

homework1文档

homework1文档

[代码结构](#)

[作业1](#)

[作业2](#)

[代码使用方法](#)

[实验结果](#)

[LDA结果](#)

[K折朴素贝叶斯结果](#)

[SVM结果](#)

[补充题](#)

[代码使用方法](#)

[实验结果](#)

[附注](#)

[实验环境](#)

代码结构

| `getData.py` 自己写的一些获取数据集的函数

| `logistic.py` 自己写的LogisticRegression类

| `part1.py` 作业2代码

| `part2.py` 补充题代码

作业1



上海交通大学

SHANGHAI JIAO TONG UNIVERSITY

理论题 1.

1) $y = e^{wx+b}$: 最小二乘法求解 w 与 b

有输入 $X = [\vec{x}_1, \vec{x}_2, \vec{x}_3, \vec{x}_4, \dots, \vec{x}_n]^T$ $X_{n \times c}$ $\vec{x}_i \in \mathbb{R}^{1 \times c}$

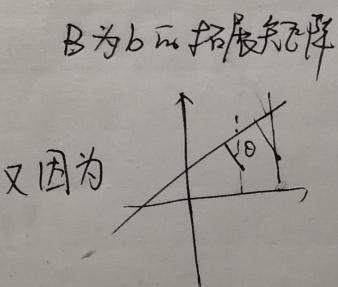
$\vec{Y}^* = [y_1, y_2, y_3, \dots, y_n]^T$ $\vec{Y}^*_{n \times 1}$ $\Rightarrow w^T_{1 \times c}$

~~$\vec{Y} = e^{w^T x + b}$~~ $\vec{Y} = e^{X \cdot w^T + b}$

$$\Rightarrow \ln Y = X \cdot w^T + B \quad \text{Target } \ln(Y^*) \quad B = \begin{bmatrix} b \\ b \\ b \\ b \end{bmatrix}_{n \times 1}$$

设 $P = \ln Y \quad P^* = \ln Y^*$

$P - Xw^T + B = 0$



距离与偏差只存在线性关系，简化为

$$\text{Error} = \arg \min \|P^* - Xw^T + B\|^2 = \arg \min (P^* - Xw^T + B)^T (P^* - Xw^T + B)$$

$$\frac{\partial E}{\partial w^T} = 2(X^T)(P^* - Xw^T + B) = 0 \quad P^* = \cancel{X}$$

$$\frac{\partial E}{\partial B} = 2(P^* - Xw^T + B)$$

把 \vec{x}_i 改为 (\vec{x}'_i) b 加入 w' , 改为 $\frac{\partial E}{\partial w'^T} = 2(-X'^T)(P^* - X'w'^T) = 0$

$$X'^T w'^T = P^* \Rightarrow w'^T = (X'^T X')^{-1} X'^T P^* \quad \text{其中 } w' \text{ 为 } [w_1, w_2, \dots, w_c, b]$$

$$\Rightarrow \cancel{X'^T} \cdot X'^T \cdot X'^T w'^T = X'^T P^* \quad X' = [\vec{x}'_1, \vec{x}'_2, \dots, \vec{x}'_n]^T$$

$$\Rightarrow W'^T = (X'^T X')^{-1} \cdot X'^T \cdot P^*$$

$$\vec{x}'_i = [x_{i1}, x_{i2}, x_{i3}, \dots, x_{iL}, 1] \\ P^* = \ln(Y^*)$$

地址：闵行东川路800号



上海交通大学

SHANGHAI JIAO TONG UNIVERSITY

2. A, B, C 工厂

	产品占比	次品率
A	$0.35 = P(A)$	$0.015 = P(\text{次} A)$
B	$0.35 = P(B)$	$0.01 = P(\text{次} B)$
C	$0.30 = P(C)$	$0.020 = P(\text{次} C)$

1) 随机抽一个为次品的概率

$$P(\text{次}) = P(A) \cdot P(\text{次}|A) + P(B) \cdot P(\text{次}|B) + P(C) \cdot P(\text{次}|C)$$

$$= 0.35 \cdot 0.015 + 0.35 \cdot 0.01 + 0.30 \cdot 0.02 = 0.01475$$

2) 抽一个次品来自

$$A: P(A|\text{次}) = \frac{P(A)P(\text{次}|A)}{P(\text{次})} = 0.3559$$

$$B: P(B|\text{次}) = 0.2373$$

$$C: P(C|\text{次}) = 0.4068$$

3. ε 的含义

- | | |
|---|--|
| $\left\{ \begin{array}{l} ① \varepsilon = 0 \\ ② 0 < \varepsilon \leq 1 \\ ③ \varepsilon > 1 \end{array} \right.$ | ① 样本可在最大边界上及外分类
② 样本可在最大间隔内分类
③ 样本会被错误分类 |
|---|--|

