

.This course focuses on **implementing tools and techniques** to effectively manage your **modeling resources** and **best serve batch and real-time inference requests**.

In this course, you'll implement **effective search strategies for the best model** that will scale for various serving needs while **constraining model complexity and hardware requirements**. You'll **optimize and manage this compute storage and IO resources** your model needs in production environments during its life cycle.

In this journey, you will continue to use the **TFX library** and rely on tools like **AutoML** for finding the **best suitable model and TensorFlow model analysis to address model fairness, explainability issues, and mitigate bottlenecks**. Along the way, you'll learn about more specialized **scenarios** like **dimensionality reduction** and **pruning** to manage your model resources wisely

exploring **pipelining** and **parallelism** and high-performance **ingestion** to make the most of your computational resources.

A Production ML system must **run nonstop, at the minimum cost** while producing the maximum performance.

## Hyperparameter Tuning

Keras Tuner.

similarities between hyperparameter tuning and neural architecture search.

**Neural architecture search, or NAS:** automating the design of neural networks -> find the optimal architecture.

Types of parameters in ML Models:

### **Trainable parameters:**

- learned by the algorithm during training. (weights bias of a neural network)

### **Hyperparameters:**

- set before launching the learning process
- not updated in each training ste (learning rate or the number of units in a dense layer)
- 

### **Hyperparameter**

- architecture options
- activation functions
- weight initialization strategy
- optimization hyperparameters (learning rate, stop condition),

### **Tuning methods:**

- Random search
- Hyperband
- Bayesian optimization
- sklearn

## Keras Autotuner

```
Get Started  autotuner.py 2, U X
autotuner.py > ...
2 import tensorflow as tf
3 from tensorflow import keras
4 mnist = tf.keras.datasets.mnist
5
6 (x_train,y_train),(x_test,y_test) = mnist.load_data()
7 x_train, x_test = x_train /255.0, x_test /255.0
8
9 model = tf.keras.models.Sequential([
10     tf.keras.layers.Flatten(input_shape=(28, 28)),
11     tf.keras.layers.Dense(512, activation='relu'),
12     tf.keras.layers.Dropout(0.2),
13     tf.keras.layers.Dense(10, activation = 'softmax')
14 ])
15
16 model.compile(optimizer='adam',
17               loss='sparse_categorical_crossentropy',
18               metrics=['accuracy'])
19 model.fit(x_train,y_train, epochs=5)
20 model.evaluate(x_test,y_test)
21
```

? `pip install -q -U keras-tuner`

Is this architecture optimal?

- Do the model need more or less hidden units to perform well?
- How does model size affect the convergence speed?
- Is there any trade off between convergence speed, model size and accuracy?
- Search automation is the natural path to take
- Keras tuner built in search functionality.

Other search strategien:

The parameters you choose will vary based on your strategy. But the important one to note is the **objective**.

In this case, our objective is `val_accuracy`, so we want to maximize on the validation accuracy. You can find details on the rest at the Keras site. Search can take a while to complete and use a lot of compute resources.

But you can configure a an **early stopping callback** that stops the search when the conditions are met. So for example, here I'm monitoring the validation loss and the patience is set to **five**, which means that it doesn't change significantly, or if it doesn't change significantly in five epics, then stop searching on this iteration. And you set the **call back as a search parameter**.

The rest of the parameter specify **how to search**, such as the **data** and **label**, the number of **epics to train** for, and the validation split. So for example how much data you will use to validate against the training set.

As it searches, you'll see the results of each trial.

In the example: searching on one parameter, the **units in the dense hidden layer**. As you can see as it updates, you can keep track of the best value so far. Which at this point that we have in the slide is 48, and in this case actually it ended up that way too when it completed. So I can go back and try my architecture again and use the results of the Keras tuner to set the number of units manually, this time with 48 neurons in the layer, which is the number that we got from from Keras tuner.

And when it's retrained, you'll see the results. It's only been five epics and the value isn't quite as good as it was before hand, but our epics are three times quicker than they were. So I can maybe train for longer to get better results knowing that at the very least I've optimized part of my architecture.

```
from gc import callbacks
from pickletools import optimize
from tracemalloc import stop

from yaml import DirectiveToken
import tensorflow as tf
from tensorflow import keras
import keras_tuner as kt

mnist = tf.keras.datasets.mnist
(x_train,y_train),(x_test,y_test) = mnist.load_data()
x_train, x_test = x_train /255.0, x_test /255.0

def model_builder(hp):
    model = keras.Sequential()
    model.add(keras.layers.Flatten(input_shape=(28,
28)))
```

```

    # hp_units=hp.Int('units', min_value=16,
max_value=512, step=16)
    # model.add(keras.layers.Dense(units=hp_units,
activation='relu'))
    model.add(keras.layers.Dense(hp.Choice('units',
[16,16,512]), activation='relu'))
    model.add(tf.keras.layers.Dropout(0.2))
    model.add(keras.layers.Dense(10))

    model.compile(optimizer='adam',
                  loss='sparse_categorical_crossentropy',
                  metrics=['accuracy'])
    return model

tuner = kt.Hyperband(model_builder,
                    objective = 'val_accuracy',
                    max_epochs=10,
                    factor=3,
                    directory='my_dir',
                    project_name = 'intro_to_kt')
# kt.RandomSearch

stop_early =
tf.keras.callbacks.EarlyStopping(monitor="val_loss",
patience=5)
# 5 steps not significantly changed

tuner.search(x_train,
            y_train,
            epochs=50,
            validation_split=0.2,
            callbacks=[stop_early])

```

```
best_model = tuner.get_best_models()[0]
```



## ▼ Ungraded Lab: Intro to Keras Tuner

Developing machine learning models is usually an iterative process. You start with an initial design then reconfigure until you get a model that can be trained efficiently in terms of time and compute resources. As you may already know, these settings that you adjust are called *hyperparameters*. These are the variables that govern the training process and the topology of an ML model. These remain constant over the training process and directly impact the performance of your ML program.

The process of finding the optimal set of hyperparameters is called *hyperparameter tuning* or *hypertuning*, and it is an essential part of a machine learning pipeline. Without it, you might end up with a model that has unnecessary parameters and take too long to train.

Hyperparameters are of two types:

1. **Model hyperparameters** which influence model selection such as the number and width of hidden layers
2. **Algorithm hyperparameters** which influence the speed and quality of the learning algorithm such as the learning rate for Stochastic Gradient Descent (SGD) and the number of nearest neighbors for a k Nearest Neighbors (KNN) classifier.

For more complex models, the number of hyperparameters can increase dramatically and tuning them manually can be quite challenging.

In this lab, you will practice hyperparameter tuning with [Keras Tuner](#), a package from the Keras team that automates this process. For comparison, you will first train a baseline model with pre-selected hyperparameters, then redo the process with tuned hyperparameters. Some of the examples and discussions here are taken from the [official tutorial provided by Tensorflow](#) but we've expounded on a few key parts for clarity.



