# STEP-BY-STEP TO A DATA SCIENTIST

Introduction of basic skills for a Data Scientist



OCTOBER 25, 2020 BY CARO_CARO

# Beginner Guide to Gaussian Process by GPy

Tweet          **Share** 0          Share          **Like** 0

Step-by-step to a Data Scientist > Blog > for beginner > Beginner Guide to Gaussian Process by GPy

The main differences between Gaussian Process and linear regression are below.

1. Modeling with nonlinear
2. Model has information on both the estimation and the uncertainty.

The first one is that since this method can handle non-linearity, it is a highly expressive model. For example, linear regression analyses assumed the linear relationship between input variables and the output. Therefore, in principle, linear regression analysis is not suitable for datasets with non-linear relationships. In contrast, in such a dataset, we can apply the Gaussian Process regression.

The second one is that the output of the model has both information about the regression values and the confidence of the output. In other words, the Gaussian model tells us the confidence of results. This is completely different from linear regression.

Let's see these differences concretely with a simple example. The complete notebook can be found on **GitHub**.

# Create Training Dataset

Here, we create the training dataset. The way is to add the noise to the base function.

$$y = -1 + x + 2x^2 + 3e^x + \varepsilon(x).$$

$\varepsilon(x)$ is the noise function and described as follows:

$$\varepsilon(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

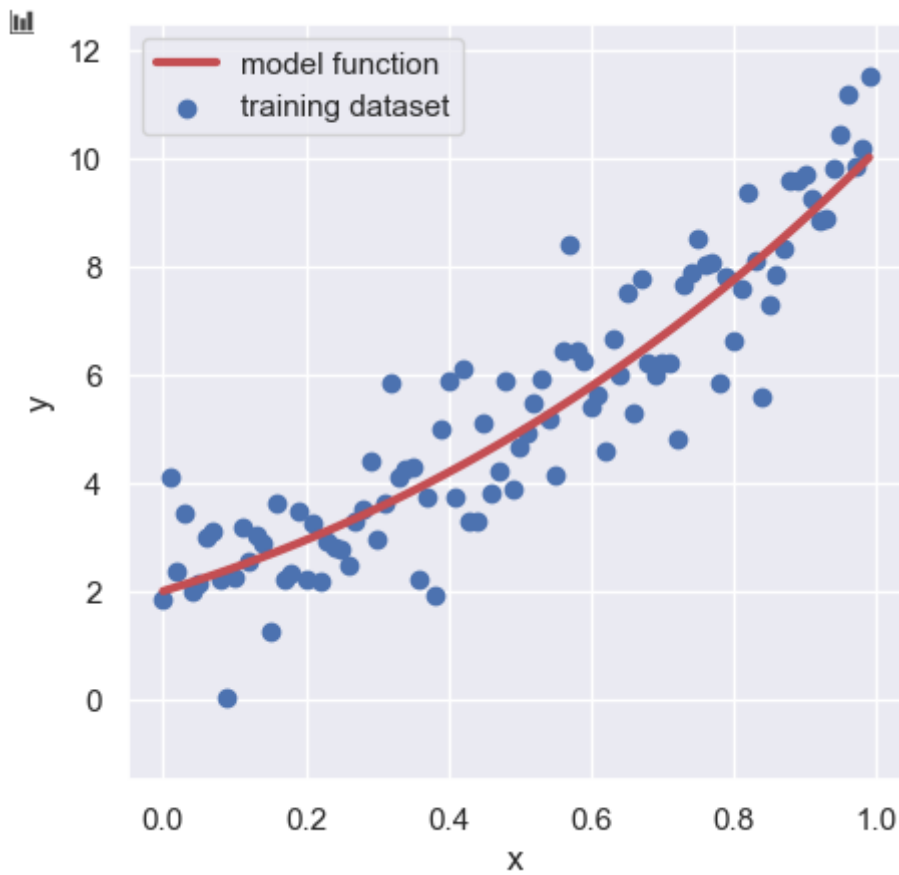where $\mu$ is the mean and $\sigma$ the standard deviation.

The code is below. The derails are introduced in another post.

