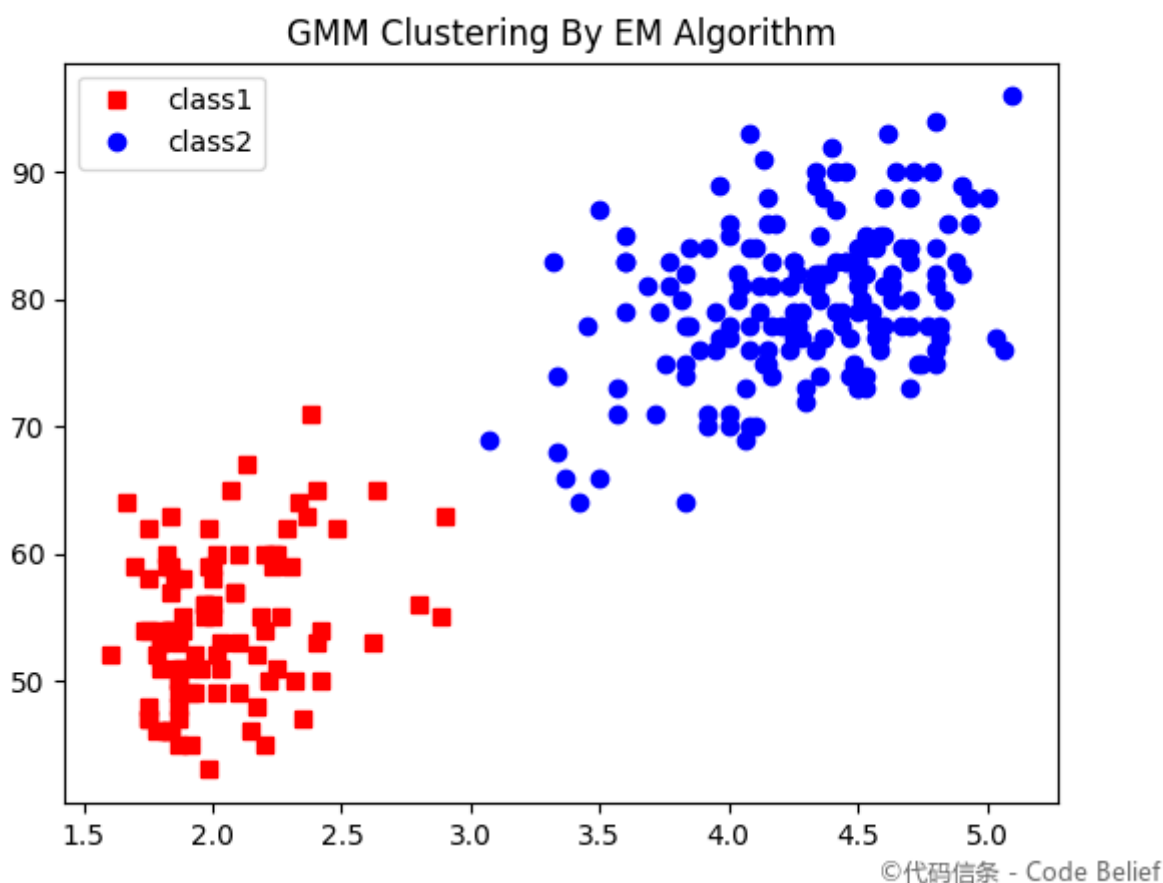




数据库管理系统索引技术概述 (<http://www.codebelief.com/article/2018/04/dbms-index-technique-overview/>)

代码信条 - Code Belief (<http://www.codebelief.com/>)



高斯混合模型 EM 算法的 Python 实现

📅 2017年11月24日 (<http://www.codebelief.com/article/2017/11/gmm-em-algorithm-implementation-by-python/>)
👤 Wray Zheng (<http://www.codebelief.com/article/author/wrayzheng/>) 1,537

