

# Learning Scikit-Learn

Qiyang Hu

UCLA IDRE/OARC Workshop

February 25<sup>th</sup>, 2021

# Scikit-learn in minutes

- A Python machine learning framework
  - Library built on numpy, scipy, matplotlib
    - Started in 2007, publicly released in 2010
    - Is currently maintained by volunteers
- Installation/Loading
  - `conda install -c intel scikit-learn`
  - On H2: `module load anaconda3`  
`conda activate sklearn`
  - Using Google Colab
- Designed for easy-to-use productions
  - Simplicity
  - Qualitative code
    - Performance
    - Elegant APIs
  - Excellent docs: <https://scikit-learn.org>

```
# 2 samples, 3 features
x = [[ 1,  2,  3],
      [11, 12, 13]]
```

Data

```
# classes of each sample
y = [0, 1]
```

Modeling

```
from sklearn.ensemble import RandomForestClassifier

clf = RandomForestClassifier(random_state=0)

clf.fit(X, y)
```

```
# predict classes of the training data
clf.predict(X)



# predict classes of new data
clf.predict([[4, 5, 6], [14, 15, 16]])
```

Predicting

# Outline

- High-level overview of scikit-learn libraries
  - According to a typical machine learning workflow
  - A lot of colab snippets as examples
- Demo for working on a fun machine learning project
  - Titanic challenge in Kaggle
- Discussion on advanced topics
  - Boosting scikit-learn performance
  - Scikit-learn extension libraries
  - High-performance machine learning

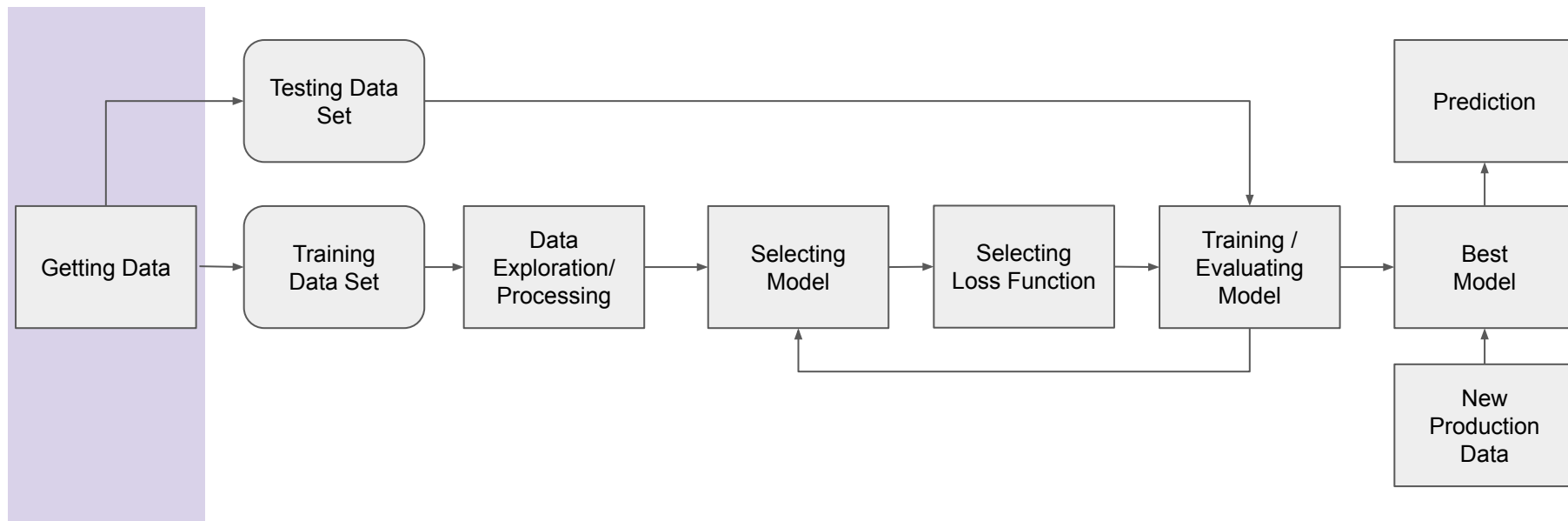
# What can/can't be expected in this class?

 CAN	 CAN'T
<ul style="list-style-type: none"><li>• Review on Machine learning workflows</li></ul>	<ul style="list-style-type: none"><li>• Introduction to various Machine Learning models</li></ul>
<ul style="list-style-type: none"><li>• A <b><i>BIG</i></b> picture on scikit-learn's features, functions &amp; components</li></ul>	<ul style="list-style-type: none"><li>• Discussions on the details of specific scikit-learn function interfaces</li></ul>
<ul style="list-style-type: none"><li>• Providing handy examples as demos (probably for studying <i>after</i> the class)</li></ul>	<ul style="list-style-type: none"><li>• Line-by-line explanation on every demo code</li></ul>
<ul style="list-style-type: none"><li>• High-level introduction on some advanced topics</li></ul>	<ul style="list-style-type: none"><li>• A lecture on high-performance machine learning</li></ul>

# Outline

- High-level overview of scikit-learn libraries
  - According to a typical machine learning workflow
  - A lot of colab snippets as examples
- Demo for working on a fun machine learning project
  - Titanic challenge in Kaggle
- Discussion on advanced topics
  - Boosting scikit-learn performance
  - Scikit-learn extension libraries
  - High-performance machine learning

