

# Contents

|    |   |    |
|----|---|----|
| 1  | <a href="#">task1/README</a>              | 2  |
| 2  | <a href="#">task1/USERS.txt</a>           | 3  |
| 3  | <a href="#">task1/src/adaboost.py</a>     | 4  |
| 4  | <a href="#">task1/src/binarysearch.py</a> | 7  |
| 5  | <a href="#">task1/src/classifier.py</a>   | 9  |
| 6  | <a href="#">task1/src/ex4 tools.py</a>    | 14 |
| 7  | <a href="#">task1/src/geo</a>             | 17 |
| 8  | <a href="#">task1/src/model.py</a>        | 18 |
| 9  | <a href="#">task1/src/models.py</a>       | 19 |
| 10 | <a href="#">task1/src/strong.py</a>       | 24 |
| 11 | <a href="#">task1/src/strong2.py</a>      | 30 |
| 12 | <a href="#">task1/src/testmodel.py</a>    | 35 |
| 13 | <a href="#">task1/src/weather.py</a>      | 36 |
| 14 | <a href="#">task1/src/wether.json</a>     | 37 |

# 1 task1/README

```
1  .:
2  README
3  src/
4  USERS.txt
5
6  ./src:
7  adaboost.py
8  BinAgent
9  binarysearch.py
10 classifier.py
11 ex4_tools.py
12 geo
13 model.py
14 models.py
15 strong.py
16 strong2.py
17 testmodel.py
18 weather.py
19 wether.json
```

## 2 task1/USERS.txt

```
1 nir.vaknin, 31151478
2 davidponar, 208504050
3 nitsanwoja, 314967837
4 amirlevy95, 204880330
```

### 3 task1/src/adaboost.py

```
1  """
2  =====
3  Introduction to Machine Learning (67577)
4  =====
5
6  Skeleton for the AdaBoost classifier.
7
8  Author: Gad Zalcberg
9  Date: February, 2019
10
11 """
12 import numpy as np
13 from ex4_tools import *
14 from matplotlib import pyplot as plt
15
16 class AdaBoost(object):
17
18     def __init__(self, WL, T, support_wights = False):
19         """
20         Parameters
21         -----
22         WL : the class of the base weak learner
23         T : the number of base learners to learn
24         """
25         WL, self.fe = WL
26         self.WL = [ WL for _ in range(T)]
27         self.T = T
28         self.h = [None]*T # list of base learners
29         self.w = np.zeros(T) # weights
30         self.support_wights = support_wights
31
32         # self._predict = np.vactorize(WL.prdeict)
33
34     def train(self, X, y):
35         """
36         Parameters
37         -----
38         X : samples, shape=(num_samples, num_features)
39         y : labels, shape=(num_samples)
40         Train this classifier over the sample (X,y)
41         After finish the training return the weights of the samples in the last iteration.
42         """
43         y = y.flatten()
44         D = np.ones( len( y ) ) * 1/len( y )
45
46         def find_loss_D(h, D, X, y):
47             _out = h.predict(X).flatten()
48             return np.sum( D[_out != y.transpose()] ), _out
49
50         def normalize(vec):
51             r = np.sum(vec )
52             return vec / r if r != 0 else 0
53
54         self.w = []
55
56         for i in range(self.T):
57             self.h[i] = self.WL[i](self.fe)
58             if self.support_wights:
59                 self.h[i].train(X, y, D)
```

```

60         else:
61             self.h[i].train(X, y)
62
63         e_t, _out = find_loss_D(self.h[i], D, X, y)
64         # print("[%]et:")
65         # print(e_t)
66         w = 0.5 * np.log( 1/e_t - 1 )
67         D = normalize(D * np.e ** ( -w * y * _out ) )
68         self.w.append(w)
69
70     self.w = np.array(self.w)
71     return D
72
73     def predict(self, X, max_t=None):
74         """
75         Parameters
76         -----
77         X : samples, shape=(num_samples, num_features)
78         :param max_t: integer < self.T: the number of classifiers to use for the classification
79         :return: y_hat : a prediction vector for X. shape=(num_samples)
80         Predict only with max_t weak learners,
81         """
82         if max_t is None:
83             max_t = self.T
84         return np.sign(
85             sum(
86                 self.w[i] * self.h[i].predict(X) for i in range(max_t)) )
87
88     def error(self, X, y, max_t=None):
89         """
90         Parameters
91         -----
92         X : samples, shape=(num_samples, num_features)
93         y : labels, shape=(num_samples)
94         :param max_t: integer < self.T: the number of classifiers to use for the classification
95         :return: error : the ratio of the correct predictions when predict only with max_t weak learners (float)
96         """
97         if max_t is None:
98             max_t = self.T
99         errors = sum(np.ones( len(y) )[ self.predict(X, max_t) != y ])
100         return errors / len(y)
101
102
103     class AdaBoostList(AdaBoost):
104         def __init__(self, WLS):
105             self.WL = WLS
106             self.T = len(WLS)
107             self.h = [None]*self.T # lisself. of base learners
108             self.w = np.zeros(self.T) # weights
109
110         def prdict(self, X):
111             return super().prdict(X, self.T)
112
113
114     def test():
115         def plot_as_function_of_agents(A, X, y, max_t):
116             yplot = [A.error(X, y, t) for t in range(max_t)]
117             plt.plot([t for t in range(max_t)], yplot)
118
119         def plot_points(A, X, y,
120             classifiers=[ 5, 10, 50, 100, 200, 500 ], D_T=None, noise=0):
121
122             for i ,T in enumerate(classifiers):
123                 plt.subplot( 320 + i + 1)
124                 decision_boundaries(A, X, y, num_classifiers=T, weights=D_T)
125
126             size_str = True if D_T is not None else False
127             plt.savefig(f"plot_points_{size_str}_{noise}.png")

```

```

128
129     for noise in [0, 0.01, 0.4]:
130         A = AdaBoost( DecisionStump, 500 )
131         X,y = generate_data(5000, noise)
132         A.train(X,y)
133
134
135         plot_as_function_of_agents(A, X, y, 500)
136         X,y = generate_data(200, noise)
137         plot_as_function_of_agents(A, X, y, 500)
138         plt.legend(["training error rate","testing error rate"])
139         plt.xlabel("agents")
140         plt.ylabel("error rate")
141         plt.ylim([0,1])
142         plt.title(f"error rate above distribution with nosie - {noise}")
143         plt.savefig(f"plot_error_{noise}.png")
144
145
146         D_T = A.w.transpose()
147         D_T = 30* D_T/max(D_T)
148         plt.figure()
149
150         for d_T in [None, D_T ]:
151             plot_points(A, X, y, D_T=d_T, noise=noise)
152             plt.figure()
153
154
155
156 if __name__ == "__main__":
157     test()