高斯混合模型简介

首先简单介绍一下，高斯混合模型（GMM, Gaussian Mixture Model）是多个高斯模型的线性叠加，高斯混合模型的概率分布可以表示如下：

文章目录 
高斯混合模型简介
GMM 的 EM 算法
E 步
M 步



$$P(x) = \sum_{k=1}^K \alpha_k \phi(x; \mu_k, \Sigma_k)$$

其中， K 表示模型的个数， α_k 是第 k 个模型的系数，表示出现该模型的概率， $\phi(x; \mu_k, \Sigma_k)$ 是第 k 个高斯模型的概率分布。

这里讨论的是多个随机变量的情况，即多元高斯分布，因此高斯分布中的参数不再是方差 σ_k ，而是协方差矩阵 Σ_k 。

我们的目标是给定一堆没有标签的样本和模型个数 K ，以此求得混合模型的参数，然后就可以用这个模型来对样本进行聚类。

GMM 的 EM 算法

我们知道，EM 算法是通过不断迭代来求得最佳参数的。在执行该算法之前，需要先给出一个初始化的模型参数。

我们让每个模型的 μ 为随机值， Σ 为单位矩阵， α 为 $\frac{1}{K}$ 即每个模型初始时都是等概率出现的。

EM 算法可以分为 E 步和 M 步。

E 步

直观理解就是我们已经知道了样本 x_i ，那么它是由哪个模型产生的呢？我们这里求的就是：样本 x_i 来自于第 k 个模型的概率，我们把这个概率称为模型 k 对样本 x_i 的“责任”，也叫“响应度”，记作 γ_{ik} ，计算公式如下：

$$\gamma_{ik} = \frac{\alpha_k \phi(x_i; \mu_k, \Sigma_k)}{\sum_{k=1}^K \alpha_k \phi(x_i; \mu_k, \Sigma_k)}$$

Normalize to 0-1!

M 步

根据样本和当前 γ 矩阵重新估计参数，注意这里 x 为列向量：

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_{ik} x_i$$

note: r_{ik} is related to sample i and model k . when compute mean, not only choose samples belong to model k , instead choose all samples, but give a prob to it.

$$\Sigma_k = \frac{1}{N_k} \sum_{i=1}^N \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^T$$

$$\alpha_k = \frac{N_k}{N}$$



Python 实现

在给出代码前，先作一些说明。

- 在对样本应用高斯混合模型的 EM 算法前，需要先进行数据预处理，即把所有样本值都缩放到 0 和 1 之间。
- 初始化模型参数时，要确保任意两个模型之间参数没有完全相同，否则迭代到最后，两个模型的参数也将完全相同，相当于一个模型。
- 模型的个数必须大于 1。当 K 等于 1 时相当于将样本聚成一类，没有任何意义。

`gmm.py` 文件:



```

1  # -*- coding: utf-8 -*-
2  # -----
3  # Copyright (c) 2017, Wray Zheng. All Rights Reserved.
4  # Distributed under the BSD License.
5  # -----
6
7  import numpy as np
8  import matplotlib.pyplot as plt
9  from scipy.stats import multivariate_normal
10
11  DEBUG = True
12
13  #####
14  # 调试输出函数
15  # 由全局变量 DEBUG 控制输出
16  #####
17  def debug(*args, **kwargs):
18      global DEBUG
19      if DEBUG:
20          print(*args, **kwargs)
21
22
23  #####
24  # 第 k 个模型的高斯分布密度函数
25  # 每 i 行表示第 i 个样本在各模型中的出现概率
26  # 返回一维列表
27  #####
28  def phi(Y, mu_k, cov_k):
29      norm = multivariate_normal(mean=mu_k, cov=cov_k)
30      return norm.pdf(Y)
31
32
33  #####
34  # E 步: 计算每个模型对样本的响应度
35  # Y 为样本矩阵, 每个样本一行, 只有一个特征时为列向量
36  # mu 为均值多维数组, 每行表示一个样本各个特征的均值
37  # cov 为协方差矩阵的数组, alpha 为模型响应度数组
38  #####
39  def getExpectation(Y, mu, cov, alpha):
40      # 样本数
41      N = Y.shape[0]
42      # 模型数
43      K = alpha.shape[0]
44
45      # 为避免使用单个高斯模型或样本, 导致返回结果的类型不一致
46      # 因此要求样本数和模型个数必须大于1
47      assert N > 1, "There must be more than one sample!"
48      assert K > 1, "There must be more than one gaussian model!"
49
50      # 响应度矩阵, 行对应样本, 列对应响应度
51      gamma = np.mat(np.zeros((N, K)))
52
53      # 计算各模型中所有样本出现的概率, 行对应样本, 列对应模型
54      prob = np.zeros((N, K))