# Workflow for a machine learning project

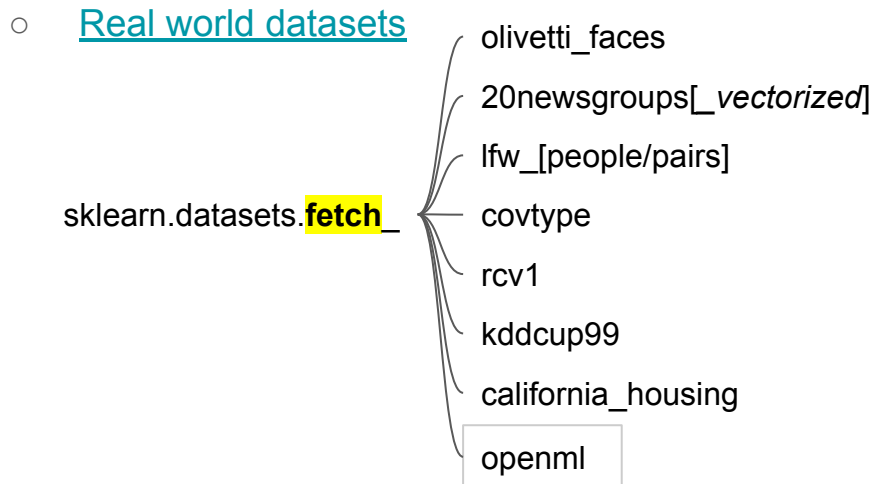
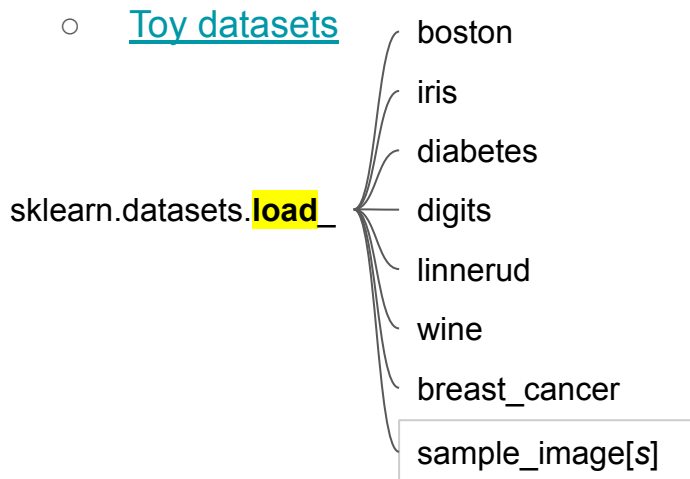


# Data input and loader

- Data format can be input directly as:
  - Dense data: `numpy.ndarray`
  - Sparse data: `scipy.sparse.matrix`
- Data can be loaded from standard datasets:

Loaders and Fetchers returns a dictionary-like **bunch** object

```
>>> b = Bunch(a=1, b=2)
>>> b['b']
2
>>> b.b
2
>>> b.a = 3
>>> b['a']
3
>>> b.c = 6
>>> b['c']
6
```



# Data Generator

sklearn.datasets.**make\_**

blob

classification

gaussian\_quantiles

hastie\_10\_2

circles

moons

multilabel\_classification

biclusters

checkerboard

regression

friedman[1/2/3]

sparse\_uncorrelated

s\_curve

swiss\_roll

low\_rank\_matrix

sparse\_coded\_signal

spd\_matrix

sparse\_spd\_matrix

```
(n_samples=100, n_features=2, *,
centers=None, cluster_std=1.0,
center_box=(- 10.0, 10.0,
shuffle=True, random_state=None,
return_centers=False)
```

