



دانشگاه صنعتی خواجه نصیرالدین طوسی
دانشکده مهندسی برق

درس یادگیری ماشین پاسخ مینی پروژه ۱

نام و نام خانوادگی	مهدی خدابنده‌لو
شماره دانشجویی	۴۰۱۰۵۳۱۴
تاریخ	بهار ۱۴۰۳



فهرست مطالب

- ۱ سوال اول ۵
- ۱.۱ فرآیند آموزش و ارزیابی یک مدل طبقه بند خطی را به صورت دیاگرامی بلوکی نمایش دهید و در مورد اجزای مختلف این دیاگرام بلوکی توضیحاتی بنویسید. تغییر نوع طبقه بندی از حالت دوکلاسه به چندکلاسه در کدام قسمت از این دیاگرام بلوکی تغییراتی ایجاد می کند؟ توضیح دهید. ۵
- ۲.۱ با استفاده از یک `sklearn.datasets`، دیتاست با ۱۰۰۰ نمونه، ۴ کلاس و ۳ ویژگی تولید کنید و آن را به صورتی مناسب نمایش دهید. آیا دیتاستی که تولید کردید چالش برانگیز است؟ چرا؟ به چه طریقی می توانید دیتاست تولیدشده خود را چالش برانگیزتر و سخت تر کنید؟ ۵
- ۳.۱ با استفاده از حداقل دو طبقه بند خطی آماده ی پایتون و در نظر گرفتن فرآیندهای مناسب چهار کلاس موجود در دیتاست قبلی را از هم تفکیک کنید. ضمن روند توضیح فرآیندها نتیجه دقت آموزش و ارزیابی را نشان دهید. برای بهبود نتایج از چه روش هایی استفاده کردید؟ ۶
- ۲ سوال دوم ۸
- ۱.۲ با مراجعه به صفحه دیتاست CWRU با Bearing یک دیتاست مربوط به حوزه تشخیص عیب آشنا شوید. با جستجوی آن در اینترنت و مقالات، توضیحاتی از اهداف، ویژگی ها و حالت های مختلف این دیتاست ارائه کنید. ۸
- ۲.۲ برای تشکیل دیتاست مراحل زیر را انجام دهید: ۸
- ۳.۲ بدون استفاده از کتابخانه های آماده پایتون، مدل طبقه بند، تابع اتلاف و الگوریتم یادگیری و ارزیابی را کدنویسی کنید تا دو کلاس موجود در دیتاست به خوبی از یکدیگر تفکیک شوند. نمودار تابع اتلاف را رسم کنید و نتیجه ارزیابی روی داده های تست را با حداقل ۲ شاخصه محاسبه کنید. نمودار تابع اتلاف را تحلیل کنید. آیا می توان از روی نمودار تابع اتلاف و قبل از مرحله ارزیابی با قطعیت در مورد عمل کرد مدل نظر داد؟ چرا و اگر نمی توان، راه حل چیست؟ ۱۲
- ۳ سوال سوم ۱۴
- ۱.۳ ابتدا هیت مپ ماتریس همبستگی و هیستوگرام پراکندگی ویژگی ها را رسم و تحلیل کنید. ۱۴



فهرست تصاویر

۶ class_sep=۱ داده های تولید شده با	۱
۷ class_sep=۰.۵ داده های تولید شده با	۲
۱۴ خروجی تابع Loss در هر مرحله	۳
۱۵ قدر مطلق وزن ها	۴
۱۶ هیت مپ	۵
۱۷ هیستوگرام	۶



فهرست جداول



فهرست برنامه‌ها

۵ (Python) data Generate	۱
۶ (Python) split test and Train	۲
۷ accuracy(Python) Calculate	۳
۸ (Python)	۴
۹ (Python)	۵
۱۰ (Python)	۶
۱۰ (Python)	۷
۱۱ (Python)	۸
۱۲ (Python)	۹
۱۳ (Python)	۱۰
۱۳ (Python)	۱۱
۱۴ (Python)	۱۲



Notebook Colab
github

۱ سوال اول

۱.۱ فرآیند آموزش و ارزیابی یک مدل طبقه بند خطی را به صورت دیاگرامی بلوکی نمایش دهید و در مورد اجزای مختلف این دیاگرام بلوکی توضیحاتی بنویسید. تغییر نوع طبقه بندی از حالت دوکلاسه به چندکلاسه در کدام قسمت از این دیاگرام بلوکی تغییراتی ایجاد می کند؟ توضیح دهید.

