

CSCE633 Homework 02

Lu Sun

September 2020

1 Data exploration:

- (a) Categorical attributes
- (b) Continuous attributes

Answer:

From the data set in file "hw2_data.csv", finally results can be get:

- Categorical attributes: *primary_strength* (Column 8)
- Continuous attributes: *stamina*, *attack_value*, *defense_value*, *capture_rate*, *flee_rate*, *spawn_chance* (Column 2-7)

2 Data exploration:

- (a) Plot 2-D scatter plots.
- (b) Compute the Pearson's correlation coefficient.
- (c) Which attributes would be the most predictive.

Answer:

The code for this part is as follow:

```
1 import matplotlib.pyplot as plt
2 import scipy.stats
3
4 y = combat_point
5 for i in range(6):
6     x = NameMap[i][1]
7     xlab = NameMap[i][0]
8     plt.scatter(x, y)
9     plt.xlabel('{0}'.format(xlab))
10    plt.ylabel('combat_point')
11    plt.show()
12
13    PeaRate = scipy.stats.pearsonr(x, y)
14    print('Pearsonr Correlation Rate for {0}: ...
        {1}'.format(xlab, PeaRate))
```

2.1 Plots

The results of plots are shown as follow:

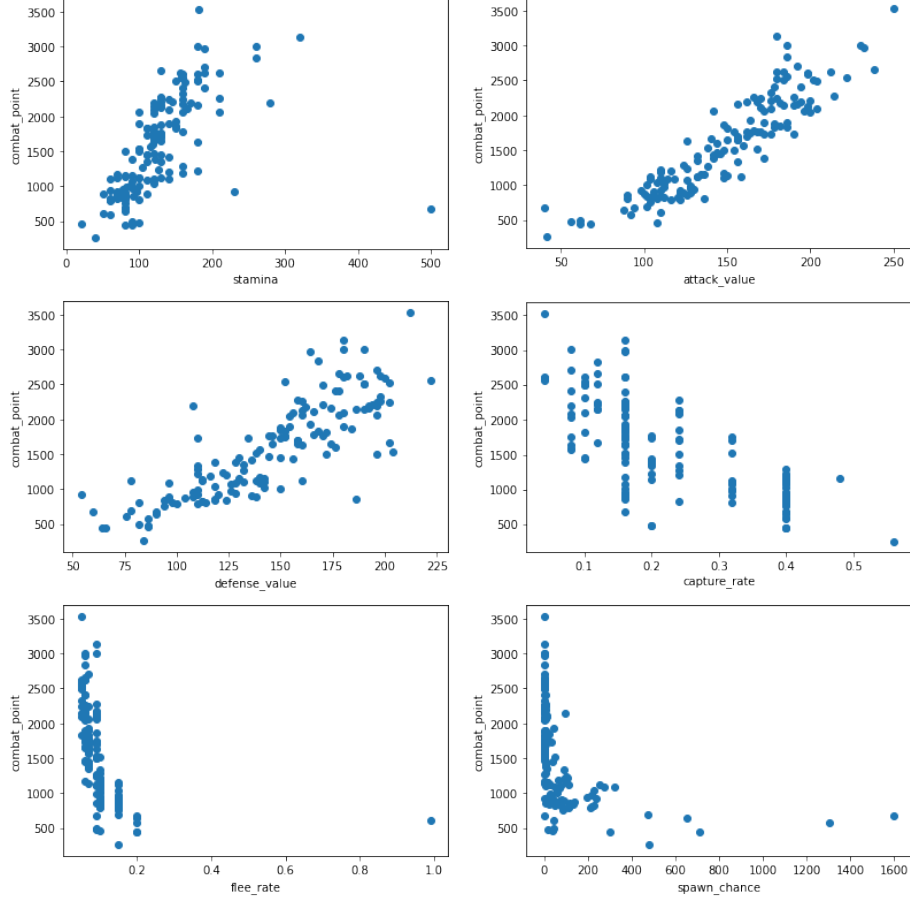


Figure 1: Scatter Plots with combat point

2.2 Pearson's correlation coefficient

The Pearson's correlation coefficients are shown in the following table:

<i>stamina</i>	<i>attack_value</i>	<i>defense_value</i>
0.582831703222926	0.9075315401042733	0.8262293053572931
<i>capture_rate</i>	<i>flee_rate</i>	<i>spawn_chance</i>
-0.7430078083529397	-0.40703421142159646	-0.42132699465983586

Table 1: Pearson's correlation coefficient with combat point

2.3 The most predictive attribute

From the table 1, *attack_value* is the most predictive attribute. The reason is as follow:

attack_value has the largest absolute value of Pearson's correlation coefficient, which is closest to 1 among all the other attributes. By the definition of Pearson's correlation coefficient, the closer distance it has to 1 or -1, the more predictable relation data set will get.

3 Data exploration:

- (a) Plot scatter plots among continuous attributes only.
- (b) Compute corresponding Pearson's correlation coefficient.
- (c) Which variables are the most correlated to each other.

