



P16 Day 2025

Skrub - Machine learning with dataframes

Riccardo Cappuzzo
Inria P16

2025-10-14



Boost your productivity with skrub!

Skrub simplifies many tedious data preparation operations



Skrub compatibility

- Skrub is fully compatible with pandas and polars
- Skrub transformers are fully compatible with scikit-learn



An example pipeline

1. Gather some data
2. Explore the data
3. Preprocess the data
4. Perform feature engineering
5. Build a scikit-learn pipeline
6. ???
7. Profit?





Exploring the data with `skrub`

```
1 from skrub import TableReport  
2 TableReport(employee_salaries)
```

TableReport Preview

Main features:

- Obtain high-level statistics about the data
- Explore the distribution of values and find outliers
- Discover highly correlated columns
- Export and share the report as an `html` file

More examples

Data cleaning with pandas/polars: setup

Pandas

Polars

```

1 import pandas as pd
2 import numpy as np
3
4 data = {
5     "Int": [2, 3, 2],    # Multiple unique values
6     "Const str": ["x", "x", "x"],    # Single unique value
7     "Str": ["foo", "bar", "baz"],    # Multiple unique values
8     "All nan": [np.nan, np.nan, np.nan],    # All missing values
9     "All empty": ["", "", ""],    # All empty strings
10    "Date": ["01 Jan 2023", "02 Jan 2023", "03 Jan 2023"],
11 }
12
13 df_pd = pd.DataFrame(data)
14 display(df_pd)

```

	Int	Const str	Str	All nan	All empty	Date
0	2	x	foo	NaN		01 Jan 2023
1	3	x	bar	NaN		02 Jan 2023
2	2	x	baz	NaN		03 Jan 2023

Nulls, datetimes, constant columns with pandas/polars

Pandas Polars

```
1 # Parse the datetime strings with a specific format
2 df_pd['Date'] = pd.to_datetime(df_pd['Date'], format='%d %b %Y')
3
4 # Drop columns with only a single unique value
5 df_pd_cleaned = df_pd.loc[:, df_pd.nunique(dropna=True) > 1]
6
7 # Function to drop columns with only missing values or empty strings
8 def drop_empty_columns(df):
9     # Drop columns with only missing values
10    df_cleaned = df.dropna(axis=1, how='all')
11    # Drop columns with only empty strings
12    empty_string_cols = df_cleaned.columns[df_cleaned.eq('').all()]
13    df_cleaned = df_cleaned.drop(columns=empty_string_cols)
14    return df_cleaned
15
16 # Apply the function to the DataFrame
17 df_pd_cleaned = drop_empty_columns(df_pd_cleaned)
```



Data cleaning with `skrub.Cleaner`

Pandas

Polars

```
1 from skrub import Cleaner
2 cleaner = Cleaner(drop_if_constant=True, datetime_format='%d %b %Y')
3 df_cleaned = cleaner.fit_transform(df_pd)
4 display(df_cleaned)
```