## ▼ Download and prepare the dataset

Let us first load the [Fashion MNIST dataset](#) into your workspace. You will use this to train a machine learning model that classifies images of clothing.

```
1 # Import keras
```

```

2 from tensorflow import keras

1 # Download the dataset and split into train and test sets
2 (img_train, label_train), (img_test, label_test) = keras.datasets.fashion_mnist

Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/fashion-mnist-10000-images.tar.gz
32768/29515 [=====] - 0s 0us/step
40960/29515 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/fashion-mnist-10000-labels.tar.gz
26427392/26421880 [=====] - 0s 0us/step
26435584/26421880 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist-10000-images.tar.gz
16384/5148 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist-10000-labels.tar.gz
4423680/4422102 [=====] - 0s 0us/step
4431872/4422102 [=====] - 0s 0us/step

```

For preprocessing, you will normalize the pixel values to make the training converge faster.

```

1 # Normalize pixel values between 0 and 1
2 img_train = img_train.astype('float32') / 255.0
3 img_test = img_test.astype('float32') / 255.0

```

## ▼ Baseline Performance

As mentioned, you will first have a baseline performance using arbitrarily handpicked parameters so you can compare the results later. In the interest of time and resource limits provided by Colab, you will just build a shallow dense neural network (DNN) as shown below. This is to demonstrate the concepts without involving huge datasets and long tuning and training times. As you'll see later, even small models can take some time to tune. You can extend the concepts here when you get to build more complex models in your own projects.

```

1 # Build the baseline model using the Sequential API
2 b_model = keras.Sequential()
3 b_model.add(keras.layers.Flatten(input_shape=(28, 28)))
4 b_model.add(keras.layers.Dense(units=512, activation='relu', name='dense_1')) #
5 b_model.add(keras.layers.Dropout(0.2))
6 b_model.add(keras.layers.Dense(10, activation='softmax'))
7
8 # Print model summary
9 b_model.summary()

```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
flatten (Flatten)	(None, 784)	0
dense_1 (Dense)	(None, 512)	401920

dropout (Dropout)	(None, 512)	0
dense (Dense)	(None, 10)	5130

```

=====
Total params: 407,050
Trainable params: 407,050
Non-trainable params: 0
=====

```

As shown, we hardcoded all the hyperparameters when declaring the layers. These include the number of hidden units, activation, and dropout. You will see how you can automatically tune some of these a bit later.

Let's then setup the loss, metrics, and the optimizer. The learning rate is also a hyperparameter you can tune automatically but for now, let's set it at 0.001.

```

1 # Setup the training parameters
2 b_model.compile(optimizer=keras.optimizers.Adam(learning_rate=0.001),
3                 loss=keras.losses.SparseCategoricalCrossentropy(),
4                 metrics=['accuracy'])

```

With all settings set, you can start training the model. We've set the number of epochs to 10 but feel free to increase it if you have more time to go through the notebook.

```

1 # Number of training epochs.
2 NUM_EPOCHS = 10
3
4 # Train the model
5 b_model.fit(img_train, label_train, epochs=NUM_EPOCHS, validation_split=0.2)

```

```

Epoch 1/10
1500/1500 [=====] - 7s 3ms/step - loss: 0.5159 - acc
Epoch 2/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3945 - acc
Epoch 3/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3587 - acc
Epoch 4/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3319 - acc
Epoch 5/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3183 - acc
Epoch 6/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3055 - acc
Epoch 7/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2911 - acc
Epoch 8/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2805 - acc
Epoch 9/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2753 - acc
Epoch 10/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2636 - acc
<keras.callbacks.History at 0x7fb8e002bf50>

```



Finally, you want to see how this baseline model performs against the test set.

```
1 # Evaluate model on the test set
2 b_eval_dict = b_model.evaluate(img_test, label_test, return_dict=True)

313/313 [=====] - 1s 3ms/step - loss: 0.3503 - accur
```

Let's define a helper function for displaying the results so it's easier to compare later.

```
1 # Define helper function
2 def print_results(model, model_name, eval_dict):
3     '''
4     Prints the values of the hyparameters to tune, and the results of model evalu
5
6     Args:
7         model (Model) - Keras model to evaluate
8         model_name (string) - arbitrary string to be used in identifying the model
9         eval_dict (dict) - results of model.evaluate
10    '''
11    print(f'\n{model_name}:')
12
13    print(f'number of units in 1st Dense layer: {model.get_layer("dense 1").units}')
14    print(f'learning rate for the optimizer: {model.optimizer.lr.numpy()}')
15
16    for key,value in eval_dict.items():
17        print(f'{key}: {value}')
18
19 # Print results for baseline model
20 print_results(b_model, 'BASELINE MODEL', b_eval_dict)
```

```
BASELINE MODEL:
number of units in 1st Dense layer: 512
learning rate for the optimizer: 0.0010000000474974513
loss: 0.35033318400382996
accuracy: 0.878000020980835
```

That's it for getting the results for a single set of hyperparameters. As you can see, this process can be tedious if you want to try different sets of parameters. For example, will your model improve if you use `learning_rate=0.00001` and `units=128`? What if `0.001` paired with `256`? The process will be even more difficult if you decide to also tune the dropout and try out other activation functions as well. Keras Tuner solves this problem by having an API to automatically search for the optimal set. You will just need to set it up once then wait for the results. You will see how this is done in the next sections.

## ▼ Keras Tuner

To perform hypertuning with Keras Tuner, you will need to:

- Define the model
- Select which hyperparameters to tune
- Define its search space
- Define the search strategy

## ▼ Install and import packages

You will start by installing and importing the required packages.

```
1 # Install Keras Tuner
2 !pip install -q -U keras-tuner
```

|██| 135 kB 4.7 MB/s

```
1 # Import required packages
2 import tensorflow as tf
3 import kerastuner as kt
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: DeprecationWarning:
  This is separate from the ipykernel package so we can avoid doing imports u
```



## ▼ Define the model

The model you set up for hypertuning is called a *hypermodel*. When you build this model, you define the hyperparameter search space in addition to the model architecture.

You can define a hypermodel through two approaches:

- By using a model builder function
- By [subclassing the HyperModel class](#) of the Keras Tuner API

In this lab, you will take the first approach: you will use a model builder function to define the image classification model. This function returns a compiled model and uses hyperparameters you define inline to hypertune the model.

The function below basically builds the same model you used earlier. The difference is there are two hyperparameters that are setup for tuning:

- the number of hidden units of the first Dense layer
- the learning rate of the Adam optimizer

You will see that this is done with a HyperParameters object which configures the hyperparameter you'd like to tune. For this exercise, you will:

- use its `Int()` method to define the search space for the Dense units. This allows you to set a minimum and maximum value, as well as the step size when incrementing between these values.
- use its `Choice()` method for the learning rate. This allows you to define discrete values to include in the search space when hypertuning.

You can view all available methods and its sample usage in the [official documentation](#)

```

1 def model_builder(hp):
2     '''
3     Builds the model and sets up the hyperparameters to tune.
4
5     Args:
6         hp - Keras tuner object
7
8     Returns:
9         model with hyperparameters to tune
10    '''
11
12    # Initialize the Sequential API and start stacking the layers
13    model = keras.Sequential()
14    model.add(keras.layers.Flatten(input_shape=(28, 28)))
15
16    # Tune the number of units in the first Dense layer
17    # Choose an optimal value between 32-512
18    hp_units = hp.Int('units', min_value=32, max_value=512, step=32)
19    model.add(keras.layers.Dense(units=hp_units, activation='relu', name='dense_1')
20
21    # Add next layers
22    model.add(keras.layers.Dropout(0.2))
23    model.add(keras.layers.Dense(10, activation='softmax'))
24
25    # Tune the learning rate for the optimizer
26    # Choose an optimal value from 0.01, 0.001, or 0.0001
27    hp_learning_rate = hp.Choice('learning_rate', values=[1e-2, 1e-3, 1e-4])
28
29    model.compile(optimizer=keras.optimizers.Adam(learning_rate=hp_learning_rate),
30                  loss=keras.losses.SparseCategoricalCrossentropy(),
31                  metrics=['accuracy'])
32
33    return model

```

## ▼ Instantiate the Tuner and perform hypertuning

Now that you have the model builder, you can then define how the tuner can find the optimal set of hyperparameters, also called the search strategy. Keras Tuner has [four tuners](#) available with built-in strategies - **RandomSearch, Hyperband, BayesianOptimization, and Sklearn.**

In this tutorial, you will use the Hyperband tuner. Hyperband is an algorithm specifically developed for hyperparameter optimization. It uses adaptive resource allocation and early-stopping to quickly converge on a high-performing model. This is done using a sports

championship style bracket wherein the algorithm trains a large number of models for a few epochs and carries forward only the top-performing half of models to the next round. You can read about the intuition behind the algorithm in section 3 of [this paper](#).

