

114-1 Machine Learning

Week 6 Assignment

Ming Hsun Wu

October 14, 2025

PROBLEM 1. Unanswered Questions

Nope!

PROBLEM 2.

Classification Dataset

Format: (Longitude, Latitude, Label)

Rules:

- If the temperature observation value is an invalid value of -999 , then label = 0.
- If the temperature observation value is valid, then label = 1.
- Use Gaussian Discriminant Analysis (GDA) to build a classification model.

NOTE OF PROBLEM 2.

Gaussian Discriminant Analysis (GDA) is a generative classification model that assumes data are generated from two Gaussian distributions:

$$p(x \mid y = 0) = \mathcal{N}(\mu_0, \Sigma), \quad p(x \mid y = 1) = \mathcal{N}(\mu_1, \Sigma)$$

with class priors

$$p(y = 1) = \phi, \quad p(y = 0) = 1 - \phi$$

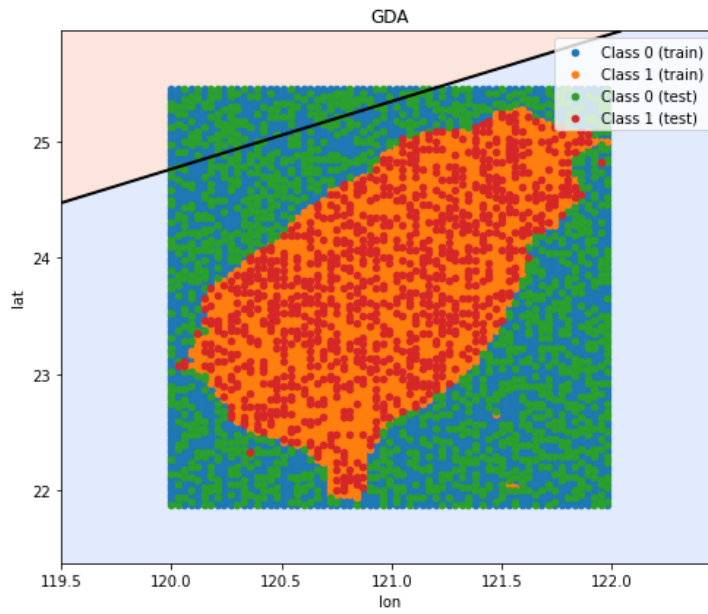
Using Bayes' theorem, we derive the posterior probability:

$$p(y = 1 \mid x) = \frac{p(x \mid y = 1) p(y = 1)}{p(x \mid y = 1) p(y = 1) + p(x \mid y = 0) p(y = 0)}$$

The idea of Gaussian Discriminant Analysis (GDA) is to model the conditional density $p(x \mid y = k)$ for each class and use Bayes' theorem to compute $p(y \mid x)$.

In this dataset, GDA learns the mean centers and spread of each class and estimates the probability that a point belongs to a class under a Gaussian assumption. If the data distribution is approximately Gaussian, GDA can classify effectively. However, if the data structure is nonlinear, a linear boundary will not fit well.

The dataset was split into 70% training and 30% testing. Model performance was evaluated by accuracy on the test set. An accuracy close to **0.5** indicates that the linear boundary cannot effectively separate the two classes. I evaluated the test set performance using the mean of the predicted and true values, i.e., $\text{mean}(\text{pred}, \text{true})$.



```
1  #%% source code
2  #%%
3  import json, re
4  import numpy as np
5  import pandas as pd
6  from sklearn.model_selection import train_test_split
7  from sklearn.metrics import accuracy_score, mean_squared_error
8  import torch
9  import torch.nn as nn
10 import torch.optim as optim
11 import matplotlib.pyplot as plt
12 #%%
13 'load data'
14 path = r"./0-A0038-003.json"
15 with open(path, "r", encoding="utf-8") as f:
16     data = json.load(f)
17
18 content = data['cwaopendata']['dataset']['Resource']['Content']
19
20 nums = re.findall(r'[-+]?\\d+\\.\\d+E[+-]?\\d+', content)
21 vals = np.array([float(x) for x in nums], dtype=float)
22
23
24 ncols, nrows = 67, 120
25 grid = vals.reshape((nrows, ncols))
26
```

```

27 bl_lon, bl_lat = 120.0, 21.88
28 res_lon, res_lat = 0.03, 0.03
29
30
31 lon = bl_lon + np.arange(ncols) * res_lon
32 lat = bl_lat + np.arange(nrows) * res_lat
33 lon_grid, lat_grid = np.meshgrid(lon, lat)
34
35
36 lon_flat = lon_grid.ravel()
37 lat_flat = lat_grid.ravel()
38 val_flat = grid.ravel()
39
40
41 labels = np.where(val_flat == -999.0, 0, 1)
42 df_classification = pd.DataFrame({"lon": lon_flat, "lat":
    lat_flat, "label": labels})
43
44
45 mask = val_flat != -999.0
46 df_regression = pd.DataFrame({
47     "lon": lon_flat[mask],
48     "lat": lat_flat[mask],
49     "value": val_flat[mask]
50 })
51
52

```

```

53  #%%
54  'data preprocessing'
55  x_cls = df_classification[["lon", "lat"]].values
56  y_cls = df_classification["label"].values
57
58  x_cls_train, x_cls_test, y_cls_train, y_cls_test =
        train_test_split(x_cls, y_cls, test_size=0.3, random_state=40)
59
60  x_reg = df_regression[["lon", "lat"]].values
61  y_reg = df_regression["value"].values
62
63  x_reg_train, x_reg_test, y_reg_train, y_reg_test =
        train_test_split(x_reg, y_reg, test_size=0.3, random_state=42)
64
65  #%%
66
67  phi = np.mean(y_cls_train)
68  mu0 = x_cls_train[y_cls_train == 0].mean(axis=0)
69  mu1 = x_cls_train[y_cls_train == 1].mean(axis=0)
70
71  # Conv matrix
72  X0 = x_cls_train[y_cls_train == 0]
73  X1 = x_cls_train[y_cls_train == 1]
74  Sigma = ((X0 - mu0).T @ (X0 - mu0) + (X1 - mu1).T @ (X1 - mu1)) /
        len(x_cls_train)
75
76  # Define predict function

```



```

77 def predict_proba(X):
78     invS = np.linalg.inv(Sigma)
79     detS = np.linalg.det(Sigma)
80     def gaussian(x, mu):
81         diff = x - mu
82         return np.exp(-0.5*diff@invS@diff.T) /
            np.sqrt((2*np.pi)**X.shape[1]*detS)
83     probs = []
84     for x in X:
85         p1 = gaussian(x, mu1)*phi
86         p0 = gaussian(x, mu0)*(1-phi)
87         probs.append(p1/(p1+p0))
88     return np.array(probs)
89
90 def predict(X):
91     return (predict_proba(X) >= 0.5).astype(int)
92
93 # Evaluate
94 y_pred = predict(x_cls_test)
95 acc = np.mean(y_pred == y_cls_test)

```