	Int	Str	Date
0	2	foo	2023-01-01
1	3	bar	2023-01-02
2	2	baz	2023-01-03



Encoding datetime features with skrub.DatetimeEncoder

```

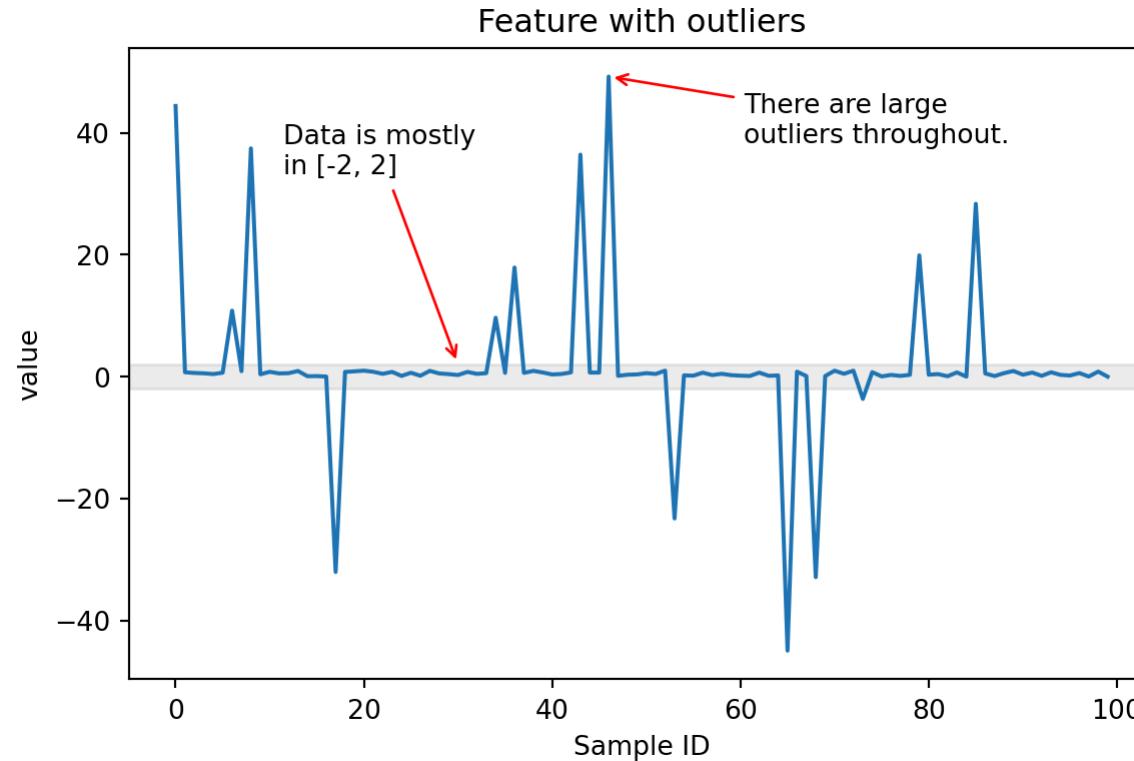
1 from skrub import DatetimeEncoder, ToDatetime
2
3 X_date = ToDatetime().fit_transform(df["date"])
4 de = DatetimeEncoder(resolution="second")
5 # de = DatetimeEncoder(periodic_encoding="spline")
6 X_enc = de.fit_transform(X_date)
7 print(X_enc)

```

shape: (3, 7)