代码使用方法

1. python part1.py -h 可以获得使用帮助

```
(base) D:\大学\大三上\机器学习\homework1>python part2.py -h
usage: part2.py [-h] [--fold FOLD] action

positional arguments:
  action      watermelon or Iris, the program will run algorithm on this two
              dataset.

optional arguments:
  -h, --help   show this help message and exit
  --fold FOLD  the num of folds in K-floid validation, default is 5.

(base) D:\大学\大三上\机器学习\homework1>python part1.py -h
usage: part1.py [-h] [--rate RATE] [--fold FOLD] action

positional arguments:
  action      type LDA or K-NB or SVM, the program will run a LDA or a Naive
              Bayes or SVM algorithm.

optional arguments:
  -h, --help   show this help message and exit
  --rate RATE  the preportion of training set, default is 0.8.
  --fold FOLD  the num of folds in K-floid validation, default is 5.

(base) D:\大学\大三上\机器学习\homework1>
```

2. python part1.py LDA (--rate 0.8) 括号内容可选且值为默认

3. python part1.py K-NB (--fold 5) 括号内容可选且值为默认

4. python part1.py SVM(--rate 0.8) 括号内容可选且值为默认

实验结果

LDA结果

```
182
183 if __name__ == "__main__":
184     parse = argparse.ArgumentParser()
185     parse.add_argument("action", type=str, help="type LDA or K-NB or SVM, the program will run a LDA o
186     parse.add_argument("--rate", type=float, default=0.8, help="the proportion of training set, defaul
187     parse.add_argument("--fold", type=int, default=5, help="the num of folds in K-fold validation, def
188
189     watermalon = pd.read_csv("watermalon.csv",engine="python")
190     iris = pd.read_csv("iris.data",header=None)
191
192     args = parse.parse_args()
193
194     if args.action=="LDA":
195         LDA(watermalon)
196
197     if args.action=="K-NB":
198         KNB(watermalon)
199
200     if args.action=="SVM":
201         SVM_func(watermalon,'rbf')
202         SVM_func(watermalon,'linear')
203         SVM_func(watermalon,'poly')
204         SVM_func(watermalon,'sigmod')
205         SVM_func(watermalon,'precomputed')
```

问题 8 输出 调试控制台 终端

```
the accuracy of SVM using kernel sigmod scores      0.25
the accuracy of SVM using kernel precomputed scores  0.25
```

```
D:\大学\大三上\机器学习\homework1>python part1.py LDA
using LDA on watermelon3 dataset,training set rate =  0.8    test accuracy : 0.5
```

```
D:\大学\大三上\机器学习\homework1>python part1.py LDA --rate 0.7
using LDA on watermelon3 dataset,training set rate =  0.7    test accuracy : 0.75
```

```
D:\大学\大三上\机器学习\homework1>python part1.py LDA --rate 0.5
using LDA on watermelon3 dataset,training set rate =  0.5    test accuracy : 1.0
```

```
D:\大学\大三上\机器学习\homework1>python part1.py LDA --rate 0.9
using LDA on watermelon3 dataset,training set rate =  0.9    test accuracy : 1.0
```

```
D:\大学\大三上\机器学习\homework1>python part1.py LDA --rate 0.4
using LDA on watermelon3 dataset,training set rate =  0.4    test accuracy : 0.75
```

```
D:\大学\大三上\机器学习\homework1>
```

K折朴素贝叶斯结果

```
183 if __name__ == "__main__":
184     parse = argparse.ArgumentParser()
185     parse.add_argument("action", type=str, help="type LDA or K-NB or SVM, the program will run a LDA or a Naive Bayes or SVM algorithm.")
186     parse.add_argument("--rate", type=float, default=0.8, help="the proportion of training set, default is 0.8.")
187     parse.add_argument("--fold", type=int, default=5, help="the num of folds in K-fold validation, default is 5.")
188
189     watermalon = pd.read_csv("watermalon.csv", engine="python")
190     iris = pd.read_csv("iris.data", header=None)
191
192     args = parse.parse_args()
193
194     if args.action=="LDA":
195         LDA(watermalon)
196
197     if args.action=="K-NB":
198         KNB(watermalon)
199
200     if args.action=="SVM":
201         SVM_func(watermalon, 'rbf')
202         SVM_func(watermalon, 'linear')
203         SVM_func(watermalon, 'poly')
204         SVM_func(watermalon, 'sigmod')
205         SVM_func(watermalon, 'precomputed')
```

问题 8 输出 调试控制台 终端

1: cmd

```
fold 2 accuracy : 0.5
fold 3 accuracy : 0.5
fold 4 accuracy : 0.5
```

