



南京邮电大学
Nanjing University of Posts and Telecommunications

联邦学习差分隐私保护方法的隐私预算方法与设计

计算机学院、软件学院、网络空间安全学院

指导老师：xx

报告人：陈瑾瑜

目录

01 课题背景

02 课题研究内容

03 实验结果分析



01 联邦学习

联邦学习 (Federated Learning) 是一种分布式训练方式，它利用分散在各参与方的数据集，通过隐私保护技术整合多方数据信息，共同构建全局模型。该设计旨在保证数据安全以及各类型数据隐私的前提下，实现多参与方或多计算节点之间高效的机器学习。其主要目标是为大数据交换提供信息安全保障，并确保在合法合规的范围内进行操作。其中使用较多的是横向联邦学习，横向联邦学习的数据特征是对齐的，即不同行的数据之间有相同的数据特征。

01

课题背景

02 差分隐私

差分隐私 (Differential Privacy)

出, 系统输出结果

出信, 真实查询结果

无法, 则其位置参数

私保: 法 $M(D) = f(D) + \text{Laplace}$

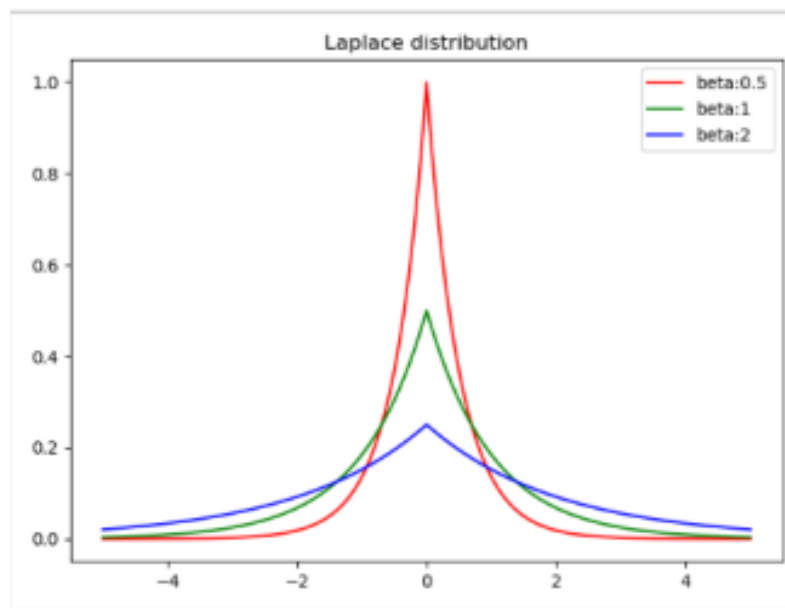


图 2.2 不同参数的 Laplace 分布图

2008年由Dwork 提

出, 系统输出结果

差分保护的方式是向

$\frac{|x|}{b}$

度为 Δf , 则随机算

在输

三方

分隐



03 聚类

聚类分析是按照某一标准将相对于其他组（聚类）中的对象而言相似度更高的对象归入同一个组（聚类）的过程。在聚类分析过程中，相似度计算所使用的评判标准与参数选择对聚类效果的衡量指标极为重要。聚类分析被视为探索性数据分析的一种方法，它能够发现数据中隐藏的模式、结构，并揭示数据中有趣的未知关系。



03 K-means算法

K-means聚类算法是目前最常用的聚类算法，其主要思想是：给定K个初始类簇中心点和K值后，分配每个数据点到距其最近的类簇中心点所代表的类簇中，将所有数据点都分配完毕后，重新计算该类簇的中心点(取平均值)，然后将分配点和更新类簇中心点的流程不断迭代，直至类簇中心点的变化很小或达到指定的迭代次数。

04 级数



设有无穷个可列实数 $x_1, x_2, \dots, x_n, \dots$ ，它们的“和”

$$x_1 + x_2 + \dots + x_n + \dots$$

记为 $\sum_{n=1}^{\infty} x_n$ ，称为无穷数项级数（也叫级数），其中 x_n 是级数的一般项或者叫做通项。

级数 $\sum_{n=1}^{\infty} x_n$ 的部分和数列 $\{S_n\}$ 的每一项定义如下所示：

$$S_1 = x_1,$$

$$S_2 = x_1 + x_2,$$

$$S_3 = x_1 + x_2 + x_3,$$

.....

$$S_n = x_1 + x_2 + \dots + x_n = \sum_{k=1}^n x_k,$$

.....