For classification and clustering

For regression

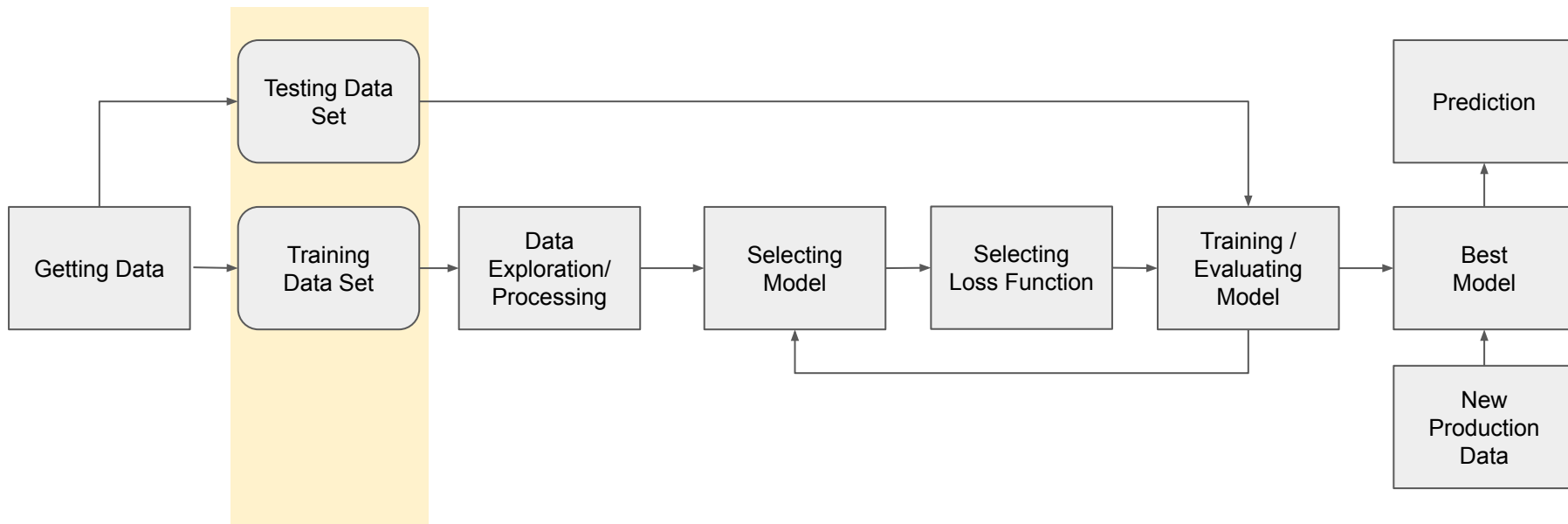
For manifold learning

For decomposition



**bit.ly/1sk1\_01**

# Workflow for a machine learning project



# Split training and testing dataset

- Essential for an unbiased evaluation of prediction performance

- Process related with model evaluation and selection

- Sklearn has helper function

`train_test_split`

to randomly split

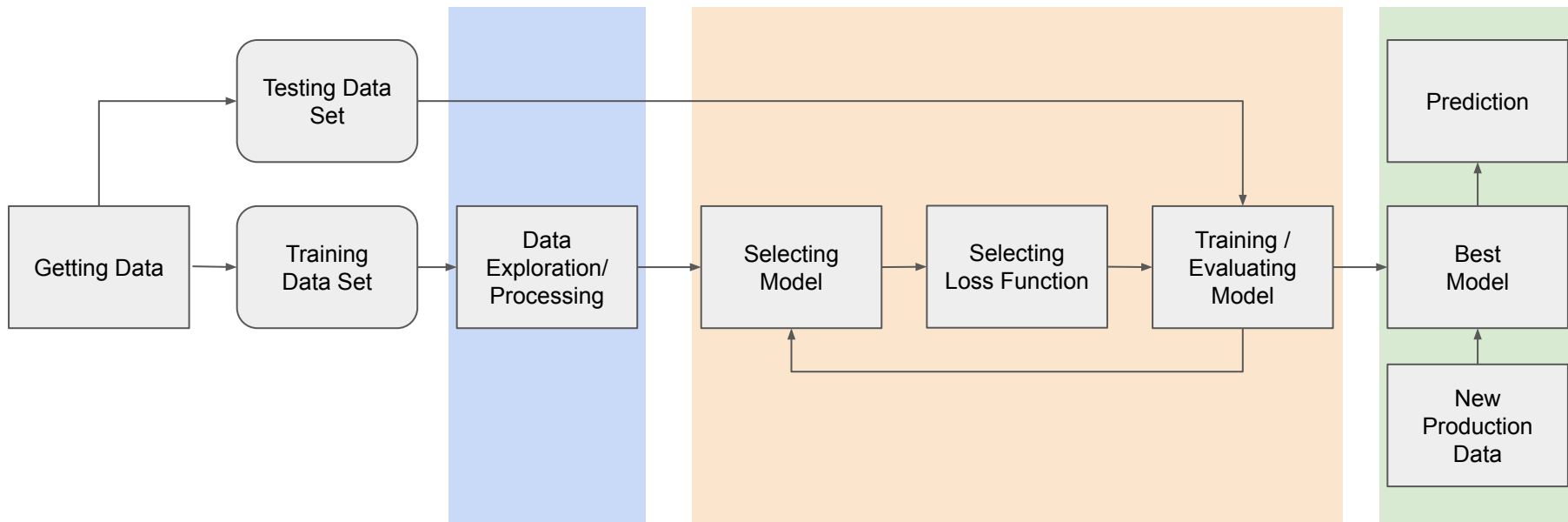
- A wrapper around ShuffleSplit
  - Only allows for stratified splitting
  - Cannot account for groups
  - As a base for the default cross-validations

```
>>> import numpy as np
>>> from sklearn.model_selection import train_test_split
>>> from sklearn import datasets
>>> from sklearn import svm

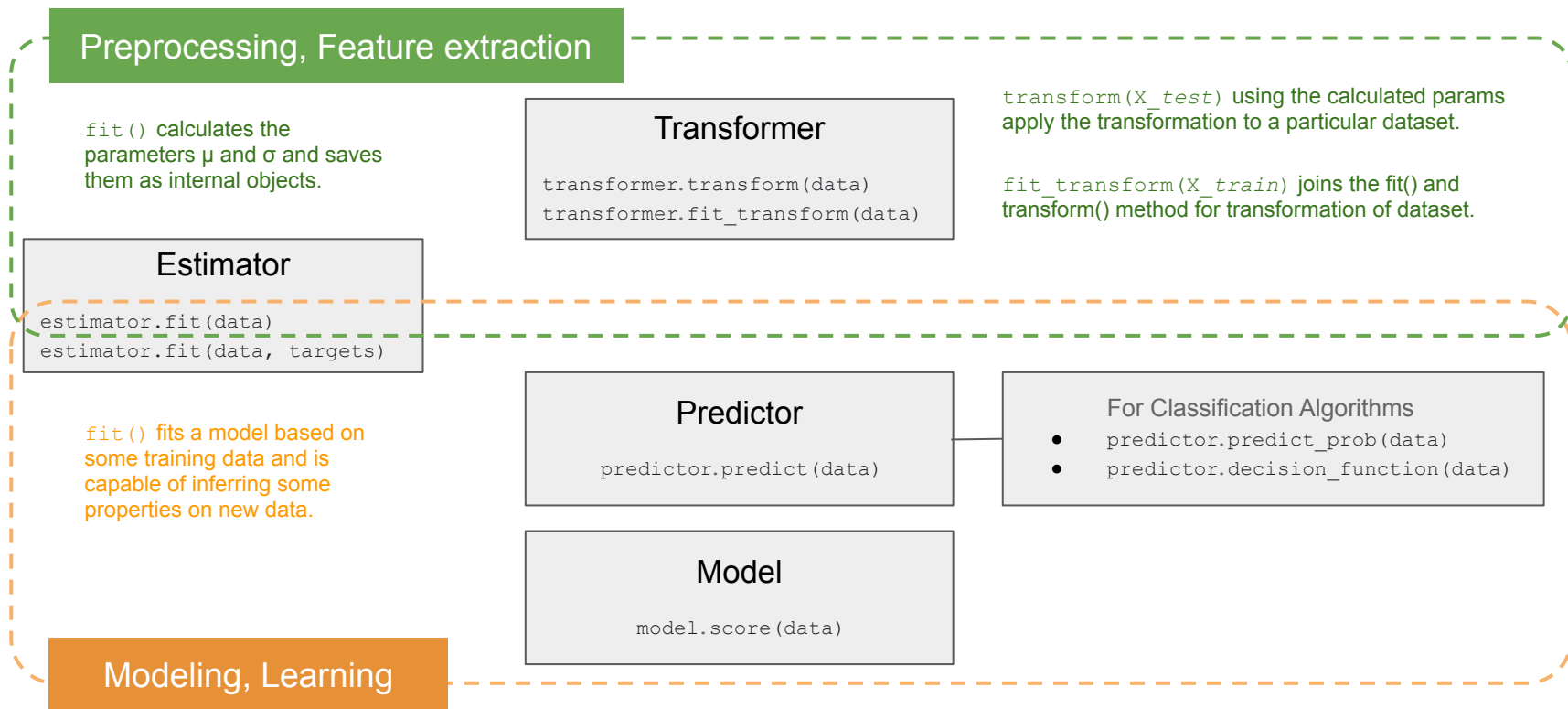
>>> X, y = datasets.load_iris(return_X_y=True)
>>> X.shape, y.shape
((150, 4), (150,))
>>> X_train, X_test, y_train, y_test = train_test_split(
...     X, y, test_size=0.4, random_state=0)

>>> X_train.shape, y_train.shape
((90, 4), (90,))
>>> X_test.shape, y_test.shape
((60, 4), (60,))
```

