

Pipelines

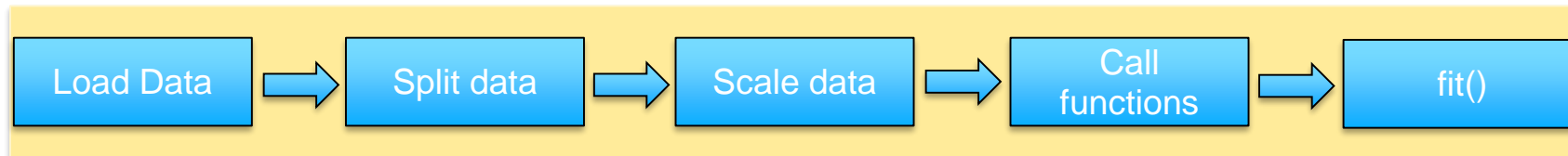
What is a pipeline-

1. Almost always, we need to tie together many different processes that we use to prepare data for machine learning based model
1. It is paramount that the stage of transformation of data represented by these processes are standardized
1. Pipeline class of sklearn helps simplify the chaining of the transformation steps and the model
1. Pipeline, along with the GridsearchCV helps search over the hyperparameter space applicable at each stage

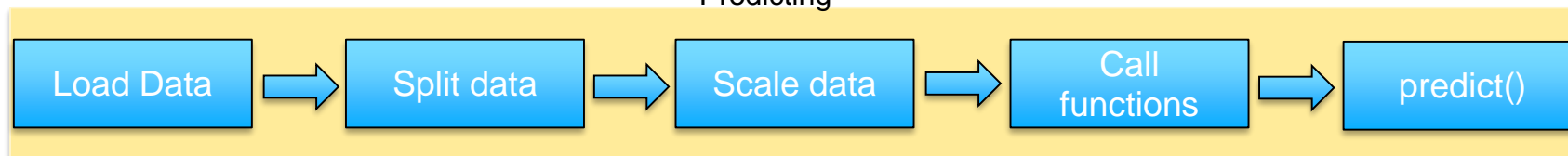
Pipelines

1. Sequentially apply a list of transforms and a final estimator.
2. Intermediate steps of the pipeline must be 'transforms', that is, they must implement fit and transform methods.
3. The final estimator only needs to implement fit
4. Helps standardize the model project by enforcing consistency in building testing and production. **Ref:** <https://scikit-learn.org/stable/modules/compose.html>

Building Model



Predicting



Ref: Pipeline_Simple.ipynb

Build a pipeline

1. Import the pipeline class
 - a. `from sklearn.pipeline import Pipeline`
2. Instantiate the class into an object by listing out the transformation steps. In the following example, a scaling function is followed by the SVC algorithm
 - a. `pipe = Pipeline([("scaler", MinMaxScaler()), ("lr", logisticregression())])`
1. Call the fit() function on the pipeline object
 - a. `pipe.fit(X_train, y_train)`
1. Call the score() function on the pipeline object or predict() function
 - a. `pipe.score(X_test, y_test)`

In the step 2b, the pipeline object is created using a dictionary of key:value pairs. The key is specified in strings for e.g. “scaler” followed by the function to be called.

The key is the name given to a step.

Pipeline_simple.ipynb

Build a pipeline (Contd...)

1. The pipeline object requires all the stages included to have a “transform()” function except for the last stage which is an estimator.
1. The transform step transforms the input data. The transformed output of a stage is the input to the next stage
1. During the call “pipeline.fit()”, the pipeline calls the fit and transform functions on each step in sequence. For the last step, only the fit function is called
1. While predicting using pipeline, similarly transform function in all the stages followed by a predict function in the last stage is performed
1. The pipeline object does not need to have a predict function. It only needs to have a fit function at least

Build a pipeline (Contd...)