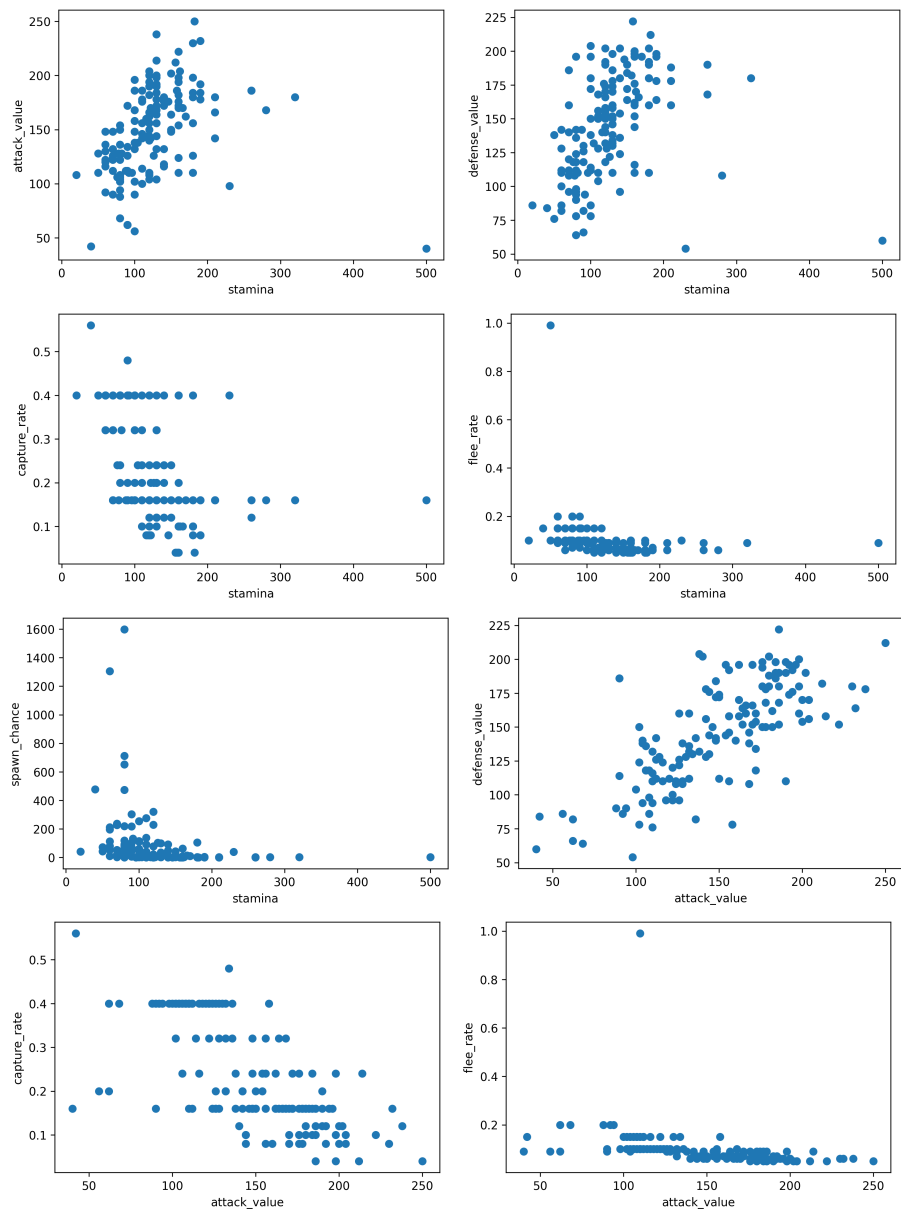
Answer:

The code for this part is as follow:

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3
4 for i in range(6):
5     j = i+1
6     while j<6:
7         filename = ','.join(['/content/drive/My ...
8             Drive/2020_FALL/CSCE633/HomeWork/HW02/Figure/...
9             prob3-'+str(i)+str(j)+'.png' ])
10
11         x = NameMap[i][1]
12         xlab = NameMap[i][0]
13         y = NameMap[j][1]
14         ylab = NameMap[j][0]
15
16         plt.scatter(x, y)
17         plt.xlabel('{0}'.format(xlab))
18         plt.ylabel('{0}'.format(ylab))
19         plt.savefig('{0}'.format(filename), dpi=300, ...
20             bbox_inches='tight') # Returns the list of datapoints
21
22         plt.show()
23
24         PeaRate = scipy.stats.pearsonr(x,y)
25         print('Pearsonr Correlation Rate for ...
26             ({0},{1}):{2}'.format(xlab,ylab,PeaRate))
27
28         j += 1
```

3.1 Plots

The scatter plots are shown as follow:



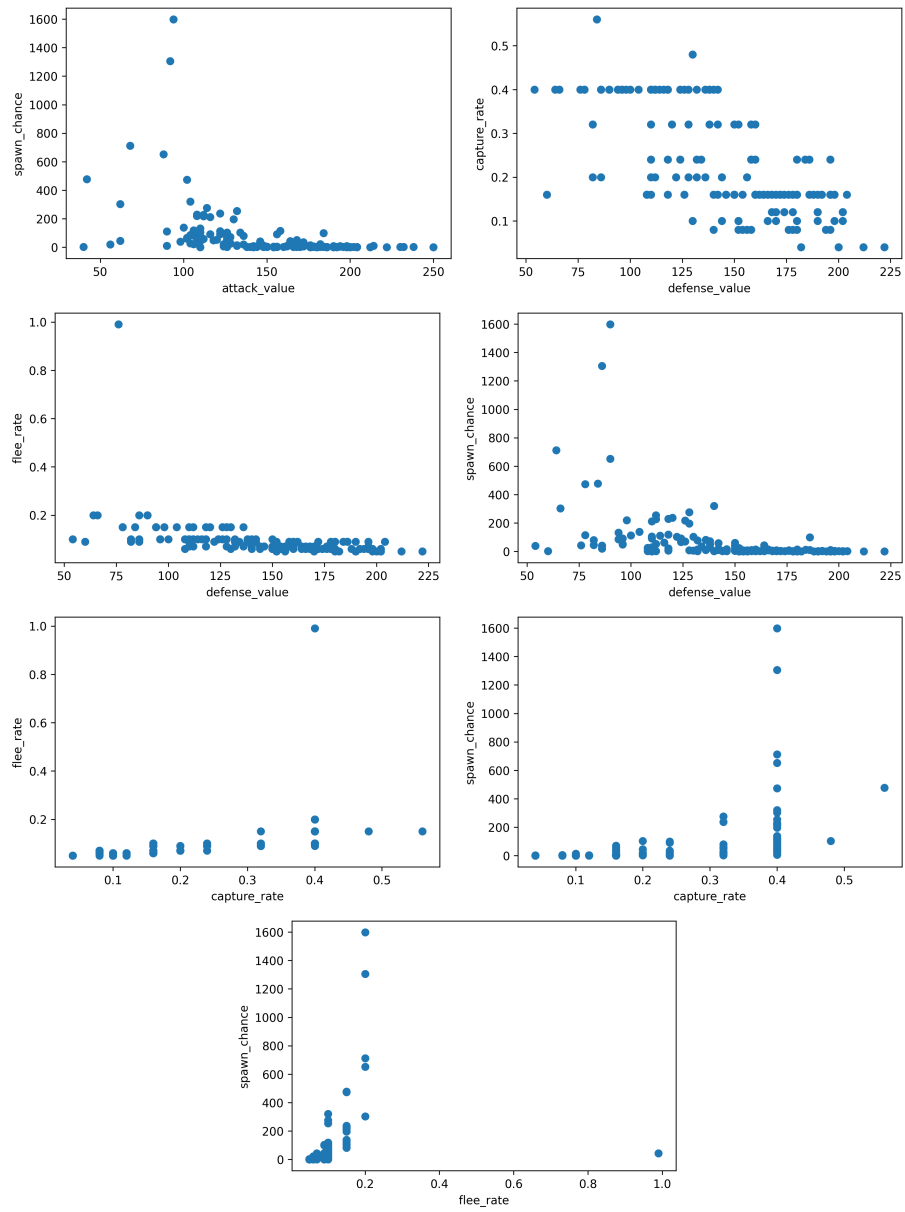


Figure 2: Scatter Plots among continuous attributes

3.2 Pearson's correlation coefficient

The Pearson's correlation coefficients among continuous attributes are shown in the following table.

$(stamina, attack_value)$ 0.3029949826738916	$(stamina, defense_value)$ 0.30266333625368935	$(stamina, capture_rate)$ -0.4468503047144601
$(stamina, flee_rate)$ -0.2710475393248393	$(stamina, spawn_chance)$ -0.2764202078836037	$(attack_value, defense_value)$ 0.7367766467515237
$(attack_value, capture_rate)$ -0.6905726716022137	$(attack_value, flee_rate)$ -0.36906414197600723	$(attack_value, spawn_chance)$ -0.4326484402010869
$(defense_value, capture_rate)$ -0.6972657162131648	$(defense_value, flee_rate)$ -0.42385975623729333	$(defense_value, spawn_chance)$ -0.43249856208332005
$(capture_rate, flee_rate)$ 0.44051150728059624	$(capture_rate, spawn_chance)$ 0.4727927266445679	$(flee_rate, spawn_chance)$ 0.29322169222082034

Table 2: Pearson's correlation coefficient among continuous attributes

3.3 The most related pair

From the above table 2, the most related pair is $(attack_value, defense_value)$ whose Pearson's correlation coefficient is 0.7367766467515237. The reason for this is total the same as that in problem 2.3.

4 Pre-processing of categorical variables:

- Count the number of different values.
- Implement the one hot encoding.

Answer:

- First, generate all the classes without duplicates by python's package *np.unique* and compute the total number of class as *onehot_length*
- Then, map each class to corresponding one hot vector as *onehot_mapTotal*
- Finally, use *combat_point_new* to record all the one hot vectors

The code for this problem is as follow:

```
1 combat_point = NameMap[6][1]
2 print(combat_point)
3
4 onehot_class = np.unique(combat_point)
5 onehot_length = len(onehot_class)
6 onehot_mapTotal = {}
```