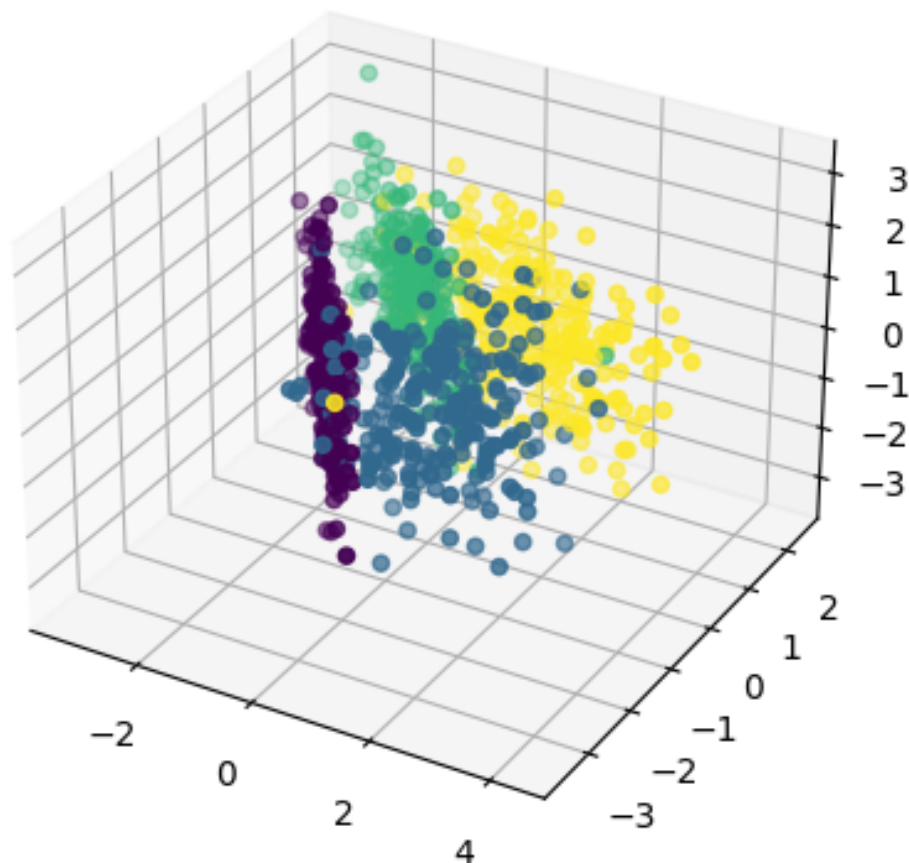
۲.۱ با استفاده از یک `sklearn.datasets`، دیتاست با ۱۰۰۰ نمونه، ۴ کلاس و ۳ ویژگی تولید کنید و آن را به صورتی مناسب نمایش دهید. آیا دیتاستی که تولید کردید چالش برانگیز است؟ چرا؟ به چه طریقی می توانید دیتاست تولیدشده خود را چالش برانگیزتر و سخت تر کنید؟

همان طور که در Code ۱ قابل مشاهده است با استفاده از دستور `make_classification` نقطه را تولید کرده ایم. و سپس آن ها را در یک نمودار سه بعدی نمایش داده ایم. شکل ۱ این نقاط را نمایش می دهد.

```
1 X,y = make_classification(n_samples=1000, n_features=3, n_classes=4,
2     n_redundant=0, random_state=14,
3     n_clusters_per_class=1, class_sep=1)
4
5 fig = plt.figure()
6 ax = fig.add_subplot(projection='3d')
7 ax.scatter(X[:,0], X[:,1], X[:,2], c=y)
```

Code 1: Generate data (Python)

همان طور که از شکل ۱ پیداست دیتاستی که تولید شده زیاد چالش برانگیز نمی باشد و می توان داده ها را از هم تفکیک کرد. برای این که داده های تولید شده را چالش بر انگیز تر کنیم می توانیم کارهای مختلفی انجام دهیم. یکی از این کارها این است که فرایارامتر `class_sep` را کاهش دهیم. شکل ۲ داده ها را در حالتی که این پارامتر 0.4 است نمایش می دهد. همان طور که مشخص است داده های کلاس های مختلف به هم نزدیک تر شده اند و تفکیک آن ها مشکل تر شده است.

شکل ۱: داده های تولید شده با `class_sep=1`

۳.۱ با استفاده از حداقل دو طبقه‌بند خطی آماده‌ی پایتون و در نظر گرفتن فرآپارامترهای مناسب چهار کلاس موجود در دیتاست قبلی را از هم تفکیک کنید. ضمن روند توضیح فرآپارامترها نتیجه دقت آموزش و ارزیابی را نشان دهید. برای بهبود نتایج از چه روش‌هایی استفاده کردید؟