Python

```python
1   import numpy as np
2   ##-- Model Function for creating the train dataset
3   def func(param, X):
4       return param[0] + param[1]*X + param[2]*np.power(X, 2) + param[3]*np.exp
5
6   ##-- Set Random Seed
7   np.random.seed(seed=99)
8
9   x = np.arange(0, 1, 0.01)
10  param = [-1.0, 1.0, 2.0, 3.0]
11
12  y_model = func(param, x)
13  y_train = func(param, x) + np.random.normal(loc=0, scale=1.0, size=Len(x))
```

## Create a Toy Dataset by the Noise Function

The toy dataset is useful when we attempt something new analysis method quickly, especially in regression analyses. In this post, we briefly see how to create a toy dataset with NumPy. Then, let's get started. Import the Libraries The code is written by Python. So firstly, we import the necessary library. Define the Model Function … Continue reading

**Step-by-step to a Data Scientist**

# Gaussian Process Regression

In this analysis, we use "*GPy*", the Gaussian Process library in Python. In this post, the version of *GPy* is assumed for "*ver. 1.9.9*".

Python

```python
1  import GPy
2  print(GPy.__version__)
3
4  >> 1.9.9
```

# Define "*the kernel function*"

First, we define the kernel function. There are many kinds of kernel functions. Here, we

Python

```
1 │ kernel = GPy.kern.Matern52(input_dim=1)
```

The option *"input_dim"* is the number of input variables. In this case, the input variable is just one *"x"*, so *"input_dim = 1"*.

# Define the Model

Second, we define the Gaussian Process model as follows:

Python

```
1 │ model = GPy.models.GPRegression(x.reshape(-1, 1), y_train.reshape(-1, 1), ke
```

There are three arguments, the input variable *(x, y)* and the kernel function.

Note that the dimensions of the input variable*(x, y)* must be two-dimensional. This is where beginners get stuck in debugging. We can easily confirm the dimensions as follows:

Python

```
1 │ x_train.ndim  # x_train.ndim
2 │
3 │ >> 1
```

Whereas,

Python

```
1 │ x_train.reshape(-1, 1).ndim  # y_train.reshape(-1, 1).ndim
2 │
3 │ >> 2
```

# Optimize the Model

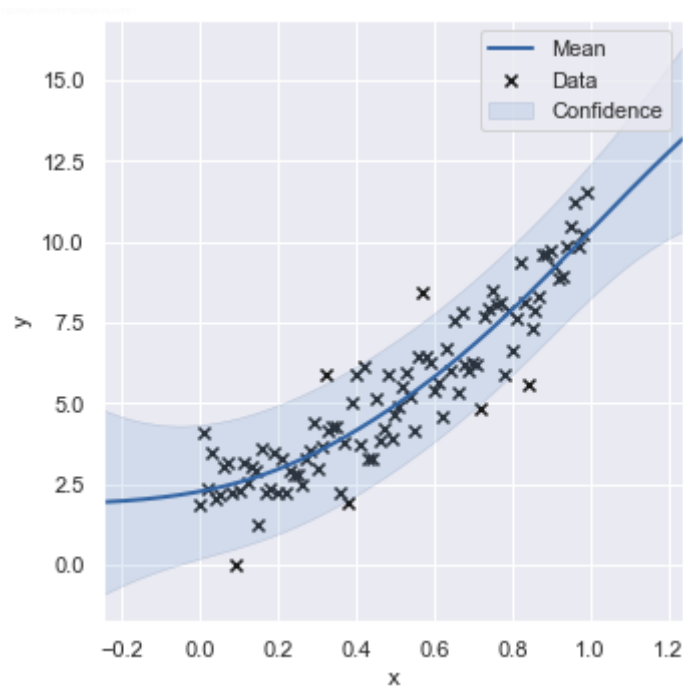Now that we have defined the model, we can optimize it.

Python

```
1 │ model.optimize()
```

*GPy* includes *matplotlib* inside, so you can quickly see the optimized model in just one line. Note that this method can be used when the number of the input variables is one. In multi variables case, the ability to plot in multiple dimensions is not supported.

Python

```
1 │ model.plot(figsize=(5, 5), dpi=100, xlabel="x", ylabel="y")
```



˅ D  I  I

Since the training data was in the range of 0 to 1, we will prepare the test data in the range of 0 to 2.