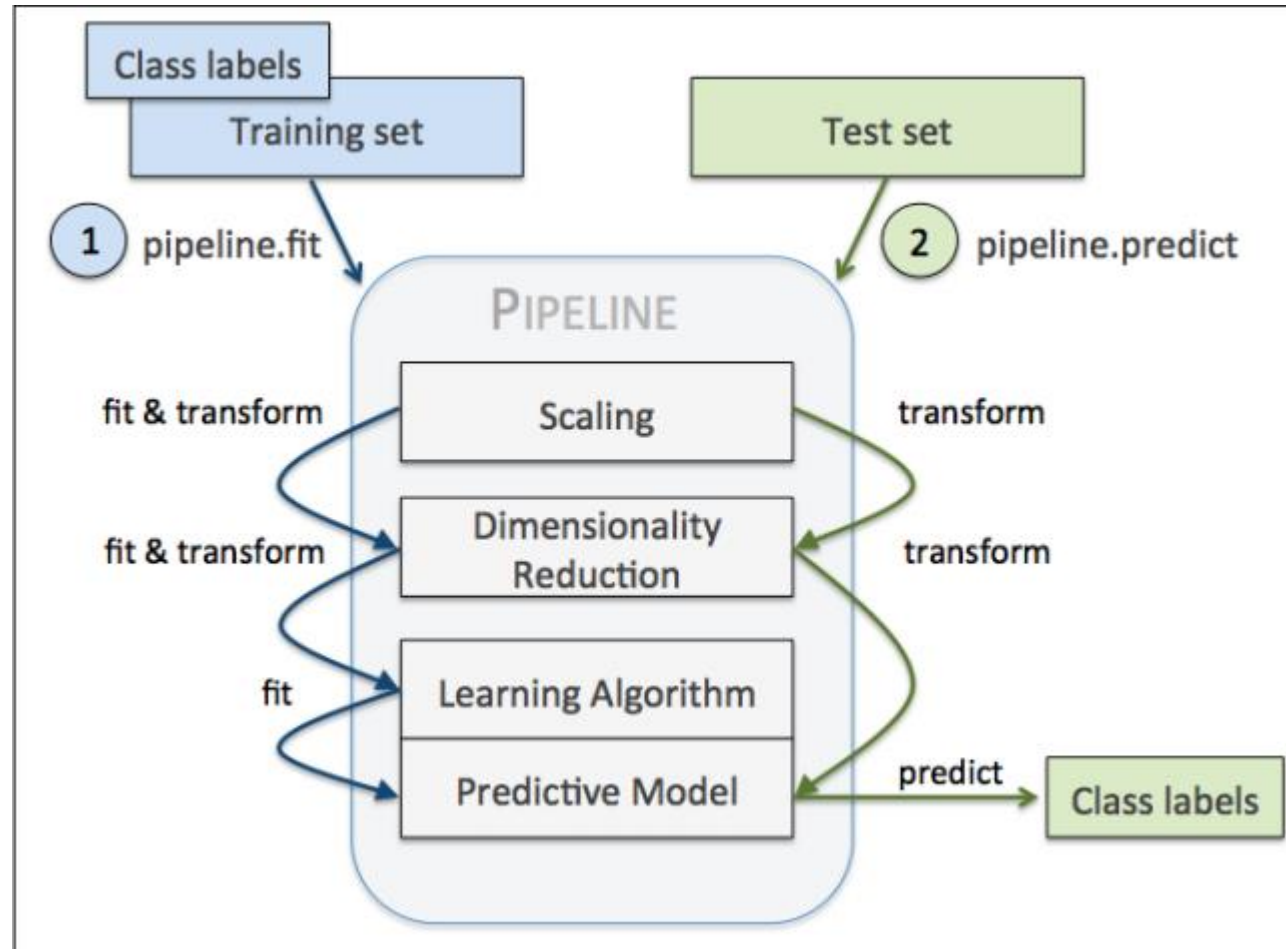


Image Source : Python Deeper Insights Into Machine Learning

make_pipeline

1. Creating the pipeline could be cumbersome. Specifying a name to each stage may not be necessary
1. Alternatively there is a “make_pipeline()” function that will create the pipeline and automatically name each step. We do not need to specify a name
 - a. `from sklearn.pipeline import make_pipeline`
 - b. `pipe = make_pipeline(MinMaxScaler(), (SVC()))`
 - c. `print(" Pipeline steps:\ n{}". format(pipe.steps))`
1. Note, we have not specified any name to the stages. The names will be automatically assigned and are usually lowercase of the class names

pipeline_introduction.ipynb

make_pipeline

1. Download Wisconsin Breast Cancer dataset from UCI
2. Create dataframe
3. Split into X and y
4. Transform y into numerical using labelencoder
5. Split data into training testing
6. Standardize dataset using standard scalar
7. Use PCA to reduce to 2 dimensions
8. Use logistic regression on the reduced dimensions
9. Create a pipeline
10. Do a fit and score