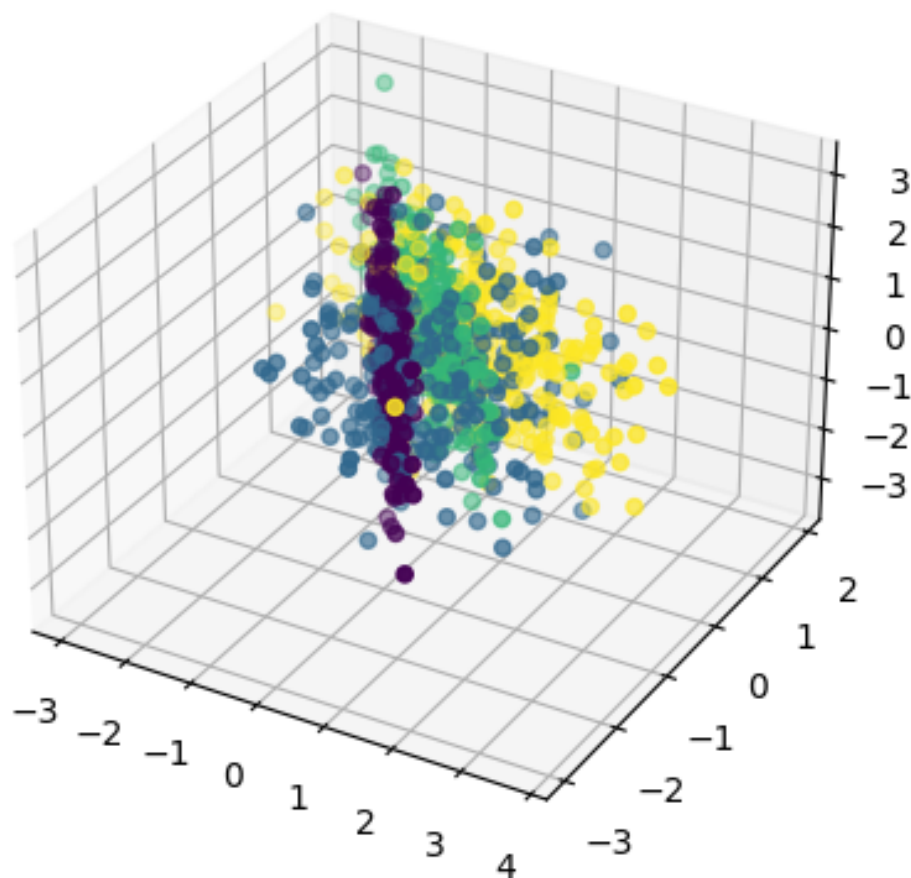
جدا کردن داده‌های آموزش و تست:

در Code ۲ داده‌های آموزش و تست را به نسبت ۸۰ به ۲۰ تقسیم کرده‌ایم.

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,  
2 random_state=14)
```

Code 2: Train and test split (Python)

Logistic_regression



شکل ۲: داده های تولید شده با $class_sep=0.5$

در Code ۳ ابتدا تعداد تکرار برای آموزش را ۵۰ در نظر گرفته ایم و پس از آموزش، با استفاده از داده های تست دقت مدل را اندازه گیری کرده ایم. با مقادیری که برای هر یک از واریاترها در نظر گرفته شده دقت برابر با ۸۹ درصد است.

```
1 logistic_regression_model = LogisticRegression(max_iter=50, random_state=14)
2 logistic_regression_model.fit(X=X_train, y=y_train)
3
4 y_hat = logistic_regression_model.predict(X_test)
5 y_hat_prob = logistic_regression_model.predict_proba(X_test)
6 accuracy = logistic_regression_model.score(X_test, y_test)
7 print('accuracy :', accuracy*100, '%')
```

Code 3: Calculate accuracy(Python)



برای بهبود نتایج باید تعداد تکرار برای آموزش را مناسب در نظر بگیریم. در صورتی که این پارامتر زیاد در نظر گرفته شده باشد باعث overfitting می‌شود، این پدیده باعث می‌شود که مدل روی داده‌های آموزش عملکرد خوبی داشته باشد ولی وقتی داده‌هایی غیر از داده‌های آموزش به مدل داده شود عملکرد آن خراب می‌شود. این پدیده مانند این است که یک دانش آموز به جای یاد گرفتن درس آن را حفظ کرده باشد. مسئله‌ی دیگری که در بهبود نتایج موثر است تعداد کافی داده‌های آموزش است، در صورتی که تعداد داده‌های آموزش کم باشد مدل به خوبی آموزش نمی‌بیند.

۲ سوال دوم

۱.۲ با مراجعه به صفحه دیتاست CWRU با Bearing یک دیتاست مربوط به حوزه تشخیص عیب آشنا شوید. با جستجوی آن در اینترنت و مقالات، توضیحاتی از اهداف، ویژگی‌ها و حالت‌های مختلف این دیتاست ارائه کنید.

برای تشخیص عیب از دو ویژگی drive end accelerometer data (که در کد x1 نامیده شده) و fan end accelerometer data (که در کد x2 نامیده شده) استفاده شده است.

۲.۲ برای تشکیل دیتاست مراحل زیر را انجام دهید:

الف: از هر کلاس M نمونه به طول N جدا کنید (M حداقل ۱۰۰ و N حداقل ۲۰۰ باشد). یک ماتریس از داده‌های هر دو کلاس به همراه برچسب آن‌ها تشکیل دهید. می‌توانید پنجره‌ای به طول N در نظر بگیرید و در نهایت یک ماتریس $M \times N$ از داده‌های هر کلاس استخراج کنید.

ما در این جا M را برابر با ۵۰۰ و N را ۲۰۰ در نظر گرفته‌ایم. همان طور که از Code ۴ مشخص است هر یک از ویژگی‌های کلاس‌های معیوب و سالم را با استفاده از دستور reshape به صورت مربعی درآورده ایم.

