

# **114-1 Machine Learning**

## **Week 4 Assignment**

Ming Hsun Wu

September 30, 2025

PROBLEM 1. Unanswered Questions

Nope!

PROBLEM 2.

(a) 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.

(b) Regression Dataset

**Format:** (Longitude, Latitude, Value)

**Rules:**

- Only valid temperature observation values are retained (all values of  $-999$  are removed).
- The Value corresponds to the temperature in degrees Celsius.

## NOTE OF PROBLEM 2.

The raw data comes from gridded observations provided by the Central Weather Bureau, consisting of  $67 \times 120$  grid points. Invalid values are marked as  $-999.0$ .

Two datasets were constructed:

- **Classification dataset:** (lon, lat, label), where label = 0 represents invalid points and label = 1 represents valid points.
- **Regression dataset:** (lon, lat, value), retaining only valid points, where *value* corresponds to the temperature in degrees Celsius.

After processing:

- The classification dataset contains 8,040 entries.
- The regression dataset contains 3,495 entries.

df_classification - DataFrame				df_regression - DataFrame			
Index	lon	lat	label	Index	lon	lat	value
0	120	21.88	0	0	120.84	21.94	26.6
1	120.03	21.88	0	1	120.72	21.97	26.7
2	120.06	21.88	0	2	120.75	21.97	26.8
3	120.09	21.88	0	3	120.78	21.97	26.1
4	120.12	21.88	0	4	120.81	21.97	24.9
5	120.15	21.88	0	5	120.84	21.97	27.5
6	120.18	21.88	0	6	120.72	22	26.5
7	120.21	21.88	0	7	120.75	22	26.6
8	120.24	21.88	0	8	120.78	22	26.3
9	120.27	21.88	0	9	120.81	22	26.9
10	120.3	21.88	0	10	120.84	22	26.9
11	120.33	21.88	0	11	120.72	22.03	26.5

PROBLEM 3.

Model Training

Using the two datasets prepared in problem 2, two simple machine learning models were trained separately:

- **Classification model:** Predicts whether a grid point is valid (1) or invalid (0) based on (longitude, latitude).
- **Regression model:** Predicts the corresponding temperature observation value based on (longitude, latitude).

### NOTE OF PROBLEM 3.

## Model design

#### (a) Classification Model

**Input:** (lon, lat)

**Output:** Binary classification

**Architecture:** A fully connected neural network (MLP) with 2 input features, followed by 4 hidden layers (128, 128, 64, 32) units. Each layer uses ReLU activation, combined with Batch Normalization and Dropout (0.3). The output layer is 2-dimensional with softmax activation.

**Loss function:** CrossEntropyLoss

**Optimizer:** Adam ( $lr = 0.001$ )

#### (b) Regression Model

**Input:** (lon, lat)

**Output:** Temperature value ( $^{\circ}\text{C}$ )

**Architecture:** A fully connected neural network with 5 hidden layers (64, 128, 256, 128, 32) units. Each layer uses ReLU activation. The final output is a single real-valued unit.