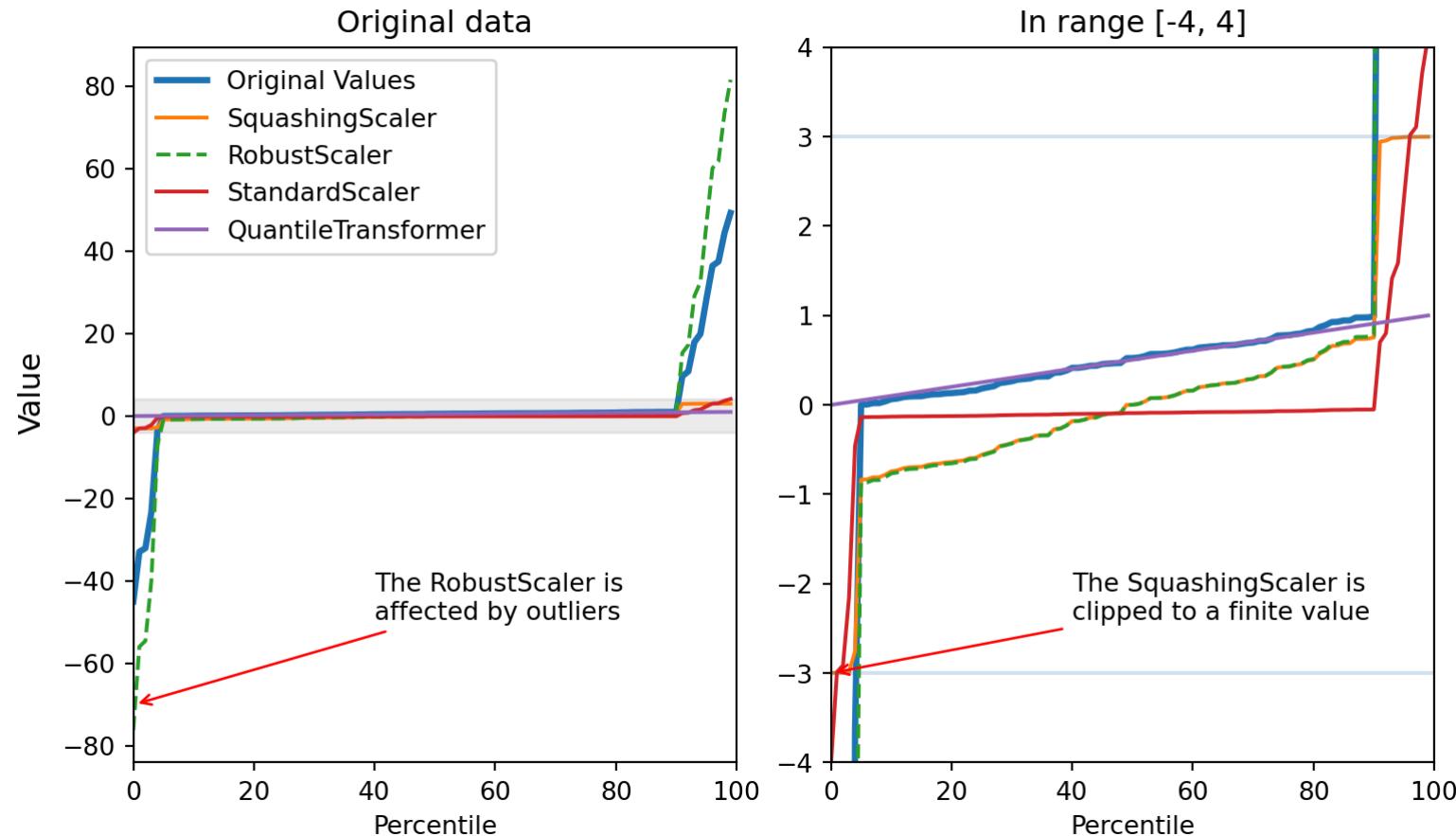
date_year	date_month	date_day	date_hour	date_minute	date_second	date_total_seconds
---	---	---	---	---	---	---
f32	f32	f32	f32	f32	f32	f32
2023.0	1.0	1.0	12.0	34.0	56.0	1.6726e9
2023.0	2.0	15.0	8.0	45.0	23.0	1.6765e9
2023.0	3.0	20.0	18.0	12.0	45.0	1.6793e9

Encoding numerical features with skrub.SquashingScaler



Encoding numerical features with skrub.SquashingScaler

Comparison of different scalers on sorted data with outliers



Encoding categorical (string/text) features

Categorical features have a “**cardinality**”: the number of unique values

- Low cardinality: `OneHotEncoder`
- High cardinality (>40 unique values):
`skrub.StringEncoder`
- Text: `skrub.TextEncoder` and pretrained models from HuggingFace Hub

Encoding *all the features*:

TableVectorizer

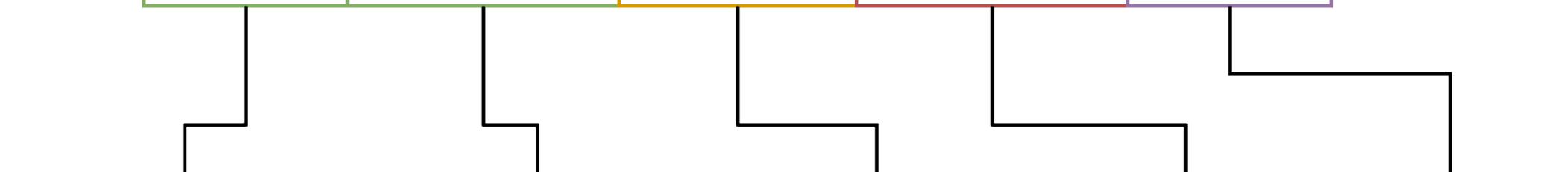
```
1 from skrub import TableVectorizer  
2  
3 table_vec = TableVectorizer()  
4 df_encoded = table_vec.fit_transform(df)
```