```
1 M = 500
2 N = 200
3 x1_normal = np.reshape(bearing_normal_dataset['X098_DE_time'][0:M*N], (M,N))
4 x2_normal = np.reshape(bearing_normal_dataset['X098_FE_time'][0:M*N], (M,N))
5 x1_fault = np.reshape(bearing_fault_dataset['X108_DE_time'][0:M*N], (M,N))
6 x2_fault = np.reshape(bearing_fault_dataset['X108_FE_time'][0:M*N], (M,N))
```

Code 4: (Python)

ب: در مورد اهمیت استخراج ویژگی در یادگیری ماشین توضیحاتی بنویسید. سپس، با استفاده از حداقل هشت عدد از روش‌های ذکر شده در جدول ۱، ویژگی‌های دیتاست قسمت «۲-آ» را استخراج کنید و یک دیتاست جدید تشکیل دهید.

Code ۵ تابعی که به منظور استخراج ویژگی نوشته شده را نشان می‌دهد. ورودی این تابع ویژگی‌های دیتاست خام می‌باشد و خروجی آن یک dictionary است که شامل ویژگی‌های استخراج شده از ویژگی‌های خام و خود ویژگی‌های خام است.



این تابع تمامی ویژگی‌های ذکر شده در جدول را محاسبه می‌کند. نکته‌ی مهمی که در این تابع وجود دارد این است که به منظور جلوگیری از تقسیم شدن بر صفر، زمانی که مقسوم علیه صفر شود، به جای صفر مقدار اپسیلون قرار می‌گیرد.

```
1 def feature_extract(x1,x2,ep):
2     x1_n = np.expand_dims(x1, axis=2)
3     x2_n = np.expand_dims(x2, axis=2)
4     xc = np.concatenate((x1_n,x2_n), axis=2)
5     xc_mean = np.mean(xc, axis=2)
6     m = np.expand_dims(xc_mean, axis=2)
7     m = np.concatenate((m,m), axis=2)
8
9     feature = {'x1': x1, 'x2': x2}
10    feature['std'] = np.std(xc, axis=2)
11    feature['p'] = np.amax(np.abs(xc), axis=2)
12    temp = np.std(xc, axis=2)**3
13    temp[temp == 0] = ep
14    feature['ske'] = np.mean((xc-m)**3, axis=2) / temp
15    temp = np.std(xc, axis=2)**4
16    temp[temp == 0] = ep
17    feature['kur'] = np.mean((xc-m)**4, axis=2) / temp
18    feature['rms'] = np.sqrt(np.mean(xc**2, axis=2))
19    temp = np.sqrt(np.mean(xc**2, axis=2))
20    temp[temp == 0] = ep
21    feature['cf'] = np.amax(np.abs(xc), axis=2) / temp
22    feature['smr'] = np.mean(np.sqrt(np.abs(xc)), axis=2) ** 2
23    temp = np.mean(np.sqrt(np.abs(xc)), axis=2) ** 2
24    temp[temp == 0] = ep
25    feature['clf'] = np.amax(np.abs(xc), axis=2) / temp
26    temp = np.mean(np.abs(xc), axis=2)
27    temp[temp == 0] = ep
28    feature['sf'] = np.sqrt(np.mean(xc**2, axis=2)) / temp
29    temp = np.mean(np.abs(xc), axis=2)
30    temp[temp == 0] = ep
31    feature['if1'] = np.amax(np.abs(xc), axis=2) / temp
32    feature['if2'] = np.max(np.abs(xc), axis=2) / temp
33    return feature
```

Code 5: (Python)



همان طور که از Code ۶ مشخص است مقدار اپسیلون را 10^{-10} در نظر گرفته‌ایم. پس از این مرحله داده‌های کلاس normal و fault را با هم ترکیب می‌کنیم. البته نیاز است که قبل از استفاده از داده‌ها آن‌ها را به خوبی بر بزنیم.