```

7  for i in range(onehot.length):
8      onehot = np.zeros(onehot.length,dtype=int)
9      onehot[i] = 1
10     onehot_map = {onehot.class[i]:onehot}
11     onehot_mapTotal.update(onehot_map)
12     #print(combat_point)
13     #print(onehot_mapTotal)
14
15     combat_point_new = []
16     for i in range(len(combat_point)):
17         combat_point_new.append(onehot_mapTotal[combat_point[i]])
18     #print(combat_point)
19     #print(combat_point_new)

```

5 Prediction combat points:

- How many parameters does OLS model have.
- Implement a linear regression model with OLS and implement RSS error computation code.
- Randomly split data and implement 5-fold cross-validation process.
- Compute square root of RSS for each fold and the average of all folds.

Answer:

5.1 Data process

From previous problems, $N_{features} = 22$ feature parameters are needed. The reason is as follow:

Here, feature parameter is represented by $\omega \in R^{N \times 1}$ and ω_0 is the parameter corresponding to constant feature 1. Then, in the previous problem, the number of continuous features are 6, namely *stamina, attack_value, defense_value, capture_rate, flee_rate, spawn_chance*. Then number for categorical feature *primary_strength*'s class are 15, namely *Bu, Dragon, Electric, Fair, Fighting, Fire, Ghost, Grass, Ground, Ice, Normal, Poison, Psychic, Rock, Water*. Finally, $22 = 1 + 6 + 15$ can be got.

5.2 Linear regression and RSS algorithms

In this problem 5 only ordinary linear regression model(OLS) is needed to generate, while in problem 6, regularity linear regression model(RLS) is needed to generate. In this part, the two methods will be generated at the same time.

First, initialize $\omega(0) = \vec{0}$ and fix hyper parameters' value, learning rate and max iteration number.

Then, implement function *score* to compute RSS within whose process function *predict* is needed for computing $X\omega$. In function *score*, finally, $(X\omega - y)^2$ can be obtained.

In the process of implementing *OLS* and *RLS* algorithm, the only difference is the way to generate gradients. In my code, basic gradient descent method is used in both algorithm. Then the formulas of gradients for *OLS* and *RLS* are shown as follow: (Here, assume $X_0[samples] = 1$) The above two formula

$$\begin{aligned} \text{OLS: } \nabla_{\omega} J_{OLS}(\omega) &= 2 \cdot X^T(X\omega - y) \\ \text{RLS: } \nabla_{\omega} J_{RLS}(\omega) &= 2 \cdot X^T(X\omega - y) + 2\lambda \cdot \omega \end{aligned}$$

can be obtained in function *OLR_gradients*, *RLR_gradients* respectively. And function *RLR_gradients* needs one more hyper parameter, regularity parameter λ , than function *OLR_gradients*

Finally, function *fit_GD_OLR* and *fit_GD_RLR* are the corresponding gradient descent method by using the above gradients respectively.

The code is shown as follow:

```

1  class LinearRegression(object):
2
3      def __init__(self, W, learning_rate, max_iter):
4          self.learning_rate = learning_rate
5          self.max_iter = max_iter
6          self.W = W
7
8      def get_params(self):
9          return self.W
10
11     def predict(self, X):
12         return X.dot(self.W)
13
14     def score(self, X, y):
15         score_tep = self.predict(X) - y
16         score = np.transpose(score_tep).dot(score_tep)
17         return np.sqrt(score)
18
19     #####Ordinary#####
20     def fit_GD_OLR(self, X, y):
21         sample_size, feature_size = X.shape
22         #self.W = np.zeros(feature_size) # Initialize model ...
23         #parameters to zeros
24
25         losses = [] # Store loss values at each epoch
26         for t in range(self.max_iter):
27             losses_tep = self.score(X, y)
28             losses.append(losses_tep)
29
30             gt = 2*self.OLR_gradients(X, y)/sample_size
31             self.W = self.W - self.learning_rate * gt
32         return self.W, losses
33
34     def OLR_gradients(self, _x, _y):
35         _g = _x.T.dot(self.predict(_x) - _y)
36         # _g = ( self.W.dot(_x) - _y ) * _x
37         return _g
38
39     #####Regularity#####

```



```

38 def fit_GD_RLR(self, X, y, labd):
39     sample_size, feature_size = X.shape
40
41     losses = [] # Store loss values at each epoch
42     for t in range(self.max_iter):
43         losses_tep = self.score(X, y)
44         losses.append(losses_tep)
45         #if losses_tep < 2900:
46         #    print("stop ahead at iteration: {0}".format(t))
47         #    break
48         gt = 2*self.RLR_gradients(X, y, labd)
49         self.W = self.W - self.learning_rate * gt
50     return self.W, losses
51
52 def RLR_gradients(self, _x, _y, labd):
53     _g = _x.T.dot(self.predict(_x) - _y) + labd * self.W
54     return _g

```

5.3 5-fold cross-validation

In this part, the 5-fold cross-validation algorithm is generated.

First, a set of random index is generated and stores in *random_index* whose length is the same as sample size.

Then, divide the training set into two part, training set and validation set according to the random index set in the first step. Namely, for example, select *random_index*[0]th, *random_index*[5]th, *random_index*[10]th, ... number in original input (*X*, *y*) as the first fold validation data set and the left as training data set. Repeat this part for 5 times to complete 5-fold cross-validation. Within each fold, use corresponding method *OLS* in this problem.

The code is as follow:

```

1
2 def cross_validation(X, y, W_ini, method, ...
3     labd=None, K_fold=5, learning_rate = 1e-7, max_iter = 15000):
4     #MethodMap = {'OrdinaryLinearR': 0,
5     #             'RegularLinearR': 1,
6     #             'l2normLogisticR': 2,
7     #             'l1normLogisticR': 3}
8     #random divide data set into 5 folds
9     random_index = np.random.RandomState(seed = ...
10                                     0).permutation(X.shape[0])
11
12     Lg = len(random_index)
13     w_tol = []
14     rss_tol = []
15     losses = None
16
17     for i in range(K_fold):
18         # get: i, i+5, i+10, i+15, ... for validation index
19         index_val = random_index[i::K_fold]
20         # get the remaining for training index
21         index_train = np.delete(random_index, np.arange(Lg) [i::K_fold])

```

```

22
23     X_val = X[index_val]
24     X_train = X[index_train]
25
26     y_val = y[index_val]
27     y_train = y[index_train]
28
29     if method == 0:
30         LRA = LinearRegression(W = W_ini, ...
31                               learning_rate = learning_rate, max_iter = max_iter)
32         w_tep, losses = LRA.fit_GD_OLR(X_train, y_train)
33         rss_tep = LRA.score(X_val, y_val)
34
35     elif method == 1:
36         LRA = LinearRegression(W = W_ini, ...
37                               learning_rate = learning_rate, max_iter = max_iter)
38         w_tep, losses = LRA.fit_GD_RLR(X_train, y_train, labd)
39         rss_tep = LRA.score(X_val, y_val)
40
41     elif method == 2:
42         clf = LogisticRegression(penalty='l2', C=labd, ...
43                                solver='saga', max_iter=max_iter)
44         clf.fit(X_train, y_train)
45         w_tep = clf.coef_
46         rss_tep = clf.score(X_val, y_val)
47
48     elif method == 3:
49         clf = LogisticRegression(penalty='l1', C=labd, ...
50                                solver='saga', max_iter=max_iter)
51         clf.fit(X_train, y_train)
52         w_tep = clf.coef_
53         rss_tep = clf.score(X_val, y_val)
54
55     if losses != None:
56         # Plot losses to make sure it converges with the fixed ...
57         # hyper parameter
58         plt.plot(losses, '+-')
59         plt.xlabel("number of epochs")
60         plt.ylabel("loss")
61         plt.show()
62
63     w_tol.append(w_tep)
64     rss_tol.append(rss_tep)
65     print('fold {0} rss:{1}'.format(i+1, rss_tep))
66     rss_ave = np.sum(rss_tol)/K_fold
67     print('average rss:', rss_ave)
68     return w_tol, rss_tol, rss_ave

```

Here, method '0' means ordinary linear regression which is used in this problem 5. The method '1' is corresponding to regularity linear regression which is needed in next problem 6. The method '2' and '3' represent logistic regression with l-2 norm and l-1 norm respectively, which will be used in the following problem 9. And finally, 10^{-7} is chosen as the learning rate which can get convergence results.

5.4 Compute RSS

The code for getting results is shown as follow:

```

1 #####Generate basic form dataset#####
2 y = np.array(NameMap[7][1])
3 X_continuous = np.array([NameMap[i][1] for i in range(6)]).T
4 X_categorical = np.array(primary_strength_new)
5
6 X = np.hstack((X_continuous,X_categorical))
7 X_basis = np.hstack((np.ones([X.shape[0],1]),X))
8
9 W_ini = np.linalg.inv(np.transpose(X_basis).dot(X_basis))...
10                                     dot(np.transpose(X_basis).dot(y))
11 #####Cross Validation#####
12 _,rss_tol,rss_ave = cross_validation(X_basis,y,0*W_ini,method = 0)
13 print('all rss:{},average rss:{}'.format(rss_tol,rss_ave))

```

The results are shown in the following table.

fold 1	1663.962670802079
fold 2	2193.803990273198
fold 3	1262.3876327445237
fold 4	2249.7536666025358
fold 5	1374.2513634282927
average	1748.8318647701258

Table 3: square root of RSS for OLS

6 Prediction combat points:

- Implement ridge regression(linear regression with regularization)
- Experiment and report the results with different λ

Answer: The corresponding method has been mentioned in problem 5 and implemented in function *fit_GD_RLR*. Initialize $\lambda = 0.001, 0.01, 1, 10, 100$. Different λ s share the same way of random dividing data set into 5 folds. The code to get results is shown as follow:

```

1 labd_total = [0.01,0.1,1,10,100]
2 rss_tol = []
3 rss_ave = []
4 for labd in labd_total:
5     _,rss_tol_tep,rss_ave_tep = ...
6         cross_validation(X_basis,y,0*W_ini,method = 1,labd = labd)
7     print('lambda:{}'.format(labd))
8     print('all rss:{},average rss:{}'.format(rss_tol_tep,rss_ave_tep))

```