若部分和数列 $\{S_n\}$ 收敛于一个有限的值 S ，则称无穷级数 $\sum_{n=1}^{\infty} x_n$ 收敛。对于收敛的无穷级数，我们记它的和为 S ，其中

$$S = \sum_{n=1}^{\infty} x_n;$$

我们称无穷级数 $\sum_{n=1}^{\infty} x_n$ 发散，当且仅当部分和数列 $\{S_n\}$ 发散。

(1) 等比级数：

$$\sum_{n=1}^{\infty} x_n = \sum_{n=1}^{\infty} q^{n-1} = 1 + q + q^2 + \dots + q^n + \dots$$

当 $|q| < 1$ 时，该级数收敛。我们可以得出它的部分和数列通项为

$$S_n = \sum_{k=1}^n q^{k-1} = \frac{1-q^n}{1-q},$$

记它的收敛值为 S ，显然：

$$S = \lim_{n \rightarrow \infty} S_n = \frac{1}{1-q}.$$

(2) 交错级数：

$$\sum_{n=1}^{\infty} x_n = \sum_{n=1}^{\infty} (-1)^{n+1} u_n = u_1 - u_2 + u_3 + \dots + (-1)^{n+1} u_n + \dots \quad (u_n > 0)$$

我们将该级数称为 Leibniz（莱布尼茨）级数，当且仅当级数项 $\{u_n\}$ 单调递减且收敛于 0。

(3) p 级数：

$$\sum_{n=1}^{\infty} x_n = \sum_{n=1}^{\infty} \frac{1}{n^p} = 1 + \frac{1}{2^p} + \frac{1}{3^p} + \dots + \frac{1}{n^p} + \dots \quad (p > 0)$$

当 $p > 1$ 时，级数收敛；当 $0 < p \leq 1$ 时，级数发散至无穷大。当 $p = 1$ 时，

我们称 $\sum_{n=1}^{\infty} \frac{1}{n}$ 为调和级数。



01 数据集的读取

数据集的读取直接采用sklearn包中的数据集，我们只需要从包中导入即可，导入后作图将数据集表示出来。

```
# 直接从sklearn中获取数据集
iris = datasets.load_iris()
X = iris.data[:, :4]    # 表示我们取特征空间中的4个维度
print(X.shape)

# 取前两个维度（萼片长度、萼片宽度），绘制数据分布图
plt.scatter(X[:, 0], X[:, 1], c="red", marker='o', label='see')
plt.xlabel('sepal length')
plt.ylabel('sepal width')
plt.legend(loc=2)
plt.show()
```




02 DPK-means算法的实现

本次毕业设计中实现的DPK-means算法实际上是以K-means算法为基础，在每一轮迭代时得到当前各个簇的质心后，给各个簇的质心加入随机噪声，该噪声符合拉普拉斯机制。当迭代次数达到阈值或聚类中心稳定时，迭代结束，并得到最终的聚类结果。

02 DPK-means算法的实现

DPK-means 工作流程：

- ①随机确定K个任意位置初始点作为中心。
- ②将数据集中的各点分配到一个簇中：具体就是找到每个点距离最近的中心, 并将该点分配到该中心所对应的簇。完成分配后, 根据该簇中所有点的平均值更新该簇的中心。更新完簇的中心后, 将更新后的簇中心值加上拉普拉斯噪声得到新的簇中心。
- ③重复上述流程, 若数据集中的各点都距离它所对应中心最近, 则算法结束。



