

FDU PRML20 Assignment 1

Dataset

Click [here](#) to visit the dataset.

The dataset is constructed using the ***create_dataset*** function, who write the samples generated by the function ***generate_gaussian_distributed_samples*** to a text file. You can load the dataset by calling the function ***load_dataset*** and split it into training set and testing set by calling the function ***split_dataset***.

Description of models

A super class called ***LinearModel*** is defined in ***source.py*** for the two models and it contains attributes and method for plotting the classification result.

The linear discriminative model is defined by Python class ***LinearDiscriminativeModel*** in ***source.py***. On the linear classification of n dimensional C classes, an instance of the model has C $n \times 1$ weight vectors, i.e. an $n \times C$ weight matrix. The weight matrix is specified on calling the method ***train*** to train the model. The method train uses **gradient descent strategy** and you can specify the training options. Once the model is trained, you can classify samples by calling the ***classify*** method, who uses softmax function to generate the probability that samples belong to each class.

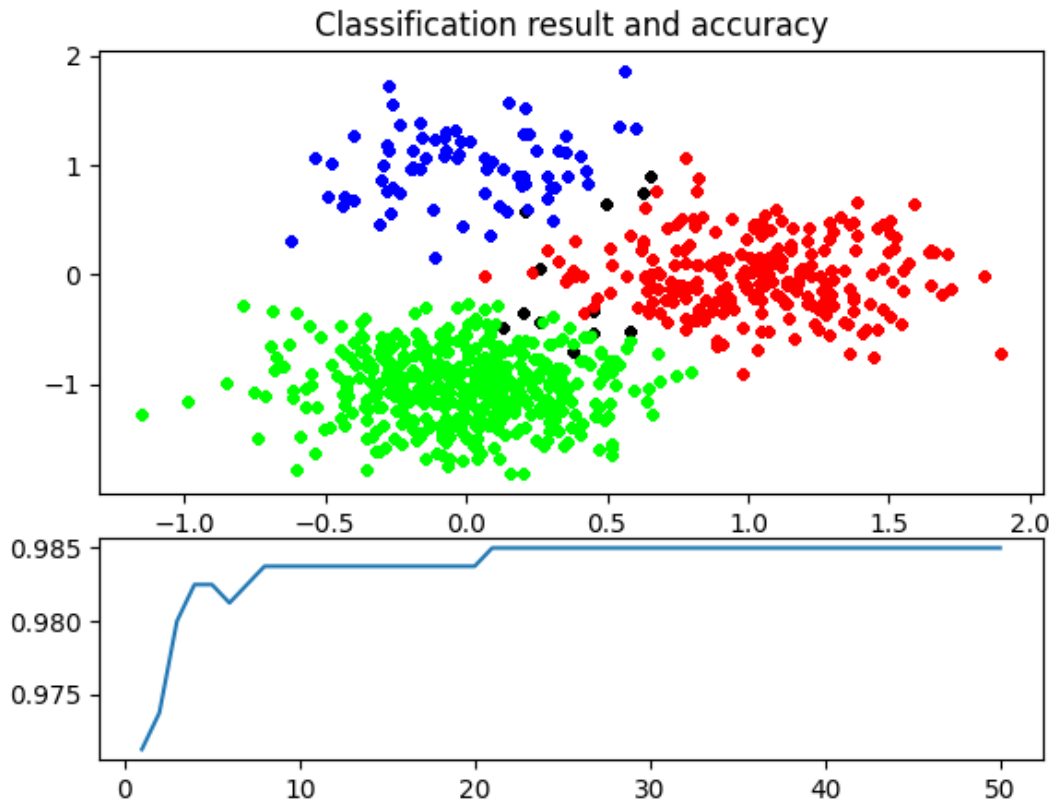
The linear generative model is defined by Python class ***LinearGenerativeModel*** in ***source.py***. On the linear classification of n dimensional C classes, an instance of the model has C $n \times 1$ weight vectors, C biases and one covariance matrix, i.e. an $n \times C$ weight matrix, a $C \times 1$ bias vector and an $n \times n$ covariance matrix. The weight matrix, bias vector and covariance matrix are specified on calling the method ***train*** to train the model. The method train calculates mathematical expectations and covariance matrices of each class using maximum likelihood estimation. Once the model is trained, you can classify samples by calling the ***classify*** method, who uses **softmax** function to generate the probability that samples belong to each class.

