

A Generative Approach to Classification

A Simple Explanation



Rahul Agarwal

Aug 25 · 5 min read



Image by Prashant Sharma from Pixabay

I always get confused whenever someone talks about generative vs. discriminative classification models.

I end up reading it again and again, yet somehow it eludes me.

So I thought of writing a post on it to improve my understanding.

This post is about understanding Generative Models and how they differ from Discriminative models.

In the end, we will create a simple generative model ourselves.

. . .

Discriminative vs. Generative Classifiers

Problem Statement: Having some input data, x we want to classify the data into labels y .

A generative model learns the **joint** probability distribution $p(x, y)$ while a discriminative model learns the **conditional** probability distribution $p(y|x)$

So really, what is the difference? They both look pretty much the same.

Suppose we have a small sample of data:

$$(x, y) : [(0, 1), (1, 0), (1, 0), (1, 1)]$$

Then $p(x, y)$ is

	$y=0$	$y=1$
$x=0$	0	0.25
$x=1$	0.50	0.25

While $p(y|x)$ is

	$y=0$	$y=1$
$x=0$	0	1
$x=1$	0.67	0.33

As you can see, they model different probabilities.

The discriminative distribution $p(y|x)$ could be used straightforward to classify an example x into a class y . An example of a discriminative classification model is Logistic regression, where we try to model $P(y|X)$.

$$P(Y=1 | X=x) = \frac{1}{1+e^{-z}}$$

Logistic Regression

Generative algorithms model $p(x,y)$. An example is the Naive Bayes model in which we try to model $P(X,y)$ and then use the Bayes equation to predict.

...

The Central Idea Behind Generative Classification

- Fit each class separately with a probability distribution.
- To classify a new point, find out which distribution is it most probable to come from.

Don't worry if you don't understand yet. You will surely get it by the end of this post.

...