```

9     rss_tol.append(rss_tol_tep)
10    rss_ave.append(rss_ave_tep)
11
12    for i in range(len(labd_total)):
13        print('lambda:{}, average rss:{}, all ...
            rss:{}'.format(labd_total[i], rss_ave[i], rss_tol[i]))

```

The results are shown in the following table:

fold	$\lambda = 0.01$	$\lambda = 0.1$	$\lambda = 1$	$\lambda = 10$	$\lambda = 100$
1	1663.96298840	1663.96584647	1663.99439964	1664.27719642	1666.85083466
2	2193.80448491	2193.80893614	2193.85340296	2194.29355673	2198.27545329
3	1262.38796634	1262.39096840	1262.42095945	1262.71793643	1265.41489576
4	2249.75391616	2249.75616197	2249.77859782	2250.00074734	2252.01674105
5	1374.25164430	1374.25417190	1374.27942298	1374.52946509	1376.80025092
Ave	1748.83220002	1748.83521697	1748.86535657	1749.16378040	1751.87163514

Table 4: square root of RSS for RLS

7 Bonus:

- (a) Use linear regression with different feature combinations.

Answer:

The code for this part is as follows:

```

1
2  #NameMap = {0:['stamina',stamina],
3  #          1:['attack.value',attack.value],
4  #          2:['defense.value',defense.value],
5  #          3:['capture.rate',capture.rate],
6  #          4:['flee.rate',flee.rate],
7  #          5:['spawn.chance',spawn.chance],
8  #          6:['primary.strength',primary.strength],
9  #          7:['combat.point',combat.point]}
10 def CombinFeature(X_new_basis,learning_rate = 1e-7,max_iter = ...
    15000):
11     y = np.array(NameMap[7][1])
12     W_ini = ...
        np.linalg.inv(np.transpose(X_new_basis).dot(X_new_basis))...
        dot(np.transpose(X_new_basis).dot(y))
13     #####Cross Validation#####
14     -,rss_tol,rss_ave = cross_validation(X_new_basis, y, 0*W_ini,...
15     method = 0,learning_rate=learning_rate,max_iter=max_iter)
16     print('all rss:{},average rss:{}'.format(rss_tol,rss_ave))
17     return rss_tol, rss_ave
18
19
20 rss_tol = []
21 rss_ave_tol = []
22 combain.features = []
23 for i in range(6):

```

```

24     j = i + 1
25     while j < 6:
26         X_new = np.array([NameMap[i][1], NameMap[j][1]]).T
27         X_new_basis = np.hstack((np.ones([X_new.shape[0], 1]), X_new))
28         name_tep = (NameMap[i][0], NameMap[j][0])
29         print(name_tep)
30         combain_features.append(name_tep)
31
32         rss, rss_ave = CombinFeature(X_new_basis, learning_rate = 1e-9)
33         rss_tol.append(rss)
34         rss_ave_tol.append(rss_ave)
35         j = j + 1
36
37     for i in range(len(combain_features)):
38         print('Combin features:{}, average rss:{}, all ...
           rss:{}'.format(combain_features[i], rss_ave_tol[i], rss_tol[i]))

```

The results is shown as follow:

Combination feature	Average RSS
(<i>stamina</i> , <i>attack_value</i>)	1761.7612002966048
(<i>stamina</i> , <i>defense_value</i>)	2105.0676869515
(<i>stamina</i> , <i>capture_rate</i>)	3312.81517445934
(<i>stamina</i> , <i>flee_rate</i>)	3312.8048894631384
(<i>stamina</i> , <i>spawn_chance</i>)	3305.896347648588
(<i>attack_value</i> , <i>defense_value</i>)	1829.5275727386877
(<i>attack_value</i> , <i>capture_rate</i>)	1865.0731029770945
(<i>attack_value</i> , <i>flee_rate</i>)	1865.1272504360481
(<i>attack_value</i> , <i>spawn_chance</i>)	1888.4128783248634
(<i>defense_value</i> , <i>capture_rate</i>)	2237.345613547558
(<i>defense_value</i> , <i>flee_rate</i>)	2237.386858013856
(<i>defense_value</i> , <i>spawn_chance</i>)	2336.7613282436782
(<i>capture_rate</i> , <i>flee_rate</i>)	2255.356812760387
(<i>capture_rate</i> , <i>spawn_chance</i>)	9394.510236534275
(<i>flee_rate</i> , <i>spawn_chance</i>)	9395.262608277268

Table 5: RSS with different combination features

In problem 2 and problem 3, the following observation can be get:

- The most related to output's features are *attack_value*, *defense_value*, *capture_rate*, *stamina* whose Pearson's correlation coefficient's absolute value larger than 0.5
- The following pairs likely have linear relation (*attack_value*, *defense_value*), (*attack_value*, *capture_rate*), (*defense_value*, *capture_rate*) whose corresponding value also larger than 0.5

From the above table 5, similar results can be explained by the results from problem 2 and problem 3. For example, combination (*stamina*, *attack_value*) has the smallest average RSS. This is corresponding to the following two things:

- *attack_value* and *stamina* value show strong relation (≈ 0.5) with output value caused by problem 2's results.
- *attack_value* and *stamina* value show no relation (≈ 0.5) with each other caused by problem 3's results.

8 Logistic Regression:

- Use the mean sample value to binarize the data.
- Split the data into training and test.
- Report the accuracy of the classifier on the test data.

Answer:

8.1 Binary the data

In this part, first *y_mean* is used for storing mean value of output. Then label *y_lab* = +1 if value *y* \geq *y_mean* and label *y_lab* = -1 if value *y* < *y_mean*. The corresponding code is as follow:

```
1 y = np.array(NameMap[7][1])
2 y_lab = np.zeros(y.shape, dtype = int)
3 y_mean = y.mean()
4
5 y_lab[np.where(y >= y_mean)] = 1
6 y_lab[np.where(y < y_mean)] = -1
7
8 print(y_mean)
9 print(y)
10 print(y_lab)
```

8.2 Split data

In this part, package *sklearn.model_selection* is needed to import function *train_test_split*. Here, set the number of test set equals 0.15 multiply the number of total data set. The code is as follow:

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X_basis, ...
    y_lab, test_size=0.15, random_state=0)
```

8.3 Accuracy

In this part, build-in function *LogisticRegression* is used. The code is as follow:

```
1 from sklearn.linear_model import LogisticRegression
2
3 clf = LogisticRegression(penalty = ...
    'none', solver='saga', max_iter=15000)
4 clf.fit(X_train, y_train)
5 print(clf.predict(X_test))
6 print(y_test)
7 print('train acc:{}, test ...
    acc:{}'.format(clf.score(X_train, y_train), clf.score(X_test, ...
    y_test)))
```

The result is as follow:

```
train acc: 0.8306451612903226
test acc: 0.9545454545454546
```

9 Regularization logistic regression:

- (a) Find optimal λ^* by 5-fold cross-validation on training data set.
- (b) Use λ^* in (a) and calculate accuracy on testing data set.

Answer:

Since there exist two kinds of norm, l-2 norm and l-1 norm, solving method '*saga*' is used to handle this problem.

9.1 Find λ^*

The codes for l-2 norm and l-1 norm are really similar. So in this part, only the results of l-1 norm will be shown. Take l-2 norm as the example. The code is as follow:

```
1 #####l2norm#####
2 labdtotal = np.linspace(0.0001,0.1,12)
3 #labdtotal = np.linspace(1000,100000,100)
4 rss_tol = []
5 rss_ave = []
6 for labd in labdtotal:
7     _,rss_tol_tep,rss_ave_tep = cross_validation(X_train,y_train, ...
    Wini=None,method = 2,labd = labd,max_iter=10000)
8     print('lambda:{}'.format(labd))
9     print('all rss:{},average rss:{}'.format(rss_tol_tep,rss_ave_tep))
10
11     rss_tol.append(rss_tol_tep)
12     rss_ave.append(rss_ave_tep)
```

```

13
14 for i in range(len(labd_total)):
15     print('l2 penalty, lambda:{}, average rss:{}, all ...
        rss:{}'.format(labd_total[i], rss_ave[i], rss_tol[i]))
16
17 plt.plot(labd_total, rss_ave, '+-')
18 plt.xlabel('lambda')
19 plt.ylabel('Average Accuracy')
20 plt.title('l-2 norm')

```

Finally, the following results can be shown:

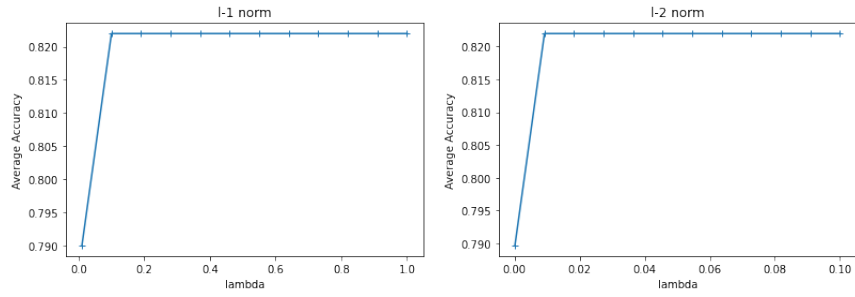


Figure 3: Find λ^*

From the above figure, choose $\lambda^* = 0.1$ for l-1 norm and $\lambda^* = 0.01$ for l-2 norm.

9.2 Calculate accuracy

The code for this part is as follow:

```

1 labd = 0.001
2 clf = LogisticRegression(penalty = ...
    'l2', C=labd, solver='saga', max_iter=15000)
3 clf.fit(X_train, y_train)
4 print(clf.predict(X_test))
5 print(y_test)
6 print('train acc:{}, test ...
    acc:{}'.format(clf.score(X_train, y_train), clf.score(X_test, ...
    y_test)))

```

The final results are shown in the following table:

norm	λ^*	Train Acc	Test Acc
l-1 norm	0.1	0.8306451612903226	0.9545454545454546
l-2 norm	0.01	0.8387096774193549	0.9545454545454546