```

## 4 task1/src/binarysearch.py

```
1  import numpy as np
2  import re
3
4
5  class binarysearch:
6
7      def __init__(self, models):
8          self.mods = models
9          self.iterations = 6
10
11
12      '''
13          for _bool in re > middle ] )
14          y = np.array( [ { False : [0] , True : [1] }[_bool ]
15      '''
16      # def load_classifiers(self, classifiers_steps):
17
18      #     pass
19
20
21      def _predict(self, x, treshold):
22          return self.mods[treshold].predict(x)
23
24      def predict(self, X, start_range, end_range):
25          times_delay = 0
26          for j in range(self.iterations):
27              middle =(start_range + end_range)/2
28              res = np.zeros( shape = start_range.shape)
29              for i, time in enumerate(start_range):
30                  res[i] = self._predict(X.iloc[[i]] , start_range[i])[0]
31              res = (res + np.ones(len(res)))/2
32              start_range, end_range = (1 - res) * start_range + res * middle , (1 - res) * middle + res * end_range
33          return middle
34
35      def dump(self, _file):
36          with open(_file, "w") as f :
37              _str = ""
38              for treshold, _mod in self._mods.items() :
39                  _str += f"{treshold}:" + _mod.dump( ) + "\n$"
40
41              _file.write(_str)
42
43      def binarysearch_read(_file):
44          mods = {}
45          for banch in re.split('$', _file.read()):
46              treshold, _strmod = re.split(':', banch)
47              mods[float(treshold)] = AdaBoost_read( )
48
49
50
51
52      # testing the idea.
53      def foo( y, s , e ):
54          middle = (s+e)/2
55          return middle * y , middle + middle * y
56
57      if __name__ == "__main__":
58          s , e = np.array( [ 0, 0 ,0]), np.array( [ 1, 1 ,1])
59          y = np.array( [1, 0, 1] )
```

```
60     print( foo(y, s, e) )
61
62
63
64
```



## 5 task1/src/classifier.py

```
1  # from strong import *
2  from models import abcModel, DecisionTree, Logistic, SVM, DecisionStumpWarper
3  from adaboost import AdaBoost, AdaBoostList
4  from itertools import combinations
5  import numpy as np
6  import pandas as pd
7  from binarysearch import binarysearch
8  from datetime import date
9  from random import shuffle
10 from sklearn.model_selection import train_test_split
11 from geopy.geocoders import Nominatim, ArcGIS
12 import json
13 import math
14 factor_reason = {'CarrierDelay': 0, 'WeatherDelay': 1, 'NASDelay': 2, 'LateAircraftDelay': 3}
15
16 #
17 # def pre_proc_class_old(_dataset, dropped_fe, categorical):
18 #     y_factor = []
19 #     _dataset = _dataset.dropna()
20 #     _dataset['month'] = pd.DatetimeIndex(_dataset['FlightDate']).month
21 #     _dataset['day'] = pd.DatetimeIndex(_dataset['FlightDate']).day
22 #     _dataset['year'] = pd.DatetimeIndex(_dataset['FlightDate']).year
23 #
24 #     for factor in _dataset["DelayFactor"]:
25 #         y_factor.append(factor_reason[factor])
26 #     y_factor = np.array(y_factor)
27 #     for index, row in _dataset.iterrows():
28 #         _dataset.loc[index, 'CRSElapsedTime'] = math.floor(row['CRSElapsedTime'] / 10)
29 #         _dataset.loc[index, 'CRSArrTime'] = math.floor(row['CRSArrTime'] / 100)
30 #         _dataset.loc[index, 'CRSDepTime'] = math.floor(row['CRSDepTime'] / 100)
31 #         _dataset.loc[index, 'Distance'] = math.floor(row['Distance'] / 10)
32 #
33 #     cat = pd.DataFrame(pd.get_dummies(_dataset[categorical].astype('category')))
34 #     _dataset = _dataset.drop(dropped_fe + categorical, axis=1)
35 #     _dataset_preproc = pd.concat([_dataset.reset_index(drop=True), cat.reset_index(drop=True)], axis=1)
36 #     return _dataset_preproc, y_factor
37
38
39
40 def proc_weather(_dataset):
41     weather = json.loads(open("wether.json").read().replace("'", "\'"))
42     for index, row in _dataset.iterrows():
43         _date = f"{row['day']:02d}-{row['month']:02d}-{row['year']:%2000}"
44         if _date in weather:
45             for key, val in weather[_date].items():
46                 if key != 'day':
47                     row[key] = val
48
49
50
51 def pre_proc_class(_dataset, categorical):
52
53
54     _dataset['month'] = pd.DatetimeIndex(_dataset['FlightDate']).month
55     _dataset['day'] = pd.DatetimeIndex(_dataset['FlightDate']).day
56     _dataset['year'] = pd.DatetimeIndex(_dataset['FlightDate']).year
57
58     proc_weather(_dataset)
59
```

```

60     for index, row in _dataset.iterrows():
61         # _dataset.loc[index, 'CRSElapsedTime'] = math.floor(row['CRSElapsedTime'] / 10)
62         _dataset.loc[index, 'CRSArrTime'] = math.floor(row['CRSArrTime'] / 100)
63         _dataset.loc[index, 'CRSDepTime'] = math.floor(row['CRSDepTime'] / 100)
64         _dataset.loc[index, 'Distance'] = math.floor(row['Distance'] / 10)
65
66
67     geolocator = ArcGIS(username="david.ponarovsky", password="mxSrWYYdSq++J7+",
68         referer="efMmfQrfh1o_Ag-MEzx5-en9lLs-m_Vu5T_JU7K45vfUhjEccY6W2cilzPnn2r9T07fpPnAeK_U8EWptaUiX0C7FosOKaovLPAXQR06PGnXrUY",
69         timeout=500000)
70
71     def DynmicGeo( df, _keys ):
72         import time
73         from pprint import pprint
74
75         #Dynmic = {}
76         Dynmic = json.loads( open("geo").read().replace("'", "\'") )
77         try :
78             for _key in _keys:
79                 df[f'lat_{_key}']=0
80                 df[f'long_{_key}']=0
81                 for x in range(len(df)):
82                     if df[_key][x] in Dynmic:
83                         location = Dynmic[df[_key][x]]
84                         if type(location) == str:
85                             df.at[x, f'lat_{_key}'] , df.at[x, f'long_{_key}'] = eval(location)
86                         else:
87                             df.at[x, f'lat_{_key}'] , df.at[x, f'long_{_key}'] = location
88                     else:
89                         location = geolocator.geocode(df[_key][x])
90                         Dynmic[df[_key][x]] = location.latitude, location.longitude
91                         # time.sleep(2)
92                         #print(location , location.latitude, location.longitude)
93                         df.at[x, f'lat_{_key}']= location.latitude
94                         df.at[x, f'long_{_key}'] = location.longitude
95         except:
96             print("except")
97
98         pprint(Dynmic)
99
100         for key, val in Dynmic.items():
101             Dynmic[key] = str(val)
102         with open('geo', 'w') as f:
103             json.dump(Dynmic, f)
104
105         return df
106
107
108     geoKeys = [ "Origin",
109                 "OriginCityName",
110                 "OriginState",
111                 "Dest",
112                 "DestCityName",
113                 "DestState"]
114
115     _dataset = DynmicGeo(_dataset, geoKeys)
116
117     if len(categorical) > 0 :
118         cat = pd.DataFrame(pd.get_dummies(_dataset[categorical].astype('category')))
119         _dataset_prepoc = pd.concat([_dataset.reset_index(drop=True), cat.reset_index(drop=True)], axis=1)
120     else:
121         _dataset_prepoc = _dataset
122
123     return _dataset_prepoc
124
125
126 def pre_proc_delay(dataset, drop_fe_delay, categorical):
127     y = []