Hyperband determines the number of models to train in a bracket by computing  $1 + \log_{\text{factor}}(\text{max\_epochs})$  and rounding it up to the nearest integer. You will see these parameters (i.e. `factor` and `max_epochs` passed into the initializer below). In addition, you will also need to define the following to instantiate the Hyperband tuner:

- the `hypermodel` (built by your model builder function)
- the `objective` to optimize (e.g. validation accuracy)
- a `directory` to save logs and checkpoints for every trial (model configuration) run during the hyperparameter search. If you re-run the hyperparameter search, the Keras Tuner uses the existing state from these logs to resume the search. To disable this behavior, pass an additional `overwrite=True` argument while instantiating the tuner.
- the `project_name` to differentiate with other runs. This will be used as a subdirectory name under the `directory`.

You can refer to the [documentation](#) for other arguments you can pass in.

```
1 # Instantiate the tuner
2 tuner = kt.Hyperband(model_builder,
3                       objective='val_accuracy',
4                       max_epochs=10,
5                       factor=3,
6                       directory='kt_dir',
7                       project_name='kt_hyperband')
```

Let's see a summary of the hyperparameters that you will tune:

```
1 # Display hypertuning settings
2 tuner.search_space_summary()
```

```
Search space summary
Default search space size: 2
units (Int)
{'default': None, 'conditions': [], 'min_value': 32, 'max_value': 512, 'step'
learning_rate (Choice)
{'default': 0.01, 'conditions': [], 'values': [0.01, 0.001, 0.0001], 'orderec
```

You can pass in a callback to stop training early when a metric is not improving. Below, we define an [EarlyStopping](#) callback to monitor the validation loss and stop training if it's not improving after 5 epochs.

```
1 stop_early = tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=5)
```

You will now run the hyperparameter search. The arguments for the search method are the same as those used for `tf.keras.model.fit` in addition to the callback above. This will take around 10 minutes to run.

```
1 # Perform hypertuning
2 tuner.search(img_train, label_train, epochs=NUM_EPOCHS, validation_split=0.2, c
```

```
Trial 30 Complete [00h 00m 46s]
val_accuracy: 0.8854166865348816
```

```
Best val_accuracy So Far: 0.8872500061988831
Total elapsed time: 00h 12m 03s
INFO:tensorflow:Oracle triggered exit
INFO:tensorflow:Oracle triggered exit
```

You can get the top performing model with the [get\\_best\\_hyperparameters\(\)](#) method.

```
1 # Get the optimal hyperparameters from the results
2 best_hps=tuner.get_best_hyperparameters()[0]
3
4 print(f"""
5 The hyperparameter search is complete. The optimal number of units in the first
6 layer is {best_hps.get('units')} and the optimal learning rate for the optimize
7 is {best_hps.get('learning_rate')}.
8 """)
9
10 best_model = tuner.get_best_models()[0]
```

```
The hyperparameter search is complete. The optimal number of units in the fir
layer is 192 and the optimal learning rate for the optimizer
is 0.001.
```

## ▼ Build and train the model

Now that you have the best set of hyperparameters, you can rebuild the hypermodel with these values and retrain it.

```
1 # Build the model with the optimal hyperparameters
2 h_model = tuner.hypermodel.build(best_hps)
3 h_model.summary()
```

```
Model: "sequential_1"
```

Layer (type)	Output Shape	Param #
flatten_1 (Flatten)	(None, 784)	0

dense_1_1 (Dense)	(None, 192)	150720
dropout_1 (Dropout)	(None, 192)	0
dense_1 (Dense)	(None, 10)	1930

```

=====
Total params: 152,650
Trainable params: 152,650
Non-trainable params: 0
=====

```

```

1 # Train the hypertuned model
2 h_model.fit(img_train, label_train, epochs=NUM_EPOCHS, validation_split=0.2)

```

```

Epoch 1/10
1500/1500 [=====] - 5s 3ms/step - loss: 0.5370 - acc: 0.1250
Epoch 2/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.4007 - acc: 0.2500
Epoch 3/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3673 - acc: 0.3750
Epoch 4/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3420 - acc: 0.5000
Epoch 5/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3281 - acc: 0.6250
Epoch 6/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3146 - acc: 0.7500
Epoch 7/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.3026 - acc: 0.8750
Epoch 8/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2920 - acc: 0.9000
Epoch 9/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2830 - acc: 0.9250
Epoch 10/10
1500/1500 [=====] - 4s 3ms/step - loss: 0.2732 - acc: 0.9500
<keras.callbacks.History at 0x7fb84d6808d0>

```

You will then get its performance against the test set.

return\_dict: If True, loss and metric results are returned as a dict, with each key being the name of the metric. If False, they are returned as a list.

```

1 # Evaluate the hypertuned model against the test set
2 h_eval_dict = h_model.evaluate(img_test, label_test, return dict=True)

```

```

313/313 [=====] - 1s 3ms/step - loss: 0.3471 - accur

```

We can compare the results we got with the baseline model we used at the start of the notebook. Results may vary but you will usually get a model that has less units in the dense layer, while having comparable loss and accuracy. This indicates that you reduced the model size and saved compute resources while still having more or less the same accuracy.

```

1 # Print results of the baseline and hypertuned model
2 print_results(b_model, 'BASELINE MODEL', b_eval_dict)
3 print_results(h_model, 'HYPERTUNED MODEL', h_eval_dict)

```

```
print_results(h_model, INTENTED_MODEL, h_eval_dict)
```

#### BASELINE MODEL:

number of units in 1st Dense layer: 512

learning rate for the optimizer: 0.0010000000474974513

loss: 0.35033318400382996

accuracy: 0.878000020980835

#### HYPERTUNED MODEL:

number of units in 1st Dense layer: 10

learning rate for the optimizer: 0.0010000000474974513

loss: 0.34706664085388184

accuracy: 0.8773000240325928

## Bonus Challenges (optional)

If you want to keep practicing with Keras Tuner in this notebook, you can do a factory reset (Runtime > Factory reset runtime) and take on any of the following:

- hypertune the dropout layer with `hp.Float()` or `hp.Choice()`
- hypertune the activation function of the 1st dense layer with `hp.Choice()`
- determine the optimal number of Dense layers you can add to improve the model. You can use the code [here](#) as reference.
- explore pre-defined `HyperModel` classes - [HyperXception](#) and [HyperResNet](#) for computer vision applications.

## Wrap Up

In this tutorial, you used Keras Tuner to conveniently tune hyperparameters. You defined which ones to tune, the search space, and search strategy to arrive at the optimal set of hyperparameters. These concepts will again be discussed in the next sections but in the context of AutoML, a package that automates the entire machine learning pipeline. On to the next!

