

# Contents

- [illegible]

1/44

# Imports and LO

## Imports

```
import os
import sys

sys.path.append(os.path.join(os.path.abspath(".."), "code"))

import matplotlib.pyplot as plt
import mglearn
import numpy as np
import pandas as pd
from plotting_functions import *
from sklearn.dummy import DummyClassifier
from sklearn.impute import SimpleImputer
from sklearn.model_selection import cross_val_score, cross_validate, train_test_split
from sklearn.pipeline import Pipeline, make_pipeline
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from utils import *

%matplotlib inline
pd.set_option("display.max_colwidth", 200)
DATA_DIR = os.path.join(os.path.abspath(".."), "data/")
```

```
from sklearn import set_config

set_config(display="diagram")
```

## Learning outcomes

From this lecture, you will be able to

- explain the need for hyperparameter optimization
- carry out hyperparameter optimization using `sklearn`'s `GridSearchCV` and `RandomizedSearchCV`

- explain different hyperparameters of `GridSearchCV`
- explain the importance of selecting a good range for the values.
- explain optimization bias
- identify and reason when to trust and not trust reported accuracies

# Hyperparameter optimization motivation

## Motivation

- Remember that the fundamental goal of supervised machine learning is to generalize beyond what we see in the training examples.
- We have been using data splitting and cross-validation to provide a framework to approximate generalization error.
- With this framework, we can improve the model's generalization performance by tuning model hyperparameters using cross-validation on the training set.

## Hyperparameters: the problem

- In order to improve the generalization performance, finding the best values for the important hyperparameters of a model is necessary for almost all models and datasets.
- Picking good hyperparameters is important because if we don't do it, we might end up with an underfit or overfit model.

## Some ways to pick hyperparameters:

- Manual or expert knowledge or heuristics based optimization
- Data-driven or automated optimization

## Manual hyperparameter optimization

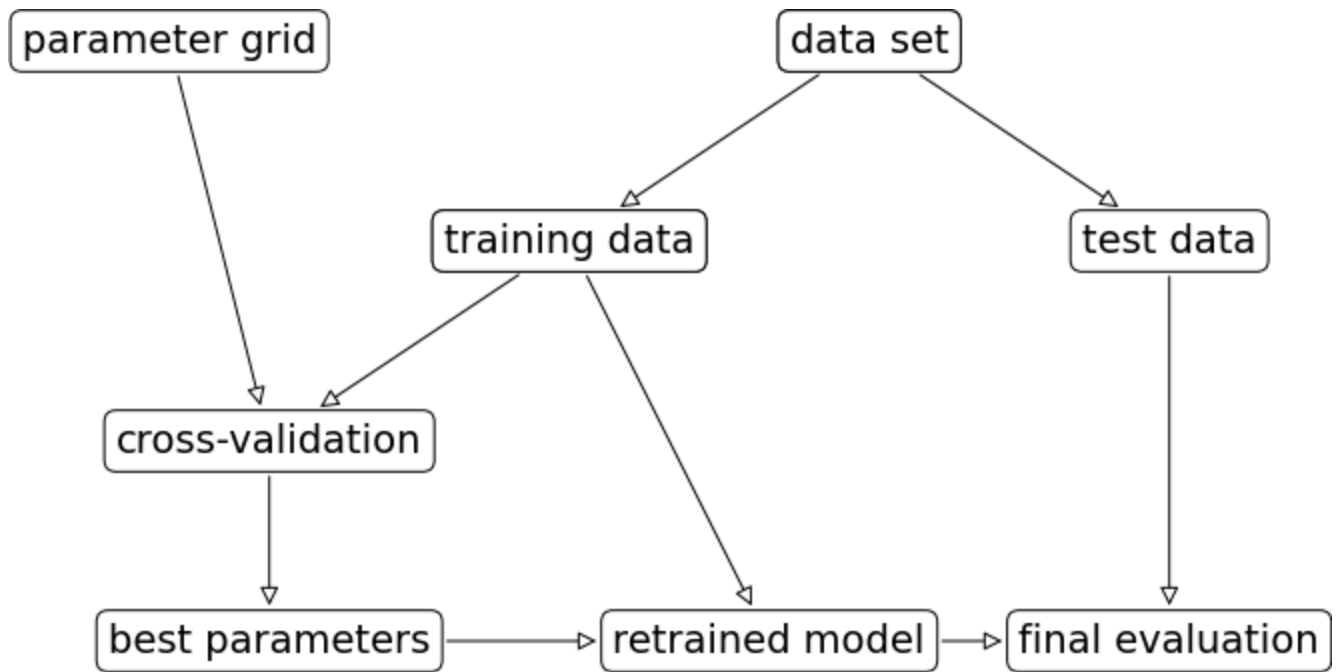
- Advantage: we may have some intuition about what might work.
  - E.g. if I'm massively overfitting, try decreasing `max_depth` or `C`.
- Disadvantages
  - it takes a lot of work
  - not reproducible
  - in very complicated cases, our intuition might be worse than a data-driven approach

## Automated hyperparameter optimization

- Formulate the hyperparameter optimization as a one big search problem.
- Often we have many hyperparameters of different types: Categorical, integer, and continuous.
- Often, the search space is quite big and systematic search for optimal values is infeasible.

In homework assignments, we have been carrying out hyperparameter search by exhaustively trying different possible combinations of the hyperparameters of interest.

```
mglearn.plots.plot_grid_search_overview()
```



Let's look at an example of tuning `max_depth` of the `DecisionTreeClassifier` on the Spotify dataset.

```

spotify_df = pd.read_csv(DATA_DIR + "spotify.csv", index_col=0)
X_spotify = spotify_df.drop(columns=["target", "artist"])
y_spotify = spotify_df["target"]
X_spotify.head()

```

	acousticness	danceability	duration_ms	energy	instrumentalness	key	live
0	0.0102	0.833	204600	0.434	0.021900	2	0.
1	0.1990	0.743	326933	0.359	0.006110	1	0.
2	0.0344	0.838	185707	0.412	0.000234	2	0.
3	0.6040	0.494	199413	0.338	0.510000	5	0.
4	0.1800	0.678	392893	0.561	0.512000	5	0.

```
X_train, X_test, y_train, y_test = train_test_split(
    X_spotify, y_spotify, test_size=0.2, random_state=123
)
```

```
numeric_feats = ['acousticness', 'danceability', 'energy',
                 'instrumentalness', 'liveness', 'loudness',
                 'speechiness', 'tempo', 'valence']
categorical_feats = ['time_signature', 'key']
passthrough_feats = ['mode']
text_feat = "song_title"
```

```
from sklearn.compose import make_column_transformer
from sklearn.feature_extraction.text import CountVectorizer

preprocessor = make_column_transformer(
    (StandardScaler(), numeric_feats),
    ("passthrough", passthrough_feats),
    (OneHotEncoder(handle_unknown = "ignore"), categorical_feats),
    (CountVectorizer(max_features=100, stop_words="english"), text_feat)
)

svc_pipe = make_pipeline(preprocessor, SVC)
```

```
best_score = 0

param_grid = {"max_depth": np.arange(1, 20, 2)}

results_dict = {"max_depth": [], "mean_cv_score": []}

for depth in param_grid["max_depth"]:
    # for each combination of parameters, train an SVC
    dt_pipe = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=depth))
    scores = cross_val_score(dt_pipe, X_train, y_train) # perform cross-validation
    mean_score = np.mean(scores) # compute mean cross-validation accuracy
    if (
        mean_score > best_score
    ): # if we got a better score, store the score and parameters
        best_score = mean_score
        best_params = {"max_depth": depth}
    results_dict["max_depth"].append(depth)
    results_dict["mean_cv_score"].append(mean_score)
```

```
best_params
```

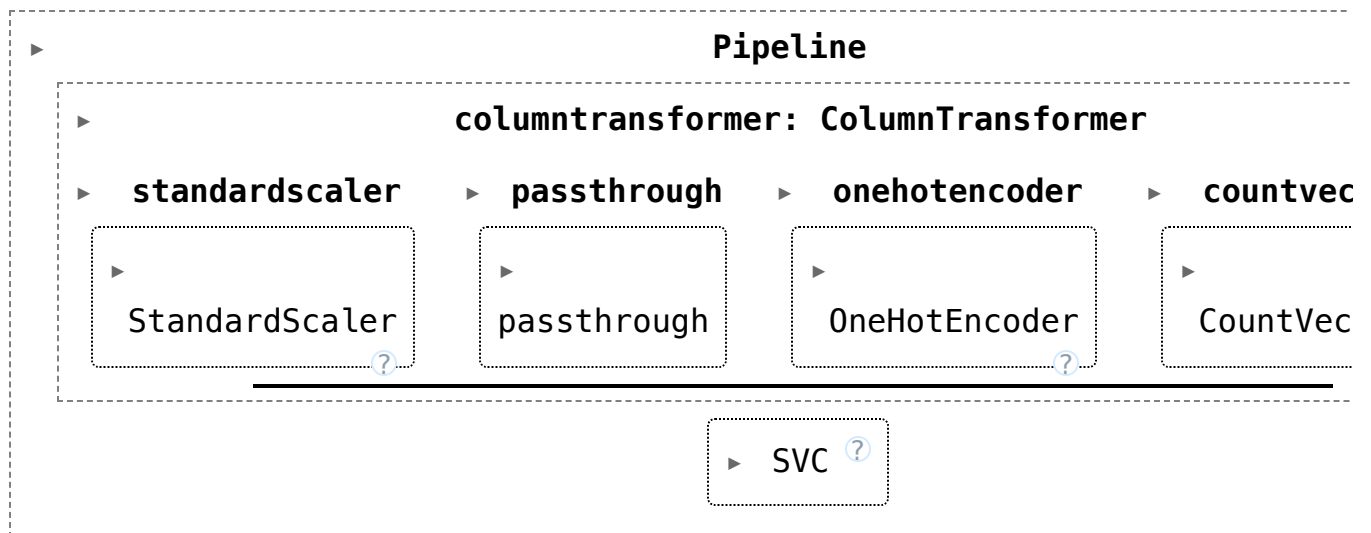
```
{'max_depth': np.int64(5)}
```

```
best_score
```

```
np.float64(0.7303155587177662)
```

Let's try SVM RBF and tuning `C` and `gamma` on the same dataset.

```
pipe_svm = make_pipeline(preprocessor, SVC()) # We need scaling for SVM RBF
pipe_svm.fit(X_train, y_train)
```



Let's try cross-validation with default hyperparameters of SVC.

```
scores = cross_validate(pipe_svm, X_train, y_train, return_train_score=True)
pd.DataFrame(scores).mean()
```

```
fit_time      0.039974
score_time    0.009185
test_score    0.734011
train_score   0.828891
dtype: float64
```

Now let's try exhaustive hyperparameter search using for loops.