```

```

128     for _bool in dataset["ArrDelay"] > 0:
129         y.append({False: [0], True: [1]}[_bool])
130     y = np.array(y)
131     dataset = dataset.drop(drop_fe_delay + categorical, axis=1)
132     return dataset, y
133
134
135 def pre_proc_factor(dataset, drop_fe_factor, categorical):
136     y_factor = []
137     dataset = dataset.dropna()
138     for factor in dataset["DelayFactor"]:
139         y_factor.append(factor_reason[factor])
140     y_factor = np.array(y_factor)
141     dataset = dataset.drop(drop_fe_factor + categorical, axis=1)
142     return dataset, y_factor
143
144
145 def pred_trees(x, y, x_test, y_test):
146     from sklearn import datasets
147     from sklearn.metrics import confusion_matrix
148     from sklearn.model_selection import train_test_split
149
150     # dividing X, y into train and test data
151
152     # training a DescisionTreeClassifier
153     from sklearn.tree import DecisionTreeClassifier
154     for depth in [4, 5, 7, 10, 13, 15, 17, 18, 19, 20]:
155         dtree_model = DecisionTreeClassifier(max_depth=depth).fit(x, y)
156         dtree_predictions = dtree_model.predict(x_test)
157         print(dtree_model.score(x_test, y_test))
158         print(depth)
159
160
161 def pred_forest(x, y, x_test, y_test):
162     from sklearn.ensemble import RandomForestClassifier
163
164     model = RandomForestClassifier(random_state=20)
165     model.fit(x, y)
166     print(model.score(x_test, y_test))
167
168
169 def identify_corelated_features(df):
170     # corr = df.corr()
171     # #print(sns.heatmap(corr))
172     # columns = np.full((corr.shape[0],), True, dtype=bool)
173     # for i in range(corr.shape[0]):
174     #     for j in range(i + 1, corr.shape[0]):
175     #         if corr.iloc[i, j] >= 0.9:
176     #             if columns[j]:
177     #                 columns[j] = False
178     # selected_columns = df.columns[columns]
179     # df = df[selected_columns]
180     return df
181
182
183 def pred(x, y, x_test, y_test):
184     from sklearn import datasets
185     from sklearn.metrics import confusion_matrix
186     from sklearn.model_selection import train_test_split
187
188     # training a KNN classifier
189     from sklearn.neighbors import KNeighborsClassifier
190     for n in [11, 13, 15, 20, 30, 35, 40, 45]:
191         knn = KNeighborsClassifier(n_neighbors=n).fit(x, y)
192
193         # accuracy on X_test
194         accuracy = knn.score(x_test, y_test)
195         print(accuracy)

```

```

196         print(n)
197
198     # creating a confusion matrix
199     # knn_predictions = knn.predict(x)
200     # cm = confusion_matrix(y_test, knn_predictions)
201
202
203     featuers = ["DayOfWeek",
204                 "FlightDate",
205                 "Reporting_Airline",
206                 "Tail_Number",
207                 "Flight_Number_Reporting_Airline",
208                 "Origin",
209                 "OriginCityName",
210                 "OriginState",
211                 "Dest",
212                 "DestCityName",
213                 "DestState",
214                 "CRSDepTime",
215                 "CRSArrTime",
216                 "CRSElapsedTime",
217                 "Distance"]
218
219     dropped_fe_factor = ['Flight_Number_Reporting_Airline',
220                         'Tail_Number',
221                         'DelayFactor',
222                         "OriginCityName",
223                         "OriginState",
224                         "DestCityName",
225                         "DestState", 'FlightDate', "CRSElapsedTime"]
226
227     dropped_fe_delay = ['Flight_Number_Reporting_Airline',
228                        'Tail_Number',
229                        "OriginCityName",
230                        "OriginState",
231                        "DestCityName",
232                        "DestState", 'FlightDate', "CRSElapsedTime", 'ArrDelay', 'DelayFactor', 'Origin', 'Dest' ]
233
234     categorical_new = ['Reporting_Airline']
235
236     def final_pre_proc(_dataset):
237         x = pre_proc_class(_dataset, categorical_new)
238         x_delay,y_delay = pre_proc_delay(x,dropped_fe_delay,categorical_new)
239         x_factor,y_factor = pre_proc_delay(x,dropped_fe_factor,categorical_new)
240
241         x_factor = identify_corelated_features(x_factor)
242         return x_delay, y_delay, x_factor,y_factor
243
244     if __name__ == '__main__':
245         _dataset = pd.read_csv("~/data/train_data.csv")
246
247         x = pre_proc_class(_dataset, categorical_new)
248
249         print(x)
250
251         x_delay,y_delay = pre_proc_delay(x,dropped_fe_delay,categorical_new)
252         x_factor,y_factor = pre_proc_delay(x,dropped_fe_factor,categorical_new)
253
254         x_factor = identify_corelated_features(x_factor)
255         x_train, x_test, y_train, y_test = train_test_split(x_factor, y_factor)
256         # print(x['month'])
257
258         # print(y_factor)
259         # pred_trees(x_train, y_train, x_test, y_test)
260         from datetime import datetime
261
262         # datetime object containing current date and time
263         print(datetime.now())

```

```
264     print(x_test.shape)
265     # pred_forest(x_train, y_train, x_test, y_test)
266     print(datetime.now())
267     # pred(x_train, y_train, x_test, y_test)
268     # print(identify_correlated_features(x, y_del))
```

## 6 task1/src/ex4 tools.py

```
1  """
2  =====
3      Introduction to Machine Learning (67577)
4  =====
5
6  This module provides some useful tools for Ex4.
7
8  Author: Gad Zalcberg
9  Date: February, 2019
10
11 """
12 import numpy as np
13 import matplotlib.pyplot as plt
14 from matplotlib.colors import ListedColormap
15 from itertools import product
16 from matplotlib.pyplot import imread
17 import os
18 from sklearn.model_selection import train_test_split
19
20
21 def find_threshold(D, X, y, sign, j):
22     """
23     Finds the best threshold.
24     D = distribution
25     S = (X, y) the data
26     """
27     # sort the data so that  $x_1 \leq x_2 \leq \dots \leq x_m$ 
28     sort_idx = np.argsort(X[:, j])
29     X, y, D = X[sort_idx], y[sort_idx], D[sort_idx]
30
31
32     thetas = np.concatenate([[ -np.inf], (X[1:, j] + X[:-1, j]) / 2, [np.inf]])
33     minimal_theta_loss = np.sum(D[y == sign]) # loss of the smallest possible theta
34     losses = np.append(minimal_theta_loss, minimal_theta_loss - np.cumsum(D * (y * sign)))
35     min_loss_idx = np.argmin(losses)
36
37     return losses[min_loss_idx], thetas[min_loss_idx]
38
39
40 class DecisionStump(object):
41     """
42     Decision stump classifier for 2D samples
43     """
44
45     def __init__(self):
46         self.theta = 0
47         self.j = 0
48         self.sign = 0
49
50
51     def fit(self, D, X, y):
52         """
53         Train the classifier over the sample (X,y) w.r.t. the weights D over X
54         Parameters
55         -----
56         D : weights over the sample
57         X : samples, shape=(num_samples, num_features)
58         y : labels, shape=(num_samples)
59         """
```

```

60     #print(f"X: {X.shape}")
61     loss_star, theta_star = np.inf, np.inf
62     for sign, j in product([-1, 1], range(X.shape[1])):
63         loss, theta = find_threshold(D, X, y, sign, j)
64         if loss < loss_star:
65             self.sign, self.theta, self.j = sign, theta, j
66             loss_star = loss
67
68     def predict(self, X):
69         """
70         Parameters
71         -----
72         X : shape=(num_samples, num_features)
73         Returns
74         -----
75         y_hat : a prediction vector for X shape=(num_samples)
76         """
77         y_hat = self.sign * ((X[:, self.j] <= self.theta) * 2 - 1)
78         return y_hat
79
80
81     def decision_boundaries(classifier, X, y, num_classifiers=1, weights=None):
82         """
83         Plot the decision boundaries of a binary classifiers over  $X \setminus \text{subseq } \mathbb{R}^2$ 
84
85         Parameters
86         -----
87         classifier : a binary classifier, implements classifier.predict(X)
88         X : samples, shape=(num_samples, 2)
89         y : labels, shape=(num_samples)
90         title_str : optional title
91         weights : weights for plotting X
92         """
93         cm = ListedColormap(['#AAAAFF', '#FFAAAA'])
94         cm_bright = ListedColormap(['#0000FF', '#FF0000'])
95         h = .003 # step size in the mesh
96         # Plot the decision boundary.
97         x_min, x_max = X[:, 0].min() - .2, X[:, 0].max() + .2
98         y_min, y_max = X[:, 1].min() - .2, X[:, 1].max() + .2
99         xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
100         Z = classifier.predict(np.c_[xx.ravel(), yy.ravel()], num_classifiers)
101         # Put the result into a color plot
102         Z = Z.reshape(xx.shape)
103         plt.pcolormesh(xx, yy, Z, cmap=cm)
104         # Plot also the training points
105         if weights is not None:
106             plt.scatter(X[:, 0], X[:, 1], c=y, s=weights, cmap=cm_bright)
107         else:
108             plt.scatter(X[:, 0], X[:, 1], c=y, cmap=cm_bright)
109         plt.xlim(xx.min(), xx.max())
110         plt.ylim(yy.min(), yy.max())
111         plt.xticks([])
112         plt.yticks([])
113         plt.title(f'num classifiers = {num_classifiers}')
114         plt.draw()
115
116
117     def generate_data(num_samples, noise_ratio):
118         """
119         generate samples X with shape: (num_samples, 2) and labels y with shape (num_samples).
120         num_samples: the number of samples to generate
121         noise_ratio: invert the label for this ratio of the samples
122         """
123         X = np.random.rand(num_samples, 2) * 2 - 1
124         radius = 0.5 ** 2
125         in_circle = np.sum(X ** 2, axis=1) < radius
126         y = np.ones(num_samples)
127         y[in_circle] = -1