---

✓ 0 s Abgeschlossen um 09:11

● ✕



# Introduction to AutoML

AutoML. A set of very versatile tools to automate the machine learning process end to end. e. You'll be learning about **search spaces** and **strategies** which are key for both hyper parameter tuning and AutoML.

We'll also discuss tools to **quantify performance estimation** on the explored models in your search.

Finally, you learn about **different AutoML offerings** available in the Cloud by major service providers.

It covers the complete pipeline from the **raw data** set to the **deployable** machine learning model.



Neural Architecture Search or **NAS** is at the **heart** of AutoML.

There are three **main parts** to Neural Architecture Search,

**a search space:**

- the range of architectures which can be represented. To **reduce the size of the search problem**, we need to limit the search space to the architectures which are best suited to the problem that we're trying to model. This helps reduce the search space, but it also means that **a human bias** will be introduced, which might prevent Neural Architecture Search from finding architectural blocks that go beyond current human knowledge.

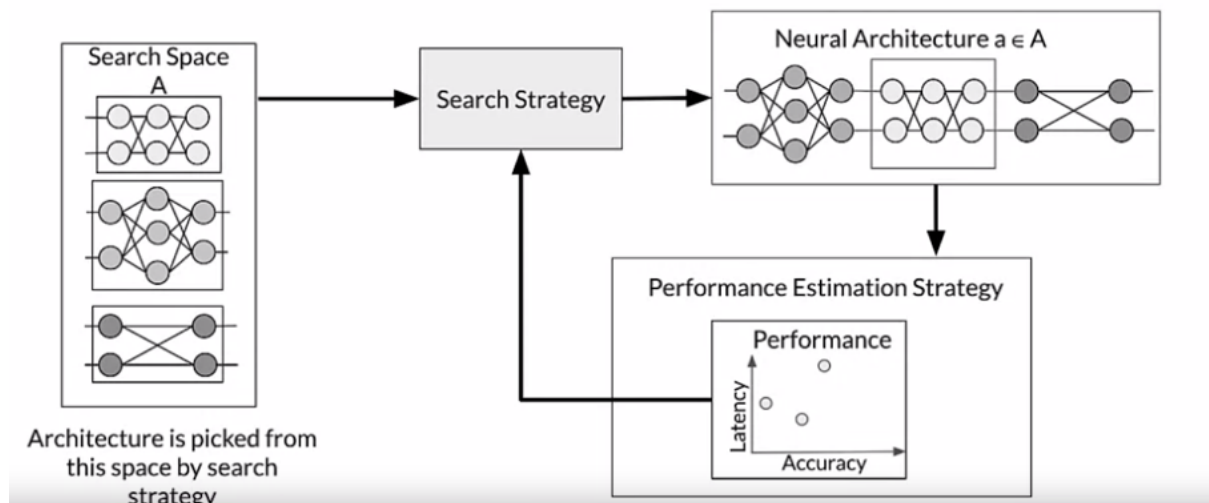
**a search strategy:**

- how we explore the search space. We want to explore the search based **quickly**, but this might lead to **premature convergence to a sub optimal region in the search space**. The objective of Neural Architecture Search is to find **architecture is that perform well on our data**.

**a performance estimation strategy:**

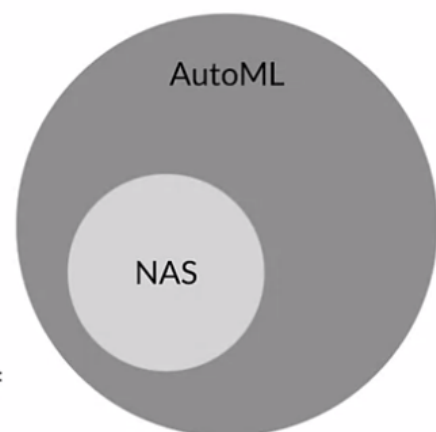
- helps in measuring and comparing the performance of various architectures
- A search strategies selects an architecture from a predefined search space of architectures. The selected architecture is passed to a performance estimation strategy, which returns it's estimated performance to the search strategy.

## Neural Architecture Search



## Neural Architecture Search

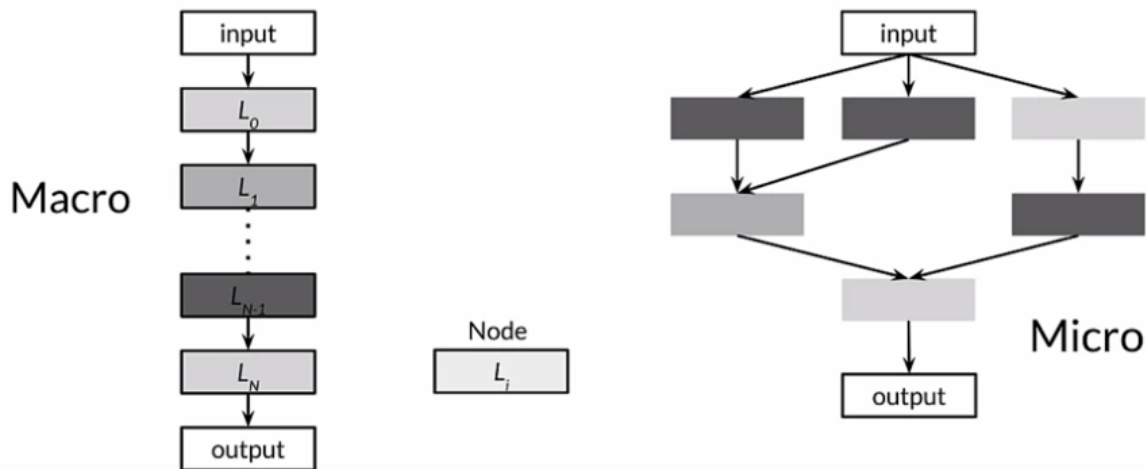
- **AutoML** automates the development of ML models
- **AutoML** is not specific to a particular type of model.
- Neural Architecture Search (**NAS**) is a subfield of AutoML
- NAS is a technique for automating the design of artificial neural networks (ANN).



## Understanding Search Spaces

search spaces, **macro** and **micro**.

## Types of Search Spaces

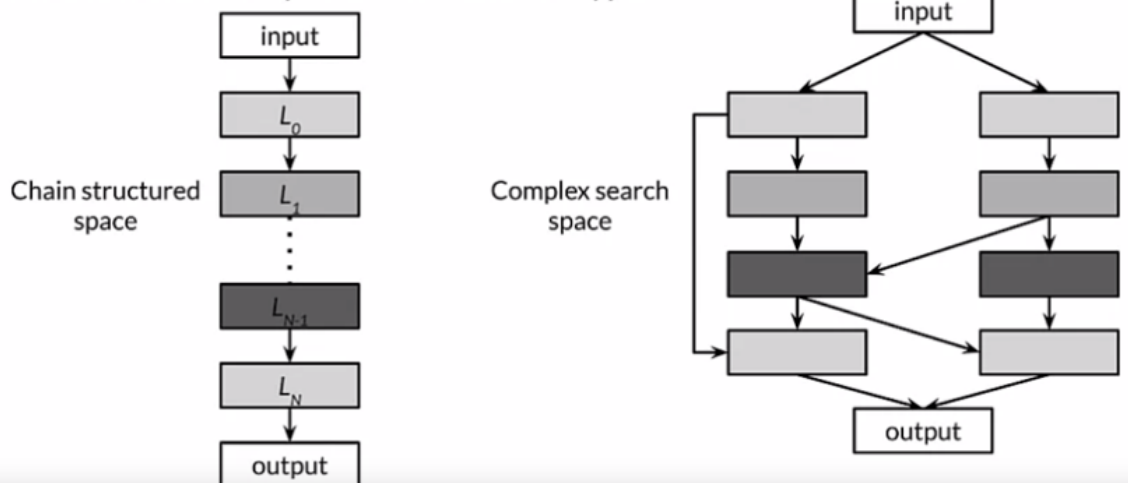


### MACRO:

- A macro search space contains the individual layers and connection types of a neural network. And neural architecture search searches within that space for the best model, building the model layer by layer
- a network can be built very simply by stacking individual layers referred to as a chain structured space or with multiple branches and skip connections in a complex space