This is what we have been doing for this:

```

for gamma in [0.01, 1, 10, 100]: # for some values of gamma
    for C in [0.01, 1, 10, 100]: # for some values of C
        for fold in folds:
            fit in training portion with the given C
            score on validation portion
            compute average score

pick hyperparameter values which yield with best average score

```

```

best_score = 0

param_grid = {
    "C": [0.001, 0.01, 0.1, 1, 10, 100],
    "gamma": [0.001, 0.01, 0.1, 1, 10, 100],
}

results_dict = {"C": [], "gamma": [], "mean_cv_score": []}

for gamma in param_grid["gamma"]:
    for C in param_grid["C"]: # for each combination of parameters, train an
        pipe_svm = make_pipeline(preprocessor, SVC(gamma=gamma, C=C))
        scores = cross_val_score(pipe_svm, X_train, y_train) # perform cross-
        mean_score = np.mean(scores) # compute mean cross-validation accuracy
        if (
            mean_score > best_score
        ): # if we got a better score, store the score and parameters
            best_score = mean_score
            best_parameters = {"C": C, "gamma": gamma}
            results_dict["C"].append(C)
            results_dict["gamma"].append(gamma)
            results_dict["mean_cv_score"].append(mean_score)

```

```
best_parameters
```

```
{'C': 1, 'gamma': 0.1}
```

```
best_score
```

```
np.float64(0.7352614272253524)
```

```
df = pd.DataFrame(results_dict)
```

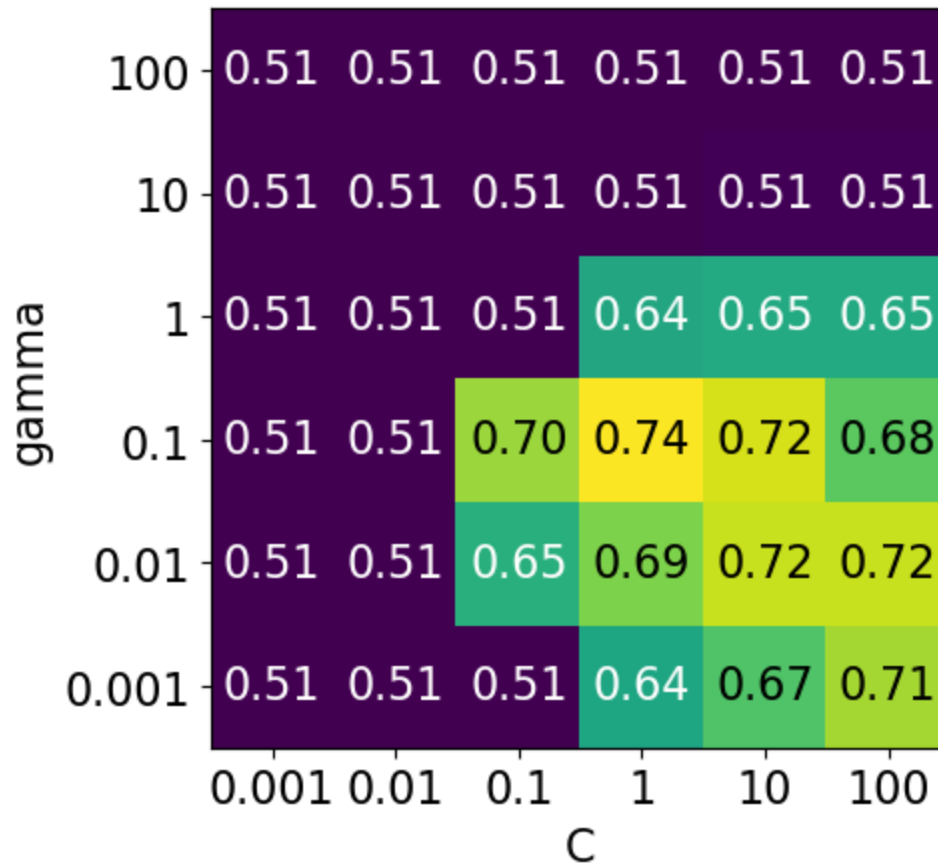


```
df.sort_values(by="mean_cv_score", ascending=False).head(10)
```

	C	gamma	mean_cv_score
15	1.0	0.100	0.735261
16	10.0	0.100	0.722249
11	100.0	0.010	0.716657
10	10.0	0.010	0.716655
5	100.0	0.001	0.705511
14	0.1	0.100	0.701173
9	1.0	0.010	0.691877
17	100.0	0.100	0.677601
4	10.0	0.001	0.673277
8	0.1	0.010	0.652828

```
scores = np.array(df.mean_cv_score).reshape(6, 6)

my_heatmap(
    scores,
    xlabel="C",
    xticklabels=param_grid["C"],
    ylabel="gamma",
    yticklabels=param_grid["gamma"],
    cmap="viridis",
    fmt="%0.2f"
);
# plot the mean cross-validation scores
```



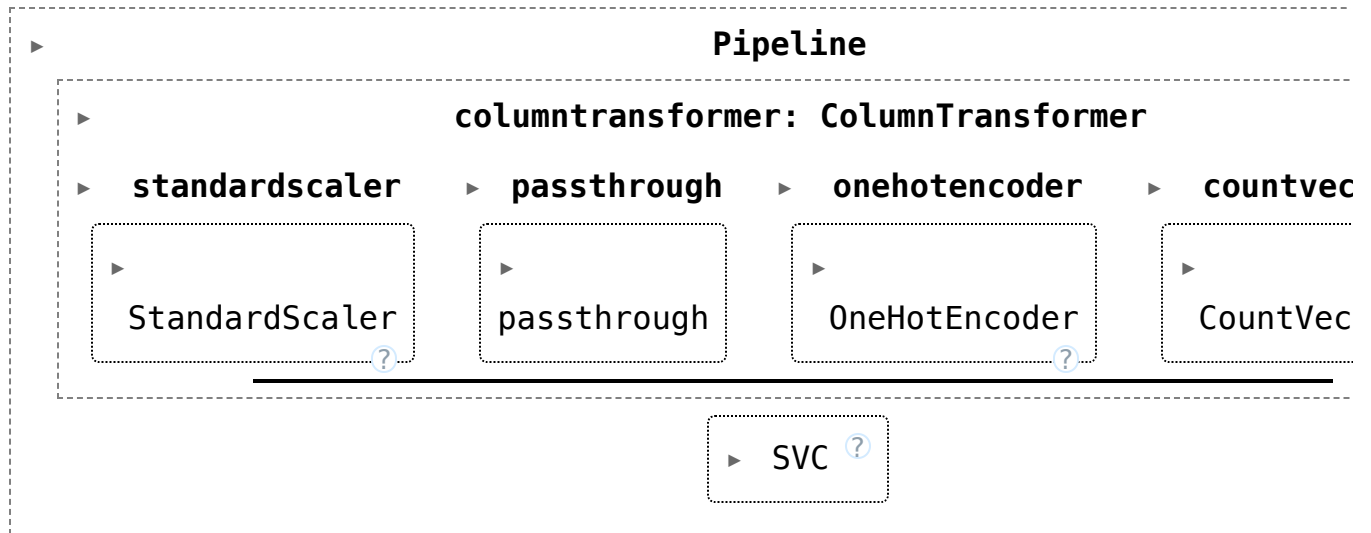
- We have 6 possible values for `C` and 6 possible values for `gamma`.
- In 5-fold cross-validation, for each combination of parameter values, five accuracies are computed.
- So to evaluate the accuracy of the SVM using 6 values of `C` and 6 values of `gamma` using five-fold cross-validation, we need to train  $36 * 5 = 180$  models!

```
np.prod(list(map(len, param_grid.values())))
```

```
np.int64(36)
```

Once we have optimized hyperparameters, we retrain a model on the full training set with these optimized hyperparameters.

```
pipe_svm = make_pipeline(preprocessor, SVC(**best_parameters))
pipe_svm.fit(
    X_train, y_train
) # Retrain a model with optimized hyperparameters on the combined training a
```



### Note

In Python, the double asterisk (`**`) followed by a variable name is used to pass a variable number of keyword arguments to a function. This allows to pass a dictionary of named arguments to a function, where keys of the dictionary become the argument names and values become the corresponding argument values.

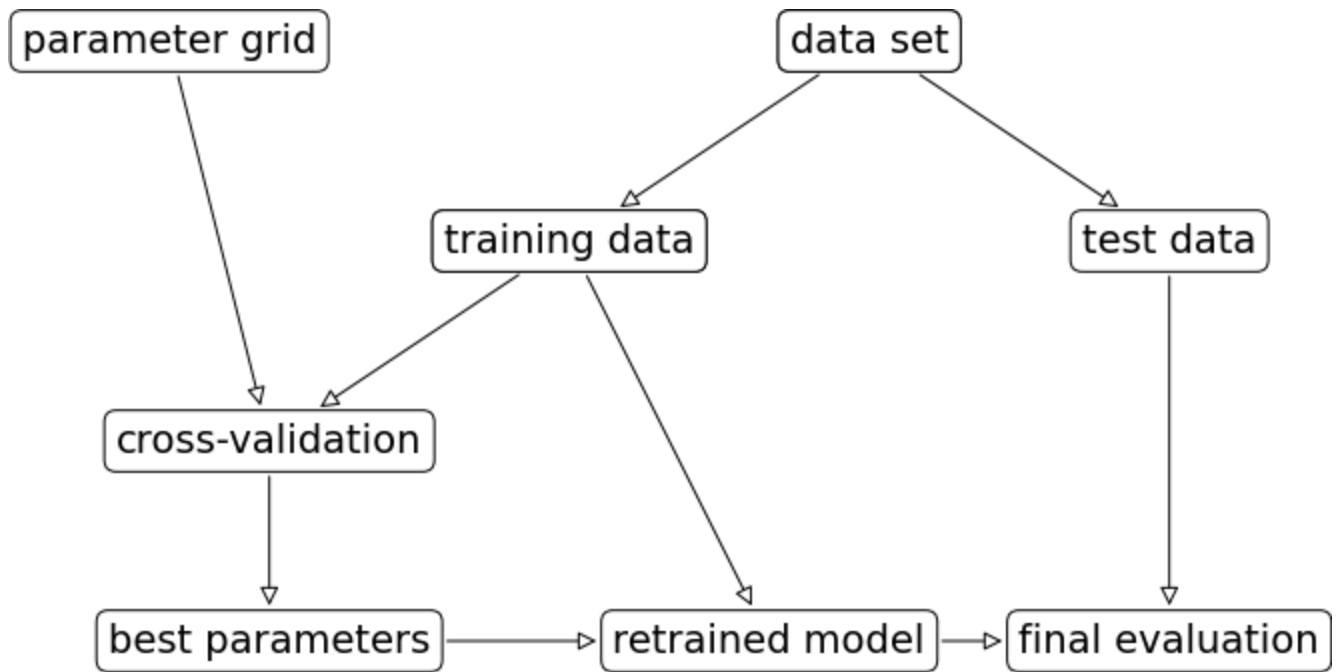
And finally evaluate the performance of this model on the test set.

```
pipe_svm.score(X_test, y_test) # Final evaluation on the test data
```

0.75

This process is so common that there are some standard methods in `scikit-learn` where we can carry out all of this in a more compact way.

```
mglearn.plots.plot_grid_search_overview()
```



In this lecture we are going to talk about two such most commonly used automated optimizations methods from `scikit-learn`.