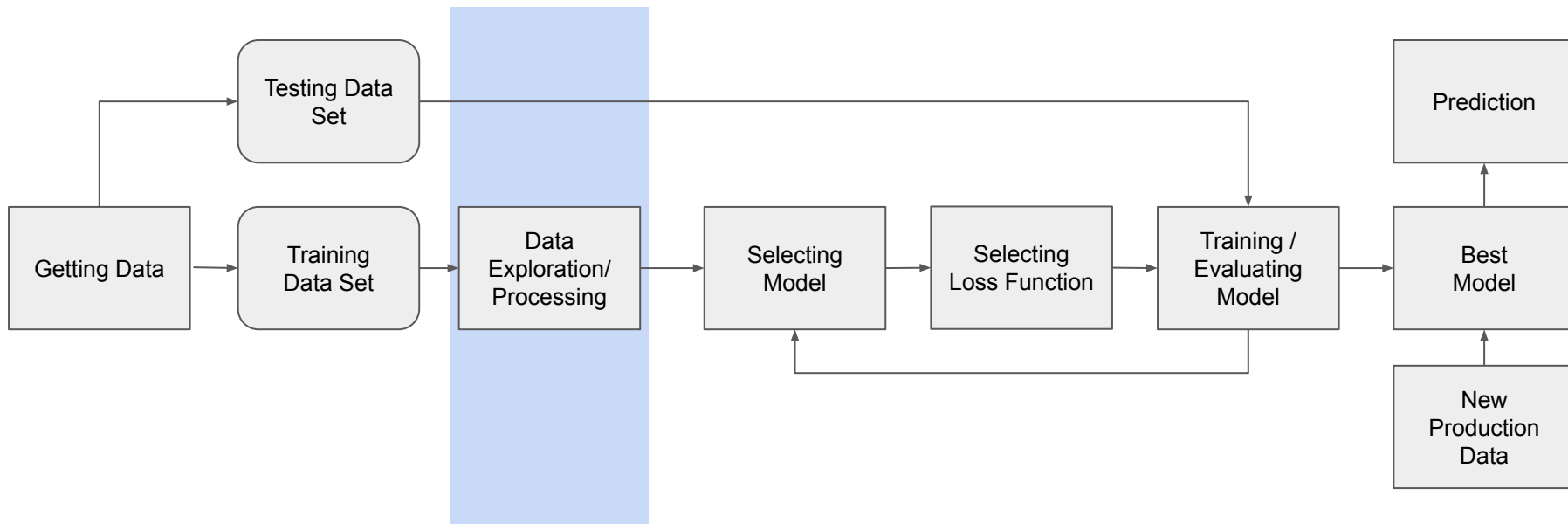
# Workflow for a machine learning project



# Core objects (estimators)



# Workflow for a machine learning project



# Preprocessing:

```
from sklearn.preprocessing import  
StandardScaler  
sc = StandardScaler()  
sc.fit_transform(X_train)  
sc.transform(X_test)
```

**sklearn.preprocessing.**

[StandardScaler](#) / RobustScaler  
MinMaxScaler / MaxAbsScaler  
KernelCenterer  
QuantileTransformer  
PowerTransformer  
normalize  
Normalizer  
OrdinalEncoder/LabelEncoder  
OneHotEncoder  
KBinsDiscretizer  
Binarizer  
FunctionTransformer  
PolynomialFeatures