A Small Example



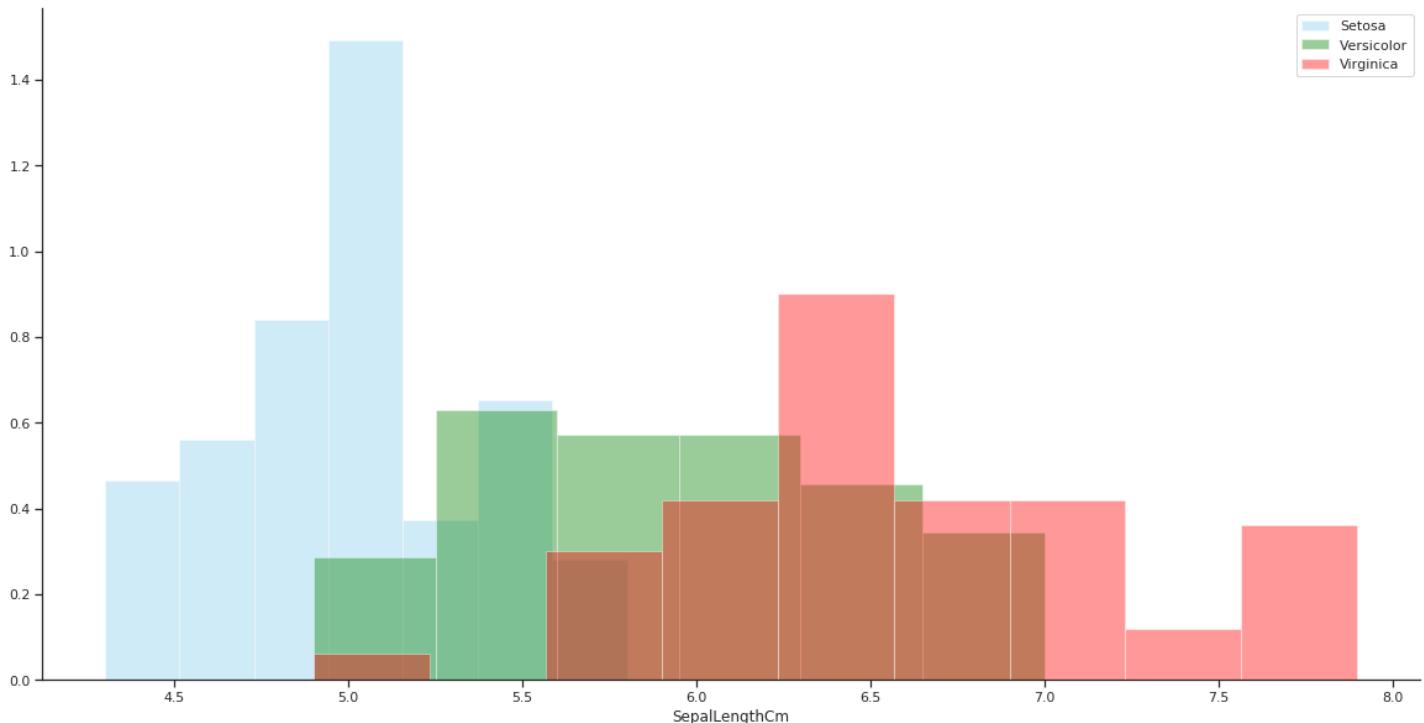


Let us work with the iris dataset.

For our simple example, we will work with a single x variable SepalLength and our target variable Species.

Let us see the distribution of sepal length with Species. I am using plotly_express for this.

```
import plotly_express as px
px.histogram(iris, x = 'SepalLengthCm', color = 'Species', nbins=20)
```



To create generative models, we need to find out two sets of values:

1. Probability of individual classes:

To get individual class probability is fairly trivial- For example, the number of instances in our dataset, which is setosa divided by the total number of cases in the dataset.

```

1 p_setosa = len(iris[iris['Species']=='Iris-setosa'])/len(iris)
2 p_versicolor = len(iris[iris['Species']=='Iris-versicolor'])/len(iris)
3 p_virginica = len(iris[iris['Species']=='Iris-virginica'])/len(iris)
4 print(p_setosa,p_versicolor,p_virginica)
```

probs.py hosted with ❤ by GitHub

[view raw](#)

0.3333333333333333 0.3333333333333333 0.3333333333333333

The iris dataset is pretty much balanced.

2. The probability distribution of x for each class:

Here we fit a probability distribution over our X. We assume here that the X data is distributed normally. And hence we can find the sample means and variance for these three distributions(As we have three classes)

