

一、今日内容

- a. 约会对象预测练习
- b. 特征预处理
- c. 特征抽取
- d. 朴素贝叶斯算法

二、昨日复习

- a. 模块
 - i. import mio
 - ii. import mio as a
 - iii. from xxx import xxx
- b. 异常
 - i. 概念：会导致程序终止
 - ii. 捕获异常

```
1 try:
2     有可能发生异常的代码
3 except Exception as e:
4     code
```

- c. 人工智能历史
 - i. 1950 ~ 1980：简单的机械运动
 - ii. 1980 ~ 2010：机器学习自然语言
 - iii. 2010 ~ ：深度学习处理图片
- d. 机器学习概念
 - i. 从历史数据中，分析获取规律，然后通过规律(模型)对结果进行预测
- e. 数据
 - i. 数据来源：BOSS 花钱购买 爬虫爬取
 - ii. 数据类型：离散型数据 连续性数据
 - iii. 数据格式：数据集：特征值 + 目标值
 - iv. 分割数据集的方法：train_test_split
- f. Knn算法
 - i. 在一个和样本空间中，选取K个最相近的点，如果这K个点中，大部分是属于同一个类别，那么我也认为我是属于这个类别
 - ii. 欧式距离公式：

iii. K值选取：选取奇数个数值。

iv. api:

1. `__init__`
2. `fit()`
3. `predict()`
4. `score()`

三、约会案例

相亲约会对象数据，这个样本是男士的数据，三个特征，玩游戏所消耗时间的百分比、每年获得的飞行常客里程数、每周消费的冰淇淋公升数。然后有一个所属类别，被女士评价的三个类别，不喜欢`didnt`、魅力一般`small`、极具魅力`large`也许也就是说飞行里程数对于结算结果或者说相亲结果影响较大，**但是统计的人觉得这三个特征同等重要。**

| 里程数 | 公升数 | 消耗时间比 | 评价 |
|-------|-----------|----------|------------|
| 14488 | 7.153469 | 1.673904 | smallDoses |
| 26052 | 1.441871 | 0.805124 | didntLike |
| 75136 | 13.147394 | 0.428964 | didntLike |
| 38344 | 1.669788 | 0.134296 | didntLike |
| 72993 | 10.141740 | 1.032955 | didntLike |
| 35948 | 6.830792 | 1.213192 | largeDoses |
| 42666 | 13.276369 | 0.543880 | largeDoses |
| 67497 | 8.631577 | 0.749278 | didntLike |
| 35483 | 12.273169 | 1.508053 | largeDoses |
| 50242 | 3.723498 | 0.831917 | didntLike |

```
1 # target: 预测 [[50000, 5, 0.5]] 女神对我的态度
2
3 # 导入pandas处理数据
4 import pandas as pd
5
6 # 定义目标值名称
7 target_names = ["didntLike", "smallDoses", "largeDoses"]
8
9
10 def main():
11
12     # 1.获取原始数据集 pd.read_csv() 读取数据之后会返回一个DataFrame数据类型
13     dating = pd.read_csv("./dating.txt")
14
15     # dataDream支持像字典一个读取数据，也支持切片操作
```

```

16     print(dating)
17
18     # 2. 确定特征值与目标值
19     x = dating[["milage", "Liters", "Consumtime"]]
20     y = dating["target"]
21
22     print("x : ",x[:1])
23
24     return 0
25
26 main()

