



13  
15  
16  
17  
18  
19 print samp

更多蓄水池抽样算

1 蓄水池抽样——

2 Reservoir Samp

3 几个随机算法 ht

最后转载下一篇蓄

数据工程师必知算

引言：众所周知，想

来考察面试者的功底。

“给出一个数据流，这

当面对这样一个问题的

器学习模型也不再适合

第二件要做的事情是：

题。所以如果像下面一

“我会首先将输入存到

下？）

第三件要做的事情是：

解，给你灵感。

蓄水池算法

如前面所说，对这个问

难一点的情况：假设数

我们读到了第一个数据

就可以满足题目要求。

接着我们继续分析有三

1和2中的一个。应该

中，如果返回数据3的

么数据1被最终留下的

- 数据1被留下
- 数据2被留下
- 数据3被留下

这个方法可以满足题目

因此，我们做一下推证

证明：假设n-1时候成

设成立。

这就是所谓的蓄水池抽

布式蓄水池抽样和加

（注：Cloudera ML

蓄水池抽样在C

分布式蓄水池抽样是C

们使用mapreduce 模

抽样的时候非常管用。

种分类的组合进行抽样

热门文章

WEKA使用教程(经典教程转载)  
阅读量：180774

概率语言模型及其变形系列(1)-PLSA及EM  
算法  
阅读量：68756

概率语言模型及其变形系列(2)-LDA及  
Gibbs Sampling  
阅读量：62173

数据挖掘-决策树ID3分类算法的C++实现  
阅读量：53544

Stanford概率图模型 ( Probabilistic  
Graphical Model ) — 第一讲 贝叶斯网络  
阅读量：50998

分类

HarryNLPIR

已关注

ACM-分类与图论

12篇

ACM-动态规划

41篇

ACM-搜索DFS/BFS

242

粉丝

2626

喜欢

50

评论

772

ACM-数学题

22篇

ACM-Hash

4篇

等级： 博客 访问：149万+

积分：1万+ 周排名：1400

勋章： 恒

最新评论

我的微博

weixin\_44315211：您好，我的信箱是  
a22352496@gmail.com，文件可以發我一份嗎，  
學習一下，謝謝！

Gradient Tree Boo...  
weixin\_41093521：误人子弟

概率语言模型及其变形系列(5)-L...  
qq\_40304158：[reply]haungzhuwei[reply] 我  
也是遇到同样的问题，请问你们解决了么

概率语言模型及其变形系列(5)-L...  
qq\_40304158：楼主请问一下迭代次数一定得100  
么

几句话弄清楚Java参数传值还是传...  
olivia77lvy：厉害了，的确一句话就懂。 “当  
Integer 做为方法参数传递进方法内时，对其的赋  
值都会导致 原...

联系我们



微信客服

QQ客服

QQ客服

kefu@csdn.net

客服论坛

400-660-0108

工作时间 8:30-22:00

归档

关于我们 招聘 广告服务 网站地图

2017年12月 1篇

百度提供站内搜索 京ICP证09002463号

2019年8月19 江苏乐知网络技术有限公司 1篇

江苏知之为计算机有限公司 北京创新乐知 1篇

2017年3月 信息技术有限公司版权所有

2016年10月 3篇

网络110报警服务 经营性网站备案信息

2016年7月 3篇

北京互联网违法和不良信息举报中心

中国互联网举报中心 展开

in the reservoir14

# with a decreasing

0 and index (inclusive)

0 dex)

= line

et/huangong\_adu/article/details/7619665

nts http://gregable.com/2007/10/reservoir-sampling.html

n-algorithm.html

载如下。

工程师——是件很难的事情。我在面试中常使用的方法是：提出即需要算法设计，又需要一些概率论

中数据只能访问一次。请写出一个随机选择算法，使得数据流中所有数据被选中的概率相等。”