Table 6: Results

Appendix

(i) Code for upload data

```
1 import numpy as np
2
3 class DataPoint(object):
4
5     def __init__(self, feats):
6         """
7         Return a Datapoint object whose attributes are given by "feats"
8         """
9         self.stamina = feats['stamina']
10        self.attack_value = feats['attack.value']
11        self.defense_value = feats['defense.value']
12        self.capture_rate = feats['capture.rate']
13        self.flee_rate = feats['flee.rate']
14        self.spawn_chance = feats['spawn.chance']
15        self.primary_strength = feats['primary.strength']
16        self.combat_point = feats['combat.point']
17
18    def feature_vector(self):
19        """
20        Return feature vector as a numpy array
21        """
22        return np.array([self.stamina, self.attack_value,
23                          self.defense_value, self.capture_rate,
24                          self.flee_rate, self.spawn_chance,
25                          self.primary_strength])
26
27
28    def full_vector(self):
29        """
30        Return feature vector as a numpy array
31        """
32        return np.array([self.stamina, self.attack_value,
33                          self.defense_value, self.capture_rate,
34                          self.flee_rate, self.spawn_chance,
35                          self.primary_strength, self.combat_point])
36
37
38 def parse_dataset(filename):
39     data_file = open(filename, 'r') # Open File "to read"
40     dataset = [] # List to hold Datapoint objects
41
42     for index, line in enumerate(data_file):
43         if index == 0: # First line: "Height Weight Label" ...
44             # describes the datapoint, it's not an actual datapoint,
45             continue # do nothing, it will skip all the following lines
46         #print(line)
47         _, stamina, attack_value, defense_value, ...
48         capture_rate, flee_rate, spawn_chance, primary_strength, ...
49         combat_point = ...
50         line.strip().split(',') # line.strip().split('\t') # ...
51         strip() removes '\n', and split('\t') splits the line at ...
52         tabs
```

```

47     dataset.append(DataPoint({'stamina':float(stamina),...
48                               'attack.value':float(attack.value),...
49                               'defense.value':float(defense.value),...
50                               'capture.rate':float(capture.rate),...
51                               'flee.rate':float(flee.rate),...
52                               'spawn.chance':float(spawn.chance),...
53                               'primary.strength':primary.strength,...
54                               'combat.point':float(combat.point)}))
55     # Create DataPoint object for the given data
56
57     print("Total Number of Data Points: {0}".format(len(dataset)))
58     print("Examples: 1 - {0} , 2 - {1}".format(dataset[0],...
59                                                dataset[-1]))
60     return dataset
61
62 datasetTrain = parse_dataset('/content/drive/My ...
63                               Drive/2020_FALL/CSCE633/HomeWork/HW02/hw2_data.csv')
64 # Returns the list of datapoints
65 Lg = len(datasetTrain)
66 stamina      =[datasetTrain[i].stamina for i in range(Lg)]
67 attack.value =[datasetTrain[i].attack.value for i in range(Lg)]
68 defense.value =[datasetTrain[i].defense.value for i in range(Lg)]
69 capture.rate  =[datasetTrain[i].capture.rate for i in range(Lg)]
70 flee.rate     =[datasetTrain[i].flee.rate for i in range(Lg)]
71 spawn.chance  =[datasetTrain[i].spawn.chance for i in range(Lg)]
72 primary.strength=[datasetTrain[i].primary.strength for i in ...
73                   range(Lg)]
74 combat.point  =[datasetTrain[i].combat.point for i in range(Lg)]
75 NameMap = {0:['stamina',stamina],
76            1:['attack.value',attack.value],
77            2:['defense.value',defense.value],
78            3:['capture.rate',capture.rate],
79            4:['flee.rate',flee.rate],
80            5:['spawn.chance',spawn.chance],
81            6:['primary.strength',primary.strength],
82            7:['combat.point',combat.point]}
83 print (NameMap[0][1])

```

(ii) Code for l-1 norm

```

1  #####l1norm#####
2  labdtotal = np.linspace(0.01,1,12)
3  #labdtotal = np.linspace(1000,100000,100)
4  rss_tol = []
5  rss_ave = []
6  for labd in labdtotal:
7      _,rss_tol_tep,rss_ave_tep = cross_validation(X_train,y_train, ...
8          W_ini=None,method = 3,labd = labd,max_iter=10000)
9      print('lambda:{0}'.format(labd))
10     print('all rss:{0},average rss:{0}'.format(rss_tol_tep,rss_ave_tep))
11
12     rss_tol.append(rss_tol_tep)
13     rss_ave.append(rss_ave_tep)
14
15 for i in range(len(labdtotal)):
16     print('12 penalty, lambd:{0}, average rss:{0}, all ...

```

```

    rss:{}'.format(labd_total[i], rss_ave[i], rss_tol[i]))
16
17 plt.plot(labd_total, rss_ave, '+-')
18 plt.xlabel('lambda')
19 plt.ylabel('Average Accuracy')
20 plt.title('l-1 norm')
21 #plt.ylim(0.82, 0.825)
22 #plt.xlim(0.09, 0.101)
23
24
25 labd = 0.1
26 clf = LogisticRegression(penalty = ...
    'l1', C=labd, solver='saga', max_iter=15000)
27 clf.fit(X_train, y_train)
28 print(clf.predict(X_test))
29 print(y_test)
30 print('train acc:{}, test ...
    acc:{}'.format(clf.score(X_train, y_train), clf.score(X_test, ...
    y_test)))

```