## Macro Architecture Search Space

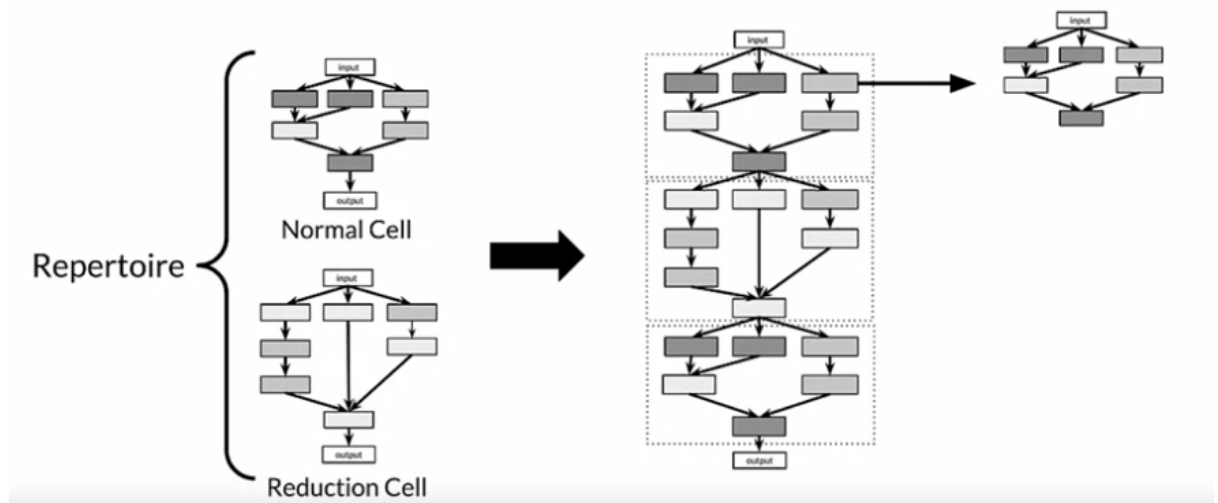
Contains individual layers and connection types



### MIRCO:

- neural architecture search builds a neural network from cells where each cell is a smaller network. Here are two different types of cells, a normal cell on the top and a reduction cell on the bottom. Cells are stacked to produce the final network. This approach has been shown to have s

## Micro Architecture Search Space



significant performance advantages compared to a macro approach.

## Search Strategies

Neural architecture search searches through the search base **for the architecture** that produces the best performance.

A variety of different **approaches** can be used for that search:

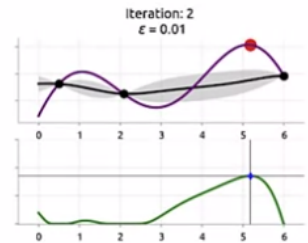
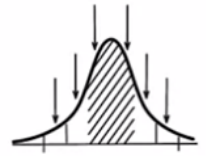
- grid search
  - search everything, that means you cover every option we have in the search base
- random search
  - you select the next option randomly within the search space

for grid and random search:

- Both of these work reasonably well **in smaller search bases**, but both also fail fairly quickly. When the search space grows beyond a certain size, which is all too common.
- Bayesian optimization
  - a little more sophisticated. It assumes that a specific probability distribution, which is typically a Gaussian distribution is underlying the performance of model architectures. So you use observations from tested architectures to constrain the probability distribution and guide the selection of the next option. This allows us to build up an architecture stochastically based on the test results and the constrained distribution.

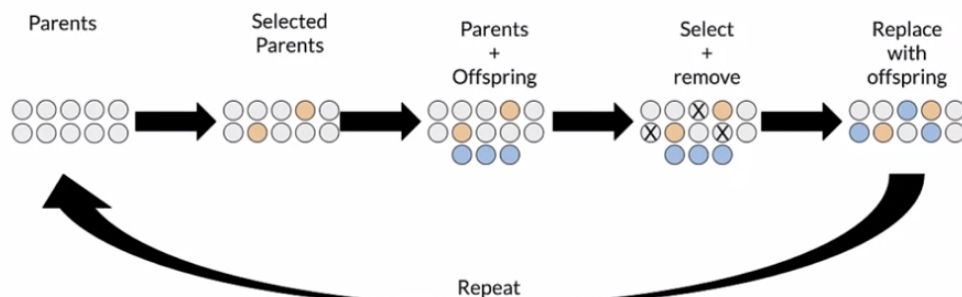
# Bayesian Optimization

- Assumes that a *specific probability distribution*, is underlying the performance.
- Tested architectures constrain the probability distribution and guide the selection of the next option.
- In this way, promising architectures can be stochastically determined and tested.



- evolutionary algorithms
  - First, an initial population of  $n$  different model architecture is **randomly generated**, the performance of each individual. In other words, each architecture is evaluated as defined by the performance estimation strategy, which we'll talk about next. Then the  $X$  highest performers are selected as parents for **a new generation**. This new generation of architectures might be copies of the respective parents with induced random **alterations** or **mutations**. Or they might arise from combinations of the parents, the performance of the offspring is assessed.
  - Again using the performance estimation strategy. The list of possible **mutations** can include operations like adding or removing a layer. Adding or removing a connection, changing the size of a layer or changing another hyper parameter.  $Y$  architectures are selected to be removed from the population. This might be the  $Y$  worst performers, the  $Y$  oldest individuals in the population. Or a selection of individuals based on a combination of these parameters. The offspring replaces the removed architectures and the process is restarted with this new population.

## Evolutionary Methods

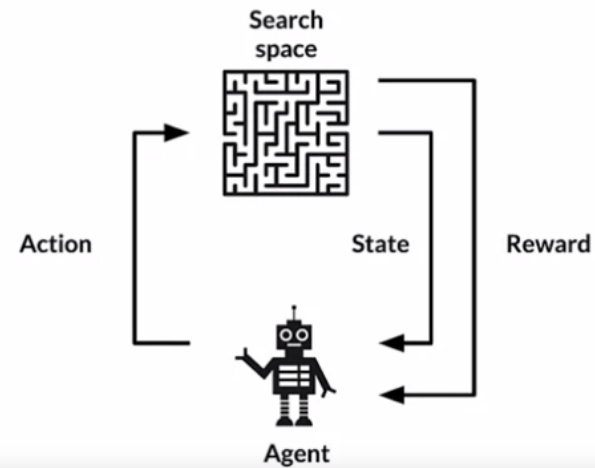


- reinforcement learning.
  - agents take actions in an environment, trying to maximize a reward. After each action, the state of the agent and the environment is updated and a reward is issued based on a performance metric. Then the range of possible next actions is