并没有玩你，相反他可能特别想雇你。他可能正在为无尽的分析请求烦恼，他的ETL流水线已经不在口

你答出来。

的面试官读过Daniel Tunkelang的关于数据工程师的面试建议，那么这个面试题很可能就是他工作中

取结束之后随机选取一个”（大哥，你没看见题目已经说了，数据流长度很大或者未知么，不怕你的

体问题（而不是抽象问题），最开始你设计的小例子可能和最后的问题之间相去甚远，但是却能启发

个数据。我们接收数据，发现数据流结束了，直接返回该数据，该数据返回的概率为1。看来很简单，

有结束。我们继续读取第二个数据，发现数据流结束了。因此我们只要保证以相同的概率返回第一个

我们就返回第一个数据，如果R大于0.5，返回第二个数据。

给流中的数据命名为1、2、3。我们陆续收到了数据1、2和前面的例子一样，我们只能保存一个数据

之一的概率淘汰一个，例如我们淘汰了2。继续读取流中的数据3，发现数据流结束了，我们知道在长

就是说，目前我们手里有1,3两个数据，我们通过一次随机选择，以1/3的概率留下数据3，以2/3的概

概率留下该数据，否则留下前n-1个数据中的一个。以这种方法选择，所有数据流中数据被选择的概

正在读取第n个数据，以1/n的概率返回它。那么前n-1个数据中数据被返回的概率为：(1/(n-1))\*((n-

。你可以在[这里](#)找到Greg写的关于蓄水池抽样的算法介绍。本文后面会介绍一下在Cloudera ML中使

见，基本的蓄水池抽样要求对数据流进行顺序读取。要进行容量为k的分布式蓄水池抽样（前面讨论的

组合中的每一个元素，都产生一个0-1的随机数，之后选取随机值最大的前k个元素。这种方法在对大数

龄，地理信息等组合。注意如果输入的数据集分布极端的不均匀，那么抽样可能不能覆盖到所有的

作纯文本或者hive中的表。

https://blog.csdn.net/yanqliuy/article/details/43924581

2/7



Stanford概率图模型 ( Probabilistic Graphical Model ) — 第一讲 贝叶斯网络  
阅读量：50998

，它们有时对于大规模数据分析和一些统计问题也特别有帮助。接下来，我推荐几篇关于算法博客《术》上面的算法吧~



白马负金羁  
367篇文章

中国互联网举报中心 [展开](#)

