



PyData Global



# Introduction to Machine Learning Pipelines

How to Prevent Data Leakage and Build Efficient Workflows



# Who am I ?



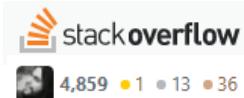
Cainã Max Couto-Silva



Data Scientist, PhD  
@ Schlumberger



cmcoutho-silva



## Hobbies

BSc. Biology



Safety Data  
Sciences Associate



Jr Data Scientist



- Statistical modeling for **Churn prediction**
- Ensemble model for clf **customer recurrence**
- Customer acquisition/recurrence **forecast**



Data Scientist

- Data engineering for **pricing automation**
- Hierarchical **Time Series Forecasting**
- Provide **training** for data scientists



Data Scientist

- Predict tool risk index
- 1 upcoming paper

7 years at academia  
3-4 years at industry

MSc. Sciences  
Ph.D. Genetics

- 3 published papers

Data Science  
& Analytics



Data Science  
Consultor / Teacher\*

- Statistics
- Data cleaning/wrangling
- Data visualization
- Clustering
- Recommendation
- Regression and classification models
- Model Deployment



# Why did I choose this topic?



**Data leakage is a flaw in machine learning that leads to overoptimistic results**

Leakage affects at least 294 papers across 17 scientific fields:

**Patterns**

 **CellPress**  
OPEN ACCESS

**Article**

## **Leakage and the reproducibility crisis in machine-learning-based science**

**Sayash Kapoor<sup>1,2,\*</sup> and Arvind Narayanan<sup>1</sup>**

<sup>1</sup>Department of Computer Science and Center for Information Technology Policy, Princeton University, Princeton, NJ 08540, USA

<sup>2</sup>Lead contact

\*Correspondence: [sayashk@princeton.edu](mailto:sayashk@princeton.edu)

<https://doi.org/10.1016/j.patter.2023.100804>



### **Anecdotal experience:**

- Code from students
- Blog & social media posts
- Courses



# What will you learn?

- What data leakage is
- How to avoid data leakage
- How to identify data leakage
  
- Learn how to use machine learning pipelines
- Understand how machine pipeline works elsewhere



**Note:** This workshop is not for you if you are highly familiar with machine learning pipelines (i.e. column transformers & pipelines) and is confident about your ability to avoid data leakage.



# Requirements

## Theory (slides)

Good to have

- Basic machine learning knowledge
- Plus: familiarity with scikit-learn

## Practice (hands-on examples)

Expected:

- Python (core functions & pandas)

Good to have

- Basic machine learning knowledge
- Basic familiarity with scikit-learn



# Agenda

01

## Machine learning lifecycle

Framework, validation & feature engineering

02

## Data Leakage

Definition and common mistakes

03

## Scikit-learn pipelines (focus on pre-processing)

Transformers, estimators, ColumnTransformers & Pipelines

04

## Hands-on!

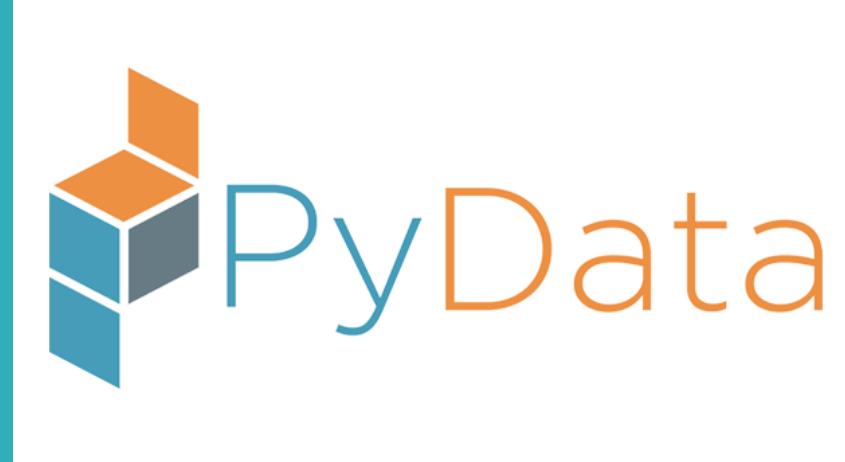
Walkthrough our code!

### Notebooks:

- 1.Pre-processing  
(manual & column transformers)
- 2.Scikit-learn and feature-engine pipelines
- 3.Imbalanced-learn pipelines
- 4.PyCaret AutoML pipelines
- 5.SparkML pipeline