The **major differences** between these two models are the **training method** and **classification method**. On the **training method**, the **discriminative model** uses gradient descent strategy, i.e. an **optimization method**, to approach an ideal weight matrix, while the **generative model** tries to generalize the statistical properties of given samples using **statistical method**. As a result, the **generative model** doesn't need to train the model by iteration and can generate new samples from its estimation of samples' mathematical expectations and covariances while the **discriminative model** have to train the model by means of iteration or other approaches and can not generate new samples. On the **classification method**, the **discriminative model** assumes the probability that samples belong to each class has the form of $\text{softmax}(w^T x)$, which comes from the **linear property** of the classifier. However, the **generative model** adopts $\text{softmax}(w^T x + b)$ as the form of probability that samples belong to each class because it assumes that **samples are compiled to Gaussian distributions**. Thus the

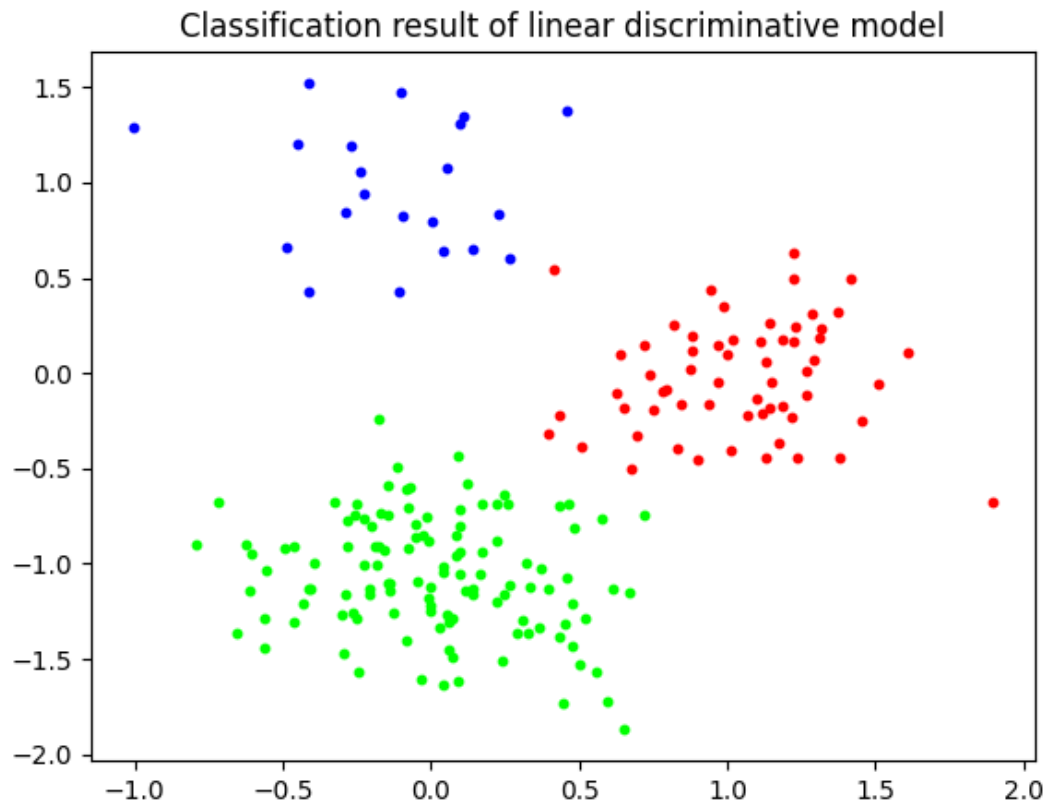
generative model is extremely relied on its assumption of samples' statistical distribution and **is sensitive to outliers or samples compiled to multi-center Gaussian distribution** while the **discriminative model** can distinguish samples compiled to any distribution as long as they are linearly separable. The performance of two models is similar.

