chapter 6
learning best practices for model
evaluation and hyperparameter tuning
전설
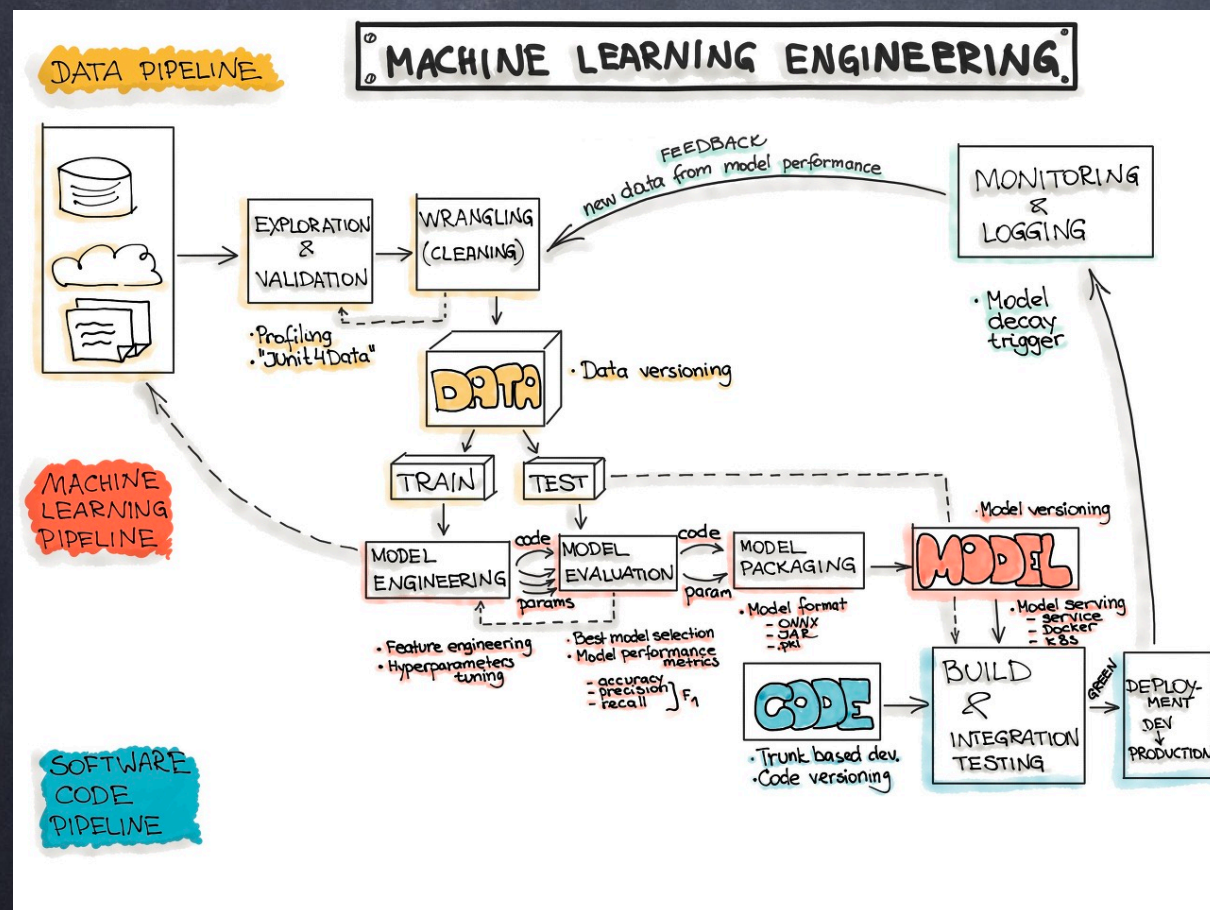
# 1. streamlining workflows with pipelines
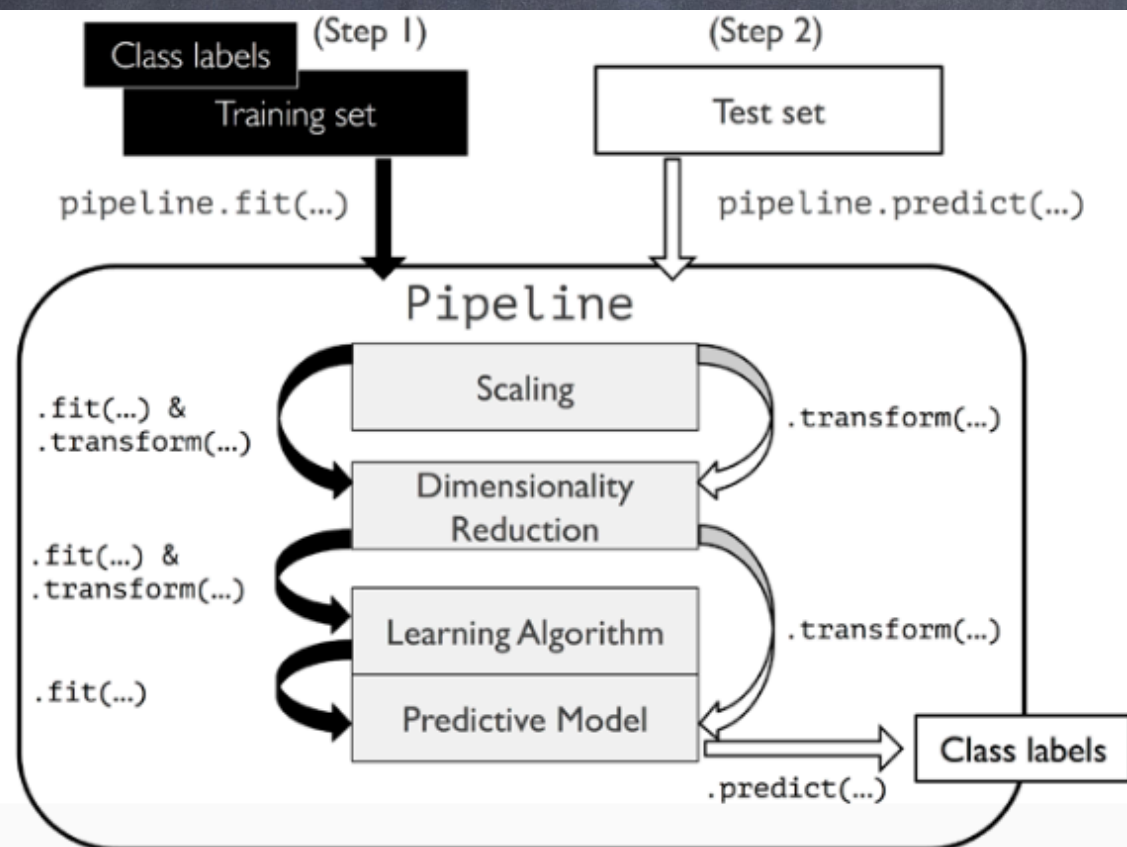
# pipelines