Python

```
1   x_test = np.arange(0, 2, 0.01)
```

Then, predict the test data by the *".predict()"* method.

Python

```
1   y_mean, y_std = model.predict(x_test.reshape(-1, 1))
```

Note that there are two returned values from *"model.predict()"*. The first one is *the regression value* and the second one is *the confidence interval*. Here, it's okay to think of the confidence interval as the standard deviation in the Gaussian function.

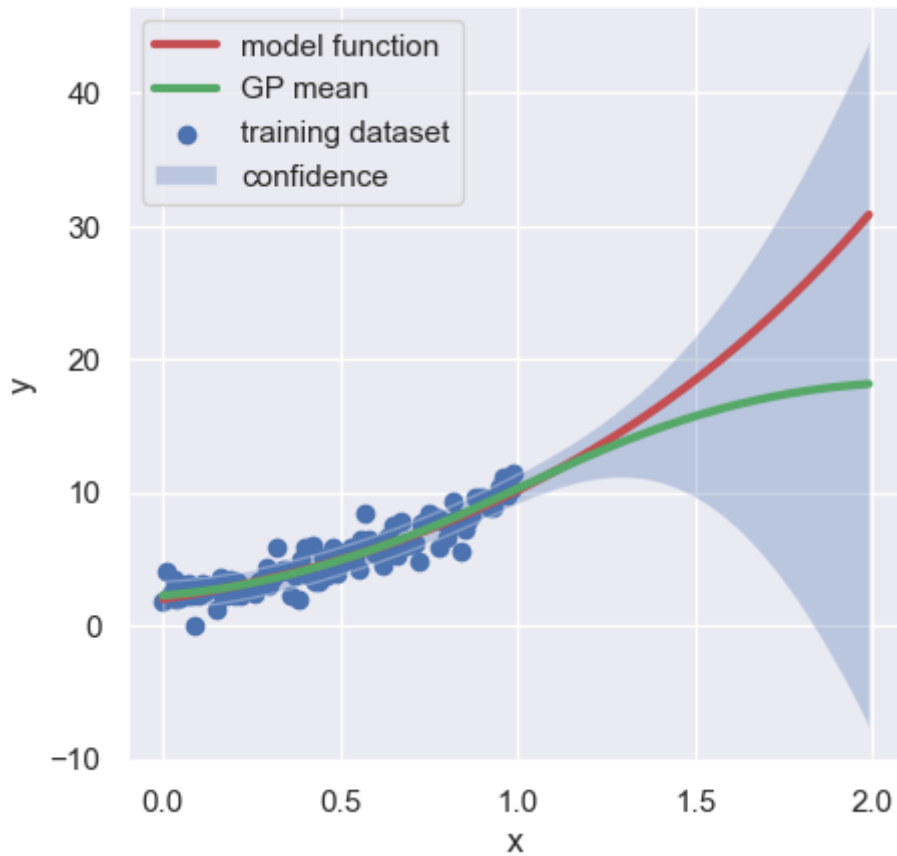Finally, let's plot the result.

Python

```
1   import matplotlib.pylab as plt
2   import seaborn as sns
3
4   plt.figure(figsize=(5, 5), dpi=100)
5   sns.set()
6   plt.xlabel("x")
7   plt.ylabel("y")
8   plt.scatter(x, y_train, lw=1, color="b", label="training dataset")
9   plt.plot(x_test, func(param, x_test), lw=3, color="r", label="model function
10  plt.plot(x_test, y_mean, lw=3, color="g", label="GP mean")
11  plt.fill_between(x_test, (y_mean + y_std).reshape(y_mean.shape[0]), (y_mean
12  plt.legend(loc="upper left")
13  plt.show()
```

Congratulations!! It can be confirmed that as the range of the test data deviates from the training data, the predicted value also deviates and the confidence interval becomes wider.

## Summary

From here, we have seen the process of Gaussian process regression with a simple example. Certainly, Gaussian process regression has the impression that it is difficult to attract attention as an analysis method because of its mathematical complexity. However, Gaussian process regression has information on confidence, making it possible to judge the validity of prediction.

The author would be happy if this post helps the reader to try Gaussian process analyses.