- Apply the `Cleaner` to all columns
- Split columns by dtype and # of unique values
- Encode each column separately

Encoding *all the features*:

TableVectorizer

Low cardinality categorical		Datetime	High cardinality cat.	Number
Gender	Contract type	Hiring date	Position	Salary
M	Fulltime	03/02/2015	Bus operator	55000
...
F	Part-time	11/11/2001	Firefighter	85000



TableVectorizer

OneHotEncoding (Gender)	OneHotEncoding (Contract type)	DatetimeEncoder (Hiring date)	StringEncoder (Position)	Passthrough (Salary)

Build a predictive pipeline

```
1 from sklearn.linear_model import Ridge  
2 model = Ridge()
```



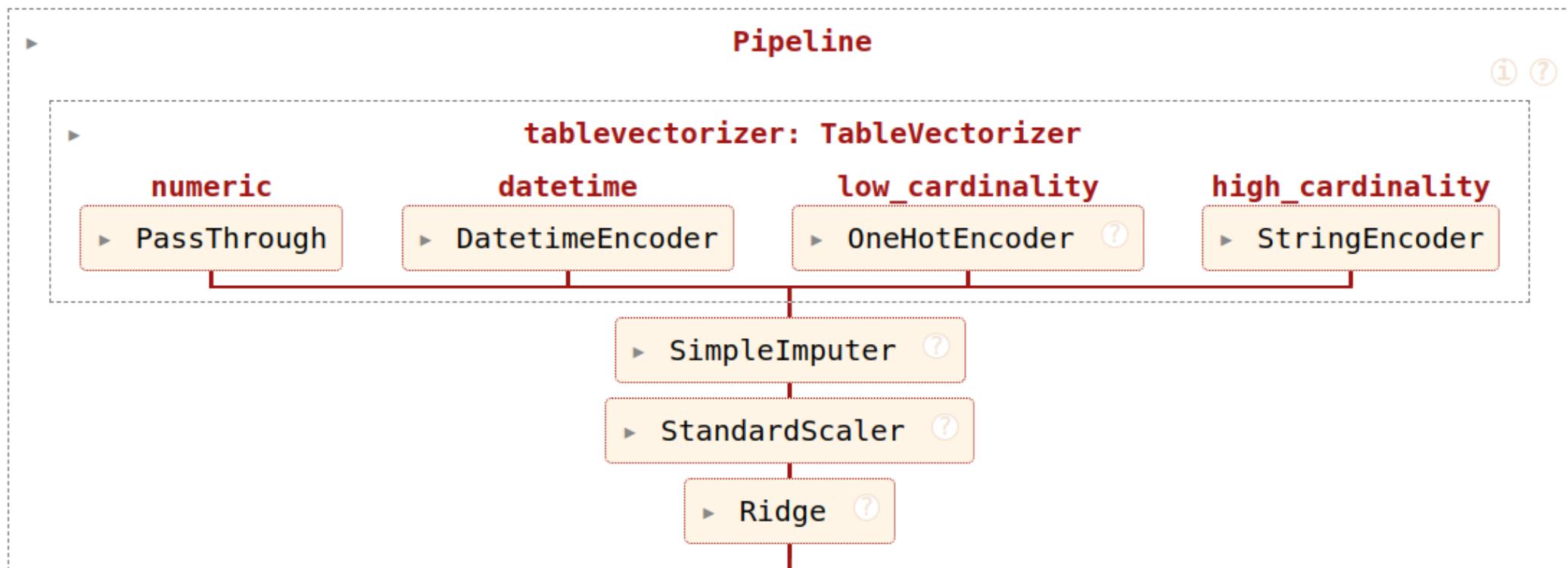
Build a predictive pipeline

```
1 from sklearn.linear_model import Ridge
2 from sklearn.pipeline import make_pipeline
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4 from sklearn.impute import SimpleImputer
5 from sklearn.compose import make_column_selector as selector
6 from sklearn.compose import make_column_transformer
7
8 categorical_columns = selector(dtype_include=object)(employees)
9 numerical_columns = selector(dtype_exclude=object)(employees)
10
11 ct = make_column_transformer(
12     (StandardScaler(),
13      numerical_columns),
14     (OneHotEncoder(handle_unknown="ignore"),
15      categorical_columns))
16
17 model = make_pipeline(ct, SimpleImputer(), Ridge())
```



