# Duck Typing, Scope, and Investigative Functions in Python

by **Adrian Tam** on February 16, 2022 in **Python for Machine Learning**

Tweet    Tweet    Share    Share

Last Updated on June 21, 2022

Python is a duck typing language. It means the data types of variables can change as long as the syntax is compatible. Python is also a dynamic programming language. Meaning we can change the program while it runs, including defining new functions and the scope of the name resolution. These give us not only a new paradigm in writing Python code but also a new set of tools for debugging. In the following, we will see what we can do in Python that cannot be done in many other languages.

After finishing this tutorial, you will know:

- How Python manages the variables you define
- How Python code uses a variable and why we don't need to define its type like in C or Java

**Kick-start your project** with my new book Python for Machine Learning, including *step-by-step tutorials* and the *Python source code* files for all examples.

Let's get started.



Duck typing, scope, and investigative functions in Python. Photo by Julissa Helmuth. Some rights reserved

# Overview

This tutorial is in three parts; they are

- Duck typing in programming languages
- Scopes and name space in Python
- Investigating the type and scope

# Duck Typing in Programming Languages

Duck typing is a feature of some modern programming languages that allow data types to be dynamic.

> A programming style which does not look at an object's type to determine if it has the right interface; instead, the method or attribute is simply called or used ("If it looks like a duck and quacks like a duck, it must be a duck.") By emphasizing interfaces rather than specific types, well-designed code improves its flexibility by allowing polymorphic substitution.

— Python Glossary

Simply speaking, the program should allow you to swap data structures as long as the same syntax still makes sense. In C, for example, you have to define functions like the following:

```
1  float fsquare(float x)
2  {
3      return x * x;
4  };
5
6  int isquare(int x)
7  {
8      return x * x;
9  };
```

While the operation x * x is identical for integers and floating-point numbers, a function taking an integer argument and a function taking a floating-point argument are not the same. Because types are static in C, we must define two functions although they perform the same logic. In Python, types are dynamic; hence we can define the corresponding function as:

```
1  def square(x):
2      return x * x
```

This feature indeed gives us tremendous power and convenience. For example, from scikit-learn, we have a function to do cross validation:

```
1   # evaluate a perceptron model on the dataset
2   from numpy import mean
3   from numpy import std
4   from sklearn.datasets import make_classification
5   from sklearn.model_selection import cross_val_score
6   from sklearn.model_selection import RepeatedStratifiedKFold
7   from sklearn.linear_model import Perceptron
8   # define dataset
9   X, y = make_classification(n_samples=1000, n_features=10, n_informative=10, n_redundant=0,
10  # define model
11  model = Perceptron()
12  # define model evaluation method
13  cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
14  # evaluate model
15  scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
16  # summarize result
17  print('Mean Accuracy: %.3f (%.3f)' % (mean(scores), std(scores)))
```

But in the above, the `model` is a variable of a scikit-learn-model object. It doesn't matter if it is a perceptron model as in the above, a decision tree, or a support vector machine model. What matters is that inside the `cross_val_score()` function, the data will be passed onto the model with its `fit()` function. Therefore, the model must implement the `fit()` member function, and the `fit()` function behaves identically. The consequence is that the `cross_val_score()` function is not expecting any particular model type as long as it looks like one. If we are using Keras to build a neural network model, we can make the Keras model look like a scikit-learn model with a wrapper:

```
1   # MLP for Pima Indians Dataset with 10-fold cross validation via sklearn
2   from keras.models import Sequential
3   from keras.layers import Dense
4   from keras.wrappers.scikit_learn import KerasClassifier
5   from sklearn.model_selection import StratifiedKFold
6   from sklearn.model_selection import cross_val_score
7   from sklearn.datasets import load_diabetes
8   import numpy
9
10  # Function to create model, required for KerasClassifier
11  def create_model():
12      # create model
13      model = Sequential()
14      model.add(Dense(12, input_dim=8, activation='relu'))
15      model.add(Dense(8, activation='relu'))
16      model.add(Dense(1, activation='sigmoid'))
17      # Compile model
18      model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
19      return model
20
21  # fix random seed for reproducibility
22  seed = 7
23  numpy.random.seed(seed)
24  # load pima indians dataset
25  dataset = numpy.loadtxt("https://raw.githubusercontent.com/jbrownlee/Datasets/master/pima-in
26  # split into input (X) and output (Y) variables
27  X = dataset[:,0:8]
28  Y = dataset[:,8]
29  # create model
30  model = KerasClassifier(build_fn=create_model, epochs=150, batch_size=10, verbose=0)
31  # evaluate using 10-fold cross validation
32  kfold = StratifiedKFold(n_splits=10, shuffle=True, random_state=seed)
33  results = cross_val_score(model, X, Y, cv=kfold)
34  print(results.mean())
```

In the above, we used the wrapper from Keras. Other wrappers exist, such as scikeras. All it does is to make sure the **interface** of the Keras model looks like a scikit-learn classifier so you can make use of the `cross_val_score()` function. If we replace the `model` above with:

```
1  model = create_model()
```

then the scikit-learn function will complain as it cannot find the `model.score()` function.