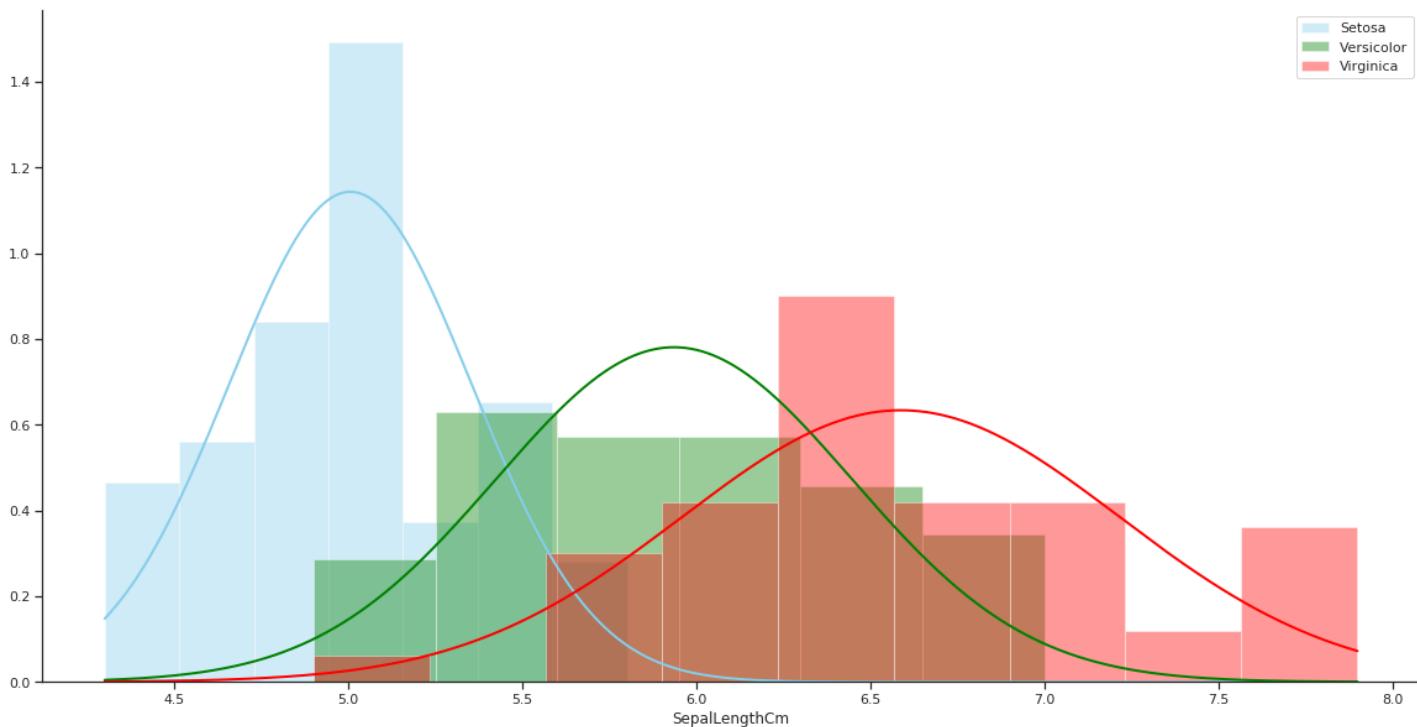
```

1 import numpy as np
2 import seaborn as sns
3 from scipy import stats
4 import matplotlib.pyplot as plt
5
6 sns.set(style="ticks")
7 # calculate the pdf over a range of values
8 xx = np.arange(min(iris['SepalLengthCm']), max(iris['SepalLengthCm']), 0.001)
9 x = iris[iris['Species']=='Iris-setosa']['SepalLengthCm']
10 sns.distplot(x, kde = False, norm_hist=True, color='skyblue', label = 'Setosa')
11 yy = stats.norm.pdf(xx, loc=np.mean(x), scale=np.std(x))
12 plt.plot(xx, yy, 'skyblue', lw=2)
13 x = iris[iris['Species']=='Iris-versicolor']['SepalLengthCm']
14 sns.distplot(x, kde = False, norm_hist=True, color='green', label = 'Versicolor')
15 yy = stats.norm.pdf(xx, loc=np.mean(x), scale=np.std(x))
16 plt.plot(xx, yy, 'green', lw=2)
17 x = iris[iris['Species']=='Iris-virginica']['SepalLengthCm']
18 g = sns.distplot(x, kde = False, norm_hist=True, color='red', label = 'Virginica')
19 yy = stats.norm.pdf(xx, loc=np.mean(x), scale=np.std(x))
20 plt.plot(xx, yy, 'red', lw=2)
21 sns.despine()
22 g.figure.set_size_inches(20,10)
```

23 g.legend()