```

四、特征预处理

4.1 概念：

将数字类型的数据，通过转换函数，将数据转换为适合机器学习的一组数据。

特征预处理： 归一化， 标准化

4.2 归一化

4.2.1 概念

概念： 将特征值数据，转换为 [0, 1]的范围之间。

4.2.2 公式

公式：

$$X' = \frac{x - \min}{\max - \min}$$

需要转换的特征值

这个特征中的最小值

特征值中的最大值

将数据转换到需要的范围之间

$$X'' = X' * (mx - mi) + mi$$

需要转换的范围的最小值。

需要转换范围的最大值

4.2.3 归一化 API接口

```
1 # 导入特征预处理的api接口
2 from sklearn.preprocessing import MinMaxScaler
3 """
4 MinMaxScaler类
5
6 __init__成员函数
7
8 函数原型:
9     def __init__(self, feature_range=(0, 1));
10
11 函数参数:
12     feature_range:需要转换的数据范围
13
14 fit成员函数
15
16 函数原型:
17     fit(self, X, y=None);
18
19 函数功能:
20     计算转换数据的标准。
21
22 函数参数:
23     X: 需要计数转换数据的标准的世俗据
24
25 transform函数
26 函数原型:
27     transform(self, X):
28
29 函数功能:
30     通过给定的数据范围,对数据进行转换处理
31
32 函数参数:
33     X: 需要转换的数据
34
35 返回值:
36     转换数据的结果
37
38 """
39
40 data = [[90, 2, 10, 40],
41         [60, 4, 15, 45],
```

```

42         [75, 3, 13, 46]]
43
44
45 convert = MinMaxScaler(feature_range=(0, 1))
46
47 # 制定转换数据的标准
48 convert.fit(data)
49 result = convert.transform(data)
50 print(result)
51
52
53 # target: 预测 [[50000, 5, 0.5]] 女神对我的态度
54
55 # 导入pandas处理数据
56 import pandas as pd
57
58 # 导入分割数据集的方法
59 from sklearn.model_selection import train_test_split
60 # 导入knn算法
61 from sklearn.neighbors import KNeighborsClassifier
62 # 导入归一化处理数据的傀儡
63 from sklearn.preprocessing import MinMaxScaler
64
65
66 # 定义目标值名称
67 target_names = [0, "didn'tLike", "smallDoses", "largeDoses"]
68
69
70 def main():
71
72     # 1.获取原始数据集 pd.read_csv() 读取数据之后会返回一个DataFrame数据类型
73     dating = pd.read_csv("./dating.txt")
74
75     # dataDream支持像字典一个读取数据，也支持切片操作
76
77     # 2. 确定特征值与目标值
78     x = dating[["milage", "Liters", "Consumtime"]]
79     y = dating["target"]
80
81     # 3. 划分数据集,将数据集划分为训练集和测试集
82     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state
83
84     # 归一化处理数据
85     convert = MinMaxScaler()

```

```

86
87     convert.fit(x_train)
88     x_train = convert.transform(x_train)
89     x_test = convert.transform(x_test)
90
91     # 4. 实例化一个算法对象
92     estimate = KNeighborsClassifier(n_neighbors=7)
93
94     # 5. 拟合数据制作模型
95     estimate.fit(x_train, y_train)
96
97     # 6. 计算模型的准确度
98     score = estimate.score(x_test, y_test)
99     print("score : ", score)
100
101     # 7. 预测模型的结果
102     data = [[80000, 5, 0.5]]
103
104     # 8 对预测数据进行归一化处理。
105
106     # 对预测数据进行归一化处理
107     data = convert.transform(data)
108     y_predict = estimate.predict(data)
109     print("模型的预测结果: ", target_names[int(y_predict)])
110     return 0
111
112
113 main()

```

五、特征抽取

5.1 概念

将任意类型的数据(图片, 字符串), 转换为适合机器学习的数字类型的特征

5.2 特征抽取原理

统计字符串中, 每一个单词出现的次数, 作为转换之后的特征值。

- i. 制定转换数据的标准: 统计字符串中出现的单词, 去除掉不重要的单词
- ii. 统计出现的次数

```
1 data = ["life is short,i like python",
```

```

2  "life is too long,i dislike python"]
3
4  # 通过每一个单词出现的次数,作为特征值。
5
6  # 制定一个转换数据的标准
7      [life, is, short, like, python, too, long, dislike]
8  数据进行转换。
9  第一篇:  1  1  1  1  1  0  0  0
10  第二篇:  1  1  0  0  1  1  1  1