Pipeline_GridSearchcv_wisc_bc_data_classwork.ipynb

HyperParameter Tuning

Hyper Parameters & Tuning

1. Hyper parameters are like handles available to the modeler to control the behavior of the algorithm used for modeling
1. Hyper parameters are supplied as arguments to the model algorithms while initializing them. For e.g. setting the criterion for decision tree building
“dt_model = DecisionTreeClassifier(criterion = 'entropy')”
1. To get a list of hyper parameters for a given algorithm, call the function `get_params()`...for e.g. to get support vector classifier hyper parameters
 1. `from sklearn.svm import SVC`
 2. `svc= SVC()`
 3. `svc.get_params()`
2. Hyper parameters are not learnt from the data as other model parameters are. For e.g. attribute coefficients in a linear model are learnt from data while cost of error is input as hyper parameter.



model_paramete
s

Hyper Parameters & Tuning

5. Fine tuning the hyper parameters is done in a sequence of steps

1. Selecting the appropriate model type (regressor or classifier such as `sklearn.svm.SVC()`)
2. Identify the corresponding parameter space
3. Decide the method for searching or sampling parameter space;
4. Decide the cross-validation scheme to ensure model will generalize
5. Decide a score function to use to evaluate the model

5. Two generic approaches to searching hyper parameter space include

1. GridSearchCV which exhaustively considers all parameter combinations
2. RandomizedSearchCV can sample a given number of candidates from a parameter space with a specified distribution.

Ref: GridSearchSimpleExample.ipynb

5. While tuning hyper parameters, the data should have been split into three parts – Training, validation and testing to **prevent data leak**

5. The testing data should be separately transformed * using the same functions that were used to transform the rest of the data for model building and hyper parameter tuning

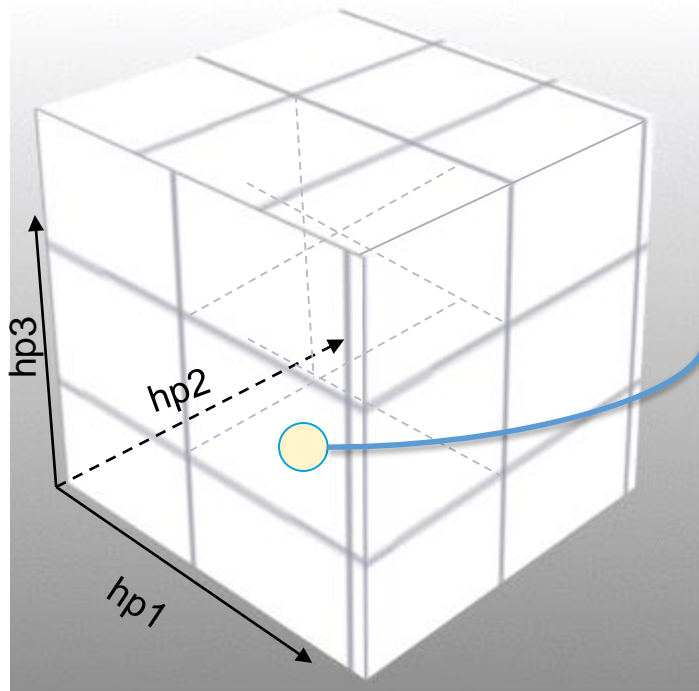
* Any transformation where rows influence each other. For e.g. using zscore. OneHotCode transformation does not come into this category. It can be done before splitting the data

Hyper Parameters & Tuning (GridsearchCV/ RandomizedSearchCv)

GridsearchCV –

1. Is at basic optimal hyperparameter tuning technique.
2. It builds a model for each permutation of all of the given hyperparameter values
3. Each such model is evaluated and ranked.
4. The combination of hyperparameter values that gives the best performing model is chosen
5. For every combination, cross validation is used and average score is calculated
6. This is an exhaustive sampling of the hyperparameter space and can be quite inefficient

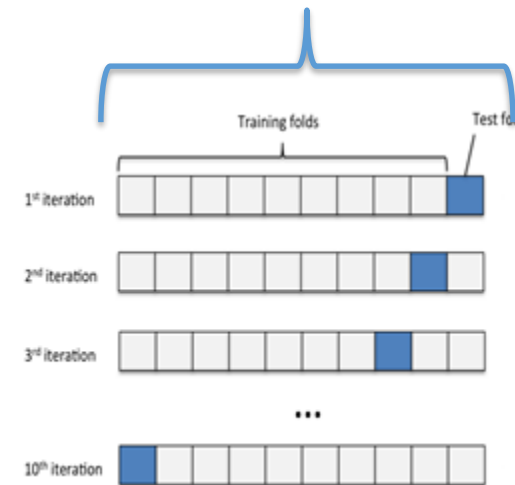
GridSearchCV



Hyper parameter space

One combination of hyper parameters used K times to train and test. The avg score of the K times is the score associated with this combination

This will repeat for all possible combinations i.e. all the cells in the space.