- Exhaustive grid search: `sklearn.model_selection.GridSearchCV`
- Randomized search: `sklearn.model_selection.RandomizedSearchCV`

The "CV" stands for cross-validation; these methods have built-in cross-validation.

## Exhaustive grid search:

`sklearn.model_selection.GridSearchCV`

- For `GridSearchCV` we need
  - an instantiated model or a pipeline
  - a parameter grid: A user specifies a set of values for each hyperparameter.

- other optional arguments

The method considers product of the sets and evaluates each combination one by one.

```
from sklearn.model_selection import GridSearchCV

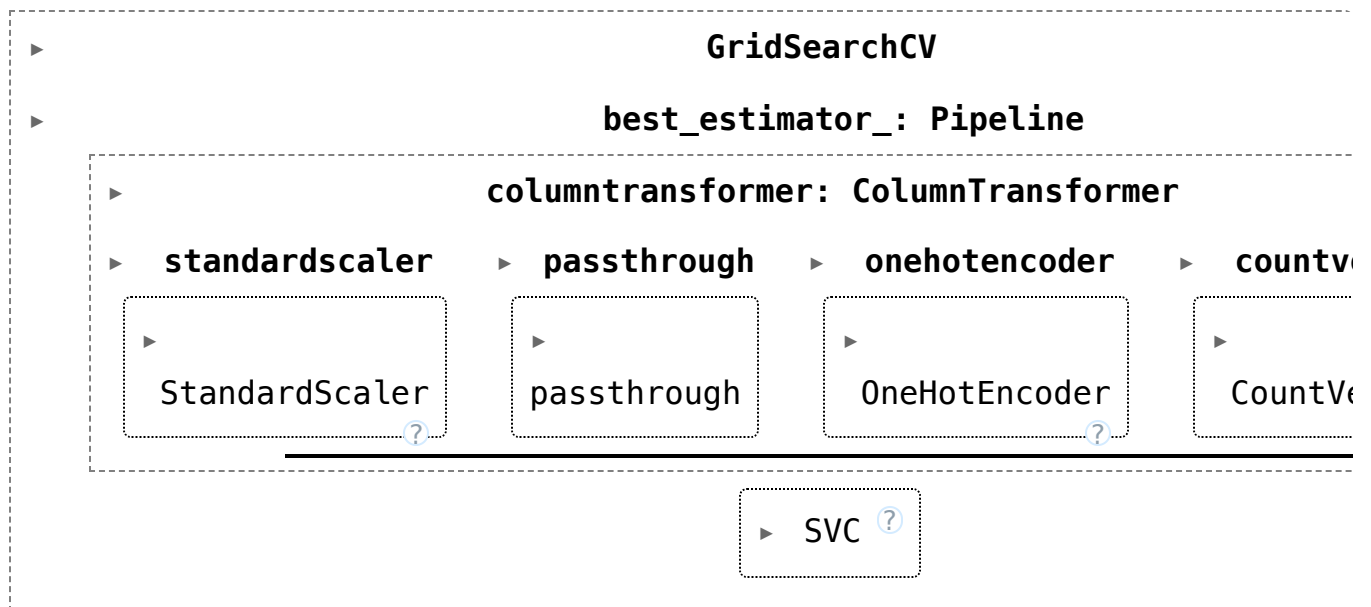
pipe_svm = make_pipeline(preprocessor, SVC())

param_grid = {
    "columntransformer__countvectorizer__max_features": [100, 200, 400, 800, 1000],
    "svc__gamma": [0.001, 0.01, 0.1, 1.0, 10, 100],
    "svc__C": [0.001, 0.01, 0.1, 1.0, 10, 100],
}

# Create a grid search object
gs = GridSearchCV(pipe_svm, param_grid=param_grid, n_jobs=-1, return_train_score=True)
```

The `GridSearchCV` object above behaves like a classifier. We can call `fit`, `predict` or `score` on it.

```
# Carry out the search
gs.fit(X_train, y_train)
```



Fitting the `GridSearchCV` object

- Searches for the best hyperparameter values
- You can access the best score and the best hyperparameters using `best_score_` and `best_params_` attributes, respectively.

```
# Get the best score  
gs.best_score_
```

```
np.float64(0.7395977155164125)
```

```
# Get the best hyperparameter values  
gs.best_params_
```

```
{'columntransformer__countvectorizer__max_features': 1000,  
 'svc__C': 1.0,  
 'svc__gamma': 0.1}
```

- It is often helpful to visualize results of all cross-validation experiments.
- You can access this information using `cv_results_` attribute of a fitted `GridSearchCV` object.

```
results = pd.DataFrame(gs.cv_results_)  
results.T
```

23 rows x 216 columns

```
results = (  
    pd.DataFrame(gs.cv_results_).set_index("rank_test_score").sort_index()  
)  
display(results.T)
```

rank_test_score	
mean_fit_time	
std_fit_time	
mean_score_time	
std_score_time	
param_columntransformer__countvectorizer__max_features	
param_svc__C	
param_svc__gamma	
params	{'columntransformer__c 1000, 's
split0_test_score	
split1_test_score	
split2_test_score	
split3_test_score	
split4_test_score	
mean_test_score	
std_test_score	
split0_train_score	
split1_train_score	
split2_train_score	
split3_train_score	
split4_train_score	
mean_train_score	
std_train_score	

22 rows × 216 columns



Let's only look at the most relevant rows.

```
pd.DataFrame(gs.cv_results_) [
    [
        "mean_test_score",
        "param_columntransformer__countvectorizer__max_features",
        "param_svc__gamma",
        "param_svc__C",
        "mean_fit_time",
        "rank_test_score",
    ]
].set_index("rank_test_score").sort_index().T
```

	rank_test_score	1	
	mean_test_score	0.739598	0.73
param_columntransformer__countvectorizer__max_features		1000.000000	2000.00
	param_svc__gamma	0.100000	0.10
	param_svc__C	1.000000	1.00
	mean_fit_time	0.076444	0.08

5 rows × 216 columns

- Other than searching for best hyperparameter values, `GridSearchCV` also fits a new model on the whole training set with the parameters that yielded the best results.
- So we can conveniently call `score` on the test set with a fitted `GridSearchCV` object.

```
gs.best_score_
```

```
np.float64(0.7395977155164125)
```

```
# Get the test scores
gs.score(X_test, y_test)
```

```
0.7574257425742574
```

Why are `best_score_` and the score above different?

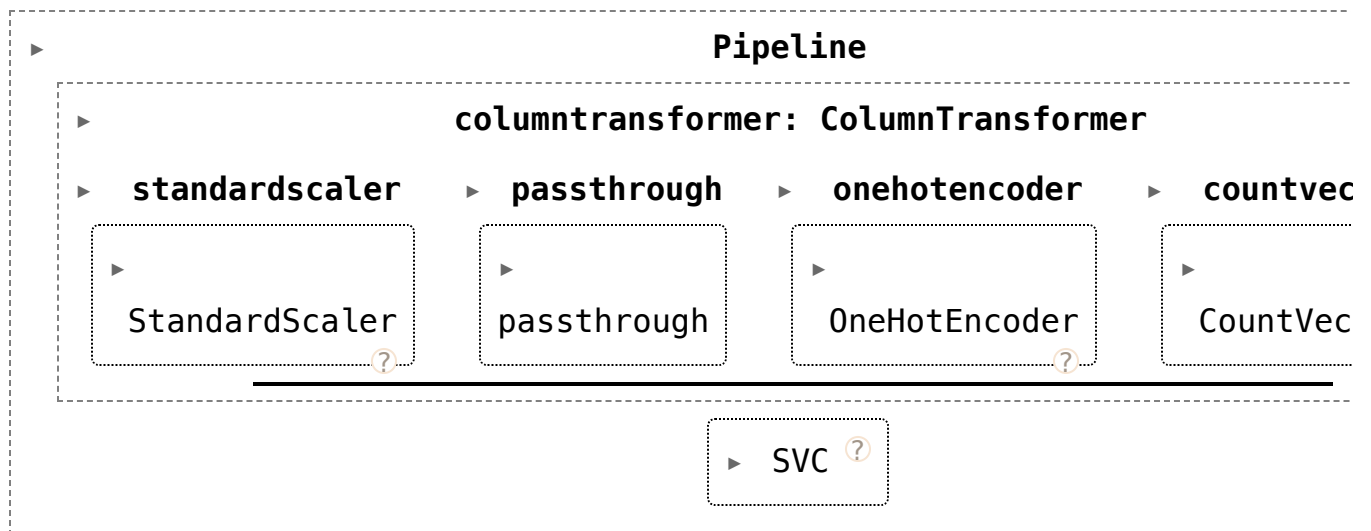
`n_jobs=-1`

- Note the `n_jobs=-1` above.
- Hyperparameter optimization can be done *in parallel* for each of the configurations.
- This is very useful when scaling up to large numbers of machines in the cloud.
- When you set `n_jobs=-1`, it means that you want to use all available CPU cores for the task.

## The `__` syntax

- Above: we have a nesting of transformers.
- We can access the parameters of the “inner” objects by using `__` to go “deeper”:
- `svc__gamma`: the `gamma` of the `svc` of the pipeline
- `svc__C`: the `C` of the `svc` of the pipeline
- `columntransformer__countvectorizer__max_features`: the `max_features` hyperparameter of `CountVectorizer` in the column transformer `preprocessor`.

```
pipe_svm
```



## Range of `C`

- Note the exponential range for `C`. This is quite common. Using this exponential range allows you to explore a wide range of values efficiently.

- There is no point trying  $C = \{1, 2, 3 \dots, 100\}$  because  $C = 1, 2, 3$  are too similar to each other.
- Often we're trying to find an order of magnitude, e.g.  $C = \{0.01, 0.1, 1, 10, 100\}$ .
- We can also write that as  $C = \{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2\}$ .
- Or, in other words,  $C$  values to try are  $10^n$  for  $n = -2, -1, 0, 1, 2$  which is basically what we have above.