03 添加服从Laplace机制的随机噪声

在完成DPK-means算法的过程中，我们需要加入Laplace噪声，接下来对添加Laplace噪声过程进行介绍。拉普拉斯分布的概率密度函数是：

$$f(x|\mu, b) = \frac{1}{2b} e^{\frac{-|x-\mu|}{b}}$$

为便于计算，我们往往令 $\mu=0$ ，即偏移量为0

拉普拉斯机制：给定数据集 D ，设有函数 $f: D \rightarrow$

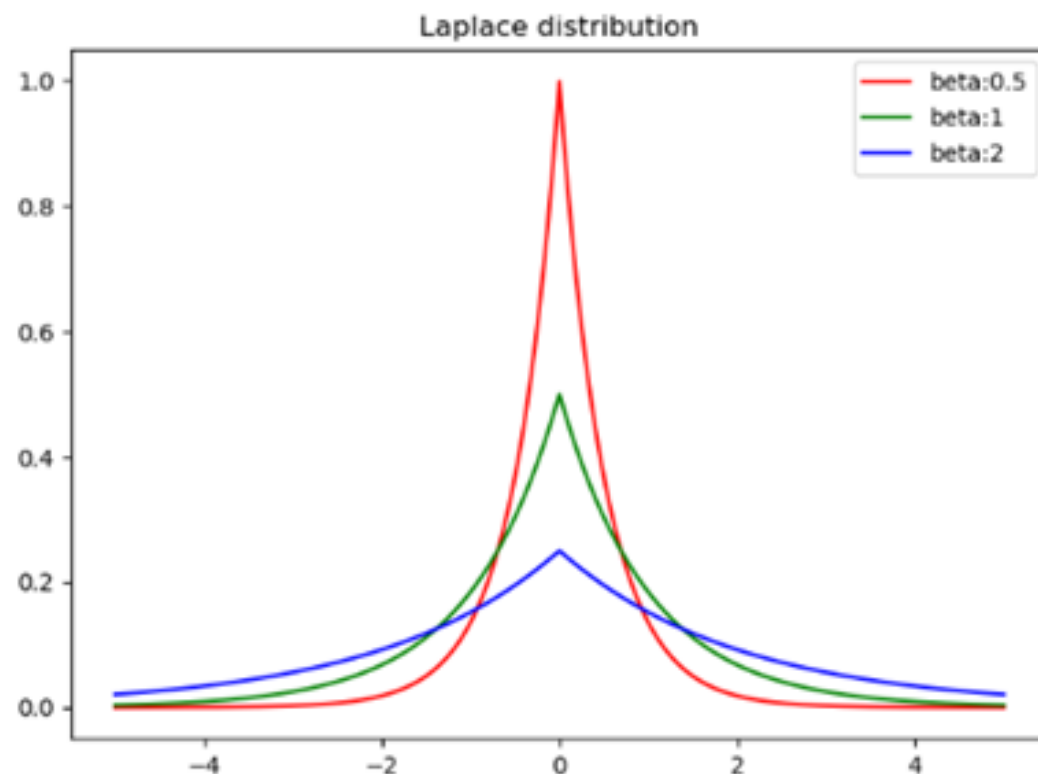
R^d ，敏感度

差分隐私保

$Y \sim L(0, \frac{\Delta f}{\epsilon})$,

出噪声值的

噪声分布如



+ Y 提供 ϵ -
噪声，

，可以看
算成反比，

03 添加服从Laplace机制的随机噪声

之后调用`np.random.laplace`函数，将生成的拉普拉斯随机噪声加入到簇中心中。

```
76     # 4.更新质心
77     for j in range(k):
78         pointsInCluster = dataSet[np.nonzero(clusterAssment[:, 0].A == j)[0]] # 获取对应簇类所有的点
79         centroids[j, :] = np.mean(pointsInCluster, axis=0) # 求均值，产生新的质心
80         # axis=0, 求的是pointsInCluster每一列的平均值，即axis是几，那就表明哪一维度被压缩成1
81     # 5.加噪声
82     epsilon = epsilon / 2
83     for i in range(len(centroids)):
84         noise = np.random.laplace(0, 1.0 / epsilon)
85         centroids[i] += 5e-3 * noise
86
87     print("cluster complete")
88     return centroids, clusterAssment
89
```





在此次毕业设计中，我们将会使用几种不同的隐私预算分配方法实现DPK-means算法。具体而言，只需要将每次迭代中添加的噪声的参数 ϵ 进行修改即可。如下为二分法、2级数法的 ϵ 参数。

```
# 5.加噪声  
  
epsilon = epsilon / 2  
for i in range(len(centroids)):  
    noise = np.random.laplace(0, 1.0 / epsilon)
```

```
# 5.加噪声  
  
epsilont = (6/math.pow(math.pi, 2)) * (epsilon/math.pow(t, 2))  
for i in range(len(centroids)):  
    noise = np.random.laplace(0, 1.0 / epsilont)  
    centroids[i] += 1e-3 * noise  
  
t = t+1
```



04 总体流程

1. 读取数据集
2. 实现K-means算法，并用数据集测试
3. 实现DPK-means算法，并用数据集测试
4. 在DPK-means算法中用不同的方法进行隐私预算的分配
5. 得到实验结果，并进行分析



03

实验结果分析

数据集1

测量数据包括：萼片长度、萼片宽度、花瓣长度、花瓣宽度。类别共分为三类：Iris Setosa,Iris Versicolour,Iris Virginica。标签0、1、2分别表示山鸢尾（Setosa）、变色鸢尾（Versicolor）、维吉尼亚鸢尾（Virginical），每个种类的鸢尾花有50份数据，共150份数据。