Standardization, or mean removal  
and variance scaling

Non-linear transformation

Normalization

Encoding categorical features

Discretization

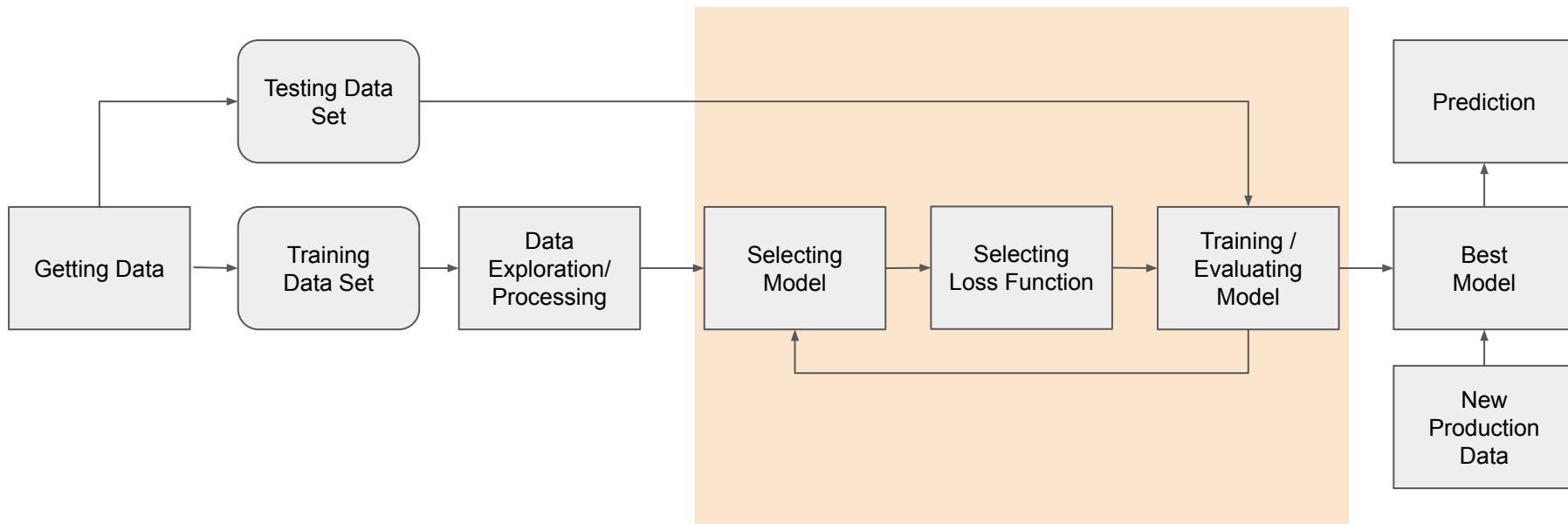
Custom transformers

**sklearn.imput.**

SimpleImputer  
IterativeImputer  
KNNImputer  
MissingIndicator

Imputation of missing values

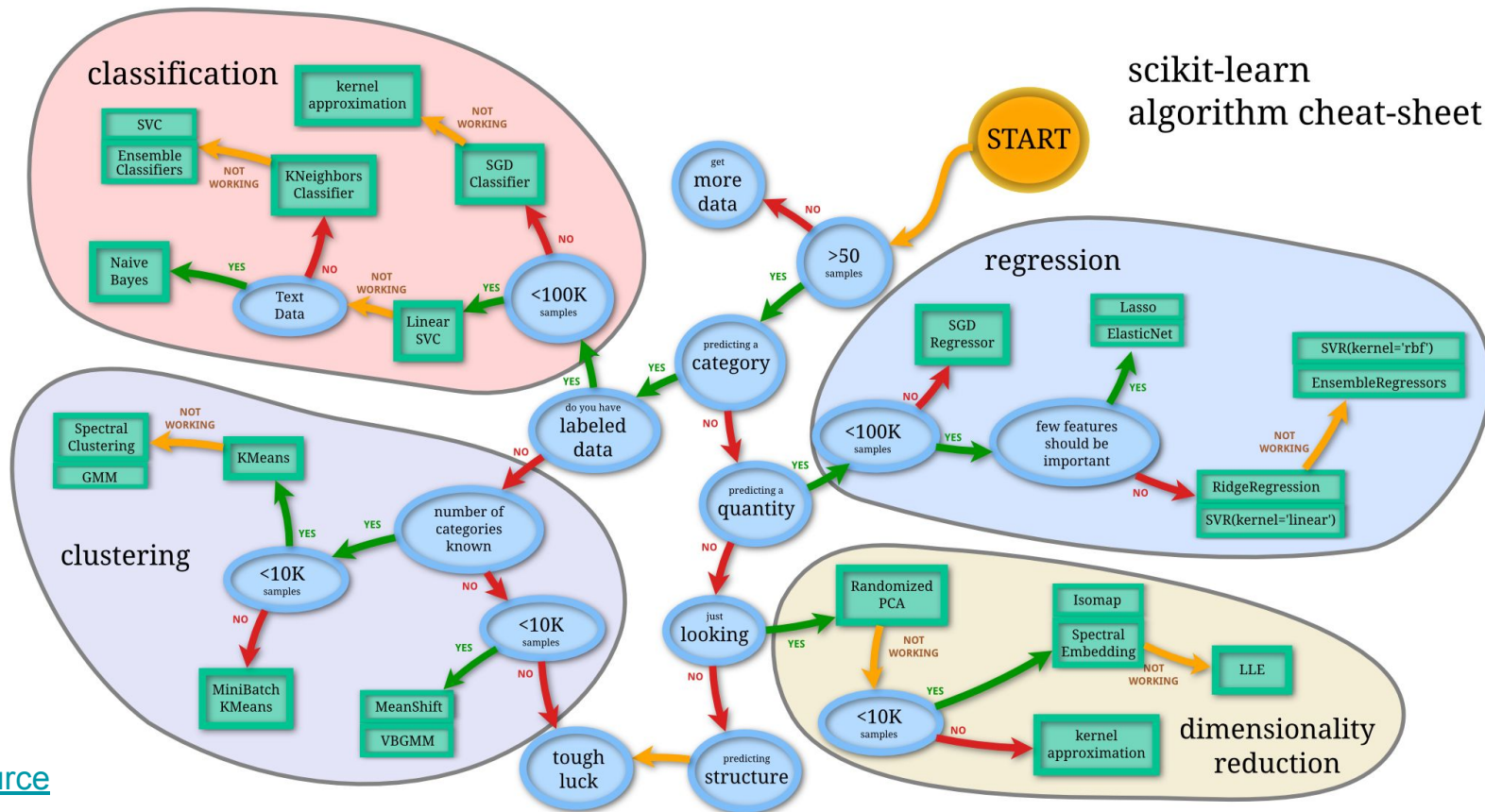
# Workflow for a machine learning project