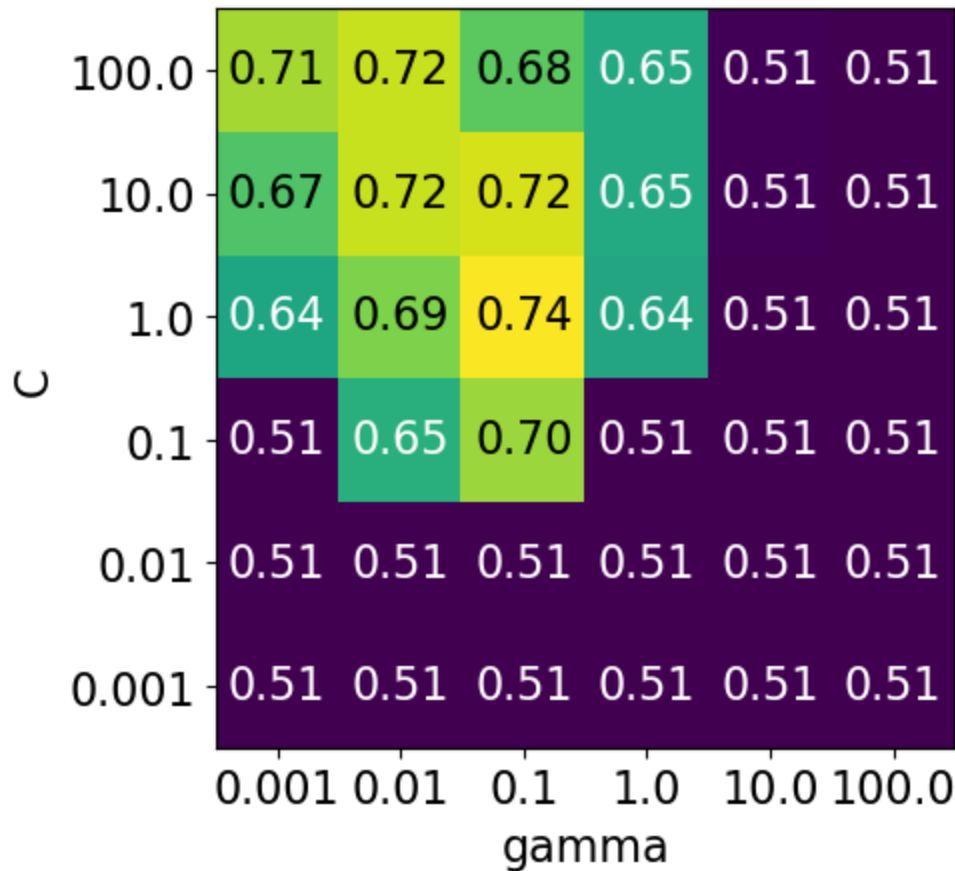
## Visualizing the parameter grid as a heatmap

```
def display_heatmap(param_grid, pipe, X_train, y_train):
    grid_search = GridSearchCV(
        pipe, param_grid, cv=5, n_jobs=-1, return_train_score=True
    )
    grid_search.fit(X_train, y_train)
    results = pd.DataFrame(grid_search.cv_results_)
    scores = np.array(results.mean_test_score).reshape(6, 6)

    # plot the mean cross-validation scores
    my_heatmap(
        scores,
        xlabel="gamma",
        xticklabels=param_grid["svc__gamma"],
        ylabel="C",
        yticklabels=param_grid["svc__C"],
        cmap="viridis",
    );
```

- Note that the range we pick for the parameters play an important role in hyperparameter optimization.
- For example, consider the following grid and the corresponding results.

```
param_grid1 = {
    "svc__gamma": 10.0**np.arange(-3, 3, 1),
    "svc__C": 10.0**np.arange(-3, 3, 1)
}
display_heatmap(param_grid1, pipe_svm, X_train, y_train)
```



- Each point in the heat map corresponds to one run of cross-validation, with a particular setting
- Colour encodes cross-validation accuracy.
  - Lighter colour means high accuracy
  - Darker colour means low accuracy
- SVC is quite sensitive to hyperparameter settings.
- Adjusting hyperparameters can change the accuracy from 0.51 to 0.74!

## Bad range for hyperparameters

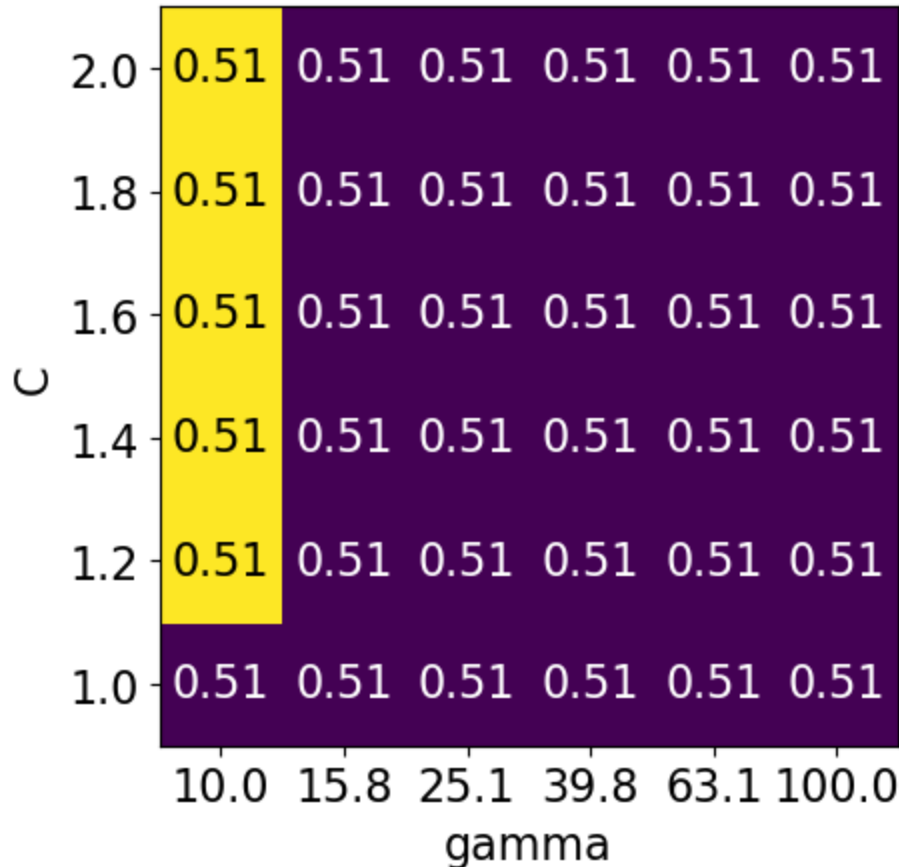
```
np.logspace(1, 2, 6)
```

```
array([ 10.          , 15.84893192, 25.11886432, 39.81071706,
        63.09573445, 100.          ])
```

```
np.linspace(1, 2, 6)
```

```
array([1. , 1.2, 1.4, 1.6, 1.8, 2. ])
```

```
param_grid2 = {"svc__gamma": np.round(np.logspace(1, 2, 6), 1), "svc__C": np.l
display_heatmap(param_grid2, pipe_svm, X_train, y_train)
```



Different range for hyperparameters yields better results!

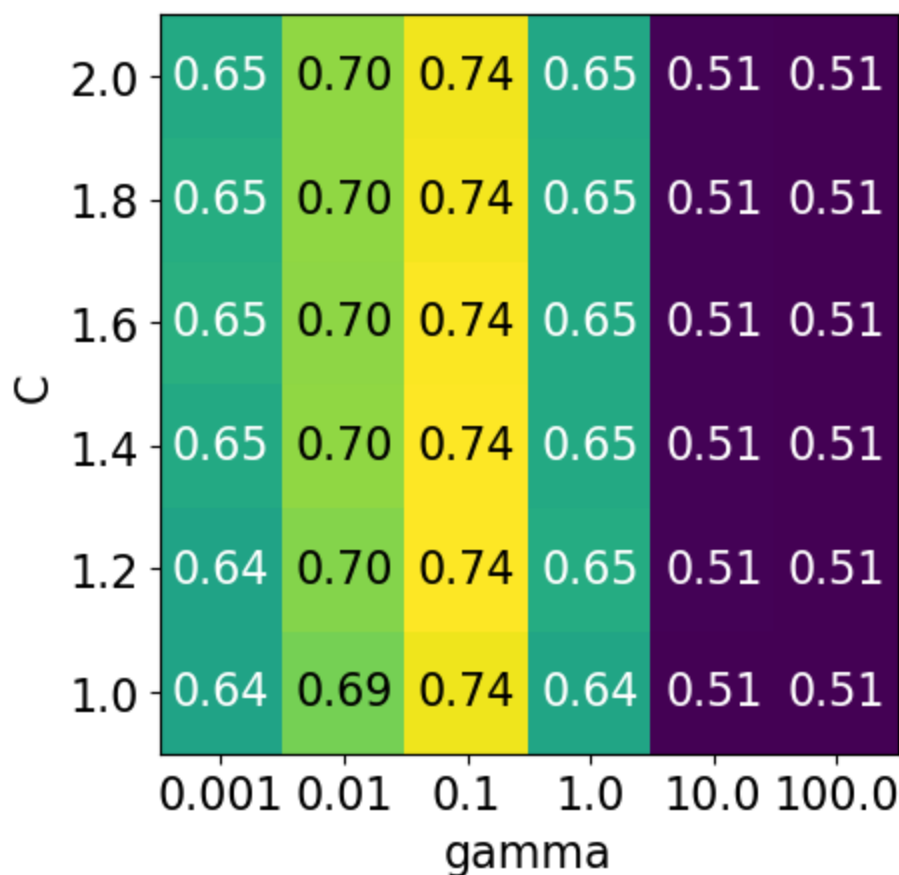
```
np.logspace(-3, 2, 6)
```

```
array([1.e-03, 1.e-02, 1.e-01, 1.e+00, 1.e+01, 1.e+02])
```

```
np.linspace(1, 2, 6)
```

```
array([1. , 1.2, 1.4, 1.6, 1.8, 2. ])
```

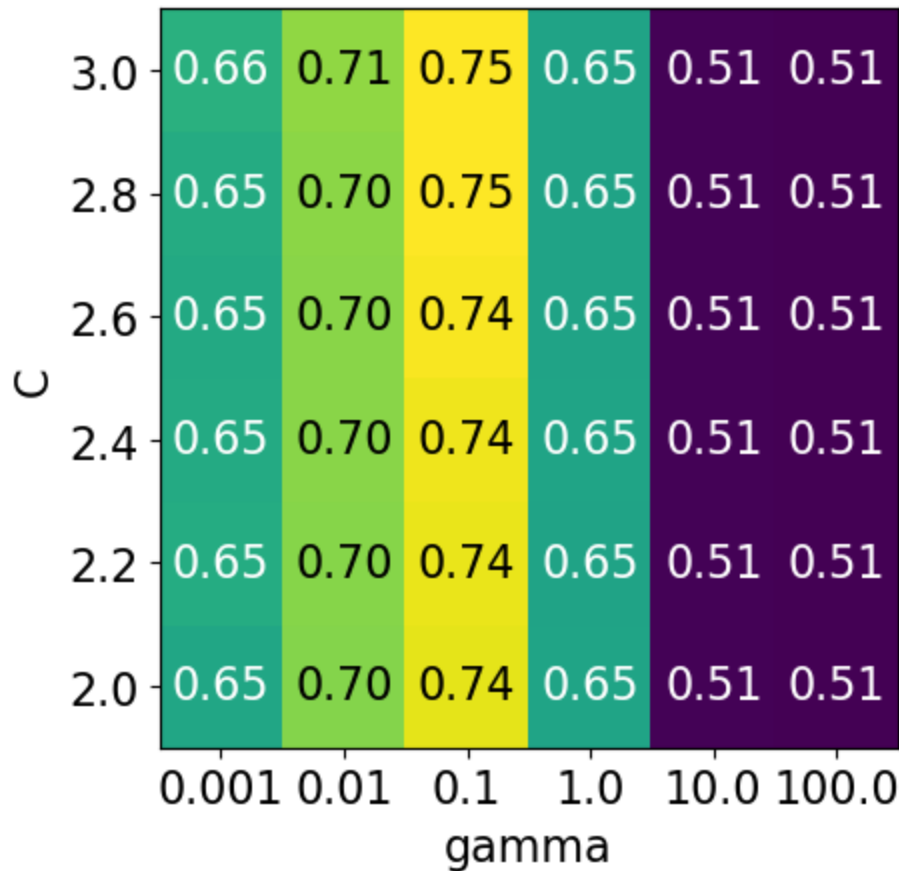
```
param_grid3 = {"svc__gamma": np.logspace(-3, 2, 6), "svc__C": np.linspace(1, 2  
display_heatmap(param_grid3, pipe_svm, X_train, y_train)
```