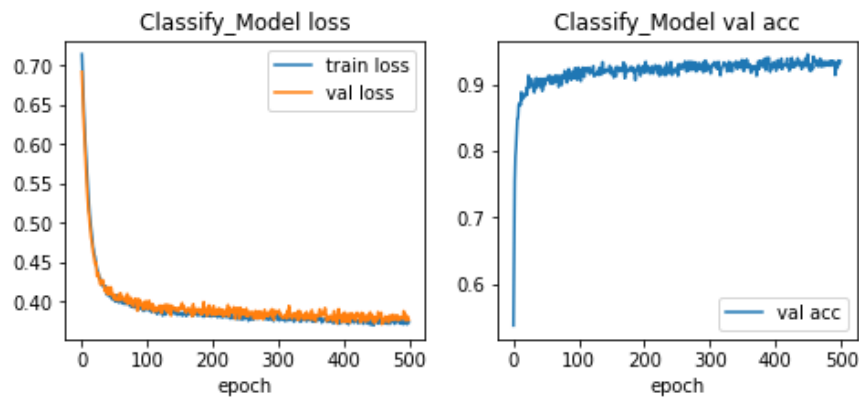
**Loss function:** MSELoss

**Optimizer:** Adam ( $lr = 0.001$ )

## Training Process and Results

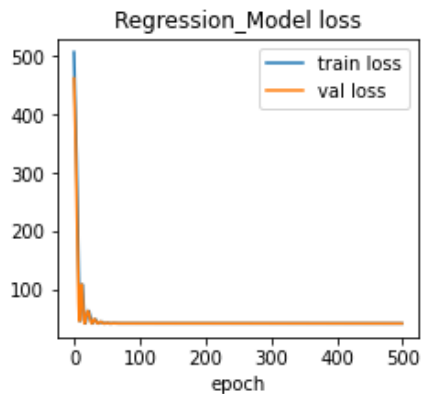
### (a) Classification Model Results

During training, the loss gradually converged, and the validation accuracy reached above 0.9 (as shown in the figure). The trained classification model achieved a prediction accuracy of approximately 0.84 on the test set.



### (b) Regression Model Results

During training, the MSE decreased from around 500 to double digits. The trained regression model achieved a mean squared error of 40.7 on the test set when comparing predictions with the ground truth.



## Discussion and Analysis

### (a) Classification Model

At the beginning, I designed the model in a relatively simple way, but the training results were very poor (the predictions on the test set all defaulted to class 0). Therefore, I adjusted the model architecture and depth, which led to the current results. I suspect the poor initial performance was mainly due to limitations and deficiencies inherent in the dataset. Despite the data imbalance, the model still maintained high accuracy, suggesting that it has learned the decision boundary effectively.

### (b) Regression Model

I believe the current regression model results are unsatisfactory. On the test set, the predictions are mostly concentrated around 21. While the training loss did converge, the model seems to have simply found a constant value that minimizes the loss, rather than capturing the true variation in the data.



---

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

```

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

```

```

53
54 #%%
55 'data preprocessing'
56 x_cls = df_classification[["lon", "lat"]].values
57 y_cls = df_classification["label"].values
58
59 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)
60 x_cls_val, x_cls_test, y_cls_val, y_cls_test =
    train_test_split(x_cls_test, y_cls_test, test_size=0.3,
    random_state=40)
61
62
63
64 x_cls_train = torch.tensor(x_cls_train, dtype=torch.float32)
65 y_cls_train = torch.tensor(y_cls_train, dtype=torch.long)
66
67 x_cls_val = torch.tensor(x_cls_val, dtype=torch.float32)
68 y_cls_val = torch.tensor(y_cls_val, dtype=torch.long)
69
70 x_cls_test = torch.tensor(x_cls_test, dtype=torch.float32)
71 y_cls_test = torch.tensor(y_cls_test, dtype=torch.long)
72
73
74
75 x_reg = df_regression[["lon", "lat"]].values
76 y_reg = df_regression["value"].values

```

```

77
78 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)
79 x_reg_val, x_reg_test, y_reg_val, y_reg_test =
    train_test_split(x_reg, y_reg, test_size=0.3, random_state=42)
80
81
82
83 x_reg_train = torch.tensor(x_reg_train, dtype=torch.float32)
84 y_reg_train = torch.tensor(y_reg_train,
    dtype=torch.float32).view(-1,1)
85 x_reg_val = torch.tensor(x_reg_val, dtype=torch.float32)
86 y_reg_val = torch.tensor(y_reg_val, dtype=torch.float32).view(-1,1)
87 x_reg_test = torch.tensor(x_reg_test, dtype=torch.float32)
88 y_reg_test = torch.tensor(y_reg_test,
    dtype=torch.float32).view(-1,1)
89
90 #%%
91 'Classify_Model Training'
92 # compile model
93 class Classify_Model(nn.Module):
94     def __init__(self):
95         super(Classify_Model, self).__init__()
96         self.layers = nn.Sequential(
97             nn.Linear(2, 128),
98             nn.ReLU(),
99             nn.BatchNorm1d(128),