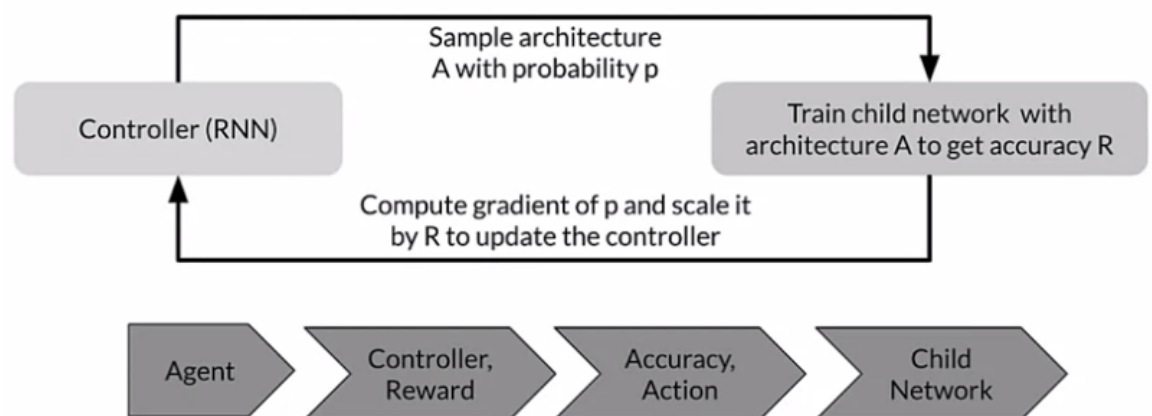
evaluated, the environment in this case is our search space. And the reward function is our performance estimation strategy.

## Reinforcement Learning

- Agents goal is to maximize a reward
- The available options are selected from the search space
- The performance estimation strategy determines the reward



- A neural network can also be specified by a variable length string where the elements of the string specify individual network layers. That enables us to use a recurrent neural network or RNN to generate that string, as we might do for an NLP model. The RNN that generates the string is referred to as the controller, after training the network referred to as the child network on real data. We can measure the accuracy on the validation set, the accuracy determines the reinforcement learning reward in this case. Based on the accuracy, we can compute the policy gradient to update the controller RNN. In the next iteration, the controller will have learned to give higher probabilities to architectures that result in higher accuracies during training.



<https://distill.pub/2020/bayesian-optimization/>

and other papers for further reading

# Measuring AutoML Efficacy

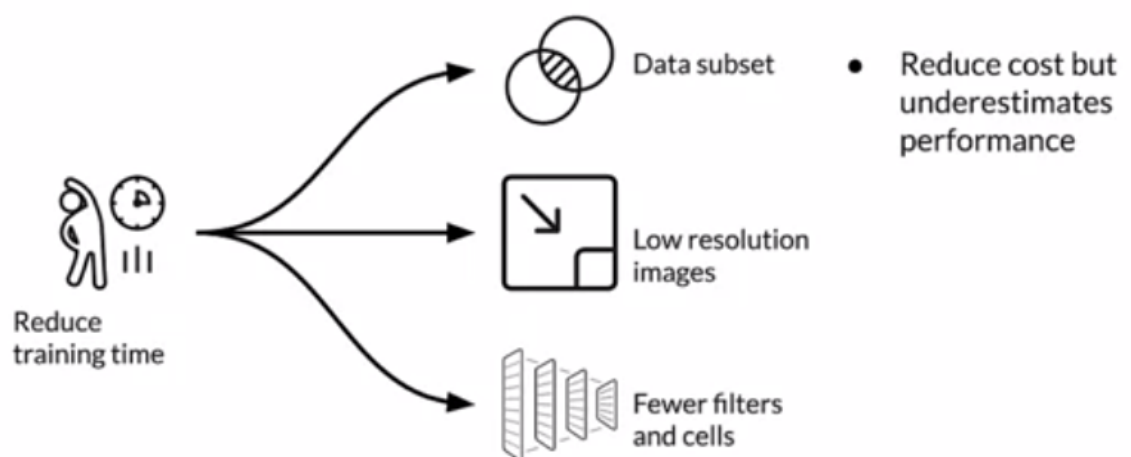
measure the accuracy or effectiveness of the different architectures that it tries. This requires a performance estimation strategy.

The simplest approach is to measure the **validation accuracy** of each architecture that is generated like we saw with the reinforcement learning approach. This becomes computational **heavy**, especially for large search spaces and complex networks. And as a result, it can take several GPU days to find the best architectures using this approach, that makes it **expensive** and **slow**. It also makes neural architecture search impractical for many use cases.

Other approaches:

- lower fidelity 逼真度 estimates
  - Lower fidelity or lower precision estimates try to **reduce the training time** by reframing the problem to make it easier to solve by training on a subset of data or using lower resolution images, for example, or using fewer filters per layer and fewer cells. It greatly reduces the computational cost, but ends up underestimating performance. That's okay if you can make sure that the relative ranking of the architectures does not change due to their lower fidelity estimates. But unfortunately, recent research has shown that this is not the case, bummer

## Lower Fidelity Estimates

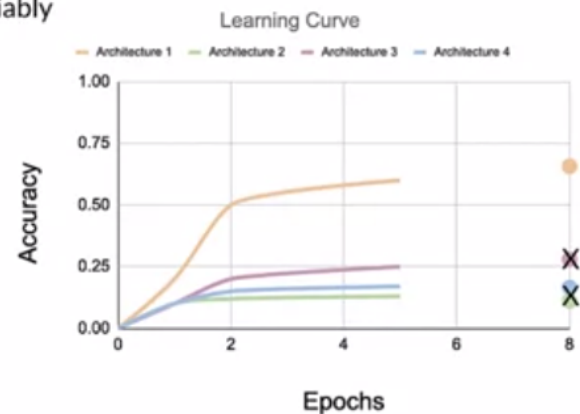


- learning curve extrapolation
  - based on the assumption that you have mechanisms to predict the learning curve reliably, and so extrapolation is a sensitive and valid choice. On the right there are learning curves for different architectures. Based on a few iterations and available knowledge, the method extrapolates the initial learning curves and terminates all architectures that perform poorly. The progressive neural architecture search algorithm, which is one of the approaches for neural architecture search, uses a similar method by training a surrogate model and using it to predict the performance using architectural properties



## Learning Curve Extrapolation

- Requires predicting the learning curve reliably
- Extrapolates based on initial learning.
- Removes poor performers



- 
- weight inheritance or network morphism
  - based on the weights of other architectures that have been trained before, similar to the way that transfer learning works. One way of achieving this is referred to as network morphism 网络射. Network morphism modifies the architecture without changing the underlying function. This is advantageous because the network inherits knowledge from the parent network, which results in methods that require only a few GPU days to design and evaluate. This allows for increasing the capacity of network successively and retaining high performance without requiring training from scratch. One advantage of this approach is that it allows for search bases that don't have an inherent upper bound on the architecture's size.