D:\大学\大三上\机器学习\homework1>python part1.py K-NB

```
fold 0 accuracy : 0.75
fold 1 accuracy : 0.75
fold 2 accuracy : 0.75
fold 3 accuracy : 0.75
fold 4 accuracy : 0.75
```

D:\大学\大三上\机器学习\homework1>python part1.py K-NB

```
fold 0 accuracy : 0.75
fold 1 accuracy : 0.75
fold 2 accuracy : 0.75
fold 3 accuracy : 0.75
fold 4 accuracy : 0.75
```

D:\大学\大三上\机器学习\homework1>

行 194, 列 27 空格

SVM结果

```

182
183 if __name__ == "__main__":
184     parse = argparse.ArgumentParser()
185     parse.add_argument("action", type=str, help="type LDA or K-NB or SVM, the program will run a LDA or a Naive Bayes or SVM algorithm.")
186     parse.add_argument("--rate", type=float, default=0.8, help="the proportion of training set, default is 0.8.")
187     parse.add_argument("--fold", type=int, default=5, help="the num of folds in K-fold validation, default is 5.")
188
189     watermalon = pd.read_csv("watermalon.csv", engine="python")
190     iris = pd.read_csv("iris.data", header=None)
191
192     args = parse.parse_args()
193
194     if args.action=="LDA":
195         LDA(watermalon)
196
197     if args.action=="K-NB":
198         KNB(watermalon)
199
200     if args.action=="SVM":
201         SVM_func(watermalon, 'rbf')
202         SVM_func(watermalon, 'linear')
203         SVM_func(watermalon, 'poly')
204         SVM_func(watermalon, 'sigmod')
205         SVM_func(watermalon, 'precomputed')
206

```

问题 8 输出 调试控制台 终端

fold 2 accuracy : 0.75
fold 3 accuracy : 0.75
fold 4 accuracy : 0.75

D:\大学\大三上\机器学习\homework1>python part1.py SVM
the accuracy of SVM using kernel rbf scores 0.5
the accuracy of SVM using kernel linear scores 0.25
the accuracy of SVM using kernel poly scores 0.5
the accuracy of SVM using kernel sigmod scores 0.25
the accuracy of SVM using kernel precomputed scores 0.25

D:\大学\大三上\机器学习\homework1>python part1.py SVM
the accuracy of SVM using kernel rbf scores 0.25
the accuracy of SVM using kernel linear scores 0.25
the accuracy of SVM using kernel poly scores 0.25
the accuracy of SVM using kernel sigmod scores 0.5
the accuracy of SVM using kernel precomputed scores 0.5

D:\大学\大三上\机器学习\homework1>

补充题

代码使用方法

1. python part2.py -h 调用使用帮助

```

D:\大学\大三上\机器学习\homework1>python part2.py -h
usage: part2.py [-h] [--fold FOLD] action

positional arguments:
  action      watermelon or Iris, the program will run algorithm on this two
              dataset.

optional arguments:
  -h, --help   show this help message and exit
  --fold FOLD  the num of folds in K-fold validation, default is 5.

```

2. python part2.py Iris(--fold 5) 括号内容可选且值为默认

3. python part2.py watermelon(--fold 5) 括号内容可选且值为默认

实例

```

D:\大学\大三上\机器学习\homework1>python part2.py Iris --fold 10
D:\大学\大三上\机器学习\homework1>logistic.py:37: RuntimeWarning: overflow encountered in exp
    return 1/(1+np.exp(X.dot(self.weight)))
using Naive Bayes on Iris with 10 folds accuracy : [1.0, 1.0, 1.0, 0.9333333333333333, 0.9333333333333333, 0.8, 0.9333333333333333, 0.86666666666666667, 1.0, 0.86666666666666667]
using LDA on Iris with 10 folds accuracy : [1.0, 0.9333333333333333, 1.0, 1.0, 0.86666666666666667, 1.0, 1.0, 1.0, 1.0, 1.0]
using SVM on Iris with 10 folds accuracy : [0.9333333333333333, 0.9333333333333333, 1.0, 1.0, 0.86666666666666667, 0.9333333333333333, 1.0, 1.0, 1.0, 0.9333333333333333]
using LogisticRegression on Iris with 10 folds accuracy : [1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]