```

```
128     y[np.random.choice(num_samples, int(noise_ratio * num_samples))] *= -1
129
130     return X, y
131
132
133
134
```



## 7 task1/src/geo

```
1  {"BNA": "(36.11958985, -86.68308703291842)", "SEA": "(47.4475673, -122.3080158569515)", "JAX": "(45.1683631, 3.6168847)", "I
```

## 8 task1/src/model.py

```
1  """
2  =====
3      Introduction to Machine Learning (67577)
4      IML HACKATHON, June 2020
5
6  Author(s):
7
8  =====
9  """
10
11 import pandas as pd
12 import pickle
13 from classifier import final_pre_proc
14 from strong import WeakTeam
15 from binarysearch import binarysearch
16 from adaboost import AdaBoost
17 from ex4_tools import DecisionStump
18
19 class FlightPredictor:
20     def __init__(self, path_to_weather=''):
21         """
22         Initialize an object from this class.
23         @param path_to_weather: The path to a csv file containing weather data.
24         """
25         #raise NotImplementedError
26         self.mod = pickle.load(open("./BinAgent" , "rb"))
27
28     def predict(self, x):
29         """
30         Recieves a pandas DataFrame of shape (m, 15) with m flight features, and predicts their
31         delay at arrival and the main factor for the delay.
32         @param x: A pandas DataFrame with shape (m, 15)
33         @return: A pandas DataFrame with shape (m, 2) with your prediction
34         """
35
36         _dataset, y, x_factor, y_factor = final_pre_proc(x)
37         self.mod.predict(x)
```

## 9 task1/src/models.py

```
1  import numpy as np
2  import pandas as pd
3  from sklearn.svm import SVC
4
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.tree import DecisionTreeClassifier
7  from sklearn.neighbors import KNeighborsClassifier
8  from sklearn.linear_model import Lasso, Ridge
9
10 import matplotlib.pyplot as plt
11 from mpl_toolkits.mplot3d import \
12     Axes3D # <--- This is important for 3d plotting
13 import pandas as pd
14 from pandas import DataFrame
15 # from plotnine import *
16
17 import matplotlib as mpl
18 import matplotlib.pyplot as plt
19
20 from random import random
21
22 from ex4_tools import DecisionStump
23
24
25 class abcModel:
26
27     def __init__(self):
28         self.mod = None
29
30     def fit(self, X, y):
31         self.mod.fit(self.bias(X), y.ravel())
32
33     def predict(self, X):
34         return self.mod.predict(self.bias(X))
35
36     def get_hyperplan(self):
37         pass
38
39     def bias(self, X):
40         return np.insert(X, 0, 1, 1)
41
42     def draw(self):
43         pass
44
45     def score(self, X, y):
46         return {
47             "num_samples": 0,
48             "error" : 0,
49             "accuracy" : 0,
50             "FPR": 0,
51             "TPR": 0,
52             "precision": 0,
53             "recall": 0
54         }
55
56
57 # [v]
58 class Perceptron(abcModel):
59     def __init__(self):
```

```

60         self.W = np.array([])
61         super().__init__()
62
63     def fit(self, _X, y):
64         X = self.bias(_X)
65
66         self.W = np.zeros(shape=X.shape[1])
67
68     def not_classifiy():
69         return [(x, y) for x in X if self.predict(x) != y]
70
71     _updated = True
72     while _updated:
73         _updated = False
74         for i in range(len(y)):
75             if np.dot(self.W, X[i]) * y[i] <= 0:
76                 self.W = self.W + (X[i] * y[i])
77             _updated = True
78     return self.W
79
80     def predict(self, X):
81         return np.sign(self.W @ self.bias(X).transpose())
82
83
84 class LDA(abcModel):
85     def __init__(self):
86         super().__init__()
87
88     def fit(self, X, y):
89         X = self.bias(X)
90
91     def gen_delta_y(X, y, y_val):
92         _X = X[y == y_val]
93         lnP = np.log(len(y == y_val) / len(y))
94
95         arithmetic_mean = np.array([np.mean(u) for u in _X.transpose()])
96         _cov = np.cov(X.transpose())
97         _inv_cov = np.linalg.pinv(_cov)
98
99     def delta(x):
100         return x.transpose() @ _inv_cov @ arithmetic_mean - \
101             0.5 * arithmetic_mean.transpose() @ _inv_cov @ arithmetic_mean + \
102             lnP
103
104     return delta
105
106     self.deltas = [gen_delta_y(X, y, y_val) for y_val in [-1, 1]]
107
108     def predict(self, U):
109         return np.array([0: -1.0, 1: 1.0][
110             np.argmax([_delta(u) for _delta in self.deltas])
111             for u in self.bias(U)])
112
113
114 class SVM(abcModel):
115     def __init__(self):
116         super().__init__()
117         self.mod = SVC(C=1e10, kernel='linear')
118
119     def coef_(self):
120         return self.mod.coef_
121
122
123 class Logistic(abcModel):
124     def __init__(self):
125         super().__init__()
126         self.mod = LogisticRegression(solver='liblinear')
127