It seems like we are getting even better cross-validation results with  $C = 2.0$  and  $\gamma = 0.1$

How about exploring different values of  $C$  close to 2.0?

```
param_grid4 = {"svc__gamma": np.logspace(-3, 2, 6), "svc__C": np.linspace(2, 3  
display_heatmap(param_grid4, pipe_svm, X_train, y_train)
```



That's good! We are finding some more options for `C` where the accuracy is 0.75. The tricky part is we do not know in advance what range of hyperparameters might work the best for the given problem, model, and the dataset.

#### **Note**

`GridSearchCV` allows the `param_grid` to be a list of dictionaries. Sometimes some hyperparameters are applicable only for certain models. For example, in the context of `SVC`, `C` and `gamma` are applicable when the kernel is `rbf` whereas only `C` is applicable for `kernel="linear"`.

## Problems with exhaustive grid search

- Required number of models to evaluate grows exponentially with the dimensionality of the configuration space.
- Example: Suppose you have
  - 5 hyperparameters
  - 10 different values for each hyperparameter

- You'll be evaluating  $10^5 = 100,000$  models! That is you'll be calling `cross_validate` 100,000 times!
- Exhaustive search may become infeasible fairly quickly.
- Other options?

## Randomized hyperparameter search

- Randomized hyperparameter optimization
  - `sklearn.model_selection.RandomizedSearchCV`
- Samples configurations at random until certain budget (e.g., time) is exhausted

```
from sklearn.model_selection import RandomizedSearchCV

param_grid = {
    "columntransformer__countvectorizer__max_features": [100, 200, 400, 800, 1000, 2000],
    "svc__gamma": [0.001, 0.01, 0.1, 1.0, 10, 100],
    "svc__C": [0.001, 0.01, 0.1, 1.0, 10, 100]
}

print("Grid size: %d" % (np.prod(list(map(len, param_grid.values()))))
      param_grid
```

Grid size: 216

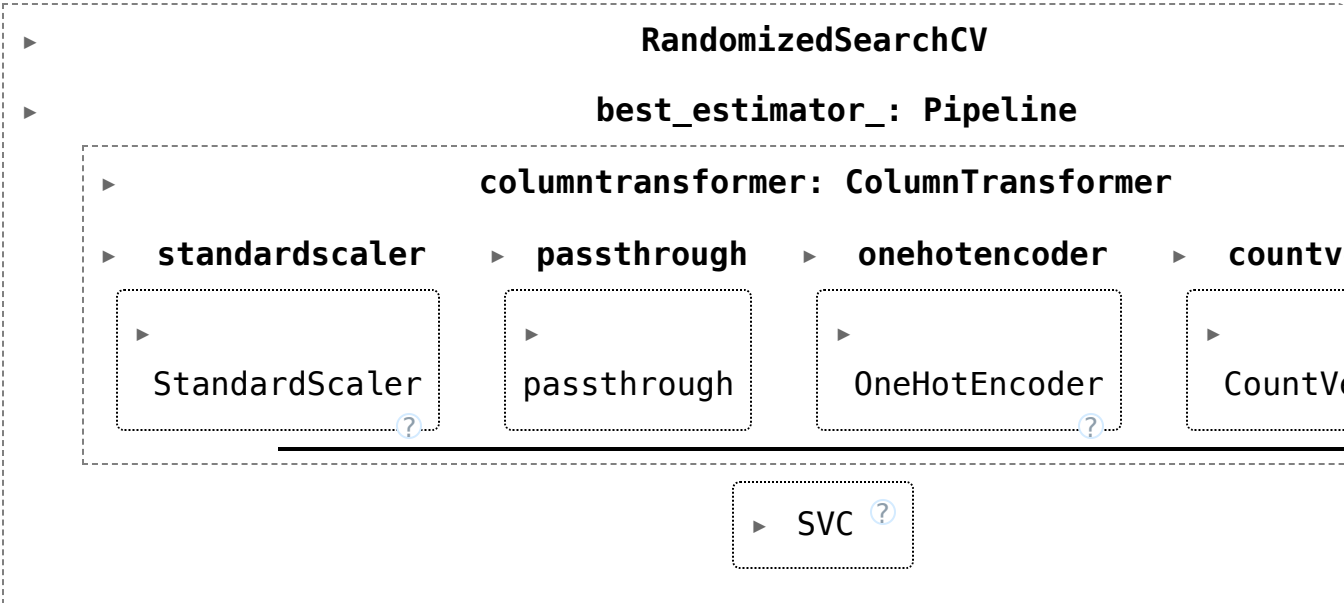
```
{'columntransformer__countvectorizer__max_features': [100,
 200,
 400,
 800,
 1000,
 2000],
'svc__gamma': [0.001, 0.01, 0.1, 1.0, 10, 100],
'svc__C': [0.001, 0.01, 0.1, 1.0, 10, 100]}
```

```
svc_pipe = make_pipeline(preprocessor, SVC())
```



```
# Create a random search object
random_search = RandomizedSearchCV(svc_pipe, param_distributions=param_grid, n
```

```
# Carry out the search
random_search.fit(X_train, y_train)
```



```
pd.DataFrame(random_search.cv_results_)[
    [
        "mean_test_score",
        "param_columntransformer__countvectorizer__max_features",
        "param_svc__gamma",
        "param_svc__C",
        "mean_fit_time",
        "rank_test_score",
    ]
].set_index("rank_test_score").sort_index().T
```

	rank_test_score	1	
	mean_test_score	0.739598	0.73
param_columntransformer__countvectorizer__max_features	1000.000000	2000.00	
	param_svc__gamma	0.100000	0.10
	param_svc__C	1.000000	1.00
	mean_fit_time	0.075229	0.08

5 rows x 100 columns

## n\_iter

- Note the `n_iter`, we didn't need this for `GridSearchCV`.
- Larger `n_iter` will take longer but it'll do more searching.
  - Remember you still need to multiply by number of folds!
- I have set the `random_state` for reproducibility but you don't have to do it.

## Passing probability distributions to random search

Another thing we can do is give probability distributions to draw from:

```
from scipy.stats import expon, lognorm, loguniform, randint, uniform, norm, ra
```

```
np.random.seed(123)
def plot_distribution(y, bins, dist_name="uniform", x_label='Value', y_label="
    plt.hist(y, bins=bins, edgecolor='blue')
    plt.title('Histogram of values from a {} distribution'.format(dist_name))
    plt.xlabel('Value')
    plt.ylabel('Frequency')
    plt.show()
```

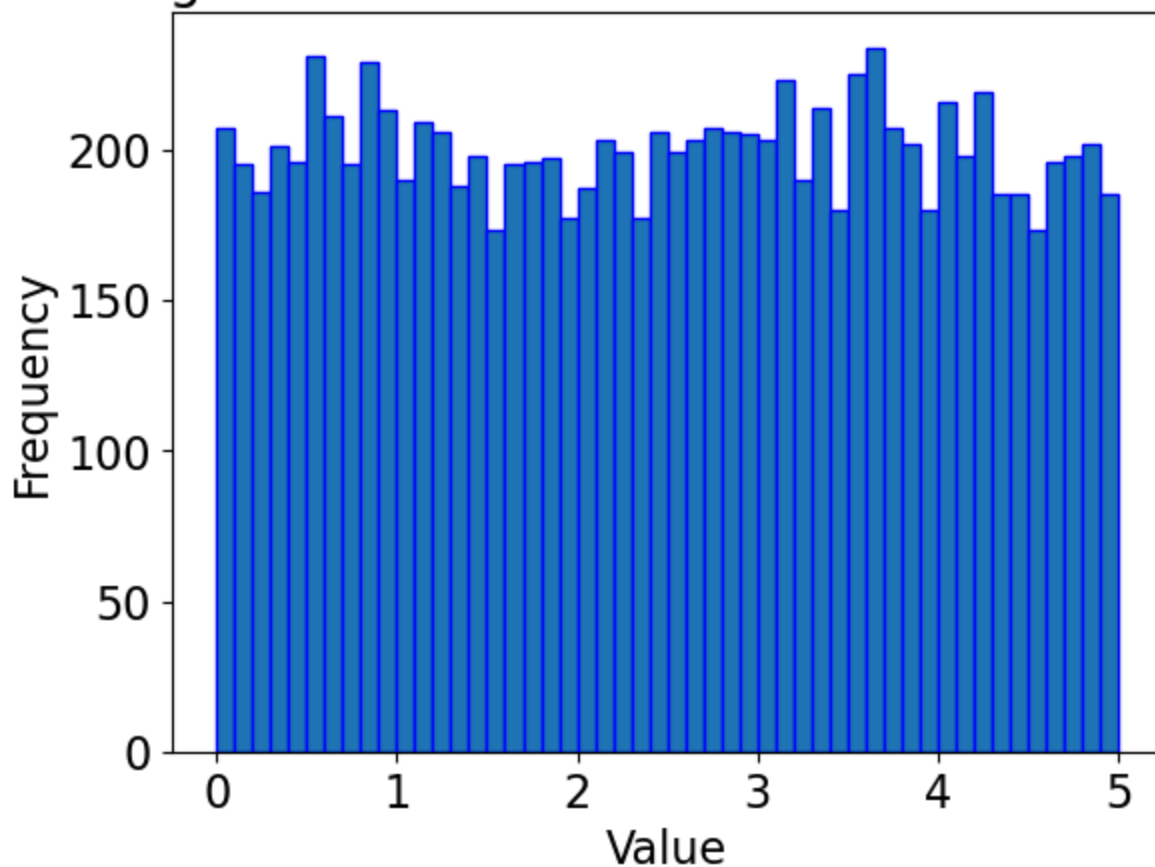
## Uniform distribution

```
# Generate random values from a uniform distribution
y = uniform.rvs(0, 5, size=10000)

# Creating bins within the range of the data
bins = np.arange(0, 5.1, 0.1) # Bins from 0 to 5, in increments of 0.1

plot_distribution(y, bins, "uniform")
```