```

```

100         nn.Dropout(0.3),
101         nn.Linear(128, 128),
102         nn.ReLU(),
103         nn.BatchNorm1d(128),
104         nn.Dropout(0.3),
105         nn.Linear(128, 64),
106         nn.ReLU(),
107         nn.Linear(64, 2),
108         nn.Softmax()
109     )
110     def forward(self, x):
111         return self.layers(x)
112
113     Classify_Model = Classify_Model()
114     optimizer_cls = optim.Adam(Classify_Model.parameters(), lr=0.001)
115
116     # train model
117     epochs = 500
118     cls_train_losses, cls_val_losses, cls_val_acc = [], [], []
119     for epoch in range(epochs):
120         Classify_Model.train()
121         optimizer_cls.zero_grad()
122         outputs = Classify_Model(x_cls_train)
123         loss = nn.CrossEntropyLoss()(outputs, y_cls_train)
124         cls_train_losses.append(loss.item())
125         loss.backward()
126         optimizer_cls.step()

```

```

127
128     with torch.no_grad():
129         y_pred_val = Classify_Model(x_cls_val)
130         val_loss = nn.CrossEntropyLoss()(y_pred_val, y_cls_val)
131         cls_val_losses.append(val_loss.item())
132
133         val_acc = accuracy_score(y_cls_val,
134                                   y_pred_val.argmax(dim=1).numpy())
135         cls_val_acc.append(val_acc)
136
137     # evaluate model
138     Classify_Model.eval()
139     with torch.no_grad():
140         Classify_preds =
141             Classify_Model(x_cls_test).argmax(dim=1).numpy()
142     acc = accuracy_score(y_cls_test, Classify_preds)
143     print("Classification accuracy:", acc)
144
145     #%%
146     'Regression_Model Training'
147     # compile model
148     class Regression_Model(nn.Module):
149         def __init__(self):
150             super(Regression_Model, self).__init__()
151             self.layers = nn.Sequential(
152                 nn.Linear(2, 64),
153                 nn.ReLU(),
154                 nn.Linear(64, 128),

```

```

152         nn.ReLU(),
153         nn.Linear(128, 256),
154         nn.ReLU(),
155         nn.Linear(256, 128),
156         nn.ReLU(),
157         nn.Linear(128, 32),
158         nn.ReLU(),
159         nn.Linear(32, 1)
160
161
162     )
163     def forward(self, x):
164         return self.layers(x)
165
166 Regression_Model = Regression_Model()
167 criterion_reg = nn.MSELoss()
168 optimizer_reg = optim.Adam(Regression_Model.parameters(), lr=0.001)
169
170 # train model
171 reg_train_losses, reg_val_losses = [], []
172 epochs = 500
173 for epoch in range(500):
174     Regression_Model.train()
175     optimizer_reg.zero_grad()
176     outputs = Regression_Model(x_reg_train)
177     loss = criterion_reg(outputs, y_reg_train)
178     reg_train_losses.append(loss.item())

```

```
179     loss.backward()
180     optimizer_reg.step()
181
182     with torch.no_grad():
183         y_pred_val = Regression_Model(x_reg_val)
184         val_loss = nn.MSELoss()(y_pred_val, y_reg_val)
185         reg_val_losses.append(val_loss.item())
186
187
188     # evaluate model
189     Regression_Model.eval()
190     with torch.no_grad():
191         Regression_preds = Regression_Model(x_reg_test).numpy()
192     mse = mean_squared_error(y_reg_test, Regression_preds)
193     print("Regression MSE:", mse)
```

---