Similarly, because of duck typing, we can reuse a function that expects a list for a NumPy array or pandas series because they all support the same indexing and slicing operation. For example, we fit a time series with ARIMA as follows:

```
 1  from statsmodels.tsa.statespace.sarimax import SARIMAX
 2  import numpy as np
 3  import pandas as pd
 4
 5  data = [266.0,145.9,183.1,119.3,180.3,168.5,231.8,224.5,192.8,122.9,336.5,185.9,
 6          194.3,149.5,210.1,273.3,191.4,287.0,226.0,303.6,289.9,421.6,264.5,342.3,
 7          339.7,440.4,315.9,439.3,401.3,437.4,575.5,407.6,682.0,475.3,581.3,646.9]
 8  model = SARIMAX(y, order=(5,1,0))
 9  res = model.fit(disp=False)
10  print("AIC = ", res.aic)
11
12  data = np.array(data)
13  model = SARIMAX(y, order=(5,1,0))
14  res = model.fit(disp=False)
15  print("AIC = ", res.aic)
16
17  data = pd.Series(data)
18  model = SARIMAX(y, order=(5,1,0))
19  res = model.fit(disp=False)
20  print("AIC = ", res.aic)
```

The above should produce the same AIC scores for each fitting.

# Scopes and Name Space in Python

In most languages, variables are defined in a limited scope. For example, a variable defined inside a function is accessible only inside that function:

```
1  from math import sqrt
2
3  def quadratic(a,b,c):
4      discrim = b*b - 4*a*c
5      x = -b/(2*a)
```

```
6        y = sqrt(discrim)/(2*a)
7        return x-y, x+y
```

The **local variable** `discrim` is in no way accessible if we are not inside the function `quadratic()`. Moreover, this may be surprising for someone:

```
1  a = 1
2
3  def f(x):
4      a = 2 * x
5      return a
6
7  b = f(3)
8  print(a, b)
```

```
1  1 6
```

We defined the variable a outside function f, but inside f, variable a is assigned to be 2 * x. However, the a inside the function and the one outside are unrelated except for the name. Therefore, as we exit from the function, the value of a is untouched. To make it modifiable inside function f, we need to declare the name a as `global` to make it clear that this name should be from the **global scope**, not the **local scope**:

```
1  a = 1
2
3  def f(x):
4      global a
5      a = 2 * x
6      return a
7
8  b = f(3)
9  print(a, b)
```

```
1  6 6
```

However, we may further complicate the issue when introducing the **nested scope** in functions. Consider the following example:

```
1   a = 1
2
3   def f(x):
4       a = x
5       def g(x):
6           return a * x
7       return g(3)
8
9   b = f(2)
10  print(b)
```

```
1  6
```

The variable a inside function f is distinct from the global one. However, when inside g, since there is never anything written to a but merely read from it, Python will see the same a from the nearest scope, i.e., from function f. The variable x, however, is defined as an argument to the function g, and it takes the value 3 when we called g(3) instead of assuming the value of x from function f.

**NOTE:** If a variable has any value assigned to it **anywhere** in the function, it is defined in the local scope. And if that variable has its value read from it before the assignment, an error is raised rather than using the value from the variable of the same name from the outer or global scope.

This property has many uses. Many implementations of memoization decorators in Python make clever use of the function scopes. Another example is the following:

```python
1   import numpy as np
2
3   def datagen(X, y, batch_size, sampling_rate=0.7):
4       """A generator to produce samples from input numpy arrays X and y
5       """
6       # Select rows from arrays X and y randomly
7       indexing = np.random.random(len(X)) < sampling_rate
8       Xsam, ysam = X[indexing], y[indexing]
9
10      # Actual logic to generate batches
11      def _gen(batch_size):
12          while True:
13              Xbatch, ybatch = [], []
14              for _ in range(batch_size):
15                  i = np.random.randint(len(Xsam))
16                  Xbatch.append(Xsam[i])
17                  ybatch.append(ysam[i])
18              yield np.array(Xbatch), np.array(ybatch)
19
20      # Create and return a generator
21      return _gen(batch_size)
```