## Histogram of values from a uniform distribution



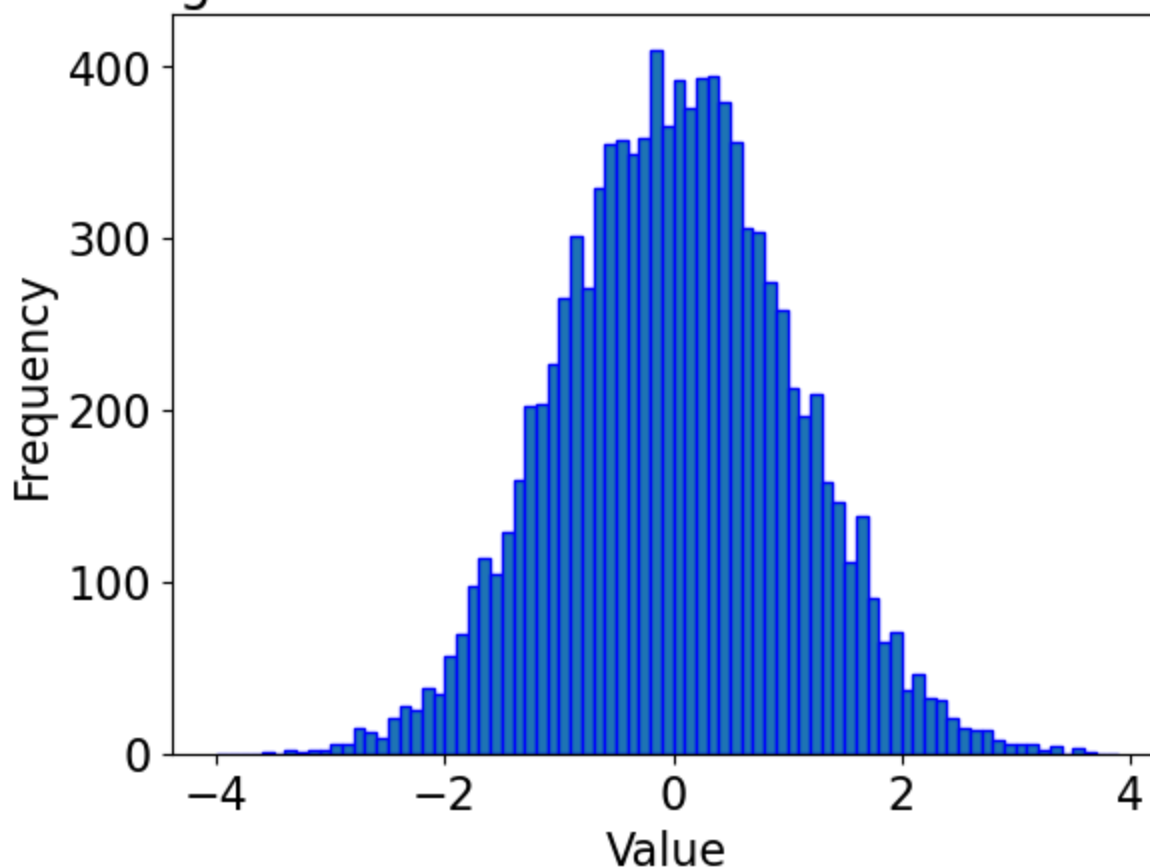
## Gaussian distribution

```
y = norm.rvs(0, 1, 10000)

# Creating bins
bins = np.arange(-4, 4, 0.1)

plot_distribution(y, bins, "normal") # Corrected distribution name
```

## Histogram of values from a normal distribution

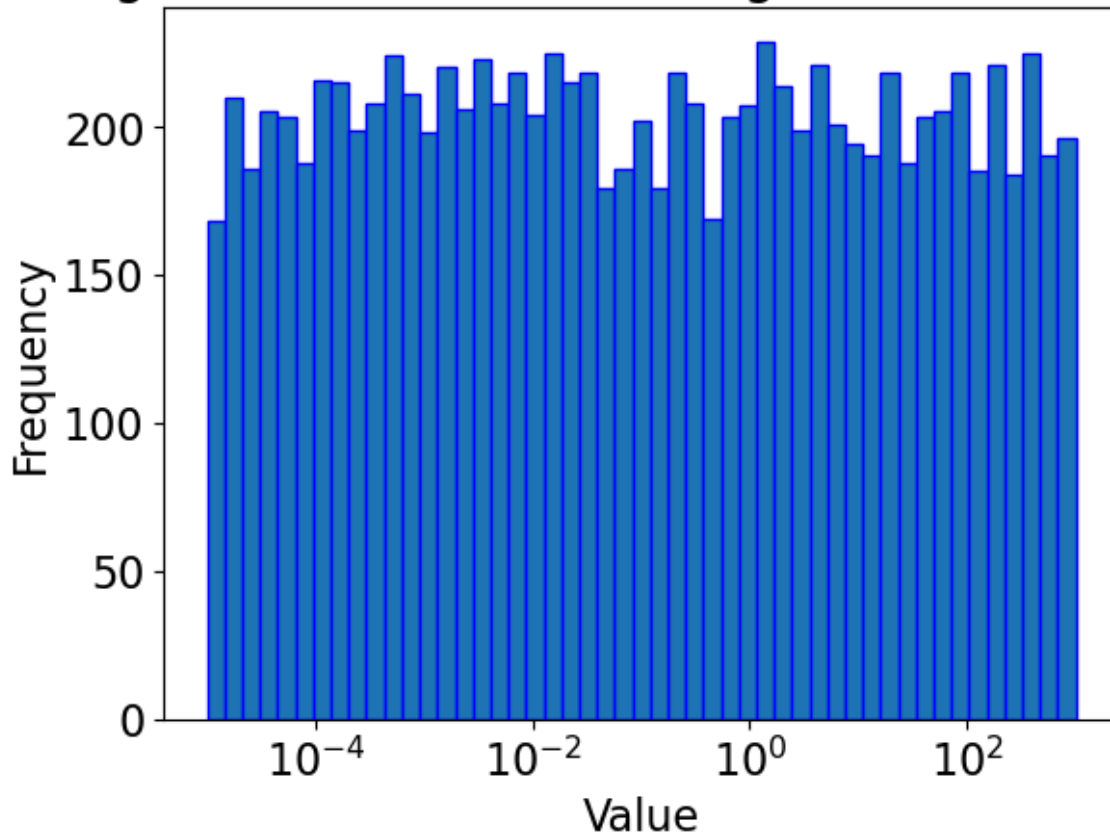


```
# Generate random values from a log-uniform distribution
# Using loguniform from scipy.stats, specifying a range using the exponents
y = loguniform.rvs(1e-5, 1e3, size=10000)

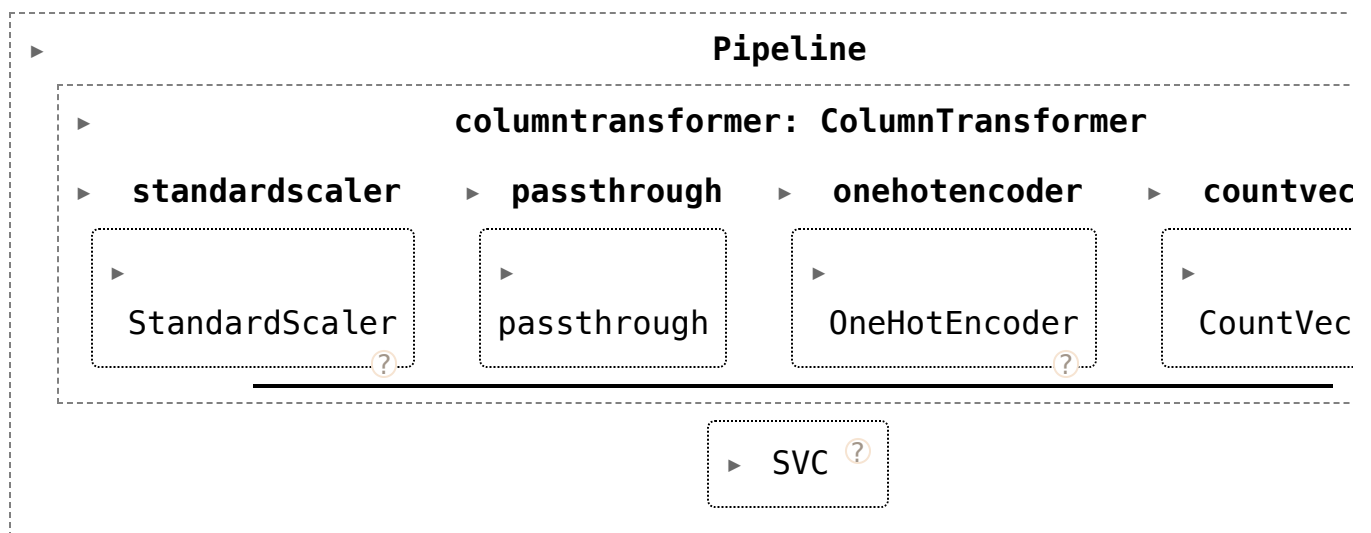
# Creating bins on a logarithmic scale for better visualization
bins = np.logspace(np.log10(1e-5), np.log10(1e3), 50)
plt.xscale('log')

plot_distribution(y, bins, "log uniform") # Corrected distribution name
```