plot.py hosted with ❤ by GitHub

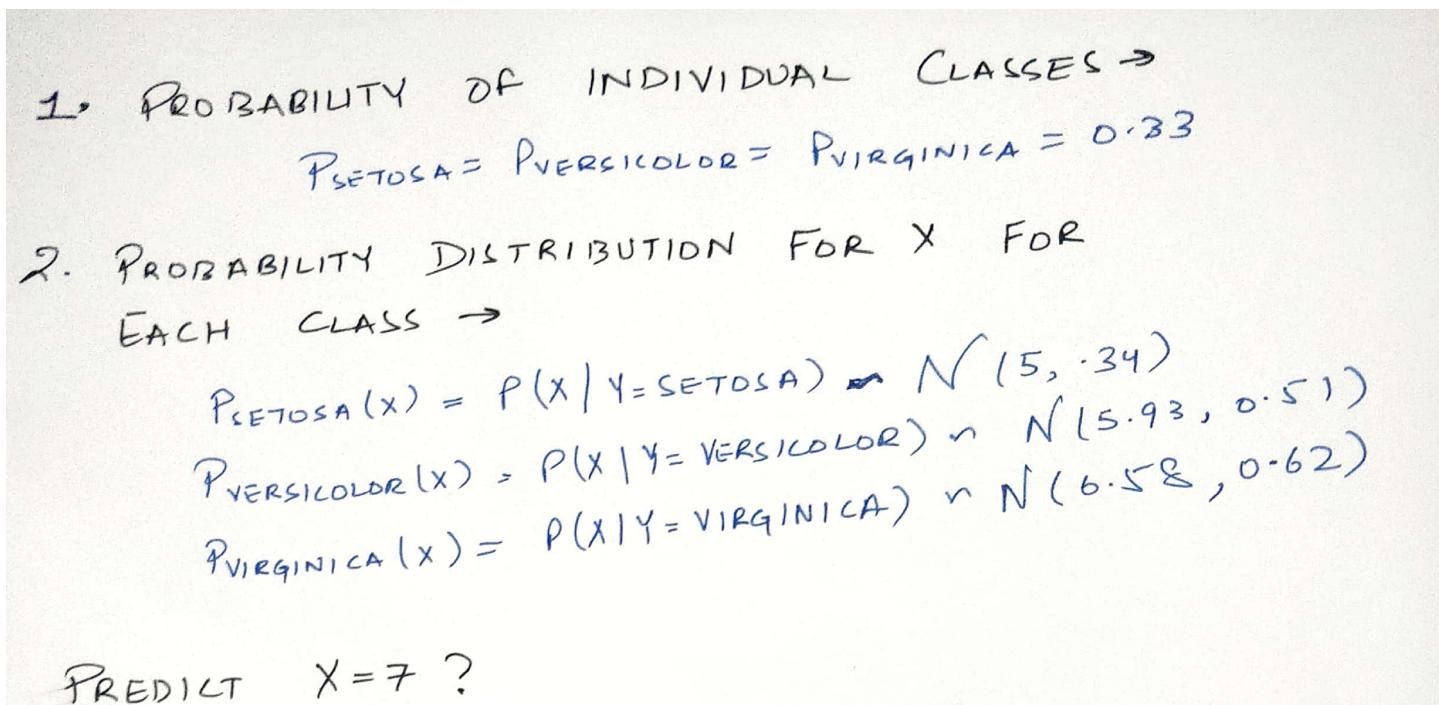
view raw



In the above graph, I have fitted three normal distributions for each of the species just using sample means and variances for each of the three species.

So, how do we predict using this?

Let us say we get a new example with SepalLength = 7 cm.



$$\text{BAYES} \rightarrow P(X, Y) = P(X|Y) \times P(Y)$$

Thus,

$$\begin{aligned} P(X=7, Y=\text{SETOSA}) &= P(X=7|Y=\text{SETOSA}) \times P(Y=\text{SETOSA}) \\ &= P_{\text{SETOSA}}(X=7) \times P_{\text{SETOSA}} \\ &= 3.062 \times 10^{-8} \end{aligned}$$

$$\begin{aligned} P(X=7, Y=\text{VERSICOLOR}) &= P_{\text{VERSICOLOR}}(X=7) \times P_{\text{VERSICOLOR}} \\ &= .0294 \end{aligned}$$

$$\begin{aligned} P(X=7, Y=\text{VIRGINICA}) &= P_{\text{VIRGINICA}}(X=7) \times P_{\text{VIRGINICA}} \\ &= \underline{\underline{.1688}} \end{aligned}$$

Since we see that the maximum probability comes for Virginica, we predict virginica for $x=7$, and based on the graph too; it looks pretty much the right choice.