天涯泪小武  
198篇文章

-柚子皮-  
607篇文章

<div><div>关注</div><div>排名:227</div></div> <div><div>热门文章</div><div>WEKA使用教程(经典教程转载) 阅读量：180774</div><div>概率语言模型及其变形系列(1)-PLSA及EM算法 阅读量：68756</div><div>概率语言模型及其变形系列(2)-LDA及Gibbs Sampling 阅读量：62173</div><div>数据挖掘-决策树ID3分类算法的C++实现 阅读量：53544</div><div>Stanford概率图模型 ( Probabilistic Graphical Model ) — 第一讲 贝叶斯网络 阅读量：50998</div></div> <div><div>解决等概率随机抽样的问题</div><div>最近在CSDN技术社区</div></div> <div><div>Apache Flink 是什么？</div><div>Apache Flink是一个面</div></div> <div><div>五大开源Web代理软件</div><div>Web代理软件转发HT</div></div> <div><div>负载均衡器部署方法</div><div>在现阶段企业网中，只</div></div> <div><div>Flink架构、原理</div><div>Apache Flink是一个面</div></div> <div><div>海量数据随机抽样</div><div>随机抽样问题表示如</div></div> <div><div>亚线性算法-水库</div><div>空间亚线性算法：由</div></div> <div><div>数据工程师必知</div><div>摘要：引言：众所周</div></div> <div><div>蓄水池采样算法</div><div>目录 问题描述分析 算</div></div> <div><div>数据挖掘：大数据</div><div>本文转载自：https://</div></div> <div><div>海量数据常见的面试题</div><div>【在海量数据中统计出</div></div> <div><div>蓄水池抽样问题</div><div>问题描述：要求从N个</div></div> <div><div>面试题：随机数生成</div><div>等概率生成rand5生成</div></div> <div><div>蓄水池抽样 ( Reservoir Sampling )</div><div>在不知道文件总行数</div></div> <div><div>python实现的四</div><div>一、单纯随机抽样 ( s</div></div>	<div><div>关注</div><div>排名:3000+</div></div> <div><div>完全随机分给4人，每人可能得...</div></div> <div><div>平台，它能够基于同一个Flink...</div></div> <div><div>透明代理，而无需客户端配置。...</div></div> <div><div>i设备，较常见是f5、redware...</div></div> <div><div>平台，它能够基于同一个Flink...</div></div> <div><div>N无法确定。这种应用的场景一...</div></div> <div><div>放入内存计算，所以一种常用...</div></div> <div><div>——数据工程师——是件很难的...</div></div> <div><div>亟描述分析 采样问题经常会被遇...</div></div> <div><div>67大数据流 即 实时收集的大...</div></div> <div><div>能够在内存中放下，比如如果...</div></div> <div><div>这种应用的场景一般是数据流的...</div></div> <div><div>杀系统</div></div> <div><div>的[0, 5)范围内的随机整数，要...</div></div> <div><div>首先想到的是我们做过类似的题...</div></div> <div><div>单位编号，再用抽签法或随机...</div></div>	<div><div>已关注</div><div>排名:167</div></div> <div><div>1661</div><div>来自：CWS_chen</div></div> <div><div>525</div><div>来自：u012361418的博客</div></div> <div><div>3330</div><div>来自：DR_eamMer的博客</div></div> <div><div>2578</div><div>来自：BrilliantEagle的专栏</div></div> <div><div>2190</div><div>来自：wtq1993的博客</div></div> <div><div>913</div><div>来自：感悟编程</div></div> <div><div>1.6万</div><div>来自：大数据技术杂谈</div></div> <div><div>6270</div><div>来自：Enweitech Software ...</div></div> <div><div>7034</div><div>来自：淡</div></div> <div><div>1164</div><div>来自：ppf19159的博客</div></div> <div><div>2.5万</div><div>来自：Hackbuteer1的专栏</div></div> <div><div>54</div><div>来自：zhangvalue的博客</div></div> <div><div>894</div><div>来自：z69183787的专栏</div></div> <div><div>2650</div><div>来自：ypfzhao</div></div> <div><div>424</div><div>来自：Miss_畅的博客</div></div> <div><div>978</div><div>来自：碎碎絮語</div></div> <div><div>1183</div><div>来自：yinjunshishui的专栏</div></div> <div><div>1639</div><div>来自：jinzhao1993的博客</div></div> <div><div>339</div><div>来自：MachileYuan的专栏</div></div> <div><div>2.4万</div><div>来自：机器学习</div></div>
---	---	--