# Histogram of values from a log uniform distribution



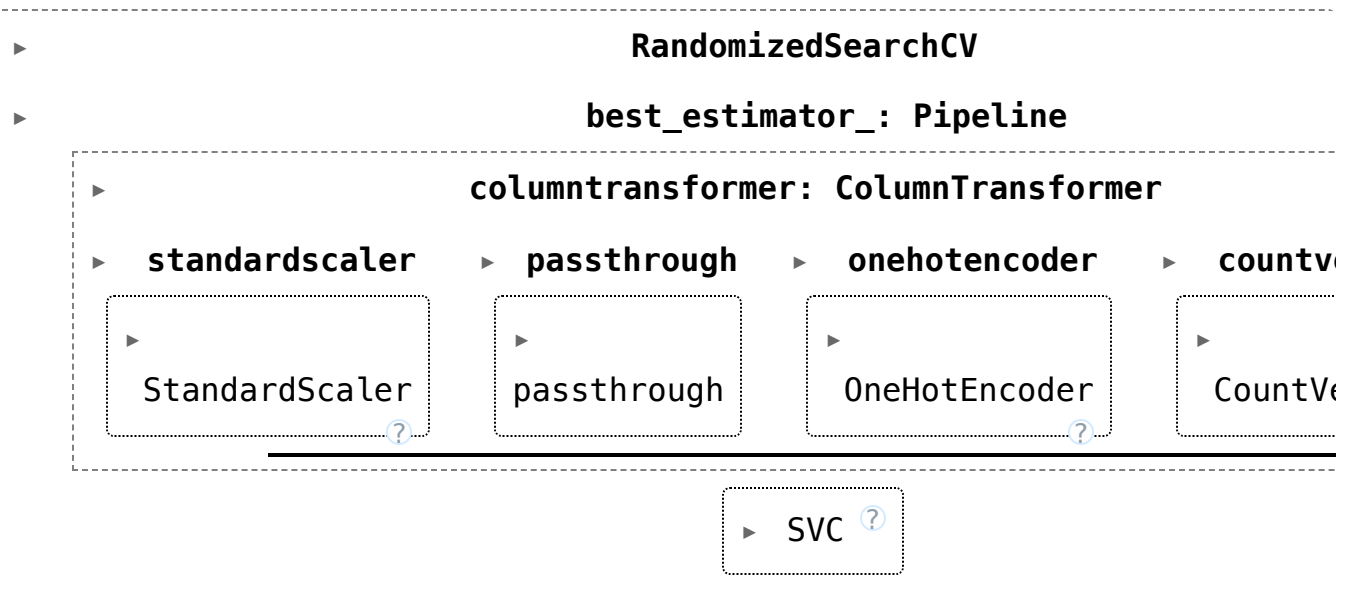
```
pipe_svm
```



```
param_dist = {
    "columntransformer__countvectorizer__max_features": randint(100, 2000),
    "svc__C": uniform(0.1, 1e4), # loguniform(1e-3, 1e3),
    "svc__gamma": loguniform(1e-5, 1e3),
}
```

```
# Create a random search object
random_search = RandomizedSearchCV(pipe_svm,
    param_distributions = param_dist,
    n_iter=100,
    n_jobs=-1,
    return_train_score=True)
```

```
# Carry out the search
random_search.fit(X_train, y_train)
```



```
random_search.best_score_
```

```
np.float64(0.7278214718381631)
```

```
pd.DataFrame(random_search.cv_results_) [
    [
        "mean_test_score",
        "param_columntransformer__countvectorizer__max_features",
        "param_svc__gamma",
        "param_svc__C",
        "mean_fit_time",
        "rank_test_score",
    ]
].set_index("rank_test_score").sort_index().T
```

	rank_test_score	1	
	mean_test_score	0.727821	0.727821
param_columntransformer__countvectorizer__max_features	1619.000000	1014.000000	1014.000000
	param_svc__gamma	0.123539	0.123539
	param_svc__C	7551.390429	477.418000
	mean_fit_time	0.104416	0.104416

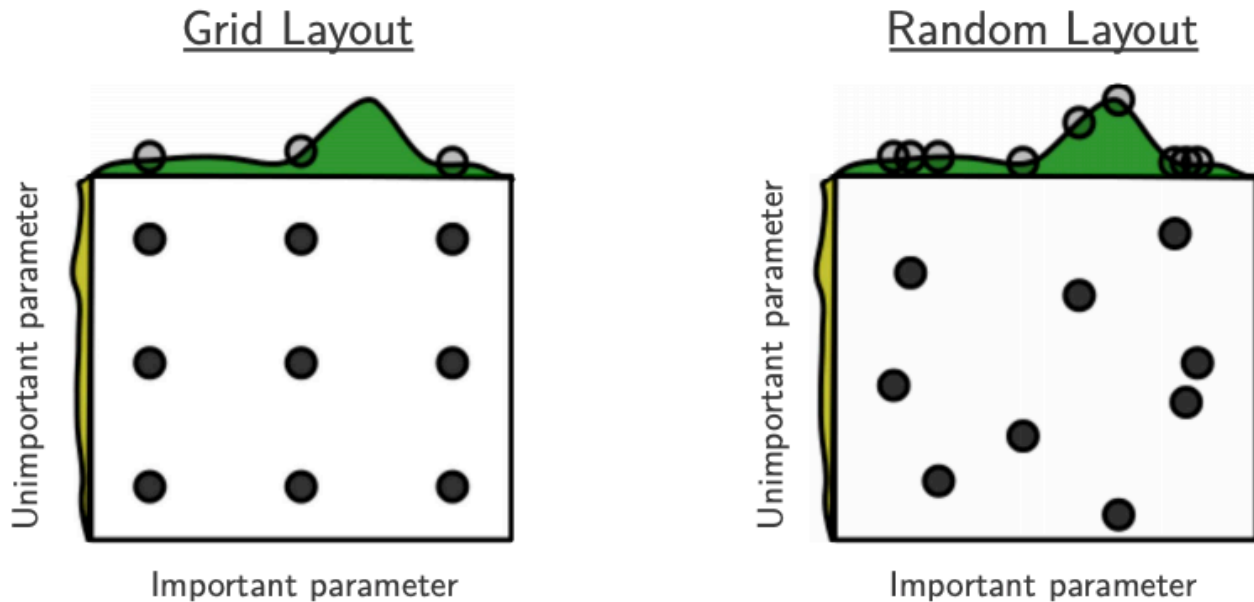
5 rows × 100 columns

- This is a bit fancy. What's nice is that you can have it concentrate more on certain values by setting the distribution.

## Advantages of RandomizedSearchCV

- Faster compared to GridSearchCV.
- Adding parameters that do not influence the performance does not affect efficiency.
- Works better when some parameters are more important than others.
- In general, I recommend using RandomizedSearchCV rather than GridSearchCV.

# Advantages of **RandomizedSearchCV**



Source: [Bergstra and Bengio, Random Search for Hyper-Parameter Optimization, JMLR 2012.](#)

- The yellow on the left shows how your scores are going to change when you vary the unimportant hyperparameter.
- The green on the top shows how your scores are going to change when you vary the important hyperparameter.
- You don't know in advance which hyperparameters are important for your problem.
- In the left figure, 6 of the 9 searches are useless because they are only varying the unimportant parameter.
- In the right figure, all 9 searches are useful.

## (Optional) Searching for optimal parameters with successive halving¶

- Successive halving is an iterative selection process where all candidates (the parameter combinations) are evaluated with a small amount of resources (e.g., small amount of training data) at the first iteration.
- Checkout [successive halving with grid search](#) and [random search](#).

```
from sklearn.experimental import enable_halving_search_cv # noqa
from sklearn.model_selection import HalvingRandomSearchCV
```

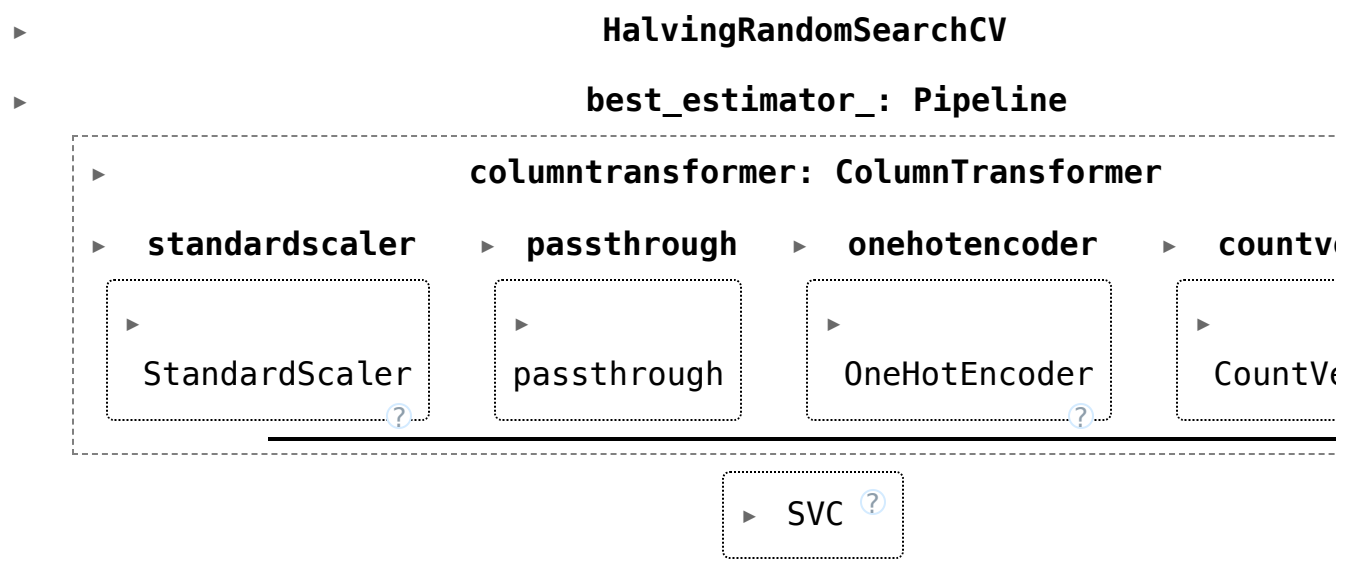


```
rsh = HalvingRandomSearchCV(
    estimator=pipe_svm, param_distributions=param_dist, factor=2, random_state=
)
rsh.fit(X_train, y_train)
```

```
/Users/kvarada/miniforge3/envs/jbook/lib/python3.12/site-packages/numpy/ma/core
_data = np.array(data, dtype=dtype, copy=copy,
```

```
/Users/kvarada/miniforge3/envs/jbook/lib/python3.12/site-packages/numpy/ma/core
_data = np.array(data, dtype=dtype, copy=copy,
```

```
/Users/kvarada/miniforge3/envs/jbook/lib/python3.12/site-packages/numpy/ma/core
_data = np.array(data, dtype=dtype, copy=copy,
```



```
results = pd.DataFrame(rsh.cv_results_)
results["params_str"] = results.params.apply(str)
results.drop_duplicates(subset=("params_str", "iter"), inplace=True)
results
```

	iter	n_resources	mean_fit_time	std_fit_time	mean_score_time	std_scor
<b>0</b>	0	20	0.003582	0.000334	0.002061	0.0
<b>1</b>	0	20	0.003032	0.000144	0.001626	0.0
<b>2</b>	0	20	0.003065	0.000167	0.001645	0.0
<b>3</b>	0	20	0.003134	0.000233	0.001638	0.0
<b>4</b>	0	20	0.003156	0.000148	0.001655	0.0
...	...	...	...	...	...	...
<b>155</b>	5	640	0.015043	0.000216	0.003529	0.0
<b>156</b>	5	640	0.014599	0.000145	0.003731	0.0
<b>157</b>	5	640	0.013772	0.000049	0.003520	0.0
<b>158</b>	6	1280	0.047535	0.000225	0.007740	0.0
<b>159</b>	6	1280	0.046873	0.000354	0.009266	0.0

160 rows × 26 columns

## (Optional) Fancier methods

- Both `GridSearchCV` and `RandomizedSearchCV` do each trial independently.
- What if you could learn from your experience, e.g. learn that `max_depth=3` is bad?