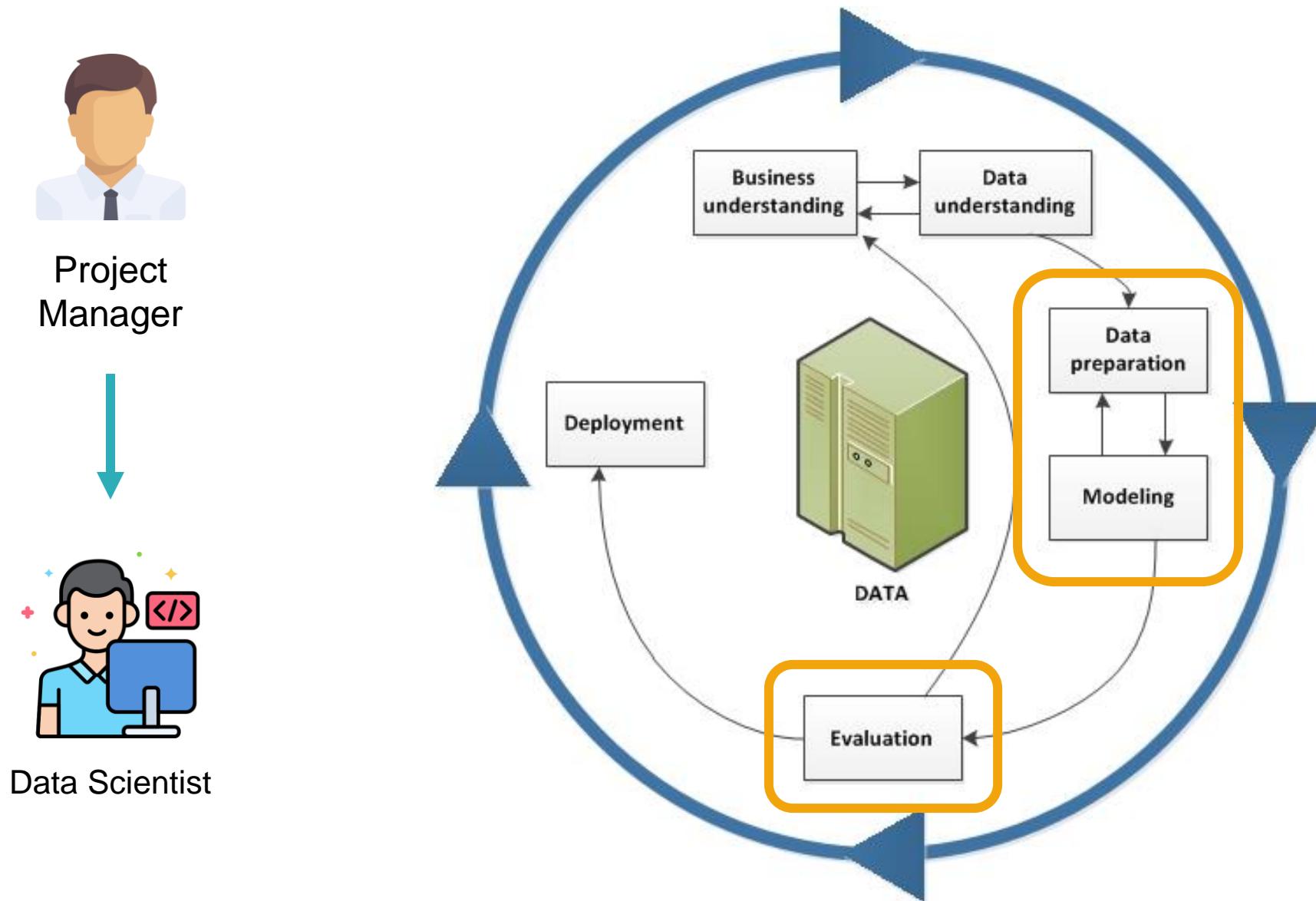


cmcouto-silva/Talks/tree/main/PyDataGlobal

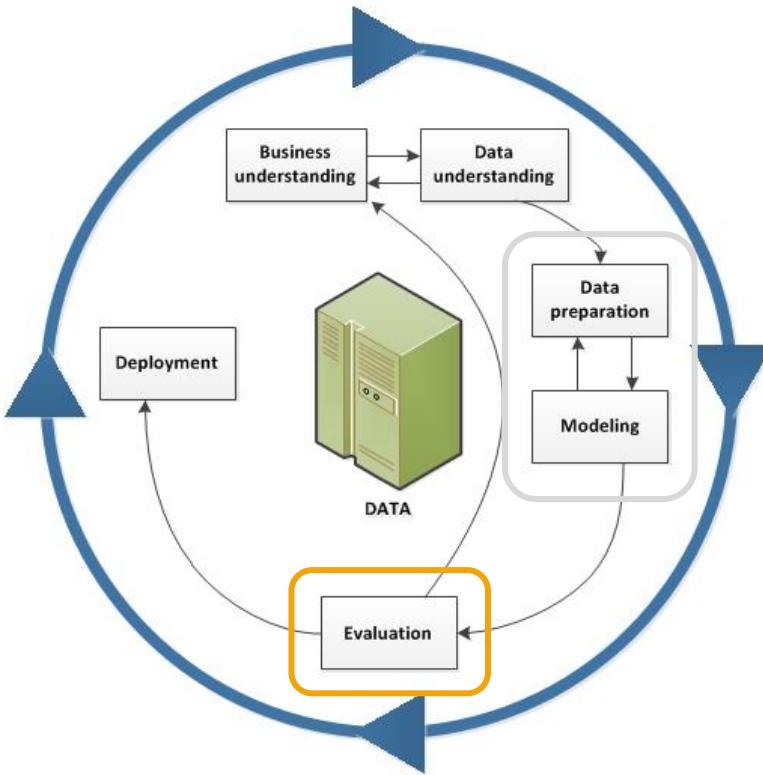


Let's begin!