```

5.3 特征抽取的API接口

```

1  from sklearn.feature_extraction.text import CountVectorizer
2  """
3  CountVectorizer类
4
5  __init__成员函数
6
7  函数原型:
8      def __init__(self);
9
10  函数功能:
11      初始化一个特征抽取的对象
12
13
14  fit成员函数
15
16  函数原型:
17      fit(self, raw_documents, y=None);
18
19  函数功能
20      统计字符串中出现的单词, 去除掉不重要的单词
21
22  函数参数:
23      raw_documents: 需要转换的文章
24
25
26  transform成员函数
27
28  函数原型:
29      transform(self, raw_documents);

```

```

30
31  函数功能：
32      对原始数据进行转换处理
33
34  函数返回值：
35      sparse类型矩阵： 方便计算机存储
36
37  get_feature_names()函数
38  函数功能：
39      返回制定的转换数据的标准
40
41  函数使用方式：
42      obj.get_feature_names()
43
44  toarray()函数
45  函数功能：
46      将数据转换为数组展示出来
47
48  使用方式：
49      obj.toarray()
50  """
51
52  data = ["life is short,i like python", "life is too long,i dislike python"]
53
54  # 实例化转换数据的对象
55  concert = CountVectorizer()
56
57  # 制定转换数据的标准
58  concert.fit(data)
59  print(concert.get_feature_names())
60
61  # 开始转换数据
62  result = concert.transform(data)
63  print(result.toarray())

```

六、中文分词 jieba 模块

6.1 jieba模块安装

- i. windows + r ---> cmd + 回车
- ii. pip install jieba

6.2 jieba模块使用

```
1 # jieba模块
2 import jieba # 导入结巴模块
3 """
4 jieba.cut函数
5
6 函数原型：
7     cut(self, sentence);
8
9 函数功能：
10    对中文文章进行分词操作
11
12 函数参数：
13    sentence：需要进行分词的中文语句
14
15 函数返回值：
16    返回值： 生成器，如果需要看内部的内容，需要转换为列表操作的。
17 """
18
19 # 原始数据
20 data = "新鲜草莓"
21
22 # 调用cut函数进行分词处理
23 result = list(jieba.cut(data))
24
25 print(result)
```

6.3 中文文章抽取

```
1 """
2 今天很残酷，明天更残酷，后天很美好，但绝大部分是死在明天晚上，所以每个人不要放弃今天。
3
4 我们看到的从很远星系来的光是在几百万年之前发出的，这样当我们看到宇宙时，我们是在看它的过去。
5
6 如果只用一种方式了解某样事物，你就不会真正了解它。了解事物真正含义的秘密取决于如何将其与我们所了
7 """
8
9 # 导入结巴模块
10 import jieba
```

```
11
12 # 导入特征抽取
13 from sklearn.feature_extraction.text import CountVectorizer
14
15
16 def cut_word():
17
18     # 获取原始数据集
19     c1 = "今天很残酷，明天更残酷，后天很美好，但绝对大部分是死在明天晚上，所以每个人不要放弃今天"
20     c2 = "我们看到的从很远星系来的光是在几百万年之前发出的，这样当我们看到宇宙时，我们是在看它的"
21     c3 = "如果只用一种方式了解某样事物，你就不会真正了解它。了解事物真正含义的秘密取决于如何将其"
22
23     # result的结果是一个列表。将列表里面的词语组成一个字符串。列表里面的字符串需要使用 space隔
24     result1 = list(jieba.cut(c1))
25     result2 = list(jieba.cut(c2))
26     result3 = list(jieba.cut(c3))
27
28     ret1 = ""
29     ret2 = ""
30     ret3 = ""
31
32     for i in result1:
33         # 1.里面的每一个值，赋值给I
34         ret1 += i
35         ret1 += " "
36     for i in result2:
37         # 1.里面的每一个值，赋值给I
38         ret2 += i
39         ret2 += " "
40     for i in result3:
41         # 1.里面的每一个值，赋值给I
42         ret3 += i
43         ret3 += " "
44     return ret1, ret2, ret3
45
46
47 # 特征抽取处理数据
48 def main():
49
50     # 1.获取原始数据集
51     ret1, ret2, ret3 = cut_word()
52
53     # 2.实例化转换器对象
54     concert = CountVectorizer()
```