```



```

55     for k in range(K):
56         prob[:, k] = phi(Y, mu[k], cov[k])
57     prob = np.mat(prob)
58
59     # 计算每个模型对每个样本的响应度
60     for k in range(K):
61         gamma[:, k] = alpha[k] * prob[:, k]
62     for i in range(N):
63         gamma[i, :] /= np.sum(gamma[i, :])
64     return gamma
65
66
67     #####
68     # M 步: 迭代模型参数
69     # Y 为样本矩阵, gamma 为响应度矩阵
70     #####
71     def maximize(Y, gamma):
72         # 样本数和特征数
73         N, D = Y.shape
74         # 模型数
75         K = gamma.shape[1]
76
77         # 初始化参数值
78         mu = np.zeros((K, D))
79         cov = []
80         alpha = np.zeros(K)
81
82         # 更新每个模型的参数
83         for k in range(K):
84             # 第 k 个模型对所有样本的响应度之和
85             Nk = np.sum(gamma[:, k])
86             # 更新 mu
87             # 对每个特征求均值
88             for d in range(D):
89                 mu[k, d] = np.sum(np.multiply(gamma[:, k], Y[:, d])) / Nk
90             # 更新 cov
91             cov_k = np.mat(np.zeros((D, D)))
92             for i in range(N):
93                 cov_k += gamma[i, k] * (Y[i] - mu[k]).T * (Y[i] - mu[k]) / Nk
94             cov.append(cov_k)
95             # 更新 alpha
96             alpha[k] = Nk / N
97         cov = np.array(cov)
98         return mu, cov, alpha
99
100
101     #####
102     # 数据预处理
103     # 将所有数据都缩放到 0 和 1 之间
104     #####
105     def scale_data(Y):
106         # 对每一维特征分别进行缩放
107         for i in range(Y.shape[1]):
108             max_ = Y[:, i].max()
109             min_ = Y[:, i].min()

```



```
110     Y[:, i] = (Y[:, i] - min_) / (max_ - min_)
111     debug("Data scaled.")
112     return Y
113
114
115     #####
116     # 初始化模型参数
117     # shape 是表示样本规模的二元组, (样本数, 特征数)
118     # K 表示模型个数
119     #####
120     def init_params(shape, K):
121         N, D = shape
122         mu = np.random.rand(K, D)
123         cov = np.array([np.eye(D)] * K)
124         alpha = np.array([1.0 / K] * K)
125         debug("Parameters initialized.")
126         debug("mu:", mu, "cov:", cov, "alpha:", alpha, sep="\n")
127         return mu, cov, alpha
128
129
130     #####
131     # 高斯混合模型 EM 算法
132     # 给定样本矩阵 Y, 计算模型参数
133     # K 为模型个数
134     # times 为迭代次数
135     #####
136     def GMM_EM(Y, K, times):
137         Y = scale_data(Y)
138         mu, cov, alpha = init_params(Y.shape, K)
139         for i in range(times):
140             gamma = getExpectation(Y, mu, cov, alpha)
141             mu, cov, alpha = maximize(Y, gamma)
142             debug("{sep} Result {sep}".format(sep="-" * 20))
143             debug("mu:", mu, "cov:", cov, "alpha:", alpha, sep="\n")
144             return mu, cov, alpha
```

main.py 文件:

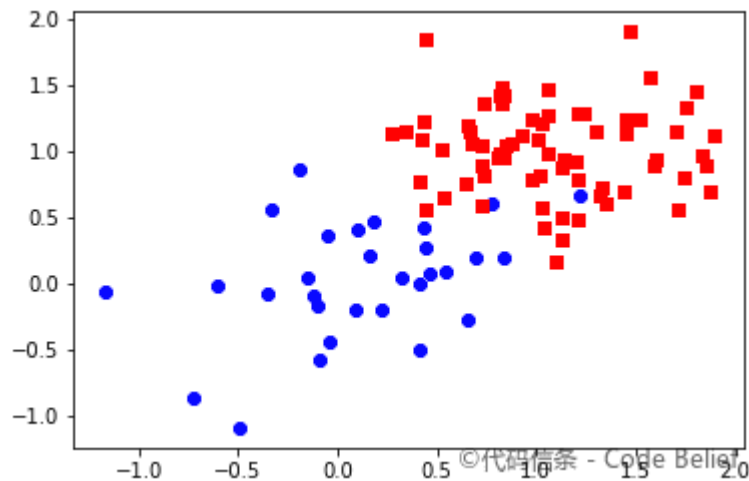


```
1  # -*- coding: utf-8 -*-
2  # -----
3  # Copyright (c) 2017, Wray Zheng. All Rights Reserved.
4  # Distributed under the BSD License.
5  # -----
6
7  import matplotlib.pyplot as plt
8  from gmm import *
9
10 # 设置调试模式
11 DEBUG = True
12
13 # 载入数据
14 Y = np.loadtxt("gmm.data")
15 matY = np.matrix(Y, copy=True)
16
17 # 模型个数，即聚类的类别个数
18 K = 2
19
20 # 计算 GMM 模型参数
21 mu, cov, alpha = GMM_EM(matY, K, 100)
22
23 # 根据 GMM 模型，对样本数据进行聚类，一个模型对应一个类别
24 N = Y.shape[0]
25 # 求当前模型参数下，各模型对样本的响应度矩阵
26 gamma = getExpectation(matY, mu, cov, alpha)
27 # 对每个样本，求响应度最大的模型下标，作为其类别标识
28 category = gamma.argmax(axis=1).flatten().tolist()[0]
29 # 将每个样本放入对应类别的列表中
30 class1 = np.array([Y[i] for i in range(N) if category[i] == 0])
31 class2 = np.array([Y[i] for i in range(N) if category[i] == 1])
32
33 # 绘制聚类结果
34 plt.plot(class1[:, 0], class1[:, 1], 'rs', label="class1")
35 plt.plot(class2[:, 0], class2[:, 1], 'bo', label="class2")
36 plt.legend(loc="best")
37 plt.title("GMM Clustering By EM Algorithm")
38 plt.show()
```