```

1 x_normal = feature_extract(x1_normal, x2_normal, 1e-10)
2 x_fault = feature_extract(x1_fault, x2_fault, 1e-10)
3
4 X = {'x1': np.concatenate((x_normal['x1'], x_fault['x1']), axis=0),
5      'x2': np.concatenate((x_normal['x2'], x_fault['x2']), axis=0),
6      'std': np.concatenate((x_normal['std'], x_fault['std']), axis=0),
7      'p': np.concatenate((x_normal['p'], x_fault['p']), axis=0),
8      'ske': np.concatenate((x_normal['ske'], x_fault['ske']), axis=0),
9      'kur': np.concatenate((x_normal['kur'], x_fault['kur']), axis=0),
10     'rms': np.concatenate((x_normal['rms'], x_fault['rms']), axis=0),
11     'cf': np.concatenate((x_normal['cf'], x_fault['cf']), axis=0),
12     'smr': np.concatenate((x_normal['smr'], x_fault['smr']), axis=0),
13     'clf': np.concatenate((x_normal['clf'], x_fault['clf']), axis=0),
14     'sf': np.concatenate((x_normal['sf'], x_fault['sf']), axis=0),
15     'if1': np.concatenate((x_normal['if1'], x_fault['if1']), axis=0),
16     'if2': np.concatenate((x_normal['if2'], x_fault['if2']), axis=0)}
17 y0 = np.zeros((M,N))#normal
18 y1 = np.ones((M,N))#fault
19 y = np.concatenate((y0,y1), axis=0)

```

Code 6: (Python)

ج: ضمن توضیح اهمیت فرآیند بر زدن، داده‌ها را در صورت امکان مخلوط کرده و با نسبت تقسیم دلخواه و معقول به دو بخش آموزش و ارزیابی تقسیم کنید.

در صورتی که داده‌های کلاس‌های مختلف به خوبی مخلوط نشوند فرایند یادگیری به درستی انجام نمی‌شود. در Code ۷ ابتدا داده‌ها را به شکل بردار درآورده‌ایم و آن‌ها را به صورت یکسان مخلوط کرده‌ایم. پس از مخلوط کردن داده‌ها را به نسبت ۸۰ به ۲۰ به داده‌های آموزش و تست تقسیم کرده‌ایم.

```

1 data_shape = y.shape
2 for key, value in X.items():
3     X[key] = X[key].flatten()
4 y = y.flatten()
5 np.random.seed(14)

```



```

6 perm = np.random.permutation(len(y))
7 for key,value in X.items():
8     X[key] = X[key][perm]
9 y = y[perm]
10 for key,value in X.items():
11     X[key] = X[key].reshape(data_shape)
12 y = y.reshape(data_shape)
13 #train and test split
14 train_size = 0.8
15 train_a = round(0.8*M)
16 X_train = {}
17 X_test = {}
18 for key, value in X.items():
19     X_train[key] = X[key][0:train_a,:]
20     X_test[key] = X[key][train_a:,:]
21 y_train = y[0:train_a,:]
22 y_test = y[train_a:,:]

```

Code 7: (Python)

د: حداقل دو روش برای نرمال سازی داده ها را با ذکر اهمیت این فرآیند توضیح دهید و با استفاده از یکی از این روش ها، داده ها را نرمال کنید. آیا از اطلاعات بخش «ارزیابی» در فرآیند نرمال سازی استفاده کردید؟ چرا؟

از آن جایی که رنج تغییرات ویژگی های مختلف متفاوت است ممکن است برخی از آنها که دامنه ی بزرگتری دارند اثر ویژگی هایی که دامنه ی کمتری دارند را خنثی کنند، از این رو بهتر است برای جلوگیری از این پدیده داده ها را نرمالیزه کنیم. Code ۸ مربوط به نرمالیزه کردن داده های تست و آموزش به صورت جداگانه است.

```

1 for key, value in X_train.items():
2     X_train[key] = (X_train[key] - np.min(X_train[key])) / ...
3         (np.max(X_train[key]) - np.min(X_train[key]))
4 for key, value in X_test.items():
5     X_test[key] = (X_test[key] - np.min(X_test[key])) / ...
6         (np.max(X_test[key]) - np.min(X_test[key]))

```

Code 8: (Python)



۳.۲ بدون استفاده از کتابخانه های آماده پایتون، مدل طبقه بند، تابع اتلاف و الگوریتم یادگیری و ارزیابی را کدنویسی کنید تا دو کلاس موجود در دیتاست به خوبی از یکدیگر تفکیک شوند. نمودار تابع اتلاف را رسم کنید و نتیجه ارزیابی رویداده های تست را با حداقل ۲ شاخصه محاسبه کنید. نمودار تابع اتلاف را تحلیل کنید. آیا می توان از روی نمودار تابع اتلاف و قبل از مرحله ارزیابی با قطعیت در مورد عمل کرد مدل نظر داد؟ چرا و اگر نمی توان، راه حل چیست؟

همان طور که در Code ۹ قابل مشاهده است به ترتیب توابع سیگموید، Logistic_regression، LOSS، گرادیان، گرادیان نزولی و دقت پیاده سازی شده اند.