```

```

128     def fit(self, X, y):
129         super().fit(X, y.flatten())
130
131     # def predict(self, X):
132     #     return self.mod.predict(self.bias(X))
133
134 class DecisionTree(abcModel):
135     def __init__(self, max_depth=2):
136         super().__init__()
137         self.mod = DecisionTreeClassifier(max_depth=2)
138
139
140 class KNearestNeighbor(abcModel):
141
142     def __init__(self):
143         super().__init__()
144         self.mod = KNeighborsClassifier(n_neighbors=40)
145
146
147 class RidgeClassifier(abcModel):
148
149     def __init__(self):
150         super().__init__()
151         self.mod = Ridge(alpha=0, normalize=True)
152
153
154 class LassoClassifier(abcModel):
155
156     def __init__(self):
157         super().__init__()
158         self.mod = Lasso(alpha=0, normalize=True)
159
160 class DecisionStumpWarper(abcModel):
161     def __init__(self):
162         super().__init__()
163         self.mod = DecisionStump( )
164
165     def fit(self, X, y, D=None):
166         self.mod.fit(D, X, y.flatten())
167
168     def predict(self, X):
169         return self.mod.predict(X)
170
171
172 def label_f(weight, bias):
173     def _label_f(x):
174         return np.sign(np.dot(weight, x) + bias)
175
176     return _label_f
177
178
179 def draw_points(m, label_function=label_f(np.array([0.3, -0.5]), 0.1)):
180     def _draw_points_n(m, n):
181         return np.random.multivariate_normal(
182             np.zeros(n), np.identity(n), m)
183
184     X = _draw_points_n(m, 2)
185     y = np.array([label_function(vec) for vec in X])
186     return X, y
187
188
189 def analyze_clsifiers():
190     # def generate_line_prec(prec):
191
192     def plot(prec, _svm, X, y):
193         plt.scatter(X[y == -1][:, 0], X[y == -1][:, 1])
194         plt.scatter(X[y == 1][:, 0], X[y == 1][:, 1])
195

```

```

196     min_x, max_x = min(X[:, 0]), max(X[:, 0])
197     min_y, max_y = min(X[:, 1]), max(X[:, 1])
198
199     _min_range, _max_range = min(min_x, min_y), max(max_x, max_y)
200     xx = [_min_range, _max_range]
201
202     def get_y(W, _x):
203         return -(W[0] + _x * W[1]) / W[2] if W[2] != 0 else -W[0]
204
205     def get_y_prep(_x):
206         return get_y(prec.W, _x)
207
208     def get_y_svm(_x):
209         print(_svm.coef_()[0])
210         return get_y(_svm.coef_()[0], _x)
211
212     def get_true_y(_x):
213         return 0.1 / 0.5 + 0.3 / 0.5 * _x
214
215     plt.xlim([_min_range, _max_range])
216     plt.ylim([_min_range, _max_range])
217
218     middle = (_min_range + _max_range) / 2
219
220     def print_line(msg, _f, _color):
221         xx = [_min_range, _max_range]
222         yy = [_f(_x) for _x in xx]
223         _x = middle + 2 * (0.5 - random())
224         _y = _f(_x)
225         plt.plot(xx, yy, color=_color)
226         plt.annotate(msg, color=_color,
227                      xy=(_x, _y), xycoords='data',
228                      xytext=(_x + 0.3, _y), textcoords='data',
229                      arrowprops=dict(arrowstyle="->"))
230
231     print_line("prep", get_y_prep, "C5")
232     print_line("svm", get_y_svm, "C4")
233     print_line("true plane", get_true_y, "C2")
234     plt.title("svm vs prep")
235     plt.xlabel("x")
236     plt.ylabel("y")
237     plt.show()
238
239     for m in [5, 10, 15, 25, 70]:
240         X, y = draw_points(m)
241         blues, reds = X[y == 1], X[y == -1]
242         _modes = [Perceptron(), SVM()]
243         for _model in _modes:
244             _model.fit(deepcopy(X), y)
245
246         plot(_modes[0], _modes[1], X, y)
247
248
249     def expanded_analyze_clssifiers():
250         times, k = 7, 1000
251         modles = []
252         models_num = 3
253
254         def genrate_real_plane(m):
255             _f = label_f(np.array([random(), random()]), random())
256             X, y = draw_points(m, label_function=_f)
257             while (-1 not in y) or (1 not in y):
258                 X, y = draw_points(m, label_function=_f)
259             return X, y, _f
260
261         def accur(mod, _f, Z):
262             _prob = 0
263             for x, y in zip(map(_f, Z), mod.predict(Z)):

```

```

264         if x == y:
265             _prob += 1
266         return _prob / len(Z)
267
268     def one_iteraion(m):
269         _modes = [Perceptron(), SVM(), LDA()]
270         X, y, _f = genrate_real_plane(m)
271         ret = []
272         for _model in _modes:
273             print(type(_model))
274             _model.fit(deepcopy(X), y)
275             Z, _ = draw_points(k)
276             ret.append(accur(_model, _f, Z))
277         return np.array(ret)
278
279     def calc_mean_performance(M=[5, 10, 15, 25, 70]):
280         ret = []
281         for m in M:
282             _mean = np.zeros(models_num)
283             for _ in range(times):
284                 _mean += one_iteraion(m)
285             ret.append(_mean / times)
286         return M, np.array(ret)
287
288     m, mean_performance = calc_mean_performance()
289     for _model_num, _name in enumerate(["perc", "svm", "lda"]):
290         print(_name)
291         print(mean_performance)
292         plt.plot(m, mean_performance[:, _model_num])
293     plt.legend(["perc", "svm", "lda"])
294     plt.title("calc_mean_performance")
295     plt.xlabel("m (size of the given training data)")
296     plt.ylabel("propability of successes")
297     plt.show()
298
299
300 if __name__ == "__main__":
301     X, y = draw_points(10)
302     # p = Perceptron()
303
304     from copy import deepcopy
305
306     models_class = [Perceptron, SVM, Logistic, DecisionTree, LDA]
307     models = []
308     for mod in models_class:
309         models.append(mod())
310         print("{} init ".format(type(mod)))
311
312     for mod in models:
313         mod.fit(deepcopy(X), y)
314         print("{} fit ".format(type(mod)))
315
316     for mod in models:
317         print(mod.predict(deepcopy(X)))
318         print("{} predict ".format(type(mod)))
319
320     # analyze_clssifiers()
321     expanded_analyze_clssifiers()