```

实验结果

多种方法iris

```
part1.py part2.py getData.py logistic.py
184     probability0 *= 1/(np.sqrt(2*np.pi)*sigma0)*np.exp(-np.square(float(item[num])-mu0)/(2*np.square(sigma0)))
185     probability1 *= 1/(np.sqrt(2*np.pi)*sigma1)*np.exp(-np.square(float(item[num])-mu1)/(2*np.square(sigma1)))
186     probability2 *= 1/(np.sqrt(2*np.pi)*sigma2)*np.exp(-np.square(float(item[num])-mu2)/(2*np.square(sigma2)))
187     porbalis = [probability0,probability1,probability2]
188     Y_hat.append((porbalis.index(max(porbalis))))
189
190     Y_hat = np.array(Y_hat)
191     Y = np.array(Y)
192
193     return np.sum(Y_hat==Y)/len(X)
194
195 if __name__ == "__main__":
196     parse = argparse.ArgumentParser()
197     parse.add_argument("action", type=str, help="watermelon or Iris, the program will run algorithm on this two dataset.")
198     # parse.add_argument("--rate", type=float, default=0.8, help="the proportion of training set, default is 0.8.")
199     parse.add_argument("--fold", type=int, default=5, help="the num of folds in K-fold validation, default is 5.")
200
201     args = parse.parse_args()
202
203     if args.action=="watermelon":
204         watermalon = pd.read_csv("watermalon.csv",engine="python")
问题 15 输出 调试控制台 终端 2:Py
(base) D:\大学\大三上\机器学习\homework1>python part2.py Iris
D:\大学\大三上\机器学习\homework1\logistic.py:37: RuntimeWarning: overflow encountered in exp
  return 1/(1+np.exp(X.dot(self.weight)))
using Naive Beyas on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.8666666666666667, 0.9333333333333333, 0.9333333333333333]
using LDA on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.9333333333333333, 0.9333333333333333, 1.0]
using SVM on Iris with 5 flocs accuracy : [0.9333333333333333, 1.0, 0.9, 1.0, 0.9666666666666667]
using LogisticRegression on Iris with 5 flocs accuracy : [1.0, 1.0, 1.0, 1.0, 1.0]

(base) D:\大学\大三上\机器学习\homework1>python part2.py
D:\大学\大三上\机器学习\homework1\logistic.py:37: RuntimeWarning: overflow encountered in exp
  return 1/(1+np.exp(X.dot(self.weight)))
using Naive Beyas on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.8666666666666667, 0.9333333333333333, 0.9333333333333333]
using LDA on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.9333333333333333, 0.9333333333333333, 1.0]
using SVM on Iris with 5 flocs accuracy : [0.9333333333333333, 1.0, 0.9, 1.0, 0.9666666666666667]
using LogisticRegression on Iris with 5 flocs accuracy : [1.0, 1.0, 1.0, 1.0, 1.0]

(base) D:\大学\大三上\机器学习\homework1>python part2.py Iris
D:\大学\大三上\机器学习\homework1\logistic.py:37: RuntimeWarning: overflow encountered in exp
  return 1/(1+np.exp(X.dot(self.weight)))
using Naive Beyas on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.8666666666666667, 0.9333333333333333, 0.9333333333333333]
using LDA on Iris with 5 flocs accuracy : [1.0, 0.9666666666666667, 0.9333333333333333, 0.9333333333333333, 1.0]
using SVM on Iris with 5 flocs accuracy : [0.9333333333333333, 1.0, 0.9, 1.0, 0.9666666666666667]
using LogisticRegression on Iris with 5 flocs accuracy : [1.0, 1.0, 1.0, 1.0, 1.0]

(base) D:\大学\大三上\机器学习\homework1>
```

多种方法西瓜

```
问题 15 输出 调试控制台 终端
(base) D:\大学\大三上\机器学习\homework1>python part2.py watermelon
using Naive Beyas on watermelon with 5 flocs accuracy : [0.75, 0.75, 0.75, 0.75, 0.75]
using LDA on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using SVM on watermelon with 5 flocs accuracy : [0.25, 0.25, 0.25, 0.25, 0.25]
using LogisticRegression on Iris with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]

(base) D:\大学\大三上\机器学习\homework1>python part2.py watermelon
using Naive Beyas on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using LDA on watermelon with 5 flocs accuracy : [0.25, 0.25, 0.25, 0.25, 0.25]
using SVM on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using LogisticRegression on Iris with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]

(base) D:\大学\大三上\机器学习\homework1>python part2.py watermelon
using Naive Beyas on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using LDA on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using SVM on watermelon with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]
using LogisticRegression on Iris with 5 flocs accuracy : [0.5, 0.5, 0.5, 0.5, 0.5]

(base) D:\大学\大三上\机器学习\homework1>
```

附注

实验环境

```
sklearn.version = '0.19.1'
```

```
D:\大学\大三上\机器学习\homework1>python
Python 3.6.4 |Anaconda, Inc.| (default, Jan 16 2018, 10:22:32) [MSC v.1900 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import sklearn
>>> sklearn.__version__
'0.19.1'
>>>
```