```
1 def sigmoid(x):
2     return 1 / (1 + np.exp(-x))
3
4 def logistic_regression(x,w):
5     u = 0
6     for key,value in x.items():
7         u += x[key] * w[key]
8     h_hat = sigmoid(u)
9     return h_hat
10
11 def bce(y, y_hat):
12     ep = 1e-10
13     loss = -(np.mean(y*np.log(y_hat + ep) + (1-y)*np.log(1-y_hat + ep)))
14     return loss
15
16 def gradient(x, y, y_hat):
17     grd = {}
18     for key, value in x.items():
19         grd[key] = np.sum(np.multiply(x[key], (y_hat-y)))
20     return grd
21
22 def gradient_decent(w, eta, grads):
23     for key, value in w.items():
24         w[key] -= eta * grads[key]
25     return w
26
27 def accuracy(y, y_hat):
```



```
28 acc = np.sum(y == np.round(y_hat)) / y.size
29 return acc
```

Code 9: (Python)

در Code ۱۰ مقادیر اولیه‌ی وزن‌ها، تعداد تکرار برای آموزش و گام آموزش تعیین شده است.

```
1 X_train['bias'] = np.ones(X_train['x1'].shape)
2 np.random.seed(14)
3 w_initial = np.random.rand(len(X)+1,1)
4 w = {'bias': w_initial[0,0], 'x1': w_initial[1,0], 'x2': w_initial[2,0],
5      'std': w_initial[3,0], 'p': w_initial[4,0], 'ske': w_initial[5,0],
6      'kur': w_initial[6,0], 'rms': w_initial[7,0], 'cf': w_initial[8,0],
7      'smr': w_initial[9,0], 'clf': w_initial[10,0], 'sf': w_initial[11,0],
8      'if1': w_initial[12,0], 'if2': w_initial[13,0]}
9 initial_eta = 0.00005
10 eta = initial_eta
11 n_epochs = 1000
```

Code 10: (Python)

پس از مراحل بالا نوبت به آموزش مدل می‌رسد. Code ۱۰ مربوط به آموزش مدل است که در هر اجرا وزن‌ها در آن آپدیت می‌شوند.

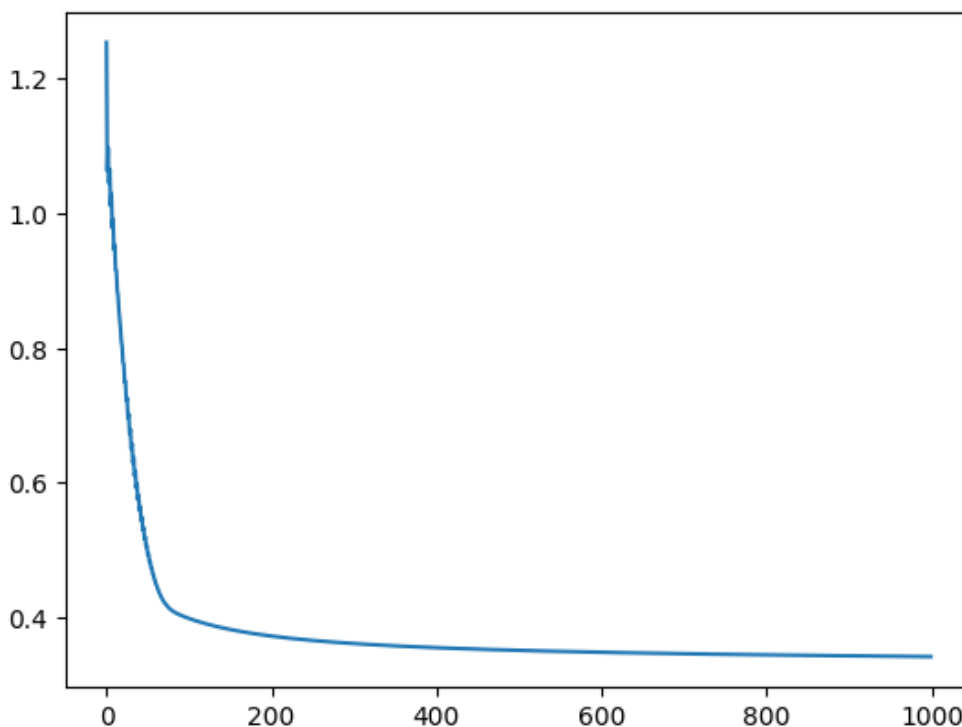
```
1 error_hist = []
2 for epoch in range(n_epochs):
3     # predictions
4     y_hat = logistic_regression(X_train, w)
5
6     # loss
7     e = bce(y_train, y_hat)
8     error_hist.append(e)
9
10    # gradients
11    grads = gradient(X_train, y_train, y_hat)
12
13    # gradient descent
```