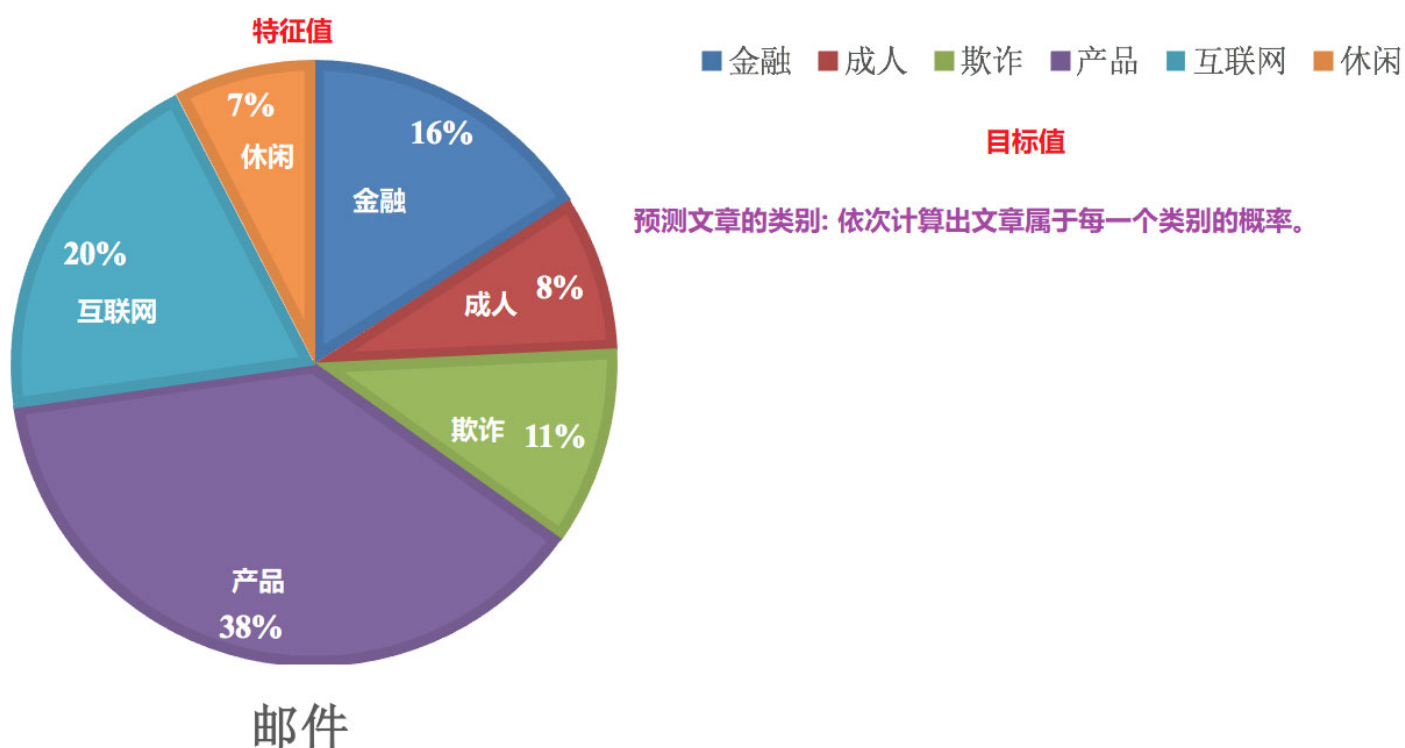
```

55 concert.fit([ret1, ret2, ret3])
56
57 # 3.对数据进行转换
58 result = concert.transform([ret1, ret2, ret3])
59 print(concert.get_feature_names())
60 print("result : ", result.toarray())
61
62
63 return 0
64
65 main()