# sklearn.pipeline



```
>>> import pandas as pd
>>> df = pd.read_csv('https://archive.ics.uci.edu/ml/'
...                   'machine-learning-databases'
...                   '/breast-cancer-wisconsin/wdbc.data',
...                   header=None)
```

```
>>> X = df.loc[:, 2:].values
>>> y = df.loc[:, 1].values
>>> le = LabelEncoder()
>>> y = le.fit_transform(y)
>>> le.classes_
array(['B', 'M'], dtype=object)
```

```
>>> le.transform(['M', 'B'])
array([1, 0])
```

```
>>> from sklearn.model_selection import train_test_split

>>> X_train, X_test, y_train, y_test = \
...     train_test_split(X, y,
...                      test_size=0.20,
...                      stratify=y,
...                      random_state=1)
```
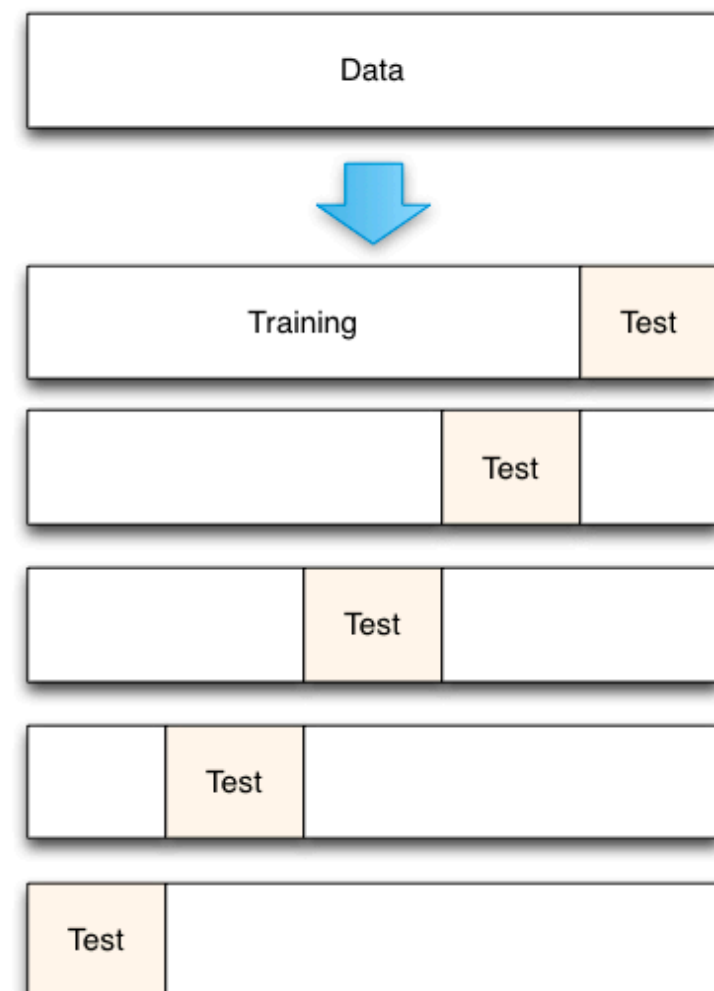
# 2. Using k-fold cross-validation to assess model performance