You can get the values using the code too.

```

1 x = iris[iris['Species']=='Iris-setosa']['SepalLengthCm']
2 print("Setosa",stats.norm.pdf(7,loc=np.mean(x),scale=np.std(x))*.33)
3 x = iris[iris['Species']=='Iris-versicolor']['SepalLengthCm']
4 print("Versicolor",stats.norm.pdf(7,loc=np.mean(x),scale=np.std(x))*.33)
5 x = iris[iris['Species']=='Iris-virginica']['SepalLengthCm']
6 print("Virginica",stats.norm.pdf(7,loc=np.mean(x),scale=np.std(x))*.33)

```

calculation.py hosted with ❤ by GitHub

[view raw](#)

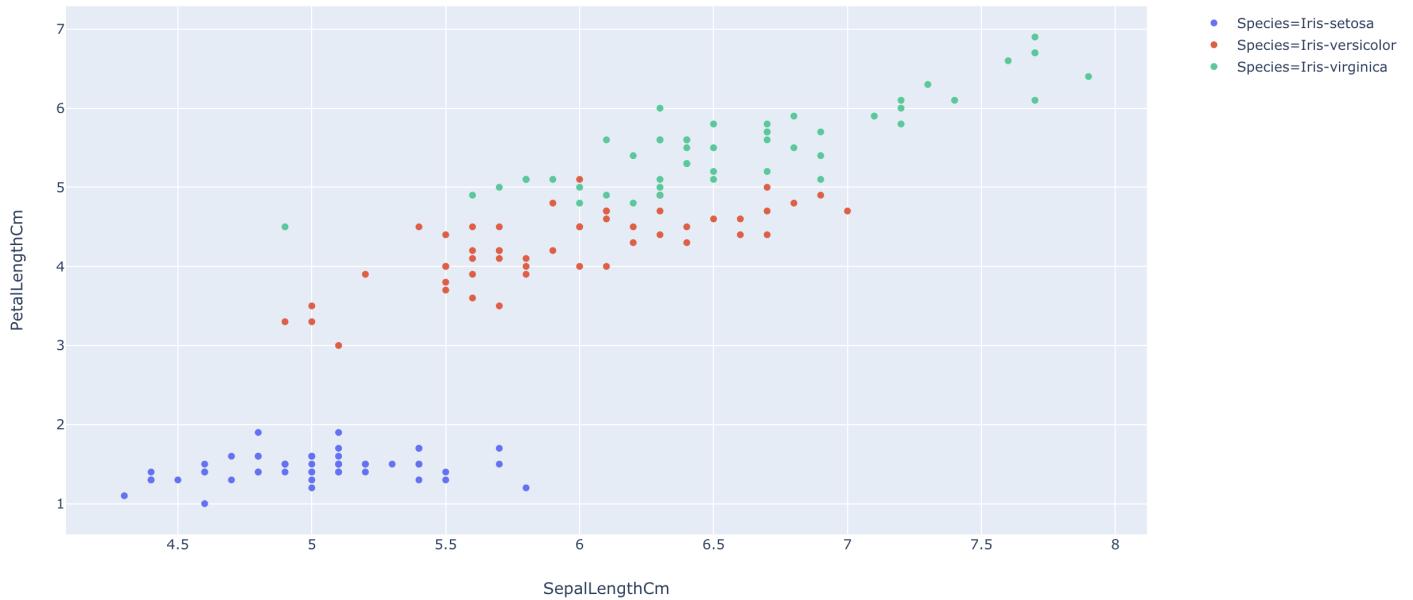
Setosa 3.062104211904799e-08
 Versicolor 0.029478757465669376
Virginica 0.16881724812694823

• • •

This is all well and good. But when do we ever work with a single variable?