```
14 w = gradient_decent(w, eta, grads)
15 acc = accuracy(y_train, y_hat)
16 print(f'Epoch={epoch}, \t E={e}, \t accuracy={acc}')
```

Code 11: (Python)

شکل ۳ خروجی تابع Loss را در هر اجرا نمایش می‌دهد.



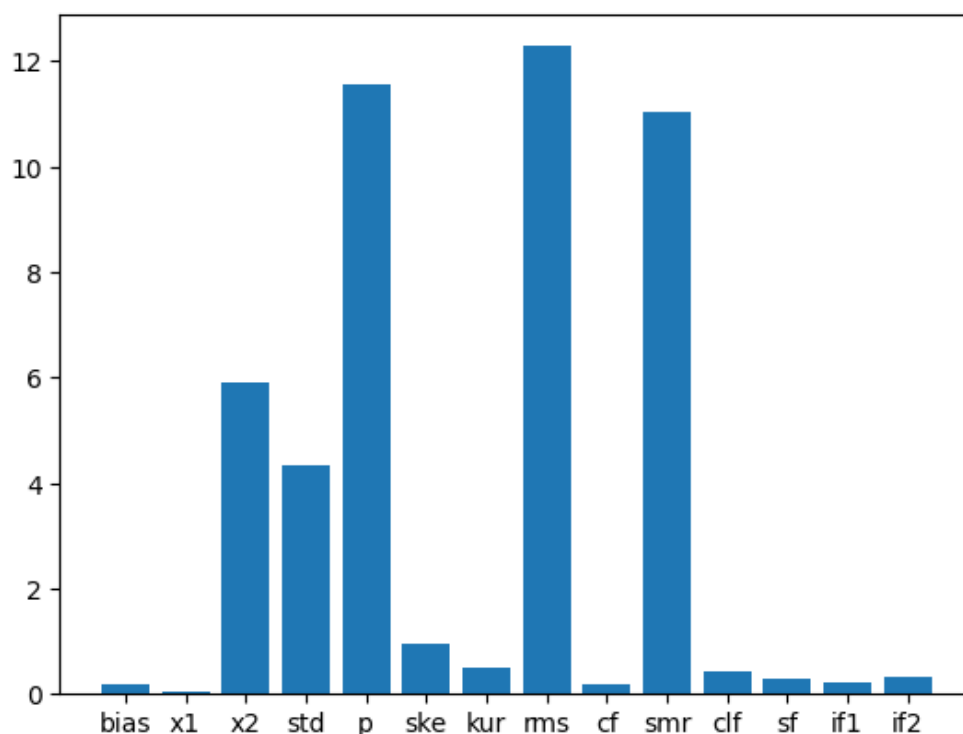
شکل ۳: خروجی تابع Loss در هر مرحله

شکل ۴ قدر مطلق وزن‌ها را نمایش می‌دهد. دقت خروجی مدل ۸۴ درصد می‌باشد.

۳ سوال سوم

۱.۳ ابتدا هیت مپ ماتریس همبستگی و هیستوگرام پراکندگی ویژگی‌ها را رسم و تحلیل کنید