```

## 10 task1/src/strong.py

```
1  from models import abcModel, DecisionTree, Logistic, SVM, DecisionStumpWarper
2  from adaboost import AdaBoost, AdaBoostList
3  from itertools import combinations
4  import numpy as np
5  import pandas as pd
6  from binarysearch import binarysearch
7  from datetime import date
8  from random import shuffle
9  import pickle
10
11  from classifier import final_pre_proc
12  from sklearn.model_selection import train_test_split
13
14
15  # class WeakFactory:
16  #     class WeakLernerByFeature(abcModel):
17  #
18  #         def __init__(self, _model):
19  #             self.mod = _model
20  #
21  #         def fit(self, X, y):
22  #             self.mod.fit(X,y)
23  #
24  #         def predict(self, X):
25  #             return self.mod.predict()
26
27
28  #     def __init__ (self):
29  #         pass
30
31  #     @staticmethod
32  #     def CreateWeaks(self):
33
34  #         return {
35
36  #             "DayOfWeek" : WeakFactory ( DecisionTree(max_depth=2) ),
37  #             "FlightDate" : WeakFactory ( DecisionTree(max_depth=2) ),
38  #             "Reporting_Airline" : WeakFactory ( DecisionTree(max_depth=2) ),
39  #             "Tail_Number" : WeakFactory ( DecisionTree(max_depth=2) ),
40  #             "Flight_Number_Reporting_Airline" : WeakFactory ( DecisionTree(max_depth=2) ),
41  #             "Origin" : WeakFactory ( DecisionTree(max_depth=2) ),
42  #             "OriginCityName" : WeakFactory ( DecisionTree(max_depth=2) ),
43  #             "OriginState" : WeakFactory ( DecisionTree(max_depth=2) ),
44  #             "Dest" : WeakFactory ( DecisionTree(max_depth=2) ),
45  #             "DestCityName" : WeakFactory ( DecisionTree(max_depth=2) ),
46  #             "DestState" : WeakFactory ( DecisionTree(max_depth=2) ),
47  #             "CRSDepTime" : WeakFactory ( DecisionTree(max_depth=2) ),
48  #             "CRSArrTime" : WeakFactory ( DecisionTree(max_depth=2) ),
49  #             "CRSElapsedTime" : WeakFactory ( DecisionTree(max_depth=2) ),
50  #             "Distance" : WeakFactory ( DecisionTree(max_depth=2) ),
51  #             "ArrDelay" : WeakFactory ( DecisionTree(max_depth=2) ),
52  #             "DelayFacto" : WeakFactory ( DecisionTree(max_depth=2) )
53
54
55  #return generateWeakClass( DecisionTree )
56  # return generateWeakClass( Logistic )
57
58  # " extracrdrd out for pickiling "
59  class WeakTeam(DecisionStumpWarper):
```



```

60
61     def __init__(self, featuers=True):
62         DecisionStumpWarper.__init__(self)
63         self.featuers = featuers
64
65     def filterX( self, X ):
66         ret = np.array( pd.DataFrame( { featuer: X[featuer] for featuer in self.featuers } ))
67         return ret
68
69     def train(self, X, y, D=None):
70         #print(f"[@] train on features : { self.featuers}")
71
72         if D is None:
73             super().fit(
74                 np.array(self.filterX(X)), y)
75         else:
76             super().fit(
77                 np.array(self.filterX(X)), y, D)
78
79     def predict(self, X):
80         return super().predict( self.filterX(X) )
81
82     def generateTeamClass(featuers):
83         return (WeakTeam, featuers)
84
85
86     #
87
88
89     # def binarysearch_read(_file):
90     #     mods = {}
91     #     for banch in re.split('$', _file.read()):
92     #         treshold, _strmod = re.split(':', banch)
93     #         mods[float(treshold)] = AdaBoost_read(banch)
94
95
96     def calc_error(model, _dataframe, y, agents):
97         _error = 0
98         z = (y.flatten() - 0.5 * np.ones(len(y))) * 2
99         for _bool in (model.predict(_dataframe, max_t=agents) != z):
100             _error += {False: 0, True: 1][_bool]
101         return _error / len(y)
102
103
104     import heapq
105
106
107     def hash_strings(featuers, _list):
108         ret = 0
109         base = 1
110         for featuer in featuers:
111             if featuer in _list:
112                 ret += base
113             base *= 10
114         return ret
115
116
117     def learn(_dataframe, y, featuers, teams=set(), depth=5, orignal=[],
118             _hased=set()):
119         if len(teams) == 0:
120             teams = [[featuer] for featuer in featuers]
121
122         agents = 1
123         _hashed = set()
124         new_team = []
125
126         def create_subgroups():
127             subgroups = []

```

```

128         for featuer in featuers:
129             for team in teams:
130                 if featuer not in team:
131                     _hash = hash_strings(original, team + [featuer])
132                     if _hash not in _hashed:
133                         # print(f"i was here {team + [featuer]}, _hash : {_hash}")
134                         subgroups.append(generateTeamClass(team + [featuer]))
135                         _hashed.add(_hash)
136         return subgroups
137
138 strongGroups = []
139 heap = []
140
141 for weak in create_subgroups():
142     _model = AdaBoost(weak, agents, support_wights=True)
143     _model.train(_dataframe, y)
144     strongGroups.append(_model)
145     heapq.heappush(heap, (
146         -calc_error(_model, _dataframe, y, agents), _model.h[0].featuers,
147         _model))
148
149     if len(heap) > 3:
150         heapq.heappop( heap )
151
152 if depth == 0:
153     while len(heap) > 1:
154         train_error, _featuers, _model = heapq.heappop( heap )
155         train_error, _featuers, _model = heapq.heappop( heap )
156         return train_error, _featuers, _model
157
158 else:
159     teams = []
160     while len(heap) > 1:
161         heapq.heappop( heap )
162         train_error, _featuers, _model = heapq.heappop( heap )
163         teams.append( _featuers )
164         return learn(_dataframe, y, featuers, teams, depth - 1,
165             original=original)
166
167
168
169
170 def generateY(_dataset, treshold=0):
171     y = []
172     for _bool in _dataset["ArrDelay"] > treshold:
173         y.append({False: [0], True: [1]}[_bool])
174     return np.array(y)
175
176
177 def generateYbyString(_dataset, s, treshold=0):
178     y = []
179     for _bool in _dataset[s] > treshold:
180         y.append({False: [0], True: [1]}[_bool])
181     return np.array(y)
182
183
184 def pairs():
185     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 0,
186                                     df_copy["DelayFactor"] == 1)]
187     df2 = df1.drop(columns=["DelayFactor"])
188     _, _, classifier1 = learn(df2, df1["DelayFactor"].as_matrix(), df2.keys(), original=df2.keys())
189
190     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 1,
191                                     df_copy["DelayFactor"] == 2)]
192     df2 = df1.drop(columns=["DelayFactor"])
193     _, _, classifier2 = learn(df2, df1["DelayFactor"].as_matrix() - 1, df2.keys(),
194                             original=df2.keys())
195

```

```

196     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 2,
197                                     df_copy["DelayFactor"] == 3)]
198     df2 = df1.drop(columns=["DelayFactor"])
199     _, _, classifier3 = learn(df2, df1["DelayFactor"].as_matrix() - 2, df2.keys(),
200                               original=df2.keys())
201
202     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 3,
203                                     df_copy["DelayFactor"] == 1)]
204     df2 = df1.drop(columns=["DelayFactor"])
205     s = df1["DelayFactor"].as_matrix()
206     s[s == 3] = 0
207     df1 = df1.drop(columns=["DelayFactor"])
208     df1["DelayFactor"] = s
209     _, _, classifier4 = learn(df2, df1["DelayFactor"].as_matrix(), df2.keys(),
210                               original=df2.keys())
211
212     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 2,
213                                     df_copy["DelayFactor"] == 0)]
214     df2 = df1.drop(columns=["DelayFactor"])
215     s = df1["DelayFactor"].as_matrix()
216     s[s == 2] = 1
217     df1 = df1.drop(columns=["DelayFactor"])
218     df1["DelayFactor"] = s
219     _, _, classifier5 = learn(df2, df1["DelayFactor"].as_matrix(), df2.keys(),
220                               original=df2.keys())
221
222     df1 = df_copy.loc[np.logical_or(df_copy["DelayFactor"] == 3,
223                                     df_copy["DelayFactor"] == 0)]
224     df2 = df1.drop(columns=["DelayFactor"])
225     s = df1["DelayFactor"].as_matrix()
226     s[s == 3] = 1
227     df1 = df1.drop(columns=["DelayFactor"])
228     df1["DelayFactor"] = s
229     _, _, classifier6 = learn(df2, df1["DelayFactor"].as_matrix(), df2.keys(),
230                               original=df2.keys())
231     classifier1.predict(df2)
232     return [classifier1, classifier2, classifier3, classifier4, classifier5, classifier6]
233
234
235
236
237 def pre_proc(_dataset, dropped_fe, categorical):
238     y = generateY(_dataset)
239
240     if len(categorical) > 0:
241         cat = pd.DataFrame(
242             pd.get_dummies(_dataset[categorical].astype('category')))
243
244     _dataset = _dataset.drop(dropped_fe + categorical, axis=1)
245     if len(categorical) > 0:
246         _dataset_prepoc = pd.concat(
247             [_dataset.reset_index(drop=True), cat.reset_index(drop=True)],
248             axis=1)
249     else:
250         _dataset_prepoc = _dataset
251
252     print(_dataset_prepoc)
253     return _dataset_prepoc, y
254
255
256 featuers = ["DayOfWeek",
257             "FlightDate",
258             "Reporting_Airline",
259             "Tail_Number",
260             "Flight_Number_Reporting_Airline",
261             "Origin",
262             "OriginCityName",
263             "OriginState",