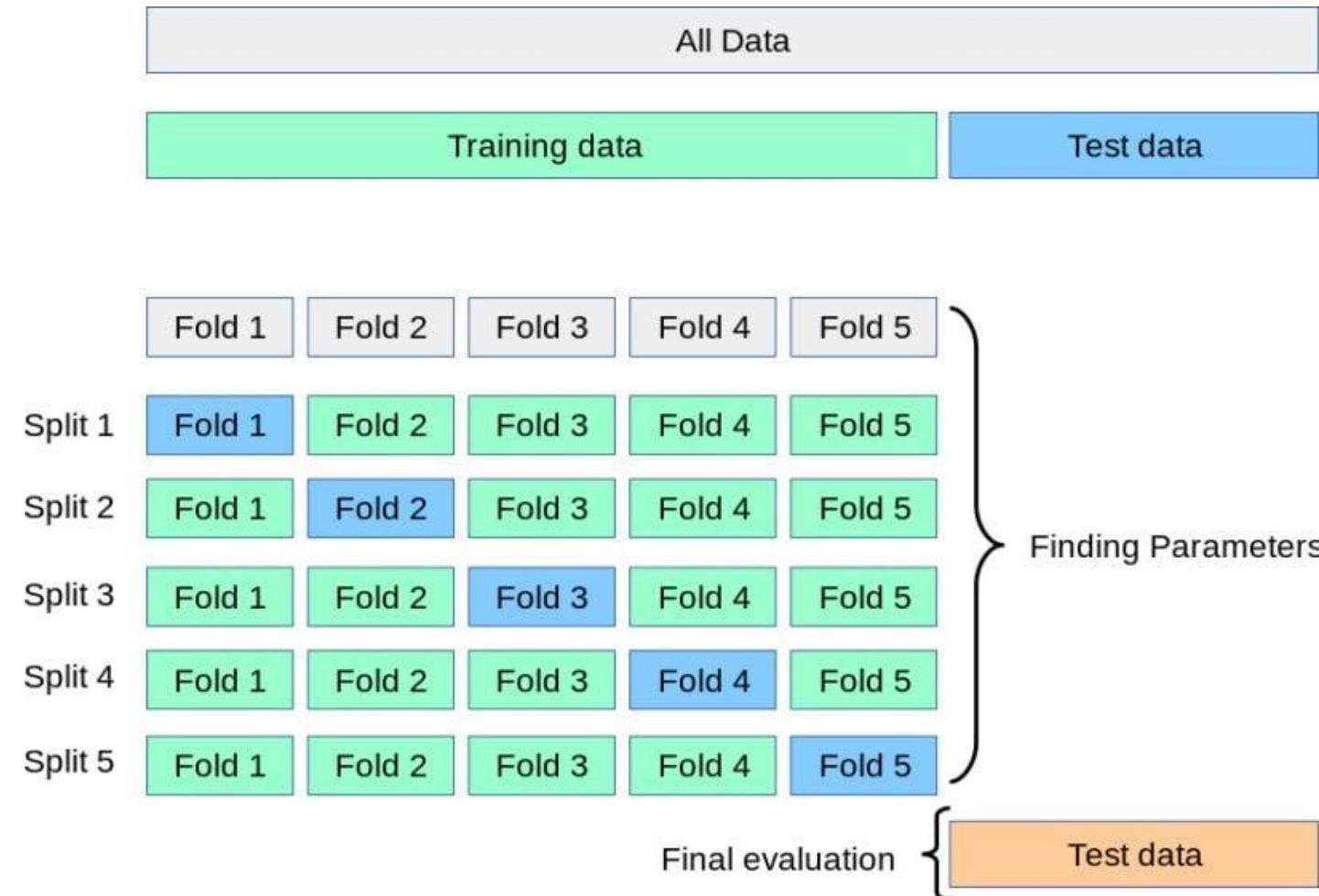
# Machine learning lifecycle



# Machine learning lifecycle

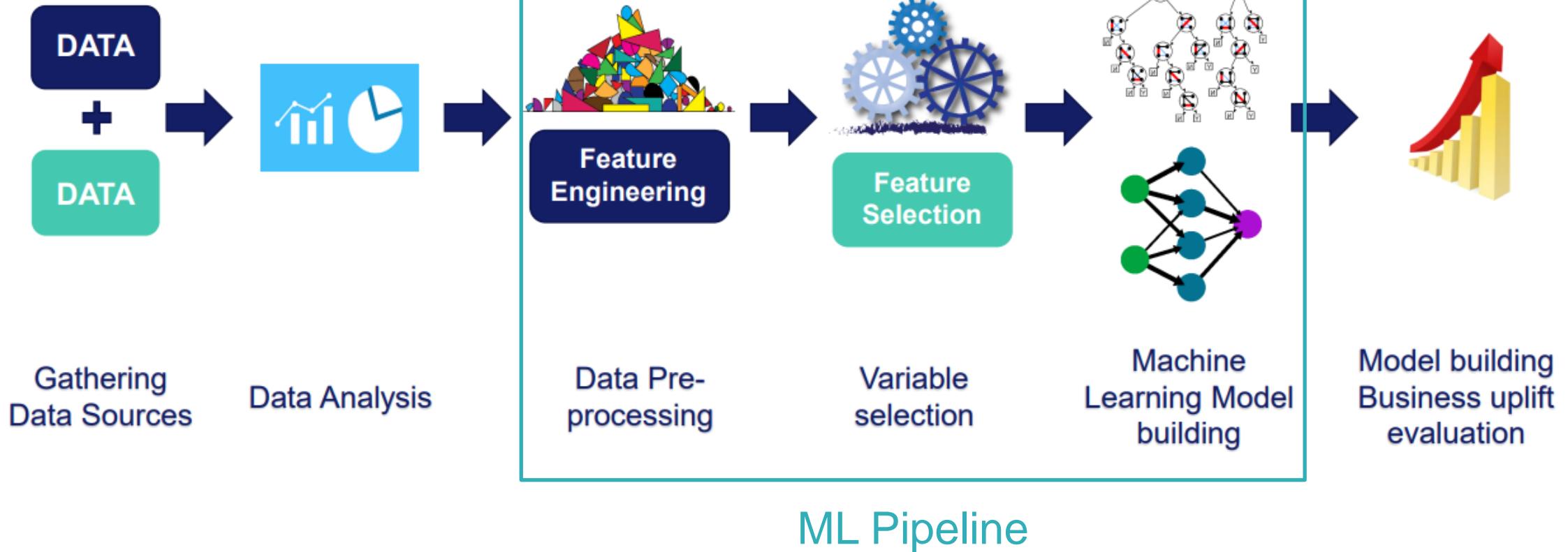
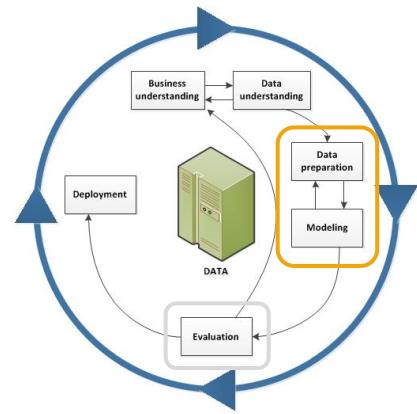