Models' performance

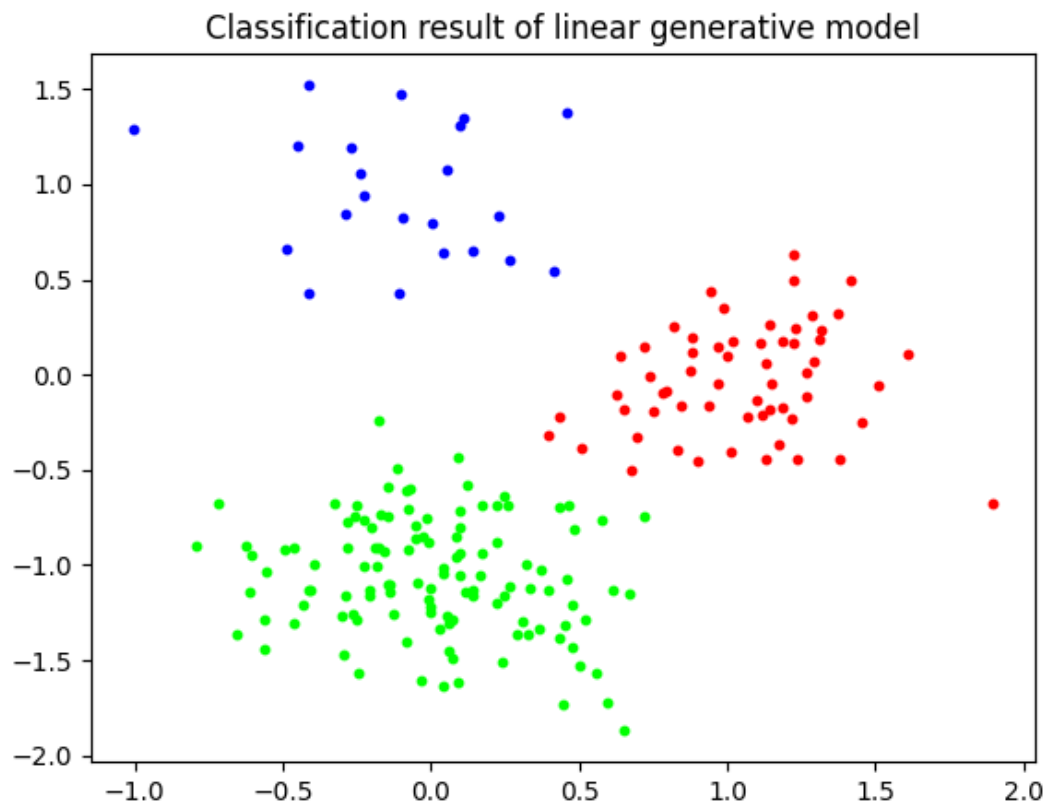
Use samples to train the **generative model**, the training process is shown in this picture. The misclassified samples are shown in black. It's evident that the model reach a high accuracy according to the lower part of the picture.



Then test the model on the testing dataset, the model can identified 99% of samples correctly. Note that the misclassified samples are **not** shown in the following picture since we won't tell the model the correct labels when calling the **classify** method.