## Weight Inheritance/Network Morphisms

- Initialize weights of new architectures based on previously trained architectures
  - Similar to transfer learning
- Uses **Network Morphism**
- Underlying function unchanged
  - New network inherits knowledge from parent network.
  - Computational speed up: only a few days of GPU usage
  - Network size not inherently bounded



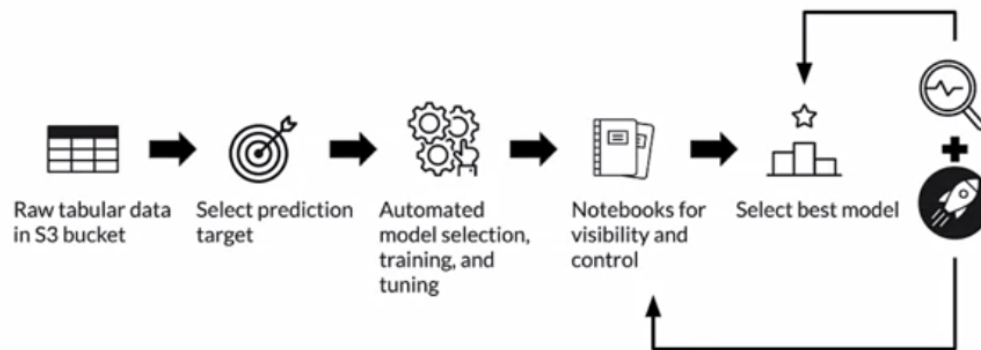
# AutoML on the Cloud

Cloud service.

## [Amazon SageMaker Autopilot](#)

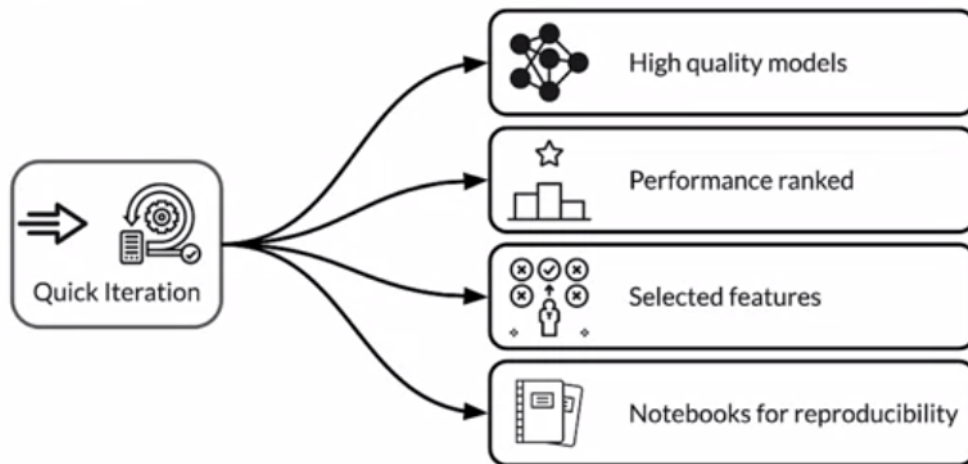
- Amazon SageMaker Autopilot automatically trains and tunes the best machine learning models for classification or regression based on your data, while allowing you to maintain full control and visibility.
- Starting with your raw data, you select a label or target. Autopilot then searches for candidate models for you to review and choose from. All of these steps are documented on executable notebooks that give you full control and reproducibility of the process. This includes a leader board of model candidates to help you select the best model for your needs. You can then deploy the model to production or iterate on the recommended solutions to further improve model quality

### Amazon SageMaker Autopilot



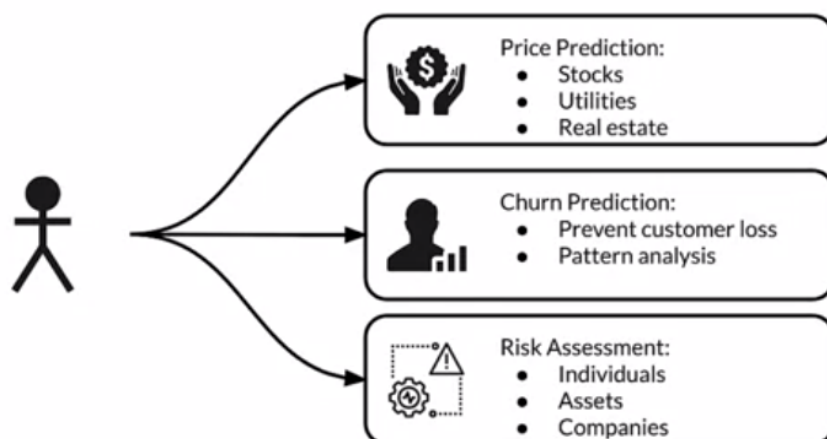
- 
- Autopilot is optimized for quick iteration. It generates **high-quality models** quickly. After the initial set of iterations, autopilot creates a **leaderboard of models ranked by performance**. You can see **which features in your dataset were selected by each model**. You can then deploy a model to production. The process of generation of the models is completely transparent. Autopilot allows you to create a SageMaker notebook from any model it created. You can then check the Notebook to dive into details of the models implementations. If need be, you can refine the model and recreate it from the Notebook at any point in time.

## Key features



- 
- usecases:
  - it can be used to predict future prices to help you make sound investment decisions based on your historical data, such as demand, seasonal trends, and the price of other commodities. Price predictions are particularly useful in the financial services sector to predict the price of stocks or real estate, to predict real estate prices, or energy and utilities to predict the prices of natural resources
  - Churn prediction can be useful in predicting the loss of customers or churn. Companies are always looking for ways to eliminate churn. Churn prediction works by learning patterns in your existing data and identifying patterns in the new datasets so that the model can predict which customers are most likely to churn.
  - Another application is risk assessment. Risk assessment requires identifying and analyzing potential events that may negatively impact individuals, assets, and your company. Risk assessment models are trained using your existing datasets so that you can optimize the models predictions for your business.

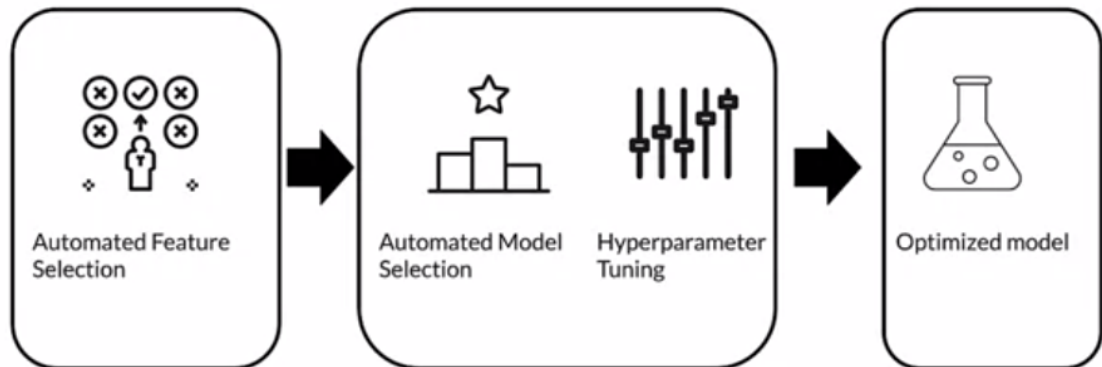
## Typical use cases



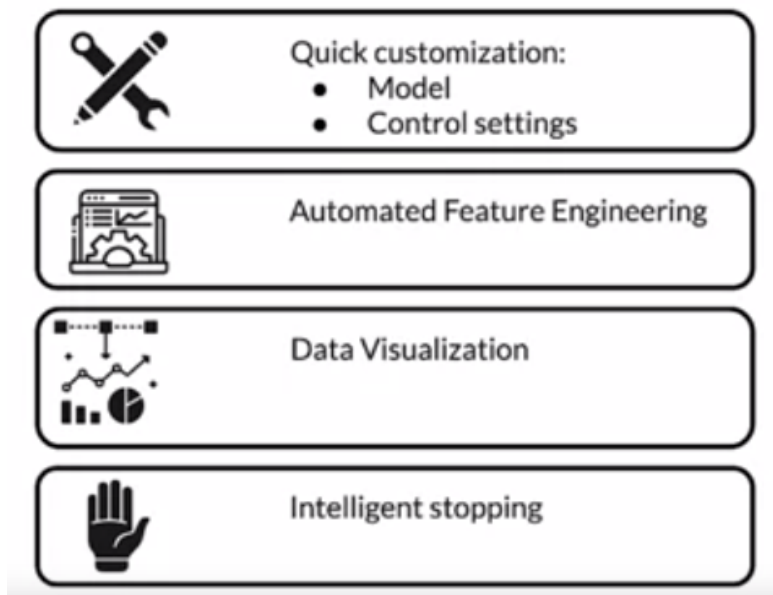
## Microsoft Azure Automated Machine Learning