Build a predictive pipeline with tabular_pipeline

```
1 import skrub
2 from sklearn.linear_model import Ridge
3 model = skrub.tabular_pipeline(Ridge())
```





```

1 from sklearn.linear_model import Ridge
2 from sklearn.pipeline import make_pipeline
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4 from sklearn.impute import SimpleImputer
5 from sklearn.compose import make_column_selector as selector
6 from sklearn.compose import make_column_transformer
7
8 categorical_columns = selector(dtype_include=object)(employees)
9 numerical_columns = selector(dtype_exclude=object)(employees)
10
11 ct = make_column_transformer(
12     (StandardScaler(), numerical_columns),
13     (OneHotEncoder(handle_unknown="ignore"),
14      categorical_columns))
15
16 model = make_pipeline(ct, SimpleImputer(), Ridge())
17

```



```

1 import skrub
2 from sklearn.linear_model import Ridge
3 tl = skrub.tabular_learner(Ridge())

```

We now have a pipeline!

1. Gather some data
2. Explore the data
 - `TableReport`
3. Pre-process the data
 - `Cleaner`, `ToDatetime` ...
4. Perform feature engineering
 - `TableVectorizer`, `SquashingScaler`, `TextEncoder`, `StringEncoder`...
5. Build a scikit-learn pipeline
 - `tabular_pipeline`
6. ???
7. Profit 

What if this is not enough??

<https://skrub-data.org/skrub-materials/>



What if...

- Your data is spread over multiple tables?
- You want to avoid data leakage?
- You want to tune more than just the hyperparameters of your model?
- You want to guarantee that your pipeline is replayed exactly on new data?

When a normal pipeline is not enough...

... the **skrub** DataOps come to the rescue 

DataOps...

- Extend the `scikit-learn` machinery to complex multi-table operations, and take care of data leakage
- Track all operations with a computational graph (a *Data Ops plan*)
- Are transparent and give direct access to the underlying object
- Allow tuning any operation in the Data Ops plan
- Guarantee that all operations are reproducible
- Can be persisted and shared easily



How do DataOps work, though?

DataOps wrap around *user operations*, where user operations are:

- any dataframe operation (e.g., merge, group by, aggregate etc.)
- scikit-learn estimators (a Random Forest, RidgeCV etc.)
- custom user code (load data from a path, fetch from an URL etc.)

 **Important**

DataOps record user operations, so that they can later be *replayed* in the same order and with the same arguments on unseen data.

Starting with the DataOps

```
1 import skrub
2 data = skrub.datasets.fetch_credit_fraud()
3
4 baskets = skrub.var("baskets", data.baskets)
5 products = skrub.var("products", data.products) # add a new variable
6
7 X = baskets[["ID"]].skb.mark_as_X()
8 y = baskets["fraud_flag"].skb.mark_as_y()
```

- `baskets` and `products` represent inputs to the pipeline.
- Skrub tracks `X` and `y` so that training and test splits are never mixed.



Applying a transformer

```
1 from skrub import selectors as s
2
3 vectorizer = skrub.TableVectorizer(
4     high_cardinality=skrub.StringEncoder()
5 )
6 vectorized_products = products(skb.apply(
7     vectorizer, cols=s.all() - "basket_ID"
8 )
```

Executing dataframe operations

```
1 aggregated_products = vectorized_products.groupby(  
2     "basket_ID"  
3 ).agg("mean").reset_index()  
4  
5 features = X.merge(  
6     aggregated_products, left_on="ID", right_on="basket_ID"  
7 )  
8 features = features.drop(columns=["ID", "basket_ID"])
```



Applying a ML model

```
1 from sklearn.ensemble import ExtraTreesClassifier  
2 predictions = features(skb.apply(  
3     ExtraTreesClassifier(n_jobs=-1), y=y  
4 ))
```