<div>android IPC通信</div> <div>android IPC通信 (上</div>	<div>热门文章</div> <div>WEKA使用教程(经典教程转载)</div> <div>阅读量：180774</div> <div>概率语言模型及其变形系列(1)-PLSA及EM算法</div> <div>阅读量：68756</div> <div>概率语言模型及其变形系列(2)-LDA及Gibbs Sampling</div> <div>阅读量：62173</div> <div>数据挖掘-决策树ID3分类算法的C++实现</div> <div>阅读量：53544</div> <div>Stanford概率图模型 ( Probabilistic Graphical Model ) — 第一讲 贝叶斯网络</div> <div>阅读量：50998</div>	<div>IPC通信 ( 中 ) - ContentProv...</div> <div>来自：Shawn_Dut的专栏</div>	<div>4643</div> <div>0</div>
<div>vmware9.0安装</div> <div>最近在vmware9.0上挂</div>		<div>ifs文件夹</div> <div>按照vmware tool，但是却无...</div> <div>来自：tankaro的专栏</div>	<div>6714</div> <div></div>
<div>在centos 中批量</div> <div>单机安装cuda，可以</div>		<div>示一步步的安装。但是当有多...</div> <div>来自：草亦花开的专栏</div>	<div>1502</div> <div></div>
<div>搭建图片服务器</div> <div>nginx是个好东西，N</div>		<div>服务器，也是一个IMAP/POP...</div> <div>来自：maoyuanming0806...</div>	<div>4449</div> <div></div>
<div>Spring Boot M</div> <div>项目地址：https://gi</div>		<div>读写分离</div> <div>DataSource/tree/dev 在 Sprin...</div> <div>来自：HelloWood</div>	<div>18943</div> <div></div>
<div>frp配置本地服务</div> <div>搭建环境：ubuntu 1</div>	<div>分类</div> <div>ACM-分类与搜索查找</div> <div>12篇</div> <div>ACM-动态规划</div> <div>41篇</div> <div>ACM-搜索DFS/BFS</div> <div>22篇</div> <div>ACM-数学题</div> <div>22篇</div> <div>ACM-Hash</div> <div>4篇</div> <div>等级：博客 访问：149万+</div> <div>积分：1万+ 周排名：1400</div> <div>勋章：恒</div> <div>最新评论</div> <div>我的微博</div> <div>程序猿搞定SVM分类-用JAVA...</div> <div>weixin_44315211：您好，我的信箱是a22352496@gmail.com，文件可以發我一份嗎，學習一下，謝謝！</div> <div>Gradient Tree Boo...</div> <div>weixin_41093521：误人子弟</div> <div>概率语言模型及其变形系列(5)-L...</div> <div>qq_40304158：[reply]haungzhuwei[reply] 我也是遇到同样的问题，请问你们解决了么</div> <div>概率语言模型及其变形系列(5)-L...</div> <div>qq_40304158：楼主请问一下迭代次数一定得100么</div> <div>几句话弄清楚Java参数传值还是传...</div> <div>olivia77livy：厉害了，的确一句话就说懂了。”当Integer 做为方法参数传递进方法内时，对其的赋值都会导致 原...</div>	<div>S(阿里云服务器) apache tomc...</div> <div>来自：Anteoy的博客</div>	<div>8290</div> <div></div>
<div>.NET和java的R</div> <div>RSA .net jva 互通 解</div>		<div>emote Forward</div> <div>1172125444948/ 实战 SSH 端...</div> <div>来自：明明</div>	<div>26677</div> <div>来自：lubiaopan的专栏</div>
<div>ODAC (odp.net</div> <div>test</div>		<div>多个倒计时. 查阅网络,基本上...</div> <div>来自：websites</div>	<div>10122</div> <div>来自：我想我是海 冬天的大...</div>
<div>SSH的端口转发:</div> <div>http://zhumeng833</div>		<div>号支付)/JSSDK的使用</div> <div>寸V3微信公众号支付PHP教程/t...</div> <div>来自：Marswill</div>	<div>62453</div> <div></div>
<div>jquery/js实现一</div> <div>jquery/js实现一个网</div>		<div>要弄懂这里面的过程，基本上...</div> <div>来自：文洲的专栏</div>	<div>16396</div> <div></div>
<div>工业相机编程模</div> <div>本文详述常见工业相</div>		<div>配置文件系统的块大小，Name...</div> <div>来自：yycdaizi的专栏</div>	<div>10311</div> <div></div>
<div>Hive小文件合并</div> <div>Hive的后端存储是HC</div>		<div>配置V3微信公众号支付PHP教程/t...</div> <div>来自：Marswill</div>	<div>19220</div> <div></div>
<div>微信支付V3微信</div> <div>扫二维码关注，获取</div>		<div>Docker了，汗汗！ Docker的...</div> <div>来自：我走小路的博客</div>	<div>3739</div> <div></div>
<div>linux上安装Doc</div> <div>最近比较有空，大四</div>		<div>：1、单击压测过程中使用过...</div> <div>来自：测试蜗牛，一步一个...</div>	<div>89059</div> <div></div>
<div>Jmeter:修改内存</div> <div>在压测过程中jmeter</div>	<div>联系我们</div> <div></div> <div>微信客服</div> <div>QQ客服</div> <div>QQ客服</div> <div>kefu@csdn.net</div> <div>客服论坛</div> <div>400-660-0108</div> <div>工作时间 8:30-22:00</div> <div>归档</div> <div>关于我们 招聘 广告服务 网站地图</div> <div>2017年12月</div> <div>百度提供站内搜索 京ICP证09002463号</div> <div>2019年8月19 江苏乐知网络技术有限公司</div> <div>江苏知之为计算机有限公司 北京创新乐知信息技术有限公司版权所有</div> <div>2017年3月</div> <div>2016年10月</div> <div>网络110报警服务 经营性网站备案信息</div> <div>2016年7月</div> <div>北京互联网违法和不良信息举报中心</div> <div>中国互联网举报中心</div> <div>展开</div>	<div>20170620</div> <div>lplus命令，基本的命令都可以...</div> <div>来自：Ape55的博客</div>	<div>13182</div> <div></div>
<div>人脸检测工具fac</div> <div>人脸检测工具face_re</div>		<div>F抗线</div> <div>] = {0,1,0,1,0,1,1,0,1,1,0,0,1,0,...</div> <div>来自：miracle的专栏</div>	<div>13477</div> <div>来自：roguesir的博客</div>
<div>plsql的命令 ( co</div> <div>command窗口是命令</div>		<div>架结构的计算过程，我们都知...</div> <div>来自：AUTO1993的博客</div>	<div>6515</div> <div></div>
<div>寻找连通线,参考</div> <div>#include using nam</div>		<div>7794</div> <div>来自：volkswageos的专栏</div>	<div>19804</div> <div></div>
<div>3D CNN框架结构</div> <div>3D CNN框架结构各层</div>			<div>21321</div> <div></div>
<div>[转]极线几何约束</div>			<div>7794</div> <div></div>