The test should know **nothing** about the training set!



→ We want a good generalist model!

# Machine learning lifecycle



# Feature engineering

## Missing values

Models do not accept missing values.

## Categorical variables

Models do not accept strings.

**Cardinality:** high number of labels  
**Rare labels:** infrequent categories

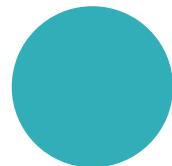
## Distribution

A better spread of values may benefit performance

## Outliers

This can lead to tremendous weight and bad generalization

# Feature engineering



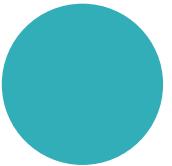
## Numeric feature scaling

Z-score and min-max normalization.



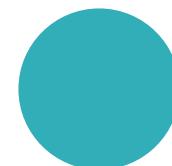
## Categorical feature encoding

One-hot, ordinal frequency or target encoding.



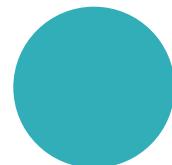
## Feature selection

Dimensionality reduction, variable importances and regularization.



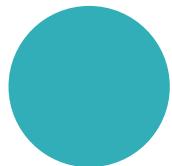
## Feature transformation

Logarithmic, polynomial, binning.



## Feature interaction

Product of two variables, custom combinations.



## Handling data imbalance

Random under- and upper-sampling, SMOTE.

**Note:** depending on the data type (text data, time-series, images), distinct pre-processing should be applied.

# When does it break?

When the training data includes information that allows the model to make predictions that it cannot make in a real-world scenario. It usually leads to overoptimistic results.

**Analogy:** It's like having the answers (or tips) to an exam in advance — it doesn't truly test your knowledge!



Data leakage is a flaw in machine learning that leads to overoptimistic results



Data leakage can be a multi-million-dollar mistake in many data science applications 💰

# Data leakage: common reasons

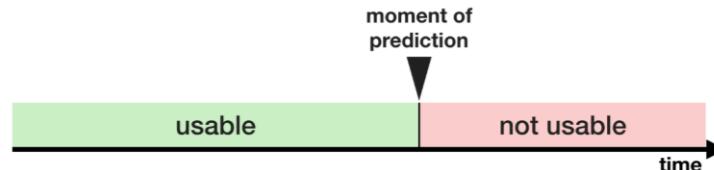
L1

## Lack of clean separation of training and test dataset

- No test set
- Pre-processing on training and test set
- Feature selection on training and test set
- Duplicates in datasets

L2

## Model uses features that are not legitimate



L3

## Test set is not drawn from the distribution of scientific interest

- Temporal leakage
- Nonindependence between training and test samples
- Sampling bias in test distribution