Inspecting the Data Ops plan

```
1 predictions(skb).full_report()
```

Execution report

Each node:

- Shows a preview of the data resulting from the operation
- Reports the location in the code where the code is defined
- Shows the run time of the node



Exporting the plan in a learner

The **Learner** is a stand-alone object that works like a scikit-learn estimator that takes a dictionary as input rather than just **X** and **y**.

```
1 learner = predictions(skb.make_learner(fitted=True))
```

Then, the **learner** can be pickled ...

```
1 import pickle
2
3 with open("learner.bin", "wb") as fp:
4     pickle.dump(learner, fp)
```

... loaded and applied to new data:

```
1 with open("learner.bin", "rb") as fp:
2     loaded_learner = pickle.load(fp)
3 data = skrub.datasets.fetch_credit_fraud(split="test")
4 new_baskets = data.baskets
5 new_products = data.products
6 loaded_learner.predict({"baskets": new_baskets, "products": new_products})
```

array([0, 0, 0, ..., 0, 0, 0], shape=(31549,))

Hyperparameter tuning in a Data Ops plan

Skrub implements four `choose_*` functions:

- `choose_from`: select from the given list of options
- `choose_int`: select an integer within a range
- `choose_float`: select a float within a range
- `choose_bool`: select a bool
- `optional`: chooses whether to execute the given operation

Tuning in scikit-learn can be complex

```
1 pipe = Pipeline([("dim_reduction", PCA()), ("regressor", Ridge())])
2 grid = [
3     {
4         "dim_reduction": [PCA()],
5         "dim_reduction_n_components": [10, 20, 30],
6         "regressor": [Ridge()],
7         "regressor_alpha": loguniform(0.1, 10.0),
8     },
9     {
10        "dim_reduction": [SelectKBest()],
11        "dim_reduction_k": [10, 20, 30],
12        "regressor": [Ridge()],
13        "regressor_alpha": loguniform(0.1, 10.0),
14    },
15    {
16        "dim_reduction": [PCA()],
17        "dim_reduction_n_components": [10, 20, 30],
18        "regressor": [RandomForestRegressor()],
19        "regressor_n_estimators": loguniform(20, 200),
20    },
21    {
22        "dim_reduction": [SelectKBest()],
23        "dim_reduction_k": [10, 20, 30],
24        "regressor": [RandomForestRegressor()],
25        "regressor_n_estimators": loguniform(20, 200),
```



Tuning with Data Ops is simple!

```
1 dim_reduction = X(skb.apply(  
2     skrub.choose_from(  
3     {  
4         "PCA": PCA(n_components=skrub.choose_int(10, 30)),  
5         "SelectKBest": SelectKBest(k=skrub.choose_int(10, 30))  
6     }, name="dim_reduction"  
7     )  
8 )  
9 regressor = dim_reduction(skb.apply(  
10    skrub.choose_from(  
11    {  
12        "Ridge": Ridge(alpha=skrub.choose_float(0.1, 10.0, log=True)),  
13        "RandomForest": RandomForestRegressor(  
14            n_estimators=skrub.choose_int(20, 200, log=True)  
15        )  
16    }, name="regressor"  
17 )  
18 )
```