Let us extend our example for two variables. This time let us use PetalLength too.

```
px.scatter(iris, 'SepalLengthCm', 'PetalLengthCm', color = 'Species')
```



So how do we proceed in this case?

The first time we had fit a Normal Distribution over our single x, this time we will fit Bivariate Normal.

```

1 import numpy as np
2 import seaborn as sns
3 from scipy import stats
4 import matplotlib.pyplot as plt
5 from matplotlib.mlab import bivariate_normal
6 sns.set(style="ticks")
7 # SETOSA
8 x1 = iris[iris['Species']=='Iris-setosa']['SepalLengthCm']
9 x2 = iris[iris['Species']=='Iris-setosa']['PetalLengthCm']
10 sns.scatterplot(x1,x2, color='skyblue',label = 'Setosa')
11 mu_x1=np.mean(x1)
12 mu_x2=np.mean(x2)
13 sigma_x1=np.std(x1)**2
14 sigma_x2=np.std(x2)**2
15 xx = np.arange(min(x1), max(x1),0.001)
16 yy = np.arange(min(x2), max(x2),0.001)
17 X, Y = np.meshgrid(xx, yy)
18 Z = bivariate_normal(X,Y, sigma_x1, sigma_x2, mu_x1, mu_x2)
```

```

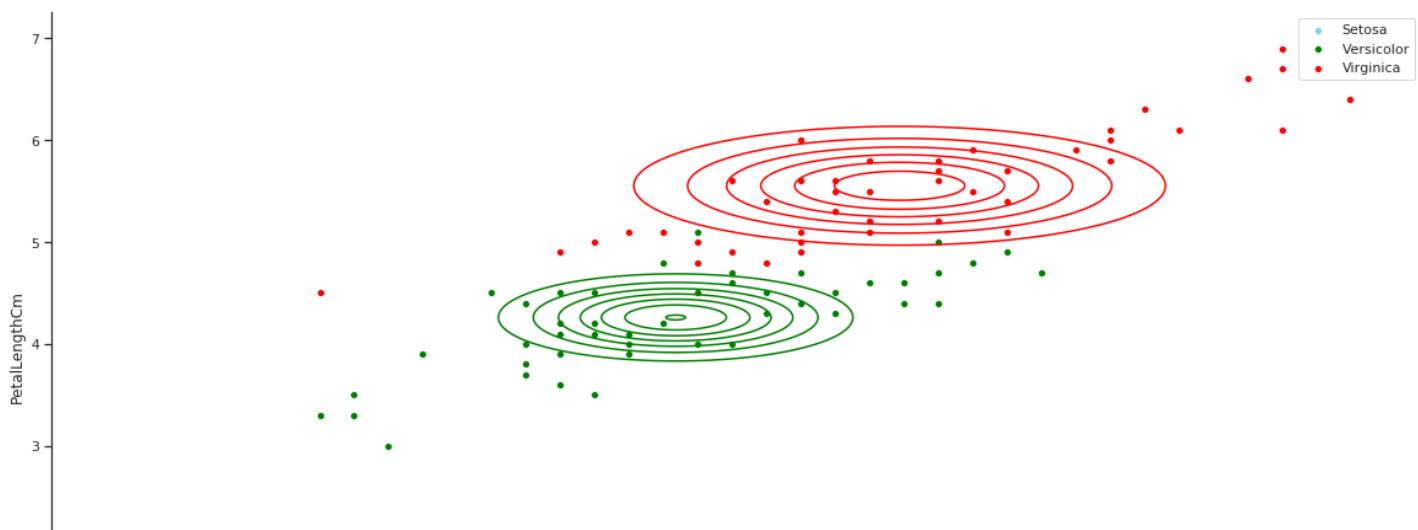
19 plt.contour(X,Y,Z,colors='skyblue')
20 # VERSICOLOR
21 x1 = iris[iris['Species']=='Iris-versicolor']['SepalLengthCm']
22 x2 = iris[iris['Species']=='Iris-versicolor']['PetalLengthCm']
23 sns.scatterplot(x1,x2,color='green',label = 'Versicolor')
24 mu_x1=np.mean(x1)
25 mu_x2=np.mean(x2)
26 sigma_x1=np.std(x1)**2
27 sigma_x2=np.std(x2)**2
28 xx = np.arange(min(x1), max(x1),0.001)
29 yy = np.arange(min(x2), max(x2),0.001)
30 X, Y = np.meshgrid(xx, yy)
31 Z = bivariate_normal(X,Y, sigma_x1, sigma_x2, mu_x1, mu_x2)
32 plt.contour(X,Y,Z,colors='green')
33 # VIRGINICA
34 x1 = iris[iris['Species']=='Iris-virginica']['SepalLengthCm']
35 x2 = iris[iris['Species']=='Iris-virginica']['PetalLengthCm']
36 g = sns.scatterplot(x1, x2, color='red',label = 'Virginica')
37 mu_x1=np.mean(x1)
38 mu_x2=np.mean(x2)
39 sigma_x1=np.std(x1)**2
40 sigma_x2=np.std(x2)**2
41 xx = np.arange(min(x1), max(x1),0.001)
42 yy = np.arange(min(x2), max(x2),0.001)
43 X, Y = np.meshgrid(xx, yy)
44 Z = bivariate_normal(X,Y, sigma_x1, sigma_x2, mu_x1, mu_x2)
45 plt.contour(X,Y,Z,colors='red')
46 sns.despine()
47 g.figure.set_size_inches(20,10)
48 g.legend()

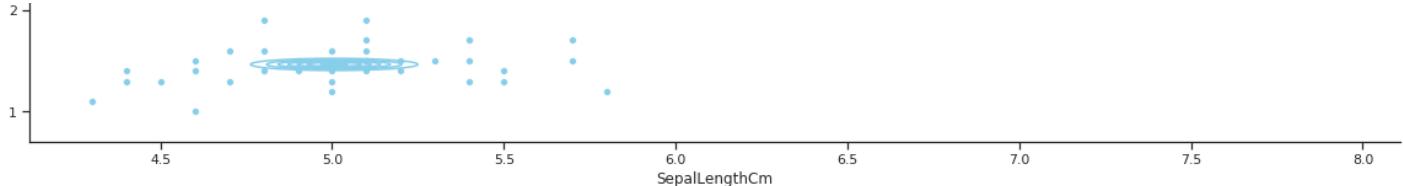
```