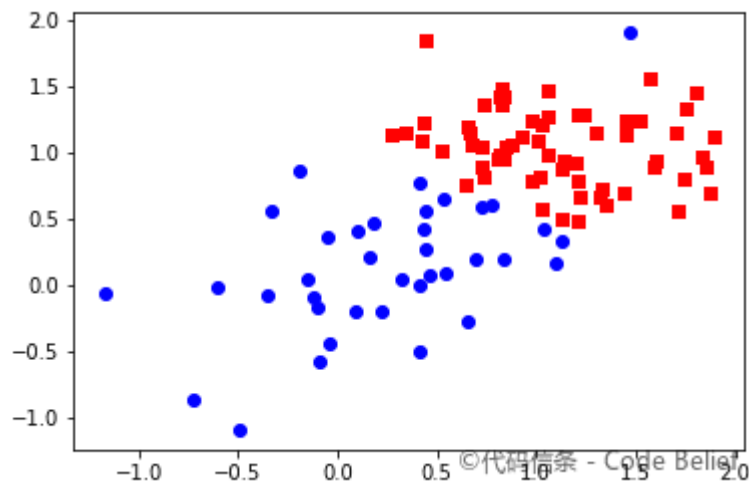
测试结果



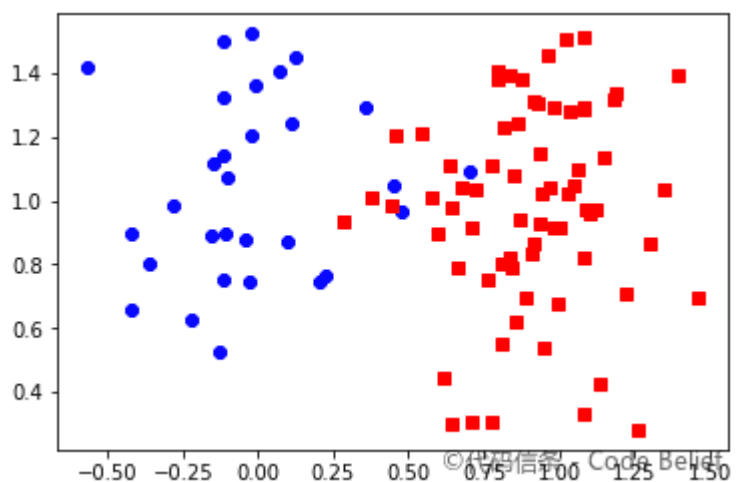
样本 1 原始类别:



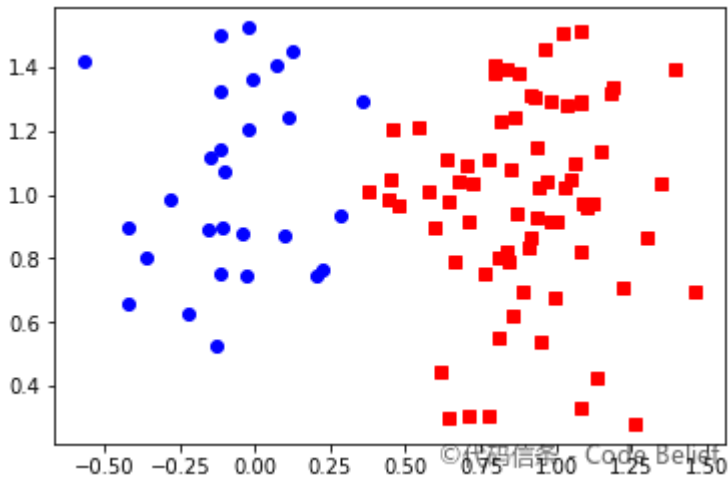
样本 1 聚类结果:



样本 2 原始类别:



样本 2 聚类结果:



测试数据

完整代码和数据放在了 Github 上, 可 点此访问 (<https://github.com/wrayzheng/gmm-em-clustering>)。

或者通过下列代码生成测试数据:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 cov1 = np.mat("0.3 0;0 0.1")
5 cov2 = np.mat("0.2 0;0 0.3")
6 mu1 = np.array([0, 1])
7 mu2 = np.array([2, 1])
8
9 sample = np.zeros((100, 2))
10 sample[:30, :] = np.random.multivariate_normal(mean=mu1, cov=cov1, size=30)
11 sample[30:, :] = np.random.multivariate_normal(mean=mu2, cov=cov2, size=70)
12 np.savetxt("sample.data", sample)
13
14 plt.plot(sample[:30, 0], sample[:30, 1], "bo")
15 plt.plot(sample[30:, 0], sample[30:, 1], "rs")
16 plt.show()
```

参考资料

- 统计学习方法, 李航
- **Pattern Recognition and Machine Learning**, Christopher M Bishop

作者: Wray Zheng (<http://www.codebelief.com>)

原文: <http://www.codebelief.com/article/2017/11/gmm-em-algorithm-implementation-by-python/>
(<http://www.codebelief.com/article/2017/11/gmm-em-algorithm-implementation-by-python/>)



相关文章

- Logistic 回归原理分析与 Python 实现 (<http://www.codebelief.com/article/2017/11/logistic-regression-analysis-and-python-realization/>)
- Python 数据可视化: matplotlib 库学习笔记 (<http://www.codebelief.com/article/2017/04/python-data-visualization-matplotlib-learning/>)
- 机器学习: scikit-learn 实现手写数字识别 (<http://www.codebelief.com/article/2017/05/python-machine-learning-scikit-learn-to-recognize-handwritten-digits/>)
- matplotlib 演示最小二乘法拟合过程 (<http://www.codebelief.com/article/2017/04/matplotlib-demonstrate-least-square-regression-process/>)
- Python 数据分析: numpy 多维数组 ndarray (<http://www.codebelief.com/article/2017/03/python-data-analysis-numpy-ndarray/>)
- [译]Python 受推崇的 super! (<http://www.codebelief.com/article/2017/03/python-super-considered-super/>)

♥ 3 人喜欢

Python (<http://www.codebelief.com/tag/python/>)、机器学习 (<http://www.codebelief.com/tag/machine-learning/>)

← 自动机: DFA转化为正则表达式
([HTTP://WWW.CODEBELIEF.COM/ARTICLE/2017/11/AUTOMATA-CONVERT-DFA-TO-REGULAR-EXPRESSION/](http://www.codebelief.com/article/2017/11/AUTOMATA-CONVERT-DFA-TO-REGULAR-EXPRESSION/))

TOMCAT WEB 应用绑定域名的几种方式 →
([HTTP://WWW.CODEBELIEF.COM/ARTICLE/2017/12/WAYS-FOR-TOMCAT-WEB-APP-TO-BIND-DOMAINS/](http://www.codebelief.com/article/2017/12/WAYS-FOR-TOMCAT-WEB-APP-TO-BIND-DOMAINS/))

发表评论

电子邮件地址不会被公开。必填项已用*标注

评论

姓名 *



电子邮件 *

站点

发表评论

文章搜索

在此输入并回车

近期文章

- ✎ 局部最优未必是全局最优 (<http://www.codebelief.com/article/2018/05/local-optimum-is-not-necessarily-global-optimum/>)
- ✎ Java 多线程下载器的设计与实现 (<http://www.codebelief.com/article/2018/04/java-multithread-downloader-design-and-implementation/>)
- ✎ 数据库管理系统索引技术概述 (<http://www.codebelief.com/article/2018/04/dbms-index-technique-overview/>)
- ✎ 彻底理解二分查找及其边界情况 (<http://www.codebelief.com/article/2018/04/completely-understand-binary-search-and-its-boundary-cases/>)
- ✎ Java logging 模块的使用 (<http://www.codebelief.com/article/2018/04/usage-of-java-logging-module/>)
- ✎ 数据库视图的概念、定义与使用 (<http://www.codebelief.com/article/2018/03/database-view-concept-define-and-usage/>)
- ✎ 理解数据库的内连接、外连接和交叉连接 (<http://www.codebelief.com/article/2018/03/understand-database-inner-join-outer-join-and-cross-join/>)
- ✎ Java Swing 编写数据库增删改查 GUI 程序 (<http://www.codebelief.com/article/2018/03/java-swing-code-a-mysql-crud-gui-program/>)
- ✎ Java 是如何利用接口避免函数回调的 (<http://www.codebelief.com/article/2018/02/java-how-to-use-interface-to-avoid-function-callback/>)
- ✎ Java 多线程的竞争条件、互斥和同步 (<http://www.codebelief.com/article/2018/02/java-multi-thread-race-condition-mutual-exclusion-and-synchronization/>)



文章归档

- 📅 2018年五月 (<http://www.codebelief.com/article/2018/05/>) (1)
- 📅 2018年四月 (<http://www.codebelief.com/article/2018/04/>) (4)
- 📅 2018年三月 (<http://www.codebelief.com/article/2018/03/>) (3)