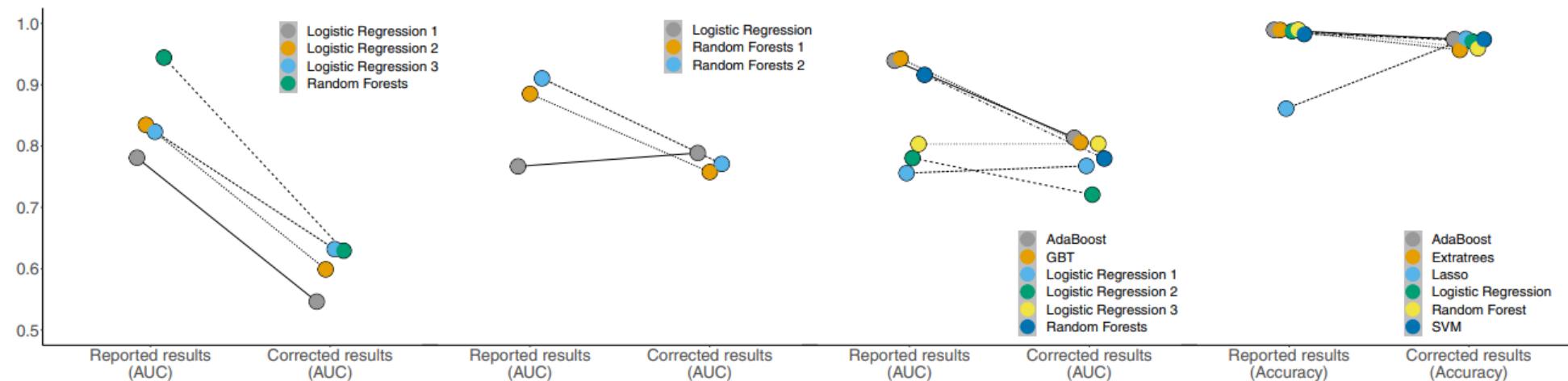
## Handling unbalanced data

- Resampling test data

## Patterns

Leakage and the reproducibility crisis in machine-learning-based science

Article



Paper	Muchlinski et al.	Colaresi and Mahmood	Wang	Kaufman et al.
<b>Claim</b>	Random Forests model drastically outperforms Logistic regression models	Random Forests models drastically outperform Logistic regression model	Adaboost and Gradient Boosted Trees (GBT) drastically outperform other models	Adaboost outperforms other models
<b>Error</b>	<b>[L1.2] Pre-proc. on train-test</b> (Incorrect imputation)	<b>[L1.2] Pre-proc. on train-test</b> (Incorrect reuse of an imputed dataset)	<b>[L1.2] Pre-proc. on train-test.</b> (Incorrect reuse of an imputed dataset) <b>[L3.1] Temporal leakage</b> ( $k$ -fold cross validation with temporal data)	<b>[L2] Illegitimate features</b> (Data leakage due to proxy variables) <b>[L3.1] Temporal leakage</b> ( $k$ -fold cross validation with temporal data)
<b>Impact</b>	Random Forests perform no better than Logistic Regression	Random Forests perform no better than Logistic Regression	Difference in AUC between Adaboost and Logistic Regression drops from 0.14 to 0.01	Adaboost no longer outperforms Logistic Regression. None of the models outperform a baseline model that predicts the outcome of the previous year
<b>Discussion</b>	Impact of the incorrect imputation is severe since 95% of the out-of-sample dataset is missing and is filled in using the incorrect imputation method	Re-use the dataset provided by Muchlinski et al., which uses an incorrect imputation method	Re-use the dataset provided by Muchlinski et al., which uses an incorrect imputation method	Use several proxy variables for the outcome as predictors (e.g., <i>colwars</i> , <i>cowwars</i> , <i>sdwars</i> , all proxies for civil war), leading to near perfect accuracy



# Solution



## Transformers & pipelines

Transformers and pipelines are specialized classes for machine learning frameworks that play a crucial role in the data preprocessing and model building process.



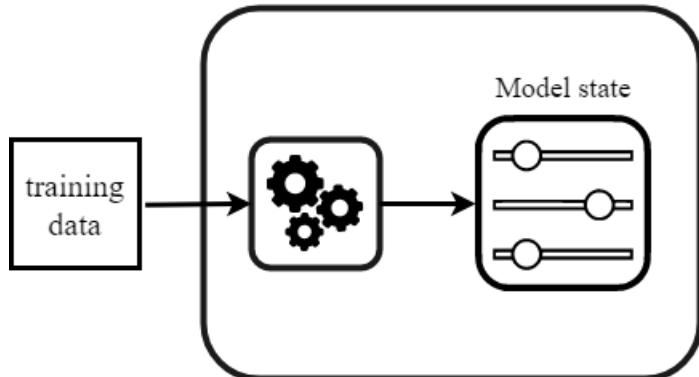
# How does scikit-learn work?

- Transformer
- Estimators
- Pipelines & column transformers

**Note:** not only scikit-learn, many libraries, inspired by the scikit-learn framework, built a similar structure.

# Scikit-learn transformer

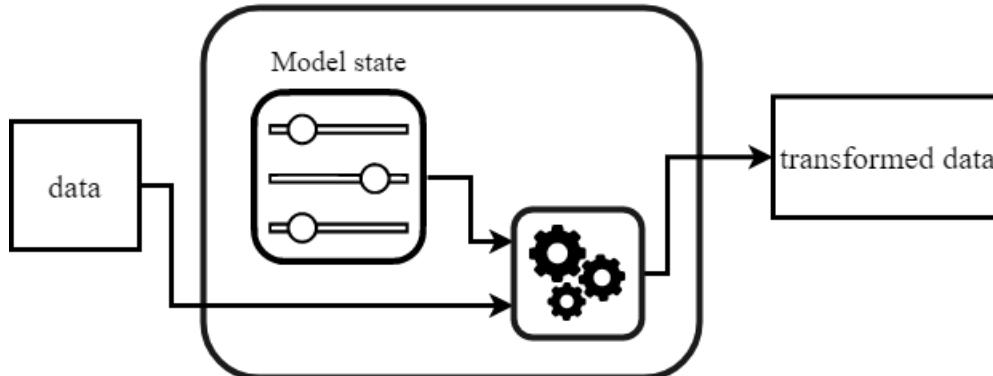
`transformer.fit(data)`



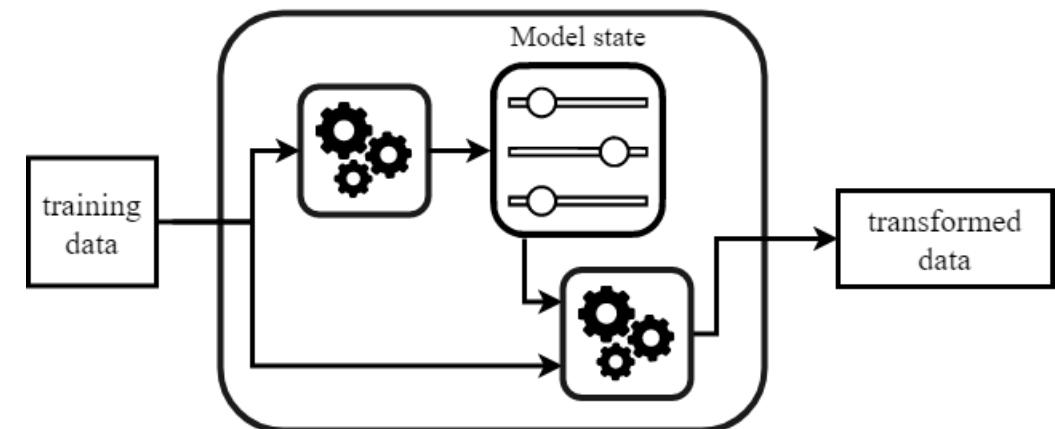
```
▶ imputer = SimpleImputer(strategy='median')  
imputer.fit(X[target_columns])  
imputer.statistics_  
  
array([29. , 70.35])
```

```
▶ scaler = StandardScaler()  
scaler.fit(X[target_columns])  
print(scaler.mean_, scaler.scale_)  
  
[32.37114866 64.76169246] [24.55773742 30.08791085]
```