```

七、文章分类 -- 朴素贝叶斯算法

7.1 前言



7.2 概率

7.2.1 概念

一件事情发生的可能性，可能发生，也可能不发生。记作 $P(A)$

7.2.2 概率计算

- 联合概率：不相关的事情同时发生的概率。
 - 记作： $P(A, B) = P(A) + P(B)$
- 条件概率：在某个条件的情况下，事件发生的概率

a. 记作: $P(A|B)$

| 样本数 | 职业 | 体型 | 女神是否喜欢 |
|-----|-----|----|--------|
| 1 | 程序员 | 超重 | 不喜欢 |
| 2 | 产品 | 匀称 | 喜欢 |
| 3 | 程序员 | 匀称 | 喜欢 |
| 4 | 程序员 | 超重 | 喜欢 |
| 5 | 美工 | 匀称 | 不喜欢 |
| 6 | 美工 | 超重 | 不喜欢 |
| 7 | 产品 | 匀称 | 喜欢 |

1、女神喜欢的概率？

$$P(\text{女神喜欢}) = \text{喜欢} / \text{全部} = 4 / 7$$

2、职业是程序员并且体型匀称的概率？

$$P(\text{程序员}, \text{匀称}) = P(\text{程序员}) * P(\text{匀称}) = 12 / 49$$

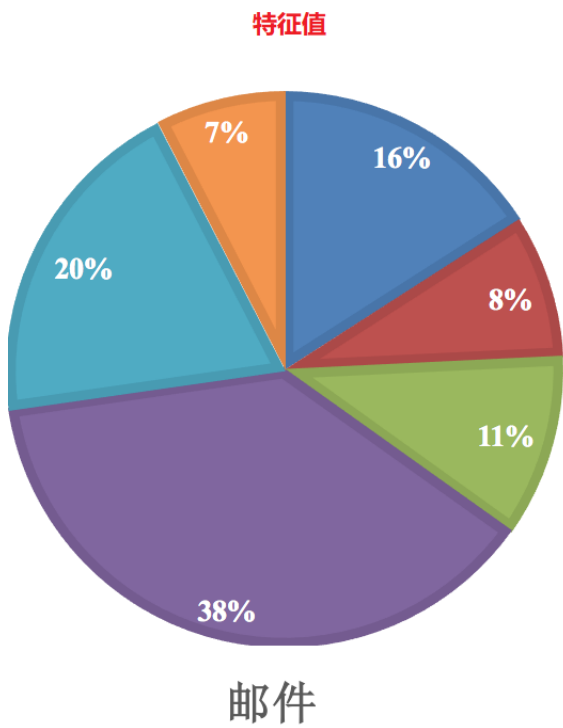
3、在女神喜欢的条件下，职业是程序员的概率？

$$P(\text{程序员}|\text{女神喜欢}) = 1 / 2$$

4、在女神喜欢的条件下，职业是产品，体重是超重的概率？

$$P(\text{产品}, \text{超重}|\text{女神喜欢}) = 1 / 2 * 1 / 4 = 1 / 8$$

7.3 文章概率如何计算？



邮件

目标值
■ 金融 ■ 成人 ■ 欺诈 ■ 产品 ■ 互联网 ■ 休闲

target: 计算邮件属于产品的概率?