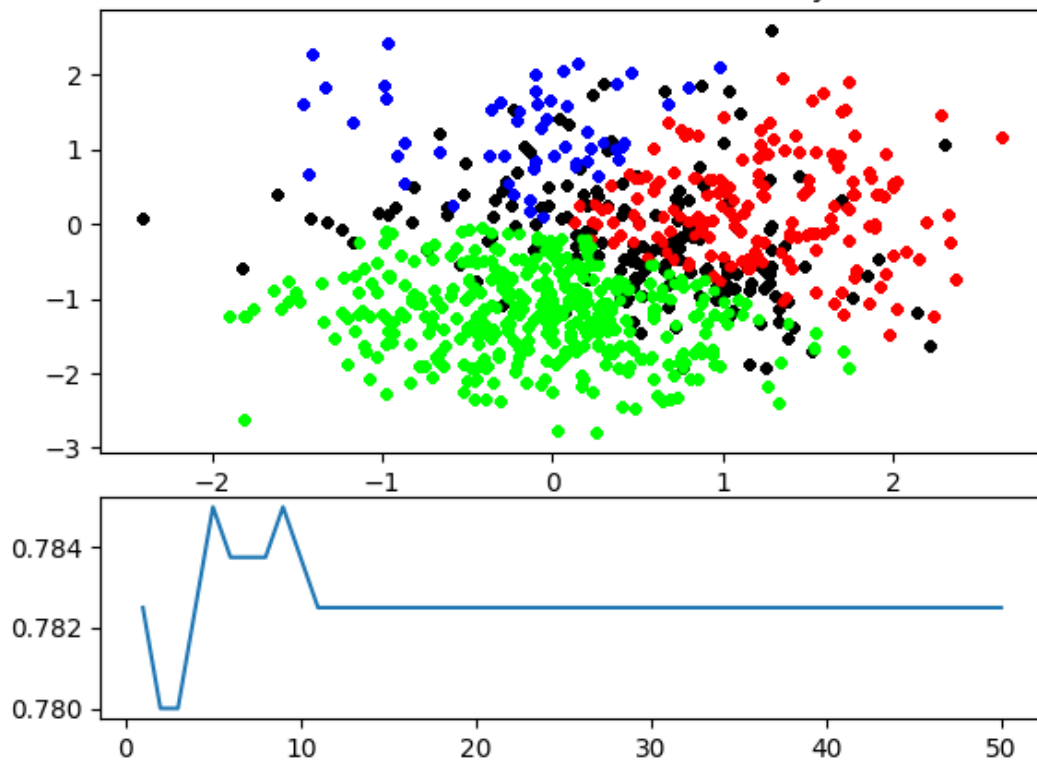


Apply the **generative model** to the same training set and testing set, the model can also identified 98.5% percent of samples correctly.