`transformer.transform(data)`



`transformer.fit_transform(data)`



# Column transformers

## Input features

CLTV Monthly Charges Contract Dependents

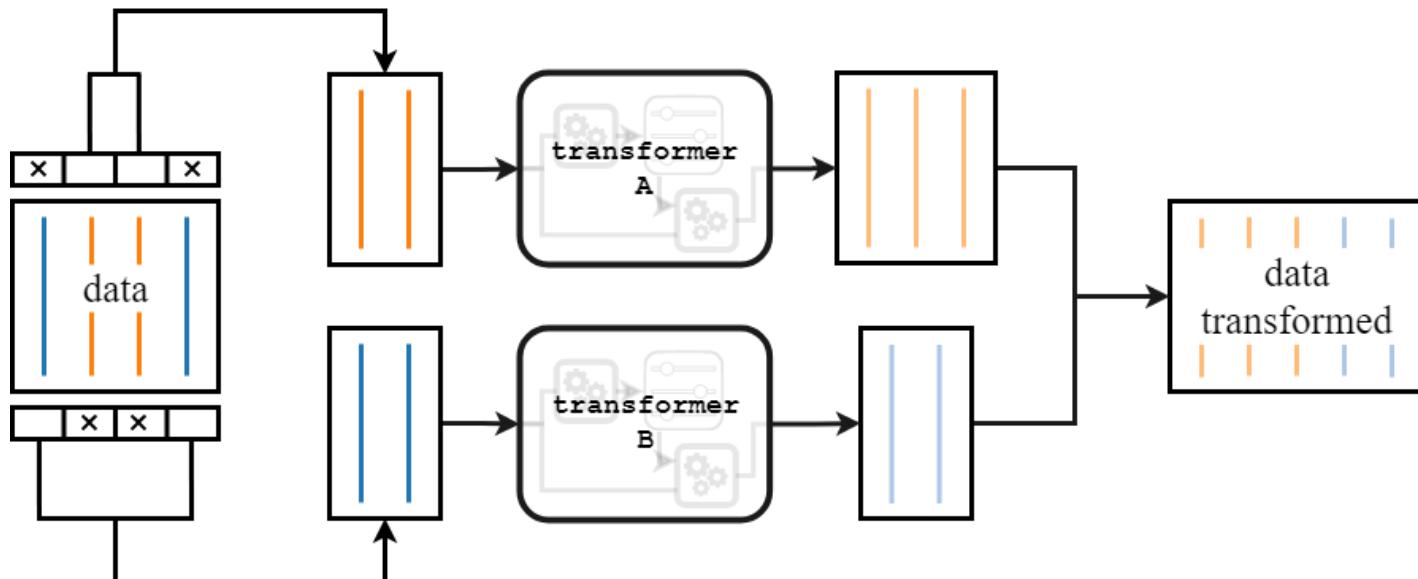
CustomerID

8695-WDYEA	4483	51.25	Month-to-month	No
3373-YZZYM	3007	19.20	Month-to-month	Yes
8748-HFWBO	5964	19.90	One year	Yes
2200-DSAAL	5914	80.65	Two year	No
0195-IESCP	3573	85.25	Month-to-month	No

Select columns

Transform

Concatenate columns



## Transformed features (output)

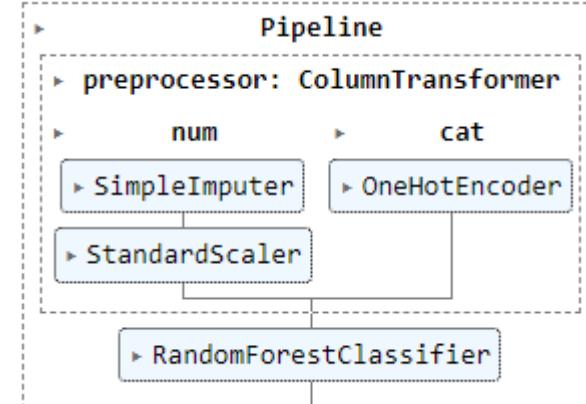
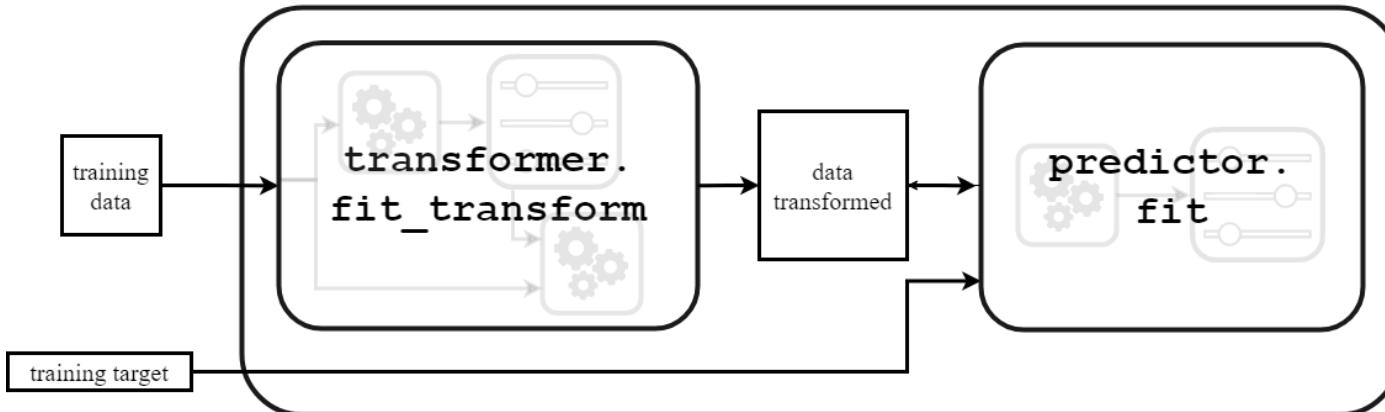
num\_CLTV num\_Monthly Charges cat\_Contract\_Month-to-month cat\_Contract\_One year cat\_Contract\_Two year cat\_Dependents\_Yes