# Choosing the right estimator (algorithm)

scikit-learn  
algorithm cheat-sheet



[Source](#)

# Pseudo-code template for modeling and learning

```
from sklearn. { linear_model  
               svm  
               tree  
               naive_bayes  
               multioutput  
               ensemble  
               cluster  
               decomposition  
               ...  
             } import SpecModel  
  
model = SpecModel( hyperparameter )  
  
model.fit( X, y )  
  
y_pred = model.predict( X_new )  
  
s = model.score( X_new )
```

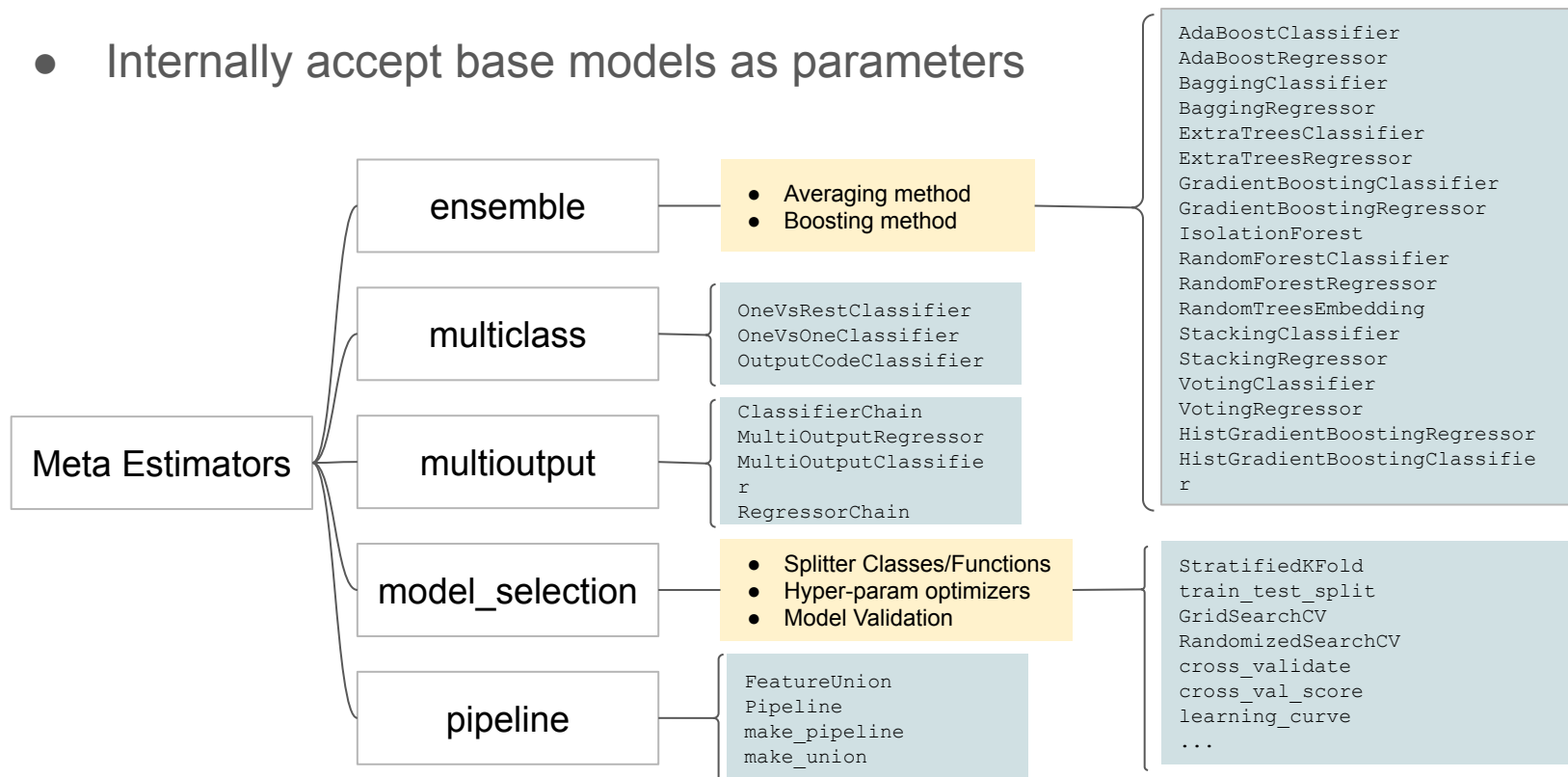
LogisticRegression  
LogisticRegressionCV  
PassiveAggressiveClassifier  
Perceptron  
RidgeClassifier  
RidgeClassifierCV  
SGDClassifier  
LinearRegression  
Ridge  
RidgeCV  
SGDRegressor  
ElasticNet  
ElasticNetCV  
Lars  
LarsCV  
Lasso  
LassoCV  
LassoLars  
LassoLarsCV  
LassoLarsIC  
OrthogonalMatchingPursuit  
OrthogonalMatchingPursuitCV  
ARDRegression  
BayesianRidge  
PoissonRegressor  
GammaRegressor  
HuberRegressor  
RANSACRegressor  
...

penalty='l2', tol=0.0001, C=0.1,  
fit\_intercept=True,  
solver='liblinear', max\_iter=100,  
multi\_class='ovr', n\_jobs=1, ...