# K-fold cross-validation

## K-Fold

- Train 데이터셋을 K개로 나눈다
- K개 중 한 개를 valid, 나머지를 training 용으로 사용하여 학습
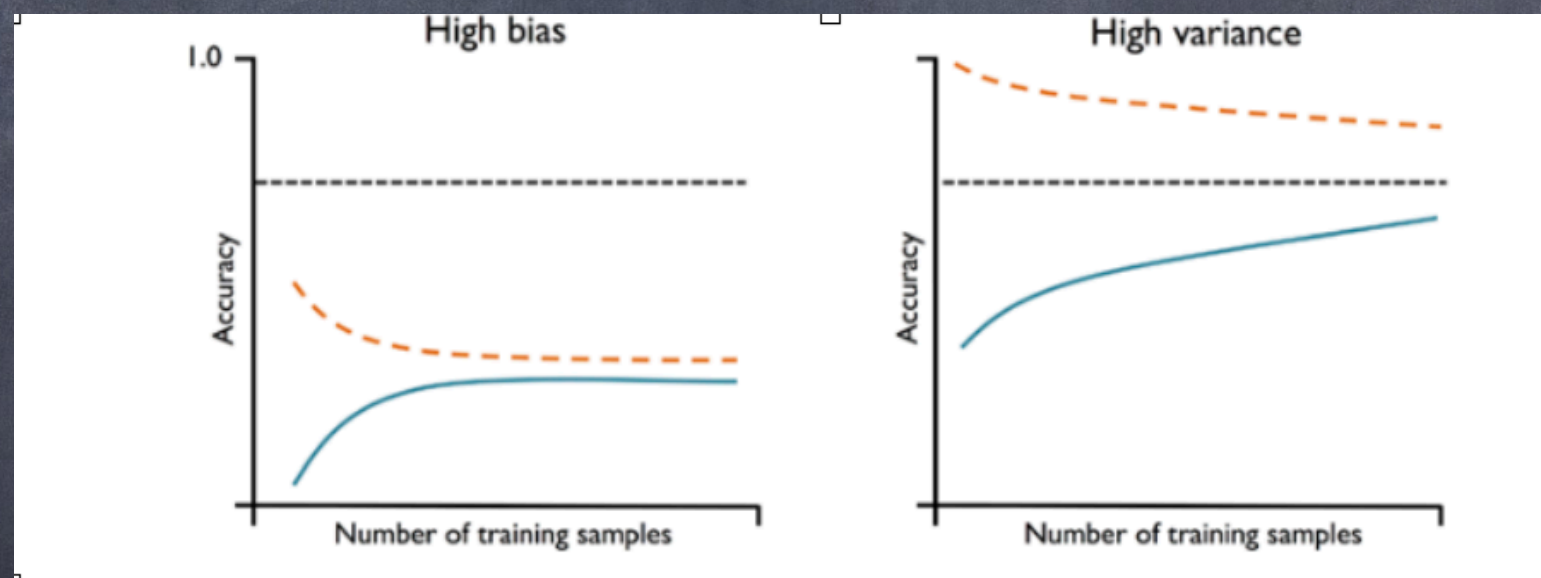- K개의 모델의 Hyper-parameter의 평균을 최종 결과로 사용한다

# 3. debugging algorithms with learning and validation curves

# Diagnosing bias and variance problems with learning curves