CustomerID

8695-WDYEA	0.062885	-0.443960	1.0	0.0	0.0	0.0
3373-YZZYM	-1.183375	-1.507622	1.0	0.0	0.0	1.0
8748-HFWBO	1.313367	-1.484390	0.0	1.0	0.0	1.0
2200-DSAAL	1.271150	0.531755	0.0	0.0	1.0	0.0
0195-IESCP	-0.705473	0.684418	1.0	0.0	0.0	0.0

# Scikit-learn pipelines

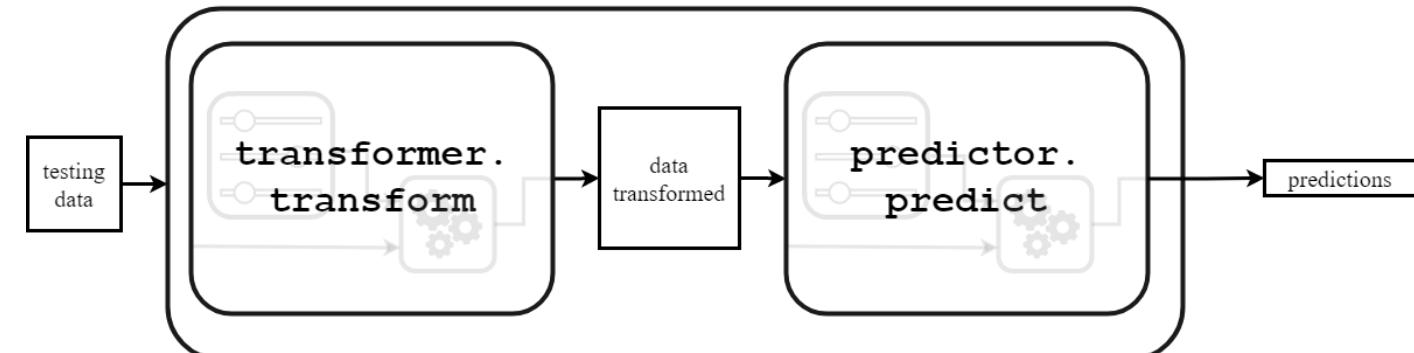
`pipeline.fit(data, target)`



`pipeline.fit(data, target):`

```
imputer.mean_
scaler.mean_, scaler.scale_
encoder.categories_, encoder.encoder_drop_idx_
classifier.estimators_
```

`pipeline.predict(data)`



## Warning: Risk of data leak

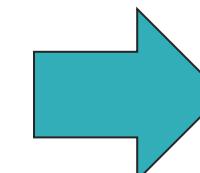
Do not use `scale` unless you know what you are doing. A common mistake is to apply it to the entire data *before* splitting into training and test sets. This will bias the model evaluation because information would have leaked from the test set to the training set. In general, we recommend using `StandardScaler` within a `Pipeline` in order to prevent most risks of data leaking: `pipe = make_pipeline(StandardScaler(), LogisticRegression())`.

## sklearn.preprocessing: Preprocessing and Normalization

The `sklearn.preprocessing` module includes scaling, centering, normalization, binarization methods.

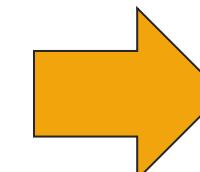
**User guide:** See the [Preprocessing data](#) section for further details.

