_feather("feature_num_features.ftr")
57
58      # (2) Categorical encoding using label encoding: 13 features
59      Pipeline([SelectCategorical(), LabelEncoder(output_suffix="")]).fit_transform(
60          pd.read_feather("train_test.ftr")
61      ).reset_index(drop=True).to_feather("feature_1way_label_encoding.ftr")
```

# Appendix: solution code (2 of 3)

```
76    # (4) 3-order combination of categorical features
77    # Use `include_cols=` kwargs to reduce the total count of combinations.
78    # 66 features (12 * 11 / 2 = 66)
79    Pipeline(
80        [
81            SelectCategorical(),
82            ConcatCombination(drop_origin=True, include_cols=["v22"], r=3),
83            LabelEncoder(output_suffix=""),
84        ]
85    ).fit_transform(pd.read_feather("train_test.ftr")).reset_index(
86        drop=True
87    ).to_feather(
88        "feature_3way_including_v22_label_encoding.ftr"
89    )
```

# Appendix: solution code (3 of 3)

```
112      # (6) 2-order Arithmetic combinations.
113      Pipeline(
114          [
115              SelectNumerical(),
116              ArithmeticCombinations(
117                  exclude_cols=["target"], drop_origin=True, operator="+", r=2,
118              ),
119          ]
120      ).fit_transform(pd.read_feather("train_test.ftr")).reset_index(
121          drop=True
122      ).to_feather(
123          "feature_arithmetic_combi2.ftr"
124      )
```