```

```

264         "Dest",
265         "DestCityName",
266         "DestState",
267         "CRSDepTime",
268         "CRSArrTime",
269         "CRSElapsedTime",
270         "Distance"]
271
272     dropped_fe = ['Flight_Number_Reporting_Airline',
273                  'Tail_Number',
274                  # 'DayOfWeek',
275                  'FlightDate',
276                  'ArrDelay',
277                  'DelayFactor']
278
279     categorical = [
280         'OriginCityName'
281         , 'OriginState'
282         , 'Origin'
283         , 'Dest'
284         , 'DestCityName'
285         , 'DestState'
286         , 'Reporting_Airline']
287
288
289
290     if __name__ == "__main__" :
291         original_dataset = pd.read_csv("~/data/train_data.csv", nrows=80)
292
293         print("#] before pre processing")
294         print(original_dataset)
295         _mods = {}
296         _minrange, _maxrange = -30, 30
297         _dataset, y, x_factor, y_factor = final_pre_proc(original_dataset)
298
299         df_copy = _dataset.copy()
300         df_copy["DelayFactor"] = y_factor
301         pairs()
302
303         start_range , end_range = np.ones(len(y)) * _minrange , np.ones(len(y)) * _maxrange
304         for i, time in enumerate( np.arange(_minrange, _maxrange, (_maxrange-_minrange)/ 2 ) ):
305             train_error, _featuers, _mods[time] = learn(_dataset,
306                                                         generateY(original_dataset,
307                                                             time),
308                                                         _dataset.keys(), teams=[
309                     ['CRSDepTime', 'DayOfWeek']], orignal=_dataset.keys())
310             print(f"{time} : { _featuers } : {train_error}")
311
312         Bagent = binarysearch( _mods )
313         # _middles = Bagent.predict( _dataset , start_range , end_range)
314         # print(_middles)
315
316         _dataset, y, x_factor, y_factor = final_pre_proc(
317             pd.read_csv("~/data/train_data.csv", nrows=10000)[9800:])
318         # _dataset, y = pre_proc(original_dataset, dropped_fe , categorical )
319         _middles = Bagent.mods[0.0].predict(_dataset)
320         _middles[_middles > 0] = 1
321         t = sum((_middles - y.flatten()) ** 2 / len(y))
322         print(f"[error]: {t}")
323         # print( )
324
325         # _dataset, y, x_factor, y_factor= final_pre_proc( pd.read_csv("~/data/train_data.csv", nrows=10000 ) [9800:] )
326         # _middles = Bagent.mods[0.0].predict( _dataset )
327         # _middles[ _middles > 0 ] = 1
328         # t = sum( (_middles- y.flatten())**2/len(y))
329         # print ( f"[error]: {t}" )
330
331         with open("./BinAgent", "wb") as f:

```

```
332         pickle.dump(Bagent, f)
333
334
335
336
337     # class StrongClassifier:
338
339     #     def __init__(self, features):
340     #         pass
341
342     #     def predict(self, X):
343     #         pass
344
345     #     def train(self, X, y):
346     #         pass
```

# 11 task1/src/strong2.py

```
1  from HUJIHACK.Task1.source.models import abcModel, DecisionTree
2  from HUJIHACK.Task1.source.adaboost import AdaBoost, AdaBoostList
3  from itertools import combinations
4  import numpy as np
5  import pandas as pd
6  from sklearn.linear_model import Ridge
7  from sklearn.model_selection import GridSearchCV
8  import matplotlib.pyplot as plt
9  from datetime import date
10 from sklearn.linear_model import LogisticRegression
11
12 import pandas as pd
13 import numpy as np
14 import seaborn as sns
15 import math
16
17 from sklearn.linear_model import Ridge
18
19 from sklearn.metrics import r2_score
20 from sklearn.metrics import mean_squared_error
21
22
23 # class WeakFactory:
24 #     class WeakLernerByFeature(abcModel):
25 #
26 #         def __init__(self, _model):
27 #             self.mod = _model
28 #
29 #         def fit(self, X, y):
30 #             self.mod.fit(X,y)
31 #
32 #         def predict(self, X):
33 #             return self.mod.predict()
34
35
36 #     def __init__ (self):
37 #         pass
38
39 #     @staticmethod
40 #     def CreateWeeks(self):
41
42 #         return {
43
44 #             "DayOfWeek" : WeakFactory ( DecisionTree(max_depth=2) ),
45 #             "FlightDate" : WeakFactory ( DecisionTree(max_depth=2) ),
46 #             "Reporting_Airline" : WeakFactory ( DecisionTree(max_depth=2) ),
47 #             "Tail_Number" : WeakFactory ( DecisionTree(max_depth=2) ),
48 #             "Flight_Number_Reporting_Airline" : WeakFactory ( DecisionTree(max_depth=2) ),
49 #             "Origin" : WeakFactory ( DecisionTree(max_depth=2) ),
50 #             "OriginCityName" : WeakFactory ( DecisionTree(max_depth=2) ),
51 #             "OriginState" : WeakFactory ( DecisionTree(max_depth=2) ),
52 #             "Dest" : WeakFactory ( DecisionTree(max_depth=2) ),
53 #             "DestCityName" : WeakFactory ( DecisionTree(max_depth=2) ),
54 #             "DestState" : WeakFactory ( DecisionTree(max_depth=2) ),
55 #             "CRSDepTime" : WeakFactory ( DecisionTree(max_depth=2) ),
56 #             "CRSArrTime" : WeakFactory ( DecisionTree(max_depth=2) ),
57 #             "CRSElapsedTime" : WeakFactory ( DecisionTree(max_depth=2) ),
58 #             "Distance" : WeakFactory ( DecisionTree(max_depth=2) ),
59 #             "ArrDelay" : WeakFactory ( DecisionTree(max_depth=2) ),
```

```

60 #             "DelayFacto" : WeakFactory ( DecisionTree(max_depth=2) )
61 #         }
62
63 def generateTeamClass(featuers):
64     class WeakTeam(DecisionTree):
65
66         def __init__(self):
67             DecisionTree.__init__(self, max_depth=3)
68             self.featuers = featuers
69
70         def filterX(self, X):
71             ret = np.array(pd.DataFrame({featurer: X[featurer] for featurer in self.featuers}))
72             return ret
73
74         def train(self, X, y):
75             print(f"[0] train on features : {self.featuers}")
76             super().fit(
77                 np.array(self.filterX(X)), y)
78
79         def predict(self, X):
80             return super().predict(self.filterX(X))
81
82     return WeakTeam
83
84
85 def pre_proc_new(_dataset, dropped_fe, categorical):
86     y_del = []
87     y_factor = []
88     for delay, factor in zip(_dataset["ArrDelay"], _dataset["DelayFactor"]):
89         y_del.append(delay)
90         y_factor.append(factor)
91     y_del = np.array(y_del)
92     y_factor = np.array(y_factor)
93     for index, row in _dataset.iterrows():
94         _dataset.loc[index, 'CRSElapsedTime'] = math.floor(row['CRSElapsedTime'] / 10)
95         _dataset.loc[index, 'CRSArrTime'] = math.floor(row['CRSArrTime'] / 100)
96         _dataset.loc[index, 'CRSDepTime'] = math.floor(row['CRSDepTime'] / 100)
97     cat = pd.DataFrame(pd.get_dummies(_dataset[categorical]).astype('category'))
98     _dataset = _dataset.drop(dropped_fe + categorical, axis=1)
99     _dataset_preproc = pd.concat([_dataset.reset_index(drop=True), cat.reset_index(drop=True)], axis=1)
100     return _dataset_preproc, y_del, y_factor
101
102
103 def ridge_reg(x, y):
104     # ridge = Ridge()
105     # parameters = {'alpha': [1, 5, 10, 20, 25, 30, 40, 50, 60, 70, 1000]}
106     # ridge_regressor = GridSearchCV(ridge, parameters, cv=5)
107     # ridge_regressor.fit(x, y)
108     # b = ridge_regressor.predict(x)
109     # print(ridge_regressor.best_params_)
110     # print(ridge_regressor.best_score_)
111     # a = range(1,10001)
112     # plt.scatter(a, y, c='b', s=1, alpha=0.6, label='Training data')
113     # plt.scatter(a, b, c='red', s=1, label='prediction')
114     # plt.plot(a, y, 'k', label='True function')
115     # plt.legend()
116     # plt.show()
117     # rr = Ridge(alpha=5)
118     rr = LogisticRegression(solver='liblinear')
119     rr.fit(x, y)
120     pred_train_rr = rr.predict(x)
121     print(rr.score(x, y))
122     # print(np.sqrt(mean_squared_error(y, pred_train_rr)))
123     # print(r2_score(y, pred_train_rr))
124
125     # pred_test_rr = rr.predict(X_test)
126     # print(np.sqrt(mean_squared_error(y_test, pred_test_rr)))
127     # print(r2_score(y_test, pred_test_rr))