This is a **generator function** that creates batches of samples from the input NumPy arrays X and y. Such a generator is acceptable by Keras models in their training. However, for reasons such as cross validation, we do not want to sample from the entire input arrays X and y but a **fixed** subset of rows from them. The way we do it is to randomly select a portion of rows at the beginning of the datagen() function and keep them in Xsam, ysam. Then in the inner function _gen(), rows are sampled from Xsam and ysam until a batch is created. While the lists Xbatch and ybatch are defined and created inside the function _gen(), the arrays Xsam and ysam are not local to _gen(). What's more interesting is when the generator is created:

```python
1   X = np.random.random((100,3))
2   y = np.random.random(100)
3
4   gen1 = datagen(X, y, 3)
5   gen2 = datagen(X, y, 4)
6   print(next(gen1))
7   print(next(gen2))
```

```
1   (array([[0.89702235, 0.97516228, 0.08893787],
2          [0.26395301, 0.37674529, 0.1439478 ],
3          [0.24859104, 0.17448628, 0.41182877]]), array([0.2821138 , 0.87590954, 0.96646776]))
4   (array([[0.62199772, 0.01442743, 0.4897467 ],
5          [0.41129379, 0.24600387, 0.53640666],
6          [0.02417213, 0.27637708, 0.65571031],
7          [0.15107433, 0.11331674, 0.67000849]]), array([0.91559533, 0.84886957, 0.30451455, 0.5
```

The function datagen() is called two times, and therefore two different sets of Xsam, yam are created. But since the inner function _gen() depends on them, these two sets of Xsam, ysam are in memory concurrently. Technically, we say that when datagen() is called, a **closure** is created with the specific Xsam, ysam defined within, and the call to _gen() is accessing that closure. In other words, the scopes of the two incarnations of datagen() calls coexist.

In summary, whenever a line of code references a name (whether it is a variable, a function, or a module), the name is resolved in the order of the LEGB rule:

1. Local scope first, i.e., those names that were defined in the same function
2. Enclosure or the "nonlocal" scope. That's the upper-level function if we are inside the nested function.
3. Global scope, i.e., those that were defined in the top level of the same script (but not across different program files)
4. Built-in scope, i.e., those created by Python automatically, such as the variable __name__ or functions `list()`

---

---

# Investigating the type and scope

Because the types are not static in Python, sometimes we would like to know what we are dealing with, but it is not trivial to tell from the code. One way to tell is using the `type()` or `isinstance()` functions. For example:

```
1  import numpy as np
2
3  X = np.random.random((100,3))
4  print(type(X))
5  print(isinstance(X, np.ndarray))
```

```
1  <class 'numpy.ndarray'>
2  True
```

The `type()` function returns a type object. The `isinstance()` function returns a Boolean that allows us to check if something matches a particular type. These are useful in case we need to know what type a

variable is. This is useful if we are debugging a code. For example, if we pass on a pandas dataframe to the `datagen()` function that we defined above:

```python
1   import pandas as pd
2   import numpy as np
3
4   def datagen(X, y, batch_size, sampling_rate=0.7):
5       """A generator to produce samples from input numpy arrays X and y
6       """
7       # Select rows from arrays X and y randomly
8       indexing = np.random.random(len(X)) < sampling_rate
9       Xsam, ysam = X[indexing], y[indexing]
10
11      # Actual logic to generate batches
12      def _gen(batch_size):
13          while True:
14              Xbatch, ybatch = [], []
15              for _ in range(batch_size):
16                  i = np.random.randint(len(Xsam))
17                  Xbatch.append(Xsam[i])
18                  ybatch.append(ysam[i])
19              yield np.array(Xbatch), np.array(ybatch)
20
21      # Create and return a generator
22      return _gen(batch_size)
23
24  X = pd.DataFrame(np.random.random((100,3)))
25  y = pd.DataFrame(np.random.random(100))
26
27  gen3 = datagen(X, y, 3)
28  print(next(gen3))
```

Running the above code under the Python's debugger `pdb` will give the following:

```
1   > /Users/MLM/ducktype.py(1)<module>()
2   -> import pandas as pd
3   (Pdb) c
4   Traceback (most recent call last):
5     File "/usr/local/lib/python3.9/site-packages/pandas/core/indexes/range.py", line 385, in
6       return self._range.index(new_key)
7   ValueError: 1 is not in range
8
9   The above exception was the direct cause of the following exception:
10
11  Traceback (most recent call last):
12    File "/usr/local/Cellar/python@3.9/3.9.9/Frameworks/Python.framework/Versions/3.9/lib/pyth
13      pdb._runscript(mainpyfile)
14    File "/usr/local/Cellar/python@3.9/3.9.9/Frameworks/Python.framework/Versions/3.9/lib/pyth
15      self.run(statement)
16    File "/usr/local/Cellar/python@3.9/3.9.9/Frameworks/Python.framework/Versions/3.9/lib/pyth
17      exec(cmd, globals, locals)
18    File "<string>", line 1, in <module>
19    File "/Users/MLM/ducktype.py", line 1, in <module>
20      import pandas as pd
21    File "/Users/MLM/ducktype.py", line 18, in _gen
22      ybatch.append(ysam[i])
23    File "/usr/local/lib/python3.9/site-packages/pandas/core/frame.py", line 3458, in __getite
24      indexer = self.columns.get_loc(key)
25    File "/usr/local/lib/python3.9/site-packages/pandas/core/indexes/range.py", line 387, in
26      raise KeyError(key) from err
27  KeyError: 1
28  Uncaught exception. Entering post mortem debugging
29  Running 'cont' or 'step' will restart the program
30  > /usr/local/lib/python3.9/site-packages/pandas/core/indexes/range.py(387)get_loc()
31  -> raise KeyError(key) from err
32  (Pdb)
```