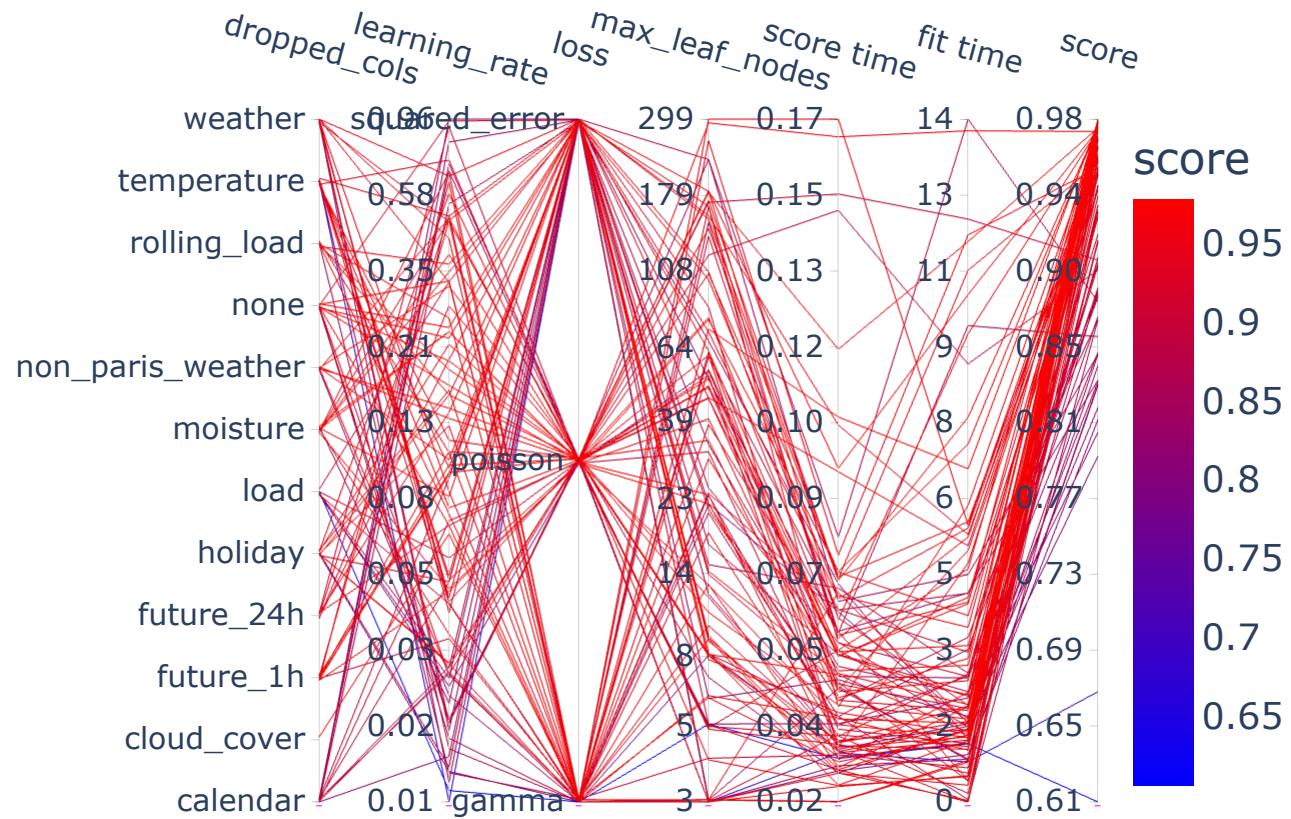


Run hyperparameter search

```
1 # fit the search
2 search = regressor(skb.make_randomized_search(
3     scoring="roc_auc", fitted=True, cv=5
4 )
5
6 # save the best learner
7 best_learner = search.best_learner_
```

A parallel coordinate plot to explore hyperparameters

```
1 search = pred(skb.get_randomized_search(fitted=True)
2 search.plot_parallel_coord()
```



More information about the Data Ops

- Skrub [example gallery](#)
- Skrub [user guide](#)
- [Tutorial](#) on timeseries forecasting at Euroscipy 2025
- [Kaggle notebook](#) on the Titanic survival challenge



Wrapping up

<https://skrub-data.org/skrub-materials/>







<https://skrub-data.org/skrub-materials/>



Getting involved

Do you want to learn more?

- [Skrub website](#)
- [Skrub materials website](#)
- [Discord server](#)

Follow skrub on:

- [Bluesky](#)
- [LinkedIn](#)

Star skrub on GitHub, or contribute directly:

- [GitHub repository](#)



tl;dw: skrub

- interactive data exploration: `TableReport`
- automated pre-processing of pandas and polars dataframes: `Cleaner`
- powerful feature engineering: `TableVectorizer`, `tabular_pipeline`
- column- and dataframe-level operations: `ApplyToCols`, `selectors`
- DataOps, plans, hyperparameter tuning, (almost) no leakage