همان طور که در Code ۱۲ مشاهده می‌شود پس از دانلود دیتاست سطرهایی که حاوی Null هستند حذف می‌شوند و سپس اطلاعات عددی از دیتاست استخراج می‌شود.

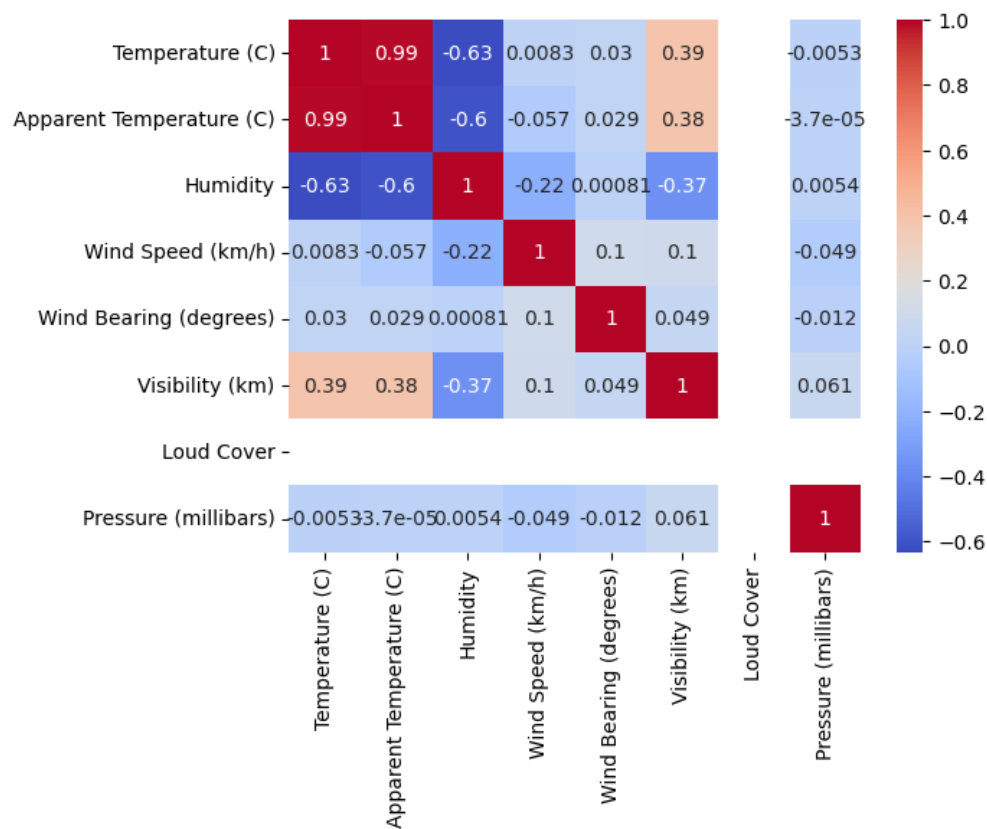


شکل ۴: قدر مطلق وزن ها

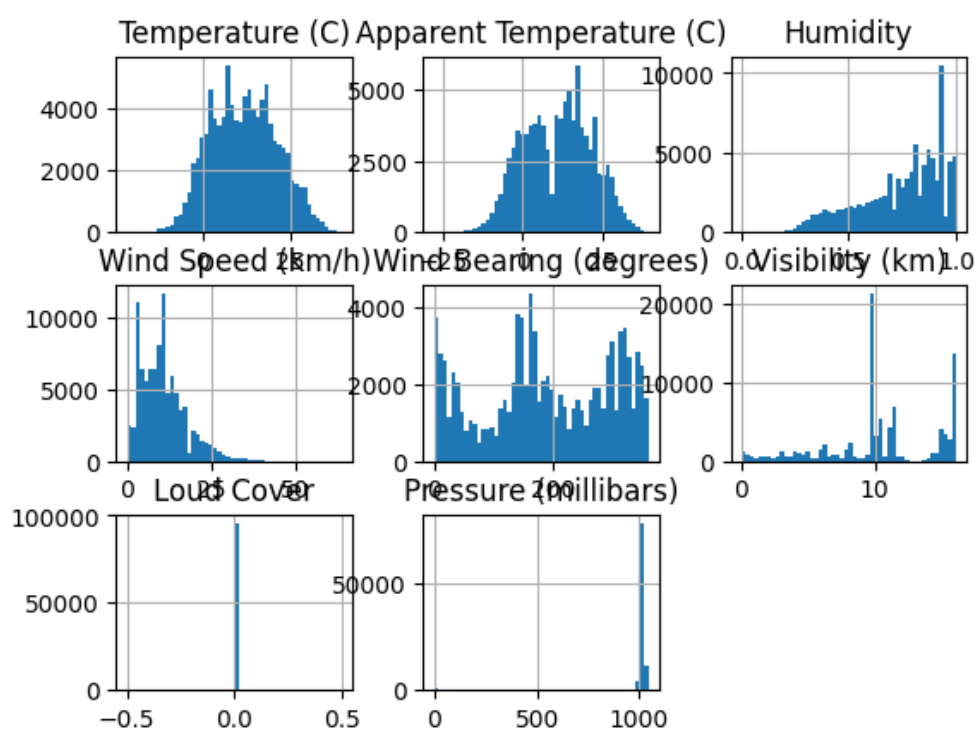
```
2 weatherHistory = pd.read_csv('weatherHistory.csv')
3
4 weatherHistory = weatherHistory.dropna()
5 numeric_weather_df = weatherHistory.select_dtypes(include=['number'])
6 print(weatherHistory.keys())
```

Code 12: (Python)

شکل ۵ هیت مپ و شکل ۶ هیستوگرام را نمایش می دهد.



شکل ۵: هیت مپ



شکل ۶: هیستوگرام