We see from the traceback that something is wrong because we cannot get `ysam[i]`. We can use the following to verify that `ysam` is indeed a Pandas DataFrame instead of a NumPy array:

```
1  (Pdb) up
2  > /usr/local/lib/python3.9/site-packages/pandas/core/frame.py(3458)__getitem__()
3  -> indexer = self.columns.get_loc(key)
4  (Pdb) up
5  > /Users/MLM/ducktype.py(18)_gen()
6  -> ybatch.append(ysam[i])
7  (Pdb) type(ysam)
8  <class 'pandas.core.frame.DataFrame'>
```

Therefore we cannot use `ysam[i]` to select row `i` from `ysam`. What can we do in the debugger to verify how we should modify our code? There are several useful functions you can use to investigate the variables and the scope:

- `dir()` to see the names defined in the scope or the attributes defined in an object
- `locals()` and `globals()` to see the names and values defined locally and globally, respectively.

For example, we can use `dir(ysam)` to see what attributes or functions are defined inside `ysam`:

```
1  (Pdb) dir(ysam)
2  ['T', '_AXIS_LEN', '_AXIS_ORDERS', '_AXIS_REVERSED', '_AXIS_TO_AXIS_NUMBER',
3  ...
4  'iat', 'idxmax', 'idxmin', 'iloc', 'index', 'infer_objects', 'info', 'insert',
5  'interpolate', 'isin', 'isna', 'isnull', 'items', 'iteritems', 'iterrows',
6  'itertuples', 'join', 'keys', 'kurt', 'kurtosis', 'last', 'last_valid_index',
7  ...
8  'transform', 'transpose', 'truediv', 'truncate', 'tz_convert', 'tz_localize',
9  'unstack', 'update', 'value_counts', 'values', 'var', 'where', 'xs']
10 (Pdb)
```

Some of these are attributes, such as `shape`, and some of these are functions, such as `describe()`. You can read the attribute or invoke the function in `pdb`. By carefully reading this output, we recalled that the way to read row `i` from a DataFrame is through `iloc`, and hence we can verify the syntax with:

```
1  (Pdb) ysam.iloc[i]
2  0     0.83794
3  Name: 2, dtype: float64
4  (Pdb)
```

If we call `dir()` without any argument, it gives you all the names defined in the current scope, e.g.,

```
1  (Pdb) dir()
2  ['Xbatch', 'Xsam', '_', 'batch_size', 'i', 'ybatch', 'ysam']
3  (Pdb) up
4  > /Users/MLM/ducktype.py(1)<module>()
5  -> import pandas as pd
6  (Pdb) dir()
7  ['X', '__builtins__', '__file__', '__name__', 'datagen', 'gen3', 'np', 'pd', 'y']
8  (Pdb)
```

where the scope changes as you move around the call stack. Similar to `dir()` without argument, we can call `locals()` to show all locally defined variables, e.g.,

```
1  (Pdb) locals()
2  {'batch_size': 3, 'Xbatch': ...,
3   'ybatch': ..., '_': 0, 'i': 1, 'Xsam': ...,
4   'ysam': ...}
5  (Pdb)
```

Indeed, `locals()` returns you a `dict` that allows you to see all the names and values. Therefore, if we need to read the variable `Xbatch`, we can get the same with `locals()["Xbatch"]`. Similarly, we can use `globals()` to get a dictionary of names defined in the global scope.

This technique is beneficial sometimes. For example, we can check if a Keras model is "compiled" or not by using `dir(model)`. In Keras, compiling a model is to set up the loss function for training and build the flow for forward and backward propagations. Therefore, a compiled model will have an extra attribute `loss` defined:

```
1   from tensorflow.keras import Sequential
2   from tensorflow.keras.layers import Dense
3
4   model = Sequential([
5       Dense(5, input_shape=(3,)),
6       Dense(1)
7   ])
8
9   has_loss = "loss" in dir(model)
10  print("Before compile, loss function defined:", has_loss)
11
12  model.compile()
13  has_loss = "loss" in dir(model)
14  print("After compile, loss function defined:", has_loss)
```

```
1  Before compile, loss function defined: False
2  After compile, loss function defined: True
```

This allows us to put an extra guard on our code before we run into an error.