- It automates time-consuming and iterative task of model development, so basically it does AutoML. It starts with automatic feature selection, followed by model selection and hyperparameter tuning on the selected model to generate the most optimized model for the task at hand

### Microsoft Azure AutoML

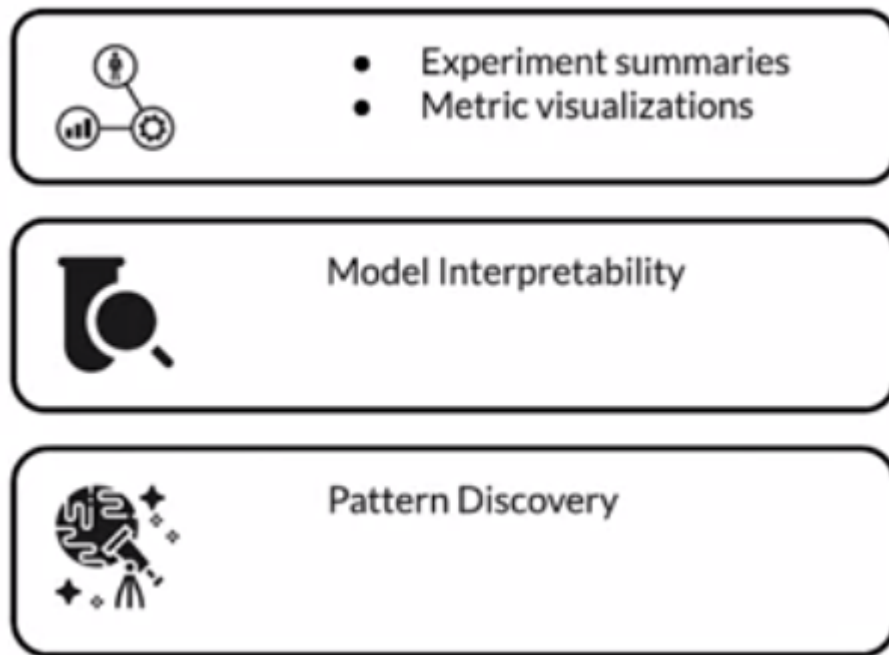


- 
- You can either create your models using a no code UI or using code first notebooks. You can quickly customize your models, apply control settings to iterations, thresholds, validations, blocked algorithms, and other experimental criteria. It also provides tools to fully automate the feature engineering process. You can also easily visualize and profile your data to spot trends, as well as discover common errors and inconsistencies in your data. This helps you better understand the recommended actions and apply them automatically. It also provides intelligent stopping to save time on computing and prioritize the primary metric and subsampling to streamline experiment runs and speed results.



- 
- .It also has built-in support for experiment runs summaries, and detailed metrics visualizations to help you understand your models and compare model performance. Model interpretability helps evaluate model fit for raw and engineered features and

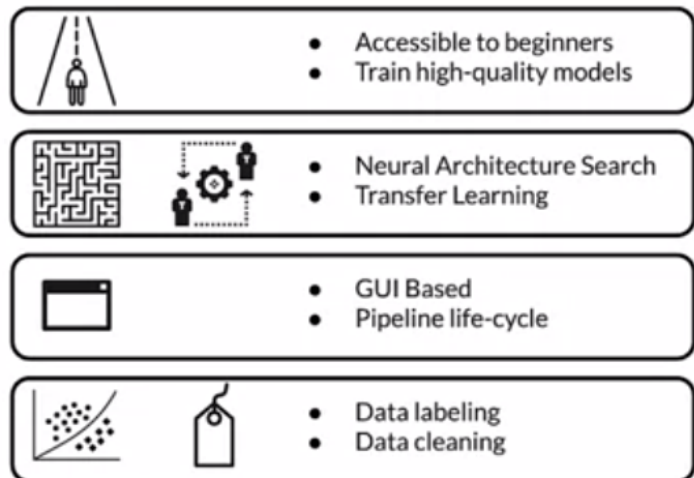
provides insights into feature importance. You can discover patterns, perform what-if analysis, and develop deeper understanding of models to support transparency and trust in your business.



### [Google Cloud AutoML.](#)

- a suite of machine learning products that enables developers with limited machine-learning expertise to train high-quality models specific to their business needs.
- It relies on Google, state of the art transfer learning and neural architecture search technologies. Cloud AutoML leverages more than 10 years of Google research to help your machine-learning models achieve faster performance and more accurate predictions.
- You can use Cloud AutoML, simple graphical user interface to train, evaluate, improve, and deploy models based on your data. You're really only a few minutes away from your own custom machine-learning model.
- Google's human labeling service can also put a team of people to work annotating or cleaning your labels to make sure that your models are being trained on high-quality data.

# Google Cloud AutoML



- 
- Because different problems and different data need to be treated differently, Cloud AutoML isn't just one thing. It's a suite of different products, each focused towards particular use cases and data types. For example, for image data, there's **AutoML Vision** and for video data, there's **AutoML Video Intelligence**, for natural language, there's **AutoML natural language** and for translation, there's **AutoML Translation**. Finally, for general structured data, there's **AutoML Tables**.

## Cloud AutoML Products

Sight	<b>Auto ML Vision</b> Derive insights from images in the cloud or at the edge.	<b>Auto ML Video Intelligence</b> Enable powerful content discovery and engaging video experiences.
Language	<b>AutoML Natural Language</b> Reveal the structure and meaning of text through machine learning.	<b>Auto ML Translation</b> Dynamically detect and translate between languages.
Structured Data	<b>AutoML Tables</b> Automatically build and deploy state-of-the-art machine learning models on structured data.	

- 
- For image data, for example, there's both vision and **classification** and Vision Object **Detection**. Then there are also Edge versions of both of these focused on optimizing for running inference at the edge in mobile applications or IoT devices.

## AutoML Vision Products

---

Auto ML Vision Classification

AutoML Vision Edge Image Classification

---

AutoML Vision Object Detection

AutoML Vision Edge Object Detection

- 
- For video, there's both **Video Intelligence classification** and **video object detection**. Well, no one knows exactly or can't say if they work for a given provider. However, it is safe to assume that the algorithms at play will be similar to the ones that we've discussed in previous lessons. Also, keep in mind that ML Engineering for Production is a rapidly changing field and new technologies and developments emerge every few months. These new advancements are typically incorporated and exploited to their greatest extent by the AutoML providers.
- Also notice that these different AutoML offerings each perform the same operations that you perform or should perform as you train your model, including feature selection and feature engineering, and cleaning your labels.
- AutoML designers understand the importance of these activities.

## So what's in the secret sauce?

### How do these Cloud offerings perform AutoML?

- We don't know (or can't say) and they're not about to tell us
- The underlying algorithms will be similar to what we've learned
- The algorithms will evolve with the state of the art

-