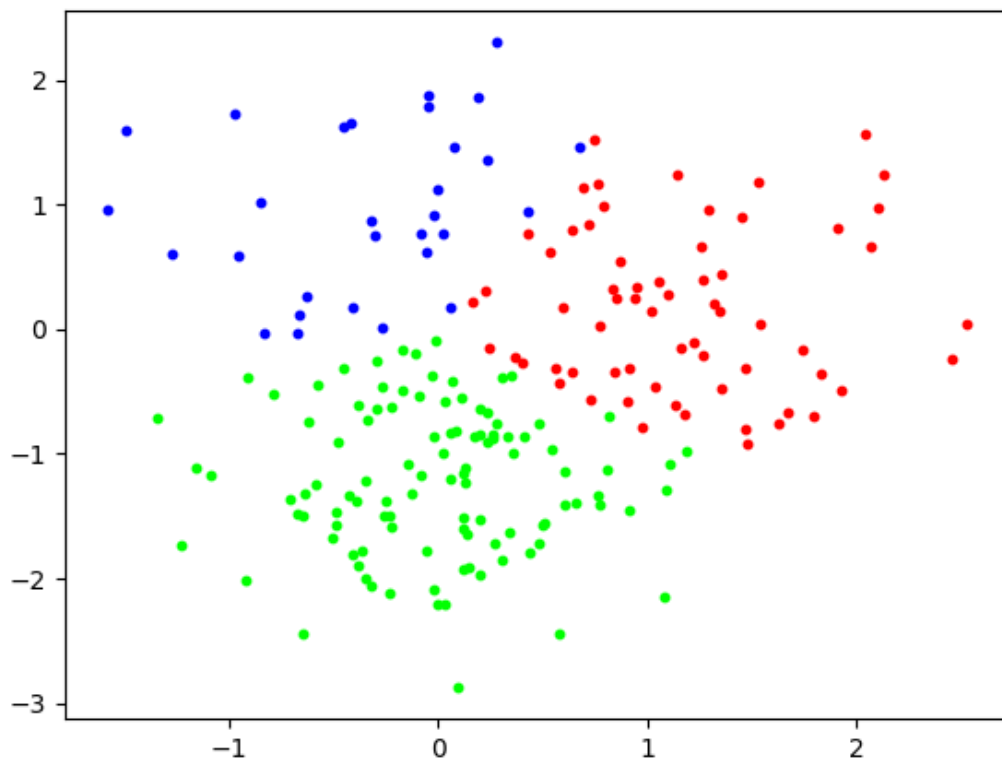


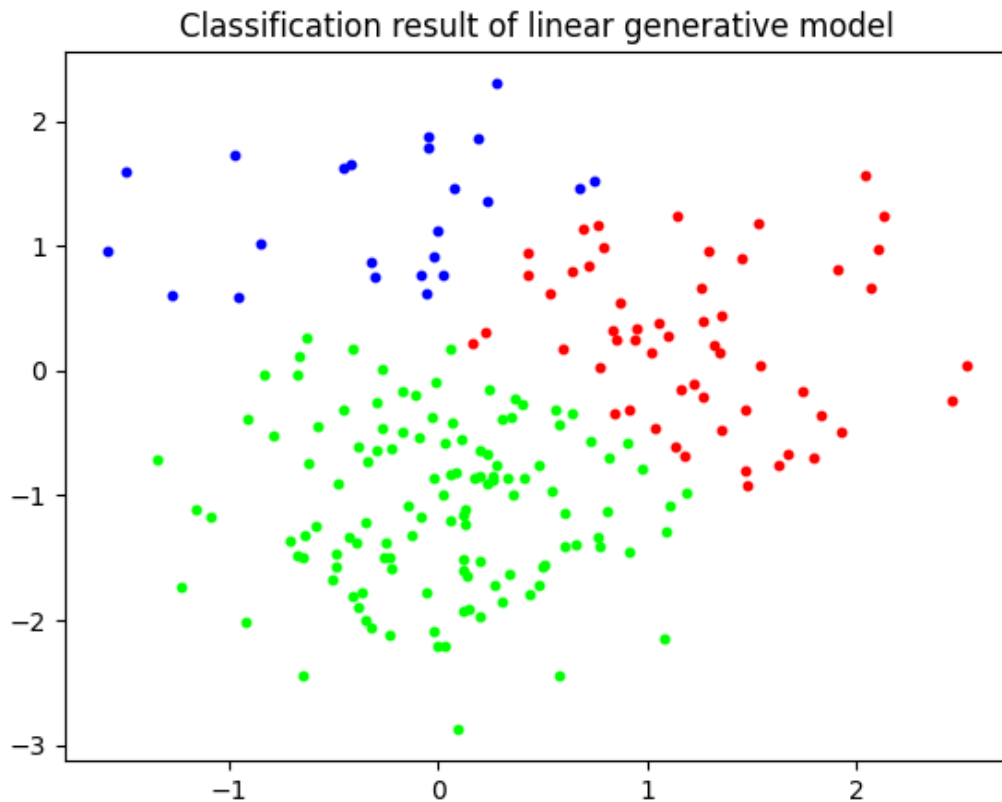
Then **increase the covariance** of Gaussian distributions and repeat above. Both of **generative model** and **discriminative model** deteriorated. They can identify 81% and 80.5% of samples respectively.

Classification result and accuracy



Classification result of linear discriminative model





Demo

This demo is also at the beginning of **source.py**.

```

1 Python 3.6.9 (default, MM dd yyyy, hh:mm:ss)
2 Type "help", "copyright", "credits" or "license" for more information.
3 >>>
4 >>> from source import *
5 >>> create_dataset()
6 >>> samples, labels = load_dataset()
7 >>> set_of_samples, set_of_labels = split_dataset(samples, labels, [0.8])
8 >>> training_samples, testing_samples = set_of_samples
9 >>> training_labels, testing_labels = set_of_labels
10 >>> model1 = LinearDiscriminativeModel()
11 >>> model1.train(training_samples, training_labels, learning_rate=.9, max_epochs=50,
    plot_training_process=False)
12 Epoch Accuracy
13 1 0.97125
14 2 0.9725
15 3 0.9775
16 4 0.98125
17 5 0.9825
18 6 0.9825
19 7 0.98125
20 8 0.9825
21 9 0.985
22 10 0.985
23 11 0.9875
24 12 0.9875
25 13 0.98875
26 14 0.98875

```

```

27 15 0.98875
28 16 0.9875
29 17 0.9875
30 18 0.98875
31 19 0.98875
32 20 0.98875
33 21 0.98875
34 22 0.98875
35 23 0.98875
36 24 0.98875
37 25 0.98875
38 26 0.98875
39 27 0.98875
40 28 0.98875
41 29 0.99
42 30 0.99
43 31 0.99
44 32 0.99
45 33 0.99
46 34 0.99
47 35 0.99
48 36 0.99
49 37 0.99
50 38 0.99
51 39 0.99
52 40 0.99
53 41 0.99
54 42 0.99
55 43 0.99
56 44 0.99
57 45 0.99
58 46 0.99
59 47 0.99
60 48 0.99
61 49 0.99
62 50 0.99
63 >>> labels1 = model1.classify(testing_samples, plot_classification_result=False)
64 >>> acc = sum(testing_labels == labels1) / testing_labels.size
65 >>> print(model1.w)
66 [[ 3.80438433 -2.04399475 -1.76038958]
67  [ 0.92932523 -5.47955759 4.55023236]]
68 >>> print(acc[0, 0])
69 0.975
70 >>> model2 = LinearGenerativeModel()
71 >>> model2.train(training_samples, training_labels)
72 >>> labels2 = model2.classify(testing_samples, plot_classification_result=False)
73 >>> acc = sum(testing_labels == labels2) / testing_labels.size
74 >>> print(model2.w)
75 [[ 9.36970937 -0.43319556 -0.04108211]
76  [ 0.08862703 -10.40075944 10.15249546]]
77 >>> print(acc[0, 0])
78 0.975

```