训练网络出现loss

粗略写一下：解决方法：

Two-pass连通域

在Two-pass连通域标

加密算法介绍及加

加密算法介绍 一. 密

JAVA结合testng

原理：1.自己构造一

web.config中的

打开某个应用程序

DirectX修复工具

最后更新：2018-12-

关于SpringBoot

问题场景描述整个项目

内连接，外链接

1.什么是连接查询呢？

强连通分量及缩点

强连通分量：简言之

热门文章

WEKA使用教程(经典教程转载)  
阅读量：180774

概率语言模型及其变形系列(1)-PLSA及EM  
算法  
阅读量：68756

概率语言模型及其变形系列(2)-LDA及  
Gibbs Sampling  
阅读量：62173

数据挖掘-决策树ID3分类算法的C++实现  
阅读量：53544

Stanford概率图模型 ( Probabilistic  
Graphical Model ) — 第一讲 贝叶斯网络  
阅读量：50998

个人分类

ACM-分治与二分查找12篇

ACM-动态规划41篇

ACM-搜索DFS/BFS32篇

ACM-数学题22篇

ACM-Hash4篇

最新评论

3行程序搞定SVM分类-用JAVA...  
weixin\_44315211：您好，我的信箱是  
a22352496@gmail.com，文件可以發我一份嗎，  
學習一下，謝謝！

Gradient Tree Boo...  
weixin\_41093521：误人子弟

概率语言模型及其变形系列(5)-L...  
qq\_40304158：[reply]haungzhuwei[/reply] 我  
也是遇到同样的问题，请问你们解决了么

概率语言模型及其变形系列(5)-L...  
qq\_40304158：楼主请问一下迭代次数一定得100  
么

几句话弄清楚Java参数传值还是传...  
olivia77livy：厉害了，的确一句话就说懂了。“当  
Integer 做为方法参数传递进方法内时，对其的赋  
值都会导致 原...

联系我们



微信客服



QQ客服

 QQ客服

 kefu@csdn.net

 客服论坛

 400-660-0108

工作时间 8:30-22:00

关于我们 招聘 广告服务 网站地图

 百度提供站内搜索 京ICP证09002463号

©1999-2019 江苏乐知网络技术有限公司

江苏知之为计算机有限公司 北京创新乐知  
信息技术有限公司版权所有

网络110报警服务 经营性网站备案信息

北京互联网违法和不良信息举报中心

中国互联网举报中心

态分布，，现在考虑下：http...  
来自：jiachen0212的博客

41860

下扫描，会将各个有效像素置...  
来自：lichengyu的专栏

5106

置换密码。1881年世界上的第...  
来自：leolewin的博客

8232

)  
2.利用testng的监听类在测试...  
来自：nicolas\_li的专栏

2680

< sessionState mode="In...  
来自：yszwn的专栏

6656

8 增强版 NEW! 版本号：V3.8...  
来自：VBcom的专栏

1741116

有关)  
框架一个module spring-boot-...  
来自：开发随笔

38665

接大总结  
从这些表中查询数据。 目的...  
来自：basycai的博客

8083

一个点也是一个连通分量 使用...  
来自：九野的博客

21170

https://blog.csdn.net/yangliuy/article/details/43924581

7/7