plot2.py hosted with ❤ by GitHub

[view raw](#)

Here is how it looks:





Now the rest of the calculations remains the same.

Just the normal gets replaced by Bivariate normal in the above equations. And as you can see, we get a pretty better separation amongst the classes by using the bivariate normal.

As an extension to this case for multiple variables (more than 2), we can use the multivariate normal distribution.

• • •

Conclusion

Generative models are good at generating data. But at the same time, creating such models that capture the underlying distribution of data is extremely hard.

Generative modeling involves a lot of assumptions, and thus, these models don't perform as well as discriminative models in the classification setting. In the above example also we assumed that the distribution is normal, which might not be correct and hence may induce a bias.

But understanding how they work is helpful all the same. One class of such models is called generative adversarial networks which are pretty useful for generating new images and are pretty interesting too.

Here is the kernel with all the code along with the visualizations.

• • •

If you want to learn more about generative models and Machine Learning, I would recommend this Machine Learning Fundamentals course from the University of San Diego. The above post is by and large inspired from content from this course in the MicroMasters from SanDiego, which I am currently working on to structure my Data Science learning.

• • •

Thanks for the read. I am going to be writing more beginner-friendly posts in the future too. Follow me up at [Medium](#) or [Subscribe to my blog](#) to be informed about them. As always, I welcome feedback and constructive criticism and can be reached on Twitter @mlwhiz.

Machine Learning

Data Science

Artificial Intelligence

Data Visualization

Programming

About Help Legal