<code>preprocessing.Binarizer(*[, threshold, copy])</code>	Binarize data (set feature values to 0 or 1) according to a threshold.
<code>preprocessing.FunctionTransformer([func, ...])</code>	Constructs a transformer from an arbitrary callable.
<code>preprocessing.KBinsDiscretizer([n_bins, ...])</code>	Bin continuous data into intervals.
<code>preprocessing.KernelCenterer()</code>	Center an arbitrary kernel matrix $K$ .
<code>preprocessing.LabelBinarizer(*[, neg_label, ...])</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.LabelEncoder()</code>	Encode target labels with value between 0 and $n_{\text{classes}}-1$ .
<code>preprocessing.MultiLabelBinarizer(*[, ...])</code>	Transform between iterable of iterables and a multilabel format.
<code>preprocessing.MaxAbsScaler(*[, copy])</code>	Scale each feature by its maximum absolute value.
<code>preprocessing.MinMaxScaler([feature_range, ...])</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.Normalizer([norm, copy])</code>	Normalize samples individually to unit norm.
<code>preprocessing.OneHotEncoder(*[, categories, ...])</code>	Encode categorical features as a one-hot numeric array.
<code>preprocessing.OrdinalEncoder(*[, ...])</code>	Encode categorical features as an integer array.
<code>preprocessing.PolynomialFeatures([degree, ...])</code>	Generate polynomial and interaction features.
<code>preprocessing.PowerTransformer([method, ...])</code>	Apply a power transform featurewise to make data more Gaussian-like.
<code>preprocessing.QuantileTransformer(*[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.RobustScaler(*[, ...])</code>	Scale features using statistics that are robust to outliers.
<code>preprocessing.SplineTransformer([n_knots, ...])</code>	Generate univariate B-spline bases for features.
<code>preprocessing.StandardScaler(*[, copy, ...])</code>	Standardize features by removing the mean and scaling to unit variance.

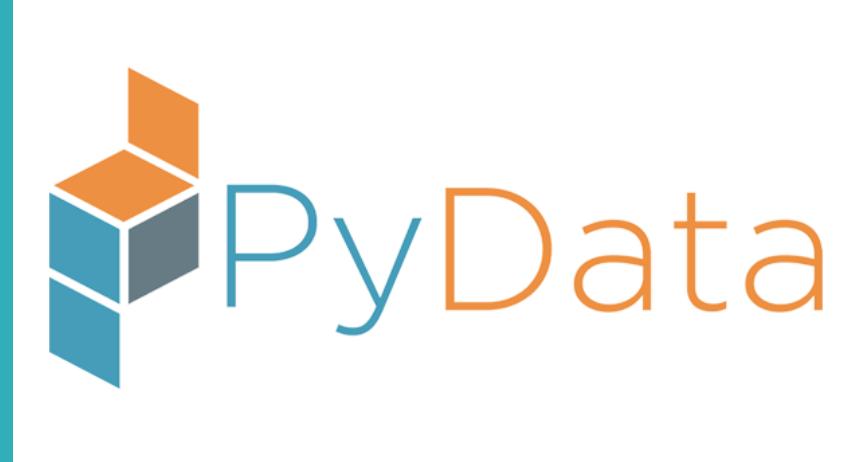


Use in pipeline

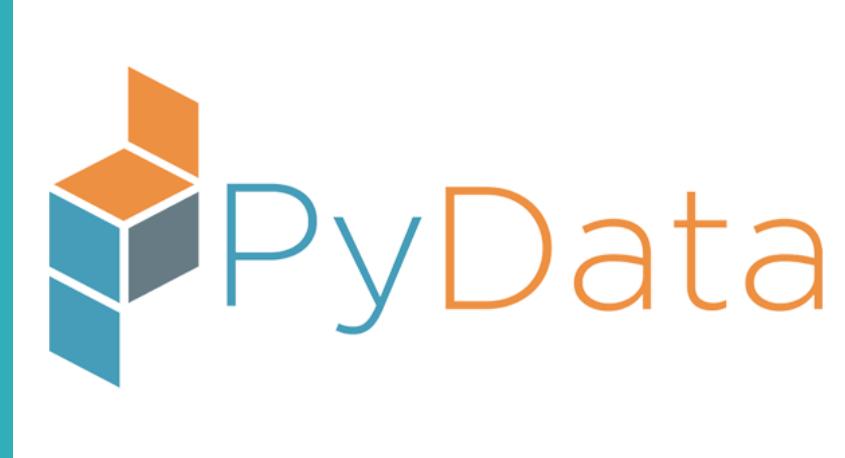
<code>preprocessing.add_dummy_feature(X[, value])</code>	Augment dataset with an additional dummy feature.
<code>preprocessing.binarize(X, *[, threshold, copy])</code>	Boolean thresholding of array-like or <code>scipy.sparse</code> matrix.
<code>preprocessing.label_binarize(y, *, classes)</code>	Binarize labels in a one-vs-all fashion.
<code>preprocessing.maxabs_scale(X, *[, axis, copy])</code>	Scale each feature to the [-1, 1] range without breaking the sparsity.
<code>preprocessing.minmax_scale(X[, ...])</code>	Transform features by scaling each feature to a given range.
<code>preprocessing.normalize(X[, norm, axis, ...])</code>	Scale input vectors individually to unit norm (vector length).
<code>preprocessing.quantile_transform(X, *[, ...])</code>	Transform features using quantiles information.
<code>preprocessing.robust_scale(X, *[, axis, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.scale(X, *[, axis, with_mean, ...])</code>	Standardize a dataset along any axis.
<code>preprocessing.power_transform(X[, method, ...])</code>	Parametric, monotonic transformation to make data more Gaussian-like.



Do not use unless you know what you are doing!



Hands-on!

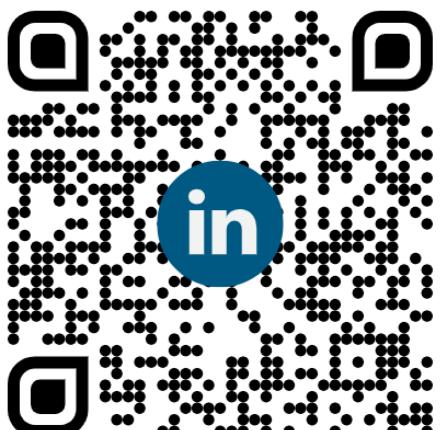


# Section Break

Time to grab a coup of coffee 



Feel free to contact me



cmcouto-silva

Thank you!