$P(\text{产品}|\text{邮件})$: 条件概率

特征抽取: 将邮件转换为 特征词

$P(\text{产品} | \text{word1, word2 ... wordN})$

$$P(C|W) = \frac{P(W|C)P(C)}{P(W)}$$

类别 文章

$P(\text{word1, word2... wordN} | \text{产品}) * P(\text{产品})$

$P(\text{word1, word2 ... wordN})$

$P(\text{word1}|\text{产品}) * P(\text{word2} | \text{产品}) * \dots * P(\text{wordN}|\text{产品}) * P(\text{产品})$

$P(\text{word1}) * P(\text{word2}) * \dots * P(\text{wordN})$

7.4 手动计算文章所属类别

训练集统计结果(指定统计词频):

| 特征\统计 | 科技(30篇) | 娱乐(60篇) | 汇总 (90篇) |
|--------|---------|---------|----------|
| “商场” | 9 | 51 | 60 |
| “影院” | 8 | 56 | 64 |
| “支付宝” | 20 | 15 | 35 |
| “云计算” | 63 | 0 | 63 |
| 汇总(求和) | 100 | 121 | 221 |

target: 计算文章属于科技的类别是多少?

$P(\text{科技}|\text{文章})$

|| 对文章进行特征抽取

$P(\text{科技} | \text{影院, 支付宝, 云计算})$

|| 贝叶斯公式

$P(\text{影院, 支付宝, 云计算} | \text{科技}) * P(\text{科技})$

$P(\text{影院, 支付宝, 云计算})$

||

$P(\text{影院}|\text{科技}) * P(\text{支付宝}|\text{科技}) * P(\text{云计算} | \text{科技}) * P(\text{科技})$

$P(\text{影院}) * P(\text{支付宝}) * P(\text{云计算})$

现有一篇被预测文档: 出现了影院, 支付宝, 云计算, 计算属于科技、娱乐的类别概率?

$P(\text{影院}|\text{科技}) * P(\text{支付宝}|\text{科技}) * P(\text{云计算} | \text{科技}) * P(\text{科技})$

$P(\text{影院}) * P(\text{支付宝}) * P(\text{云计算})$

$(8 / 100) * (20 / 100) * (63 / 100) * (30 / 90) = 0.00335664$

$(64/221) * (35 / 221) * (63/ 221) \approx 0.0130740983 \approx 25.67\%$

7.5 拉普拉斯平滑系数

i. 目的:防止计算的概率为 0

ii. 拉普拉斯平滑系数: 添加在条件概率中的.

iii. 拉普拉斯平滑系数公式:

条件下出现的次数

拉普拉斯平滑系数: 1

$$P(F1|C) = \frac{Ni + \alpha}{N + \alpha m}$$

条件概率的计算

总的出现次数

训练集的特征值个数

7.6 朴素贝叶斯算法API接口

```

1 # 朴素贝叶斯算法API接口
2 from sklearn.naive_bayes import MultinomialNB
3 """
4 MultinomialNB成员函数
5
6 __init__成员函数:
7
8 函数原型:
9     __init__(self, *, alpha=1.0);
10
11 函数参数:
12     alpha: 拉普拉斯平滑系数
13
14
15 fit成员函数
16
17 函数原型:
18     fit(self, X, y);
19
20 函数功能:
21     通过训练集数据,制作模型
22
23 函数参数:
24     X: 训练集的特征值
25     y: 训练集的目标值
26
27
28 predict(self, x): 通过数据预测模型的标签值

```

```
29 score(self, x, y):计算模型的准确度
30 """
```

八、新闻分类案例

20个新闻组

20个新闻组数据集

20个新闻组数据集是大约20,000个新闻组文档的集合，平均分布在20个不同的新闻组中。据我所知，它最初是由Ken Lang收集的，可能是因为他的[Newsweeder: 学习过滤网络新闻纸](#)，尽管他没有明确地提到这个集合。这20个新闻组集合已经成为机器学习技术的文本应用中的实验流行数据集，例如文本分类和文本聚类。