- That could save time because you wouldn't try combinations involving `max_depth=3` in the future.
- We can do this with `scikit-optimize`, which is a completely different package from `scikit-learn`
- It uses a technique called "model-based optimization" and we'll specifically use "Bayesian optimization".
  - In short, it uses machine learning to predict what hyperparameters will be good.
  - Machine learning on machine learning!
- This is an active research area and there are sophisticated packages for this.

Here are some examples

- [hyperopt-sklearn](#)
- [auto-sklearn](#)
- [SigOptSearchCV](#)
- [TPOT](#)
- [hyperopt](#)
- [hyperband](#)
- [SMAC](#)
- [MOE](#)
- [pybo](#)
- [spearmin](#)
- [BayesOpt](#)

## ? ? Questions for you

### (iClicker) Exercise 6.1

Select all of the following statements which are TRUE.

- (A) If you get best results at the edges of your parameter grid, it might be a good idea to adjust the range of values in your parameter grid.
- (B) Grid search is guaranteed to find the best hyperparameter values.

- (C) It is possible to get different hyperparameters in different runs of `RandomizedSearchCV`.



V's Solutions!



## Questions for class discussion (hyperparameter optimization)

- Suppose you have 10 hyperparameters, each with 4 possible values. If you run `GridSearchCV` with this parameter grid, how many cross-validation experiments will be carried out?
- Suppose you have 10 hyperparameters and each takes 4 values. If you run `RandomizedSearchCV` with this parameter grid with `n_iter=20`, how many cross-validation experiments will be carried out?

## Optimization bias/Overfitting of the validation set

### Overfitting of the validation error

- Why do we need to evaluate the model on the test set in the end?
- Why not just use cross-validation on the whole dataset?
- While carrying out hyperparameter optimization, we usually try over many possibilities.
- If our dataset is small and if your validation set is hit too many times, we suffer from **optimization bias** or **overfitting the validation set**.

# Optimization bias of parameter learning

- Overfitting of the training error
- An example:
  - During training, we could search over tons of different decision trees.
  - So we can get “lucky” and find a tree with low training error by chance.

## Optimization bias of hyper-parameter learning

- Overfitting of the validation error
- An example:
  - Here, we might optimize the validation error over 1000 values of `max_depth`.
  - One of the 1000 trees might have low validation error by chance.

## (Optional) Example 1: Optimization bias

Consider a multiple-choice (a,b,c,d) “test” with 10 questions:

- If you choose answers randomly, expected grade is 25% (no bias).
- If you fill out two tests randomly and pick the best, expected grade is 33%.
  - Optimization bias of ~8%.
- If you take the best among 10 random tests, expected grade is ~47%.
- If you take the best among 100, expected grade is ~62%.
- If you take the best among 1000, expected grade is ~73%.
- If you take the best among 10000, expected grade is ~82%.
  - You have so many “chances” that you expect to do well.

**But on new questions the “random choice” accuracy is still 25%.**

```
# (Optional) Code attribution: Rodolfo Lourenzutti
number_tests = [1, 2, 10, 100, 1000, 10000]
for ntests in number_tests:
    y = np.zeros(10000)
    for i in range(10000):
        y[i] = np.max(np.random.binomial(10.0, 0.25, ntests))
    print(
        "The expected grade among the best of %d tests is : %0.2f"
        % (ntests, np.mean(y) / 10.0)
    )
```

The expected grade among the best of 1 tests is : 0.25  
The expected grade among the best of 2 tests is : 0.33  
The expected grade among the best of 10 tests is : 0.47  
The expected grade among the best of 100 tests is : 0.62

The expected grade among the best of 1000 tests is : 0.73

The expected grade among the best of 10000 tests is : 0.83

## (Optional) Example 2: Optimization bias

- If we instead used a 100-question test then:
  - Expected grade from best over 1 randomly-filled test is 25%.
  - Expected grade from best over 2 randomly-filled test is ~27%.
  - Expected grade from best over 10 randomly-filled test is ~32%.
  - Expected grade from best over 100 randomly-filled test is ~36%.
  - Expected grade from best over 1000 randomly-filled test is ~40%.
  - Expected grade from best over 10000 randomly-filled test is ~43%.
- The optimization bias **grows with the number of things we try**.
  - "Complexity" of the set of models we search over.
- But, optimization bias **shrinks quickly with the number of examples**.
  - But it's still non-zero and growing if you over-use your validation set!

```
# (Optional) Code attribution: Rodolfo Lourenzutti
number_tests = [1, 2, 10, 100, 1000, 10000]
for ntests in number_tests:
    y = np.zeros(10000)
    for i in range(10000):
        y[i] = np.max(np.random.binomial(100.0, 0.25, ntests))
    print(
        "The expected grade among the best of %d tests is : %0.2f"
        % (ntests, np.mean(y) / 100.0)
    )
```

The expected grade among the best of 1 tests is : 0.25  
The expected grade among the best of 2 tests is : 0.27  
The expected grade among the best of 10 tests is : 0.32

The expected grade among the best of 100 tests is : 0.36

The expected grade among the best of 1000 tests is : 0.40

The expected grade among the best of 10000 tests is : 0.43

## Optimization bias on the Spotify dataset

```
X_train_tiny, X_test_big, y_train_tiny, y_test_big = train_test_split(
    X_spotify, y_spotify, test_size=0.99, random_state=42
)
```

```
X_train_tiny.shape
```

```
(20, 14)
```

```
X_train_tiny.head()
```

	acousticness	danceability	duration_ms	energy	instrumentalness	key
<b>130</b>	0.055100	0.547	251093	0.643	0.000000	1
<b>1687</b>	0.000353	0.420	210240	0.929	0.000747	7
<b>871</b>	0.314000	0.430	193427	0.734	0.000286	9
<b>1123</b>	0.082100	0.725	246653	0.711	0.000000	10
<b>1396</b>	0.286000	0.616	236960	0.387	0.000000	9

```
pipe = make_pipeline(preprocessor, SVC())
```

```
from sklearn.model_selection import RandomizedSearchCV

param_grid = {
    "svc__gamma": 10.0 ** np.arange(-20, 10),
    "svc__C": 10.0 ** np.arange(-20, 10),
}
print("Grid size: %d" % (np.prod(list(map(len, param_grid.values()))))
param_grid
```

Grid size: 900

```
{'svc__gamma': array([1.e-20, 1.e-19, 1.e-18, 1.e-17, 1.e-16, 1.e-15, 1.e-14, 1.e-13, 1.e-12, 1.e-11, 1.e-10, 1.e-09, 1.e-08, 1.e-07, 1.e-06, 1.e-05, 1.e-04, 1.e-03, 1.e-02, 1.e-01, 1.e+00, 1.e+01, 1.e+02, 1.e+03, 1.e+04, 1.e+05, 1.e+06, 1.e+07, 1.e+08, 1.e+09]),
 'svc__C': array([1.e-20, 1.e-19, 1.e-18, 1.e-17, 1.e-16, 1.e-15, 1.e-14, 1.e-13, 1.e-12, 1.e-11, 1.e-10, 1.e-09, 1.e-08, 1.e-07, 1.e-06, 1.e-05, 1.e-04, 1.e-03, 1.e-02, 1.e-01, 1.e+00, 1.e+01, 1.e+02, 1.e+03, 1.e+04, 1.e+05, 1.e+06, 1.e+07, 1.e+08, 1.e+09])}
```



```
random_search = RandomizedSearchCV(
    pipe, param_distributions=param_grid, n_jobs=-1, n_iter=900, cv=5, random_
)
random_search.fit(X_train_tiny, y_train_tiny);
```

```
pd.DataFrame(random_search.cv_results_) [
    [
        "mean_test_score",
        "param_svc__gamma",
        "param_svc__C",
        "mean_fit_time",
        "rank_test_score",
    ]
].set_index("rank_test_score").sort_index().T
```

rank_test_score	1	1	1	1	1	
mean_test_score	6.500000e-01	0.650000	0.650000	0.650000	0.650000	0
param_svc__gamma	1.000000e-20	0.010000	0.100000	1.000000	10.000000	100
param_svc__C	1.000000e-20	0.010000	0.010000	0.010000	0.010000	C
mean_fit_time	2.063718e-02	0.006592	0.004226	0.005245	0.005446	C

4 rows × 900 columns

Given the results: one might claim that we found a model that performs with 0.8 accuracy on our dataset.

- Do we really believe that 0.80 is a good estimate of our test data?
- Do we really believe that `gamma`=0.0 and C=1\_000\_000\_000 are the best hyperparameters?
- Let's find out the test score with this best model.

```
random_search.score(X_test, y_test)
```

```
0.5024752475247525
```

- The results above are overly optimistic.
  - because our training data is very small and so our validation splits in cross validation would be small.
  - because of the small dataset and the fact that we hit the small validation set 900 times and it's possible that we got lucky on the validation set!
- As we suspected, the best cross-validation score is not a good estimate of our test data; it is overly optimistic.
- We can trust this test score because the test set is of good size.

```
X_test_big.shape
```

```
(1997, 14)
```

## Overfitting of the validation data

The following plot demonstrates what happens during overfitting of the validation data.



[Source](#)

- Thus, not only can we not trust the cv scores, we also cannot trust cv's ability to choose of the best hyperparameters.

## Why do we need a test set?

- This is why we need a test set.
- The frustrating part is that if our dataset is small then our test set is also small 😞.
- But we don't have a lot of better alternatives, unfortunately, if we have a small dataset.

## When test score is much lower than CV score

- What to do if your test score is much lower than your cross-validation score:

- Try simpler models and use the test set a couple of times; it's not the end of the world.
- Communicate this clearly when you report the results.

## Large datasets solve many of these problems

- With infinite amounts of training data, overfitting would not be a problem and you could have your test score = your train score.
  - Overfitting happens because you only see a bit of data and you learn patterns that are overly specific to your sample.
  - If you saw "all" the data, then the notion of "overly specific" would not apply.
- So, more data will make your test score better and robust.

## ? ? Questions for you

Attribution: From Mark Schmidt's notes

### Exercise 6.2

Would you trust the model?

- You have a dataset and you give me 1/10th of it. The dataset given to me is rather small and so I split it into 96% train and 4% validation split. I carry out hyperparameter optimization using a single 4% validation split and report validation accuracy of 0.97. Would it classify the rest of the data with similar accuracy?

1. Probably
2. Probably not

# Final comments and summary

## Automated hyperparameter optimization

- Advantages
  - reduce human effort
  - less prone to error and improve reproducibility
  - data-driven approaches may be effective
- Disadvantages
  - may be hard to incorporate intuition
  - be careful about overfitting on the validation set

Often, especially on typical datasets, we get back `scikit-learn`'s default hyperparameter values. This means that the defaults are well chosen by `scikit-learn` developers!

- The problem of finding the best values for the important hyperparameters is tricky because
  - You may have a lot of them (e.g. deep learning).
  - You may have multiple hyperparameters which may interact with each other in unexpected ways.
- The best settings depend on the specific data/problem.

## Optional readings and resources

- [Preventing "overfitting" of cross-validation data](#) by Andrew Ng