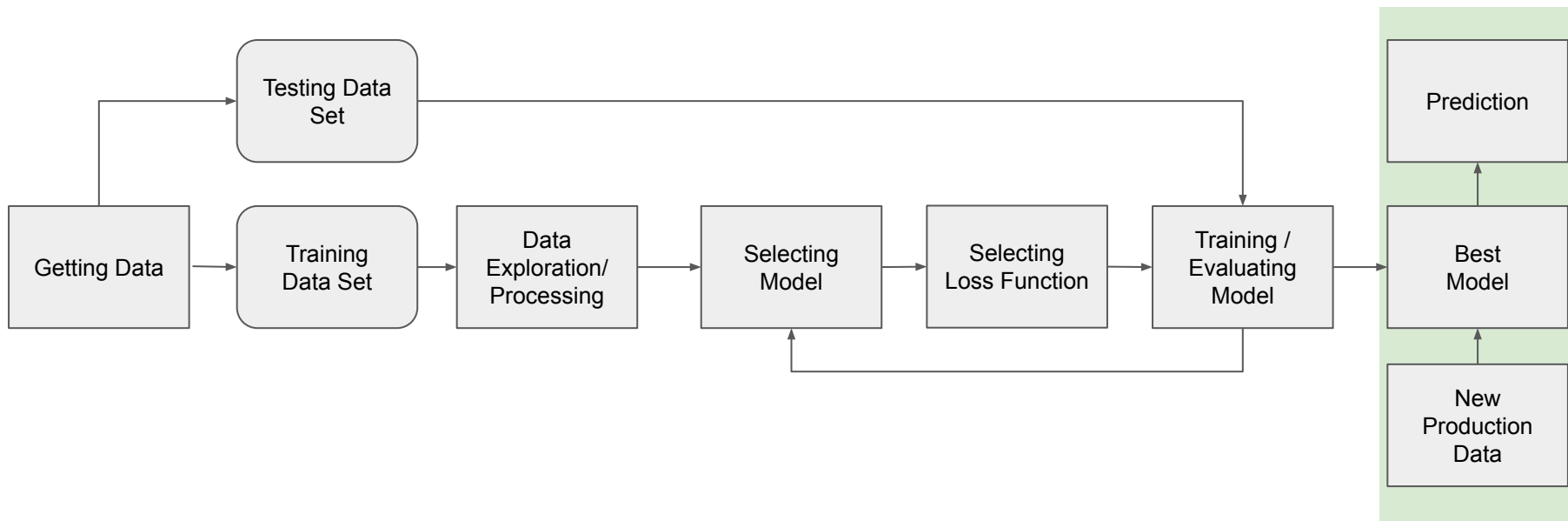
The diagram illustrates the flow of data in a pseudo-code template for modeling and learning. It shows the import of various modules from sklearn, the creation of a SpecModel object with hyperparameters, and the subsequent fitting and prediction steps. A list of models and hyperparameters is provided on the right, and a specific set of hyperparameters is highlighted in a box.

# Meta-estimator: as an assembly of base estimators

- Internally accept base models as parameters



# Workflow for a machine learning project



# Model persistence (saving/restoring a trained model)

- Python's built-in serialization:
  - Using `pickle` or `joblib`: dump and load
  - Custom transformers in Pipeline cannot be serialized by pickle or joblib
    - Consider using Neuralxle's module to [save custom pipeline](#) in step wise
  - Pickled model better to be deployed using containers to avoid portability issues
- Other exporting formats
  - Open Neural Network Exchange (ONNX)
    - [sklearn-onnx](#)
  - Predictive Model Markup Language (PMML)
    - [sklearn2pmml](#)

# Outline

- High-level overview of scikit-learn libraries
  - According to a typical machine learning workflow
  - A lot of colab snippets as examples
- Demo for working on a fun machine learning project
  - Titanic challenge in Kaggle
- Discussion on advanced topics
  - Boosting scikit-learn performance
  - Scikit-learn extension libraries
  - High-performance machine learning

# Titanic Kaggle Challenge

To predict the survival or the death of a given passenger in Titanic

- Titanic Facts:
  - Survivors
    - 492 passengers
    - 214 crews
  - Victims
    - 832 passengers
    - 685 crews
    - Death causes: drowning, hypothermia, injury, suicide, ...
    - [List of deaths](#)



**bit.ly/1sk1\_02**



# Outline

- High-level overview of scikit-learn libraries
  - According to a typical machine learning workflow
  - A lot of colab snippets as examples
- Demo for working on a fun machine learning project
  - Titanic challenge in Kaggle
- Discussion on advanced topics
  - Boosting scikit-learn performance
  - Scikit-learn extension libraries
  - High-performance machine learning

# Basic performance tips for scikit-learn projects

- Profiling your Python code
  - Using `line_profiler`, `memory_profiler`
- Using Numpy and Scipy as much as possible to replace nested for loops
  - [Performance tips for numpy and scipy](#)
- Writing cython wrappers for well-maintained C/C++ algorithm implementation
  - example: `liblinear`, `libsvm` used in Logistic Regression and SVM models.
- Multiple-core parallelism:
  - Specifying `n_jobs` for model training, evaluation, hyperparameter tuning
  - Using OpenMP functions through Cython imported by `cimport openmp`
  - Using `joblib.Parallel`
    - [joblib](#) is a set of tools to provide lightweight pipelining in Python.
    - Already used in some sklearn classes (`ElasticNet`, `SGDClassifier`, ...)
    - Various parallel backends: `locky`, `multiprocessing`, `threading`, `dask`, `ray`, ...