组织

数据被组织成20个不同的新闻组，每个新闻组对应不同的主题。一些新闻组彼此之间关系密切（例如comp.sys.ibm.pc.hardware / comp.sys.mac.hardware），而另一些新闻组则非常不相关（例如misc.forsale / soc.religion.christian）。以下是根据主题划分的20个新闻组列表（或多或少）：

| | | |
|---|--|---|
| comp.graphics comp.os.ms-windows.misc comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x | rec.autos rec.motorcycles rec.sport.baseball rec.sport.hockey | sci.crypt sci.electronics sci.med sci.space |
| misc.forsale | talk.politics.misc talk.politics.guns talk.politics.mideast | talk.religion.misc alt.atheism soc.religion.christian |

```
1 # 导入新闻网站的数据
2 from sklearn.datasets import fetch_20newsgroups
3 # 导入分割数据集的方法
4 from sklearn.model_selection import train_test_split
5 # 导入特征抽取，对数据进行处理
6 from sklearn.feature_extraction.text import CountVectorizer
7 # 导入朴素贝叶斯算法
8 from sklearn.naive_bayes import MultinomialNB
9
10
11
12 def main():
13
14     # 1.获取到原始数据集
15     news = fetch_20newsgroups(subset="all") # 训练集数据与测试集数据全都获取
16
17     # 2.确定特征值与目标值
18     x = news.data
19     y = news.target
20
21     # 3.划分数据集,将数据集划分为训练集和测试集
22     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3)
23
24     # 4.实例化一个转换数据的对象
25     concert = CountVectorizer()
```

```

26
27     # 5.制定转换数据的标准
28     concert.fit(x_train)
29
30     # 6.开始转换数据
31     x_train = concert.transform(x_train)
32     x_test = concert.transform(x_test)
33
34     # 7.实例化朴素贝叶斯算法对象
35     estimate = MultinomialNB()
36     estimate.fit(x_train, y_train)
37
38     score = estimate.score(x_test, y_test)
39     print(score)
40
41     # 8.对数据进行预测
42     data = ""Amid the rapprochement momentum in the Middle East and the US' strategic c
43
44 Blinken's trip, from June 6 to 8, comes after a visit by White House national security a
45
46 According to the US State Department, Blinken will meet with Saudi officials to "discuss
47
48 The top US diplomat will also attend the Gulf Cooperation Council (GCC) talks during his
49
50 Before Blinken's visit, Saudi Arabia's foreign minister Faisal bin Farhan met with his J
51
52 In April, Saudi Arabia and Iran formally restored diplomatic ties after a seven-year rift
53
54 With the rapprochement process in the Middle East accelerating, observers are not optimi
55
56 Liu Zhongmin, a professor with the Middle East Studies Institute of Shanghai Internation
57
58 On the one hand, the US is shifting its strategic focus to the Asia-Pacific region to co
59
60 The US cannot stop the trend of increasing autonomy in the Middle East, so it is trying
61
62 Li Haidong, a professor at the China Foreign Affairs University, told the Global Times c
63
64 In the process of adapting to the new situation in the Middle East, the US would try to
65
66 While the US is less likely to play a constructive role in the Middle East, its disrupti
67
68 The US is likely to take stock of the situation and become more involved in the region i
69

```



```
70 The oil-producing countries of the Middle East seem to have woken up to the risks, espec
71
72 On June 2, Saudi minister of energy Prince Abdulaziz bin Salman met Zhang Jianhua, admir
73
74 Middle East countries faced a historic opportunity to liquidate the power of the US and
75
76 Chinese Foreign Ministry spokesperson Wang Wenbin said on Monday that China supports reg
77
78 As a good friend of regional countries, China will continue to play an active and constr
79
80     data = concert.transform([data])
81     result = estimate.predict(data)
82
83     print("文章的类别: ", news.target_names[int(result)])
84     return 0
85
86
87 main()
88
89
90
91
```