- 📅 2018年二月 (<http://www.codebelief.com/article/2018/02/>) (2)
- 📅 2018年一月 (<http://www.codebelief.com/article/2018/01/>) (4)
- 📅 2017年十二月 (<http://www.codebelief.com/article/2017/12/>) (3)
- 📅 2017年十一月 (<http://www.codebelief.com/article/2017/11/>) (5)
- 📅 2017年十月 (<http://www.codebelief.com/article/2017/10/>) (3)
- 📅 2017年九月 (<http://www.codebelief.com/article/2017/09/>) (5)
- 📅 2017年八月 (<http://www.codebelief.com/article/2017/08/>) (1)
- 📅 2017年七月 (<http://www.codebelief.com/article/2017/07/>) (1)
- 📅 2017年六月 (<http://www.codebelief.com/article/2017/06/>) (5)
- 📅 2017年五月 (<http://www.codebelief.com/article/2017/05/>) (7)
- 📅 2017年四月 (<http://www.codebelief.com/article/2017/04/>) (8)
- 📅 2017年三月 (<http://www.codebelief.com/article/2017/03/>) (8)
- 📅 2017年二月 (<http://www.codebelief.com/article/2017/02/>) (7)
- 📅 2017年一月 (<http://www.codebelief.com/article/2017/01/>) (5)
- 📅 2016年十二月 (<http://www.codebelief.com/article/2016/12/>) (4)
- 📅 2016年十一月 (<http://www.codebelief.com/article/2016/11/>) (2)
- 📅 2016年十月 (<http://www.codebelief.com/article/2016/10/>) (1)

标签

Python (23) (<http://www.codebelief.com/tag/python/>)

Java (13) (<http://www.codebelief.com/tag/java/>)

Linux (12) (<http://www.codebelief.com/tag/linux/>)

工具资源 (12) (<http://www.codebelief.com/tag/tool-resources/>)

数据库 (9) (<http://www.codebelief.com/tag/database/>)

算法 (8) (<http://www.codebelief.com/tag/algorithm/>)

命令行 (8) (<http://www.codebelief.com/tag/command-line/>)

SQL (7) (<http://www.codebelief.com/tag/sql/>)

GUI (6) (<http://www.codebelief.com/tag/gui/>)

C++ (5) (<http://www.codebelief.com/tag/cpp/>)

操作系统 (4) (<http://www.codebelief.com/tag/operating-system/>)

翻译 (4) (<http://www.codebelief.com/tag/translation/>)

PyQt (4) (<http://www.codebelief.com/tag/pyqt/>)



数据结构 (3) (<http://www.codebelief.com/tag/data-structure/>)

Python标准库 (3) (<http://www.codebelief.com/tag/python-standard-library/>)

正则表达式 (3) (<http://www.codebelief.com/tag/regular-expression/>)

Vim (3) (<http://www.codebelief.com/tag/vim/>)

Qt (3) (<http://www.codebelief.com/tag/qt/>)

机器学习 (3) (<http://www.codebelief.com/tag/machine-learning/>)

Web (3) (<http://www.codebelief.com/tag/web/>)

自动机 (3) (<http://www.codebelief.com/tag/automata/>)

Node.js (2) (<http://www.codebelief.com/tag/node-js/>)

计算机网络 (2) (<http://www.codebelief.com/tag/network/>)

matplotlib (2) (<http://www.codebelief.com/tag/matplotlib/>)

Jupyter Notebook (2) (<http://www.codebelief.com/tag/jupyter-notebook/>)

感想 (2) (<http://www.codebelief.com/tag/sentiment/>)

Swing (2) (<http://www.codebelief.com/tag/swing/>)

Maven (2) (<http://www.codebelief.com/tag/maven/>)

软件工程 (2) (<http://www.codebelief.com/tag/software-engineering/>)

并发 (2) (<http://www.codebelief.com/tag/concurrency/>)

Git (2) (<http://www.codebelief.com/tag/git/>)

Awt (1) (<http://www.codebelief.com/tag/awt/>)

sed (1) (<http://www.codebelief.com/tag/sed/>)

shell (1) (<http://www.codebelief.com/tag/shell/>)

Ubuntu (1) (<http://www.codebelief.com/tag/ubuntu/>)

VPS (1) (<http://www.codebelief.com/tag/vps/>)

Boost (1) (<http://www.codebelief.com/tag/boost/>)

cURL (1) (<http://www.codebelief.com/tag/curl/>)

scikit-learn (1) (<http://www.codebelief.com/tag/scikit-learn/>)

排序 (1) (<http://www.codebelief.com/tag/sort/>)

科学计算 (1) (<http://www.codebelief.com/tag/scientific-computing/>)

二叉树 (1) (<http://www.codebelief.com/tag/binary-tree/>)

Express (1) (<http://www.codebelief.com/tag/express/>)

并查集 (1) (<http://www.codebelief.com/tag/union-find/>)

Eclipse (1) (<http://www.codebelief.com/tag/eclipse/>)



Copyright © 2016-2018 · 代码信条 - Code Belief (<http://www.codebelief.com/>) · All rights reserved.



黑公网安备 23010302000341号