# Scikit-learn extension libraries

- Libraries adopting scikit-learn functionalities
  - Data formats: [sklearn\\_pandas](#), [sklear\\_xarray](#), ...
  - Auto-ML: [auto-sklearn](#), [Featuretools](#), [Neuraxle](#), ...
  - Model visualization: [dtreeviz](#), [eli5](#), ...
  - Model selection: [scikit-optimize](#), [sklearn-deap](#), ...
  - Model export: [onnxmltools](#), [sklearn2pmmml](#), ...
  - Parallelization: [sk-dist](#)
  - Plotting: [scikit-plot](#)
- Libraries compatible with scikit-learn interfaces
  - Time-series models: [tslearn](#), [sktime](#), [seglearn](#), ...
  - Deep learning: [keras](#), [skorch](#), ...
  - Other regression/classification: [xgboost](#), [ML\\_Ensemble](#), [gplearn](#), ...
  - Decomposition and clustering: [lda](#), [hdbscan](#), ...

Libraries *nothing* related  
with scikit-learn

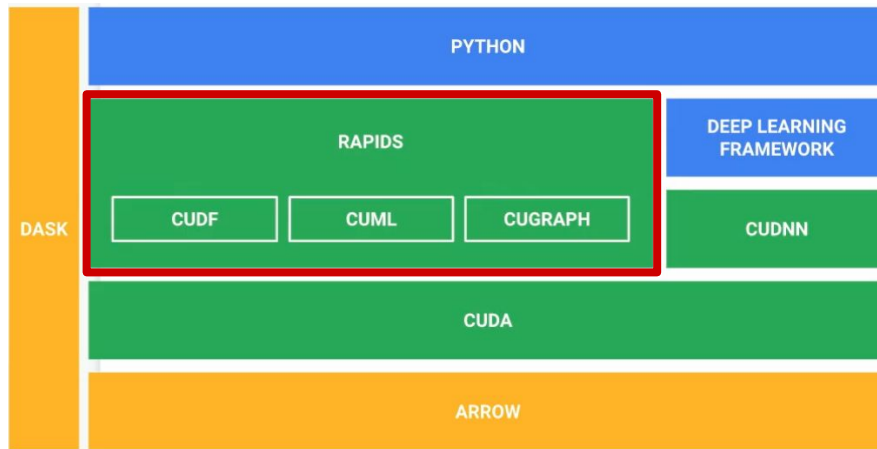
- [scikit-opt](#)

# High-performance machine learning

- Scikit-learn alone doesn't have support to GPU or TPU.

- Option 1: [RAPIDS](#) software libraries from nvidia

- Exactly same APIs
- Kernels rewritten by CUDA
  - *Numpy* -> *CuPy*
  - *pandas* -> *cuDF*
  - *scikit-learn* -> *cuML*
  - *networkx* -> *cuGraph*
- Dask + RAPIDS + blazingSQL, RAPIDS + Spark, xgboost, ...



- Option 2: JAX

- [JAX](#): compile NumPy using XLA, on GPUs and TPUs for high-performance machine learning.
- [sklearn-jax-kernel](#) is in an early stage.

# High-performance machine learning (cont'd)

- Option 3: using [daal4py](#)
  - Python APIs to Intel's OneDAL
  - Provide alternative estimators
  - Acceleration without code change:

- Intel CPU optimization patching

```
from daal4py.sklearn import patch_sklearn
patch_sklearn()
```

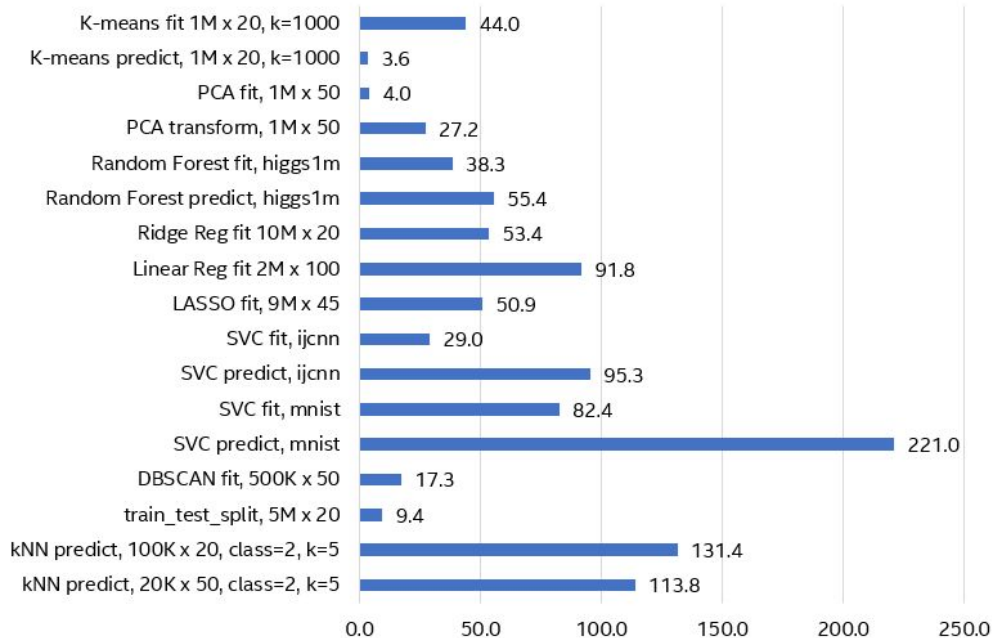
- Intel CPU/GPU optimizations patching

```
from daal4py.sklearn import patch_sklearn
from daal4py.oneapi import sycl_context
patch_sklearn()
with sycl_context("gpu"):
```

- Installation:

```
conda install daal4py -c intel
```

Speedups of daal4py-powered Scikit-learn over the original Scikit-learn



**bit.ly/1sk1\_03**

# Key Takeaways

- Learning Scikit-learn means:
  - Understanding the machine learning workflows
  - Learning its API conventions
  - Familiarizing its documentations
  - Inspecting the example codes
- Making scikit-learn more productive means:
  - Making sure everything work fine without performance optimization
  - Carefully profiling the code to address the performance bottleneck
  - Considering the scikit-learn extensive libraries
  - Trying RAPIDS, daal4py, watching JAX