数据集2

样本数据个数	442
特征个数（数据维度）	10
特征意义	年龄、性别、BMI指数、平均血压、S1、S2、S3、S4、S5、S6
特征取值范围	-0.2至0.2

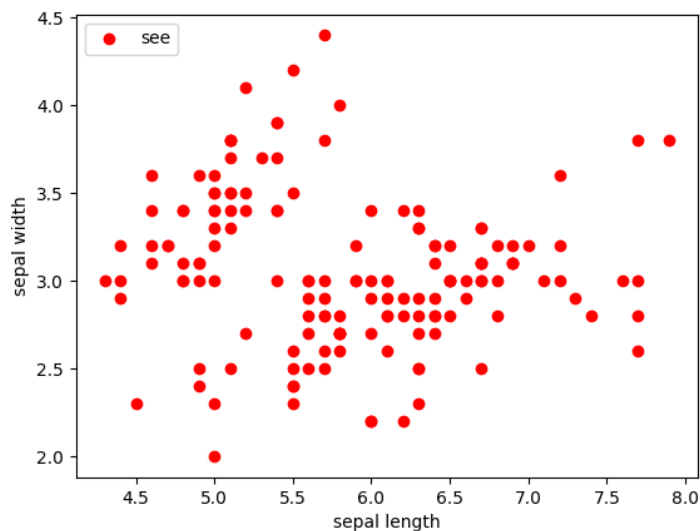
03

实验结果

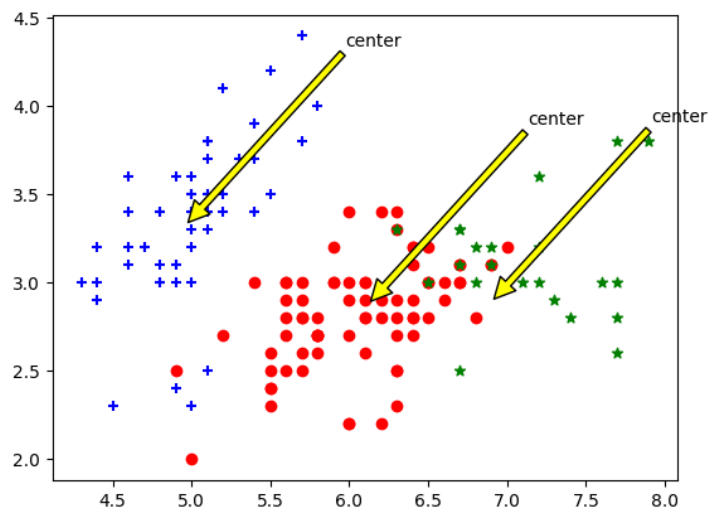


实验1 用不同的隐私预算分配方式对隐私预算进行分配，得出DPK-means聚类结果

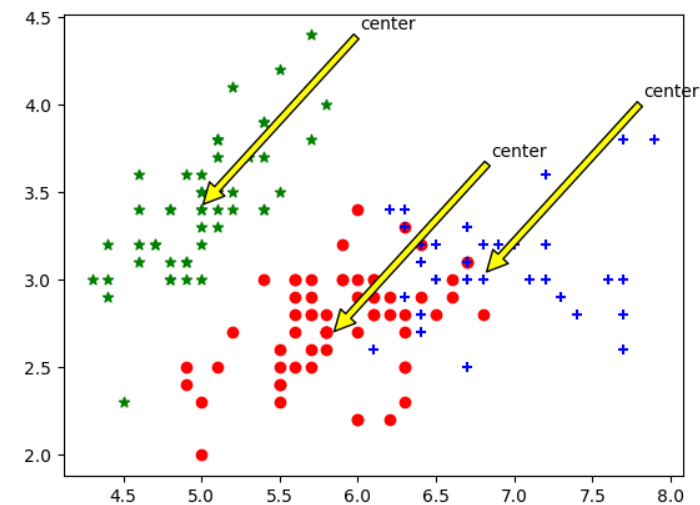
初始数据集



DPK-means算法结果（二分法）



DPK-means算法结果（二级数法）





对鸢尾花数据集的实验结果进行分析。首先我们计算噪声评估因子，根据噪声评估因子的复杂度我们可以得出，特殊级数法与 p 级数法的隐私预算分配方式更为优秀。接着我们使用F-measure值对结果进行分析和评估，具体的计算结果显示，二分法、特殊级数法、2级数法的F-measure值在DPK-means算法的迭代次数小于等于10时比较接近；而当迭代的次数大于等于20时，特殊级数