```

```

128
129
130 def identify_corelated_features(df):
131     corr = df.corr()
132     print(sns.heatmap(corr))
133     columns = np.full((corr.shape[0],), True, dtype=bool)
134     for i in range(corr.shape[0]):
135         for j in range(i + 1, corr.shape[0]):
136             if corr.iloc[i, j] >= 0.9:
137                 if columns[j]:
138                     columns[j] = False
139     selected_columns = df.columns[columns]
140     df = df[selected_columns]
141     return df
142
143
144 def learn(_dataframe, y, featuers):
145     agents = 3
146     group_size = 1
147     subgroups = [generateTeamClass(team) for team in combinations(featuers, group_size)]
148
149     def calc_error(model, _dataframe, y):
150         _error = 0
151         for _bool in (model.predict(_dataframe, max_t=1) != y).flatten():
152             _error += {False: 0, True: 1][_bool]
153         return _error / len(y)
154
155     strongGroups = []
156     for weak in subgroups:
157         _model = AdaBoost(weak, agents)
158         _model.train(_dataframe, y)
159         strongGroups.append(_model)
160
161     return strongGroups[np.argmin([calc_error(_model, _dataframe, y) for _model in strongGroups])]
162
163
164 if __name__ == "__main__":
165     _dataset = pd.read_csv("../data/train_data.csv", nrows=10000)
166
167
168     # categorical = ['view', 'waterfront', 'bedrooms', 'grade', 'floors', 'condition', 'bathrooms']
169     # cat = pd.DataFrame( { 'id' : _dataset_prepoc['id'] }, pd.get_dummies(_dataset_prepoc[categorical]).astype('category')
170
171     def pre_proc(_dataset, dropped_fe, categorical):
172
173         def generateY(_dataset):
174             y = []
175             for _bool in _dataset["ArrDelay"] > 0:
176                 y.append({False: [0], True: [1]}[_bool])
177             return np.array(y)
178
179         # def remove_end_cases(_frame):
180         #     return _frame[ _frame['price'] > 0 ]
181         y = generateY(_dataset)
182
183         cat = pd.DataFrame(pd.get_dummies(_dataset[categorical]).astype('category'))
184         _dataset = _dataset.drop(dropped_fe + categorical, axis=1)
185         _dataset_prepoc = pd.concat([_dataset.reset_index(drop=True), cat.reset_index(drop=True)], axis=1)
186         return _dataset_prepoc, y
187
188
189     featuers = ["DayOfWeek",
190                "FlightDate",
191                "Reporting_Airline",
192                "Tail_Number",
193                "Flight_Number_Reporting_Airline",
194                "Origin",
195                "OriginCityName",

```



```

196         "OriginState",
197         "Dest",
198         "DestCityName",
199         "DestState",
200         "CRSDepTime",
201         "CRSArrTime",
202         "CRSElapsedTime",
203         "Distance"]
204
205     dropped_fe = ['Flight_Number_Reporting_Airline',
206                  'Tail_Number',
207                  'FlightDate',
208                  'ArrDelay',
209                  'DelayFactor']
210
211     dropped_fe_new = ['Flight_Number_Reporting_Airline',
212                      'Tail_Number',
213                      'FlightDate',
214                      'ArrDelay',
215                      'DelayFactor',
216                      "OriginCityName",
217                      "OriginState",
218                      "DestCityName",
219                      "DestState", ]
220
221     categorical_new = ['Reporting_Airline', 'Origin', 'Dest']
222
223     categorical = [
224         'OriginCityName'
225         , 'OriginState'
226         , 'Origin'
227         , 'Dest'
228         , 'DestCityName'
229         , 'DestState'
230         , 'Reporting_Airline']
231
232     # "ArrDelay",
233     # "DelayFactor"]
234
235     # print("#] before pre processing")
236     # print(_dataset)
237     # _dataset, y = pre_proc(_dataset, dropped_fe, categorical)
238     # print("#] after pre processing")
239     # print(_dataset)
240     # print("#] y' vector ")
241     # print(y)
242     #
243     # print(_dataset.keys())
244     # _mod = learn(_dataset, y, _dataset.keys())
245     # print("#] best featuers:")
246     # print(_mod.h[0].featuers)
247     x, y_del, y_factor = pre_proc_new(_dataset, dropped_fe_new, categorical_new)
248     print(x)
249     x = identify_corelated_features(x)
250     print(x)
251     # print(x)
252     # print(x.shape)
253     # print(y_del)
254     # print(y_del.shape)
255     # print(y_factor)
256     # ridge_reg(x, y_del)
257     # print(x)
258     # print("after corr\n")
259     # print(identify_corelated_features(x,y_del))
260
261     # class StrongClassifier:
262
263     #     def __init__(self, features):

```

```
264 #         pass
265
266 #     def predict(self, X):
267 #         pass
268
269 #     def train(self, X, y):
270 #         pass
```

## 12 task1/src/testmodel.py

```
1
2
3 from model import FlightPredictor
4 import pandas as pd
5 from classifier import final_pre_proc
6 from strong import WeakTeam
7 from binarysearch import binarysearch
8 from adaboost import AdaBoost
9 from ex4_tools import DecisionStump
10
11 if __name__ == "__main__":
12     print("[@] create FlightPredictor")
13     fp = FlightPredictor( )
14     print("[@] predict")
15     ret = fp.predict( pd.read_csv("~/data/train_data.csv", nrows=10000) )
16     print(ret)
```

## 13 task1/src/weather.py

```
1  import json
2  import pandas as pd
3  import pickle
4
5
6  if __name__ == "__main__" :
7      df = pd.read_csv("~/data/all_weather_data.csv")
8      _json = { }
9      for index, row in df.iterrows():
10         for key in df.keys():
11             _json[ row['day'] ] = {}
12             if key != "day":
13                 _json[ row['day'] ][key] = row[key]
14         json.dump( _json , open("./wether.json", "w") )
15
16
```

## 14 task1/src/wether.json

```
1  {"01-01-10": {"max_wind_gust_kts": "0"}, "02-01-10": {"max_wind_gust_kts": "0"}, "03-01-10": {"max_wind_gust_kts": "0"}, "04
```