Design GridSearch with Pipeline

1. Use the earlier wisc_bc_dataset
2. Standardize dataset using standard scalar
3. Use PCA
4. Use SVC
5. Create a pipeline for scaling, pca and svc
6. Create param grid for PCA and SVC stage separately
7. Use Gridsearchcv with 5 CV
8. Print the best parameters
9. Print the best score

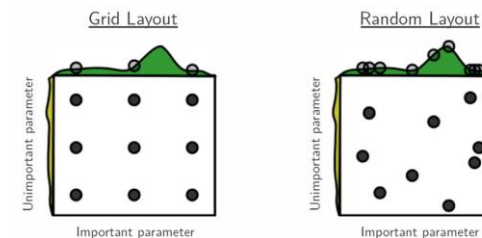
Pipeline_GridSearchcv_wisc_bc_data_classwork.ipynb

Hyper Parameters & Tuning (GridsearchCV/ RandomizedSearchCv)

RandomizedSearchCV –

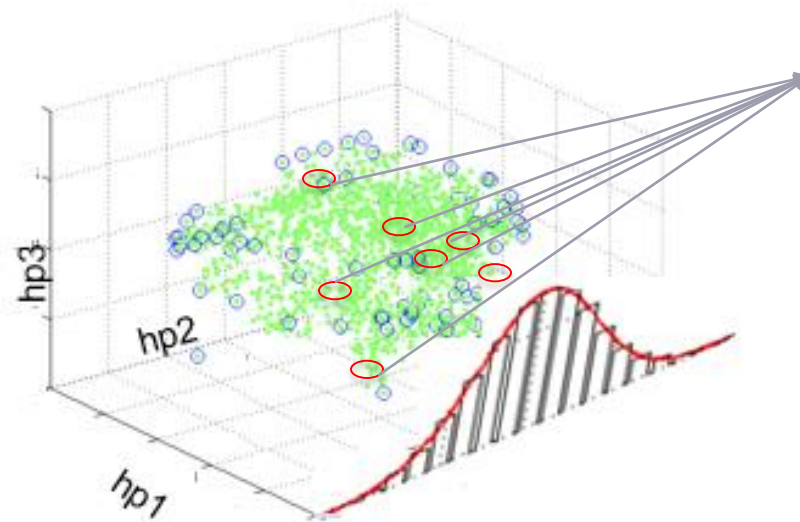
1. Random search differs from grid search. Instead of providing a discrete set of values to explore on each hyperparameter (parameter grid), we provide a statistical distribution.
1. Values for the different hyper parameters are picked up at random from this combine distribution
1. The motivation to use random search in place of grid search is that for many cases, hyperparameters are not *equally* important.

*A Gaussian process analysis of the function from hyper-parameters to validation set performance reveals that **for most data sets only a few of the hyper-parameters really matter, but that different hyper-parameters are important on different data sets.** This phenomenon makes grid search a poor choice for configuring algorithms for new data sets. - [Bergstra, 2012](#)*



bergstra12a.pdf

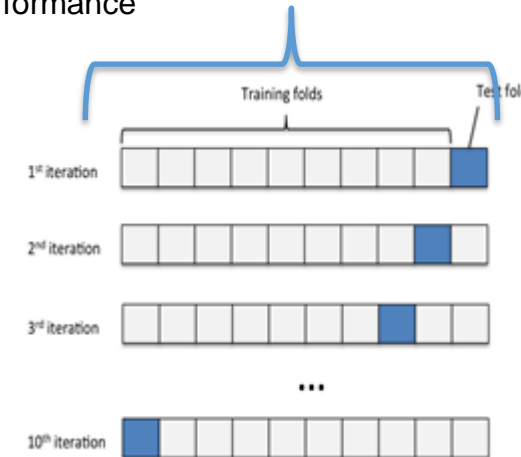
Picture by [Bergstra, 2012](#)



Hyper parameter space

RandomizedSearchCV

Randomly pick up n-iter samples from the hyper parameter distribution as sample, Use it K times and find avg performance



Ref: RandomizedSearchCV_GridSearchCV.ipynb

4. In contrast to GridSearchCV, not all combinations are evaluated. A fixed number of parameter settings is sampled from the specified distributions.
5. The number of parameter settings that are tried is given by n_iter
6. If all parameters are presented as a list, sampling without replacement is performed. If at least one parameter is given as a distribution, sampling with replacement is used. It is highly recommended to use continuous distributions for continuous parameters
7. Randomsearch has higher chance of hitting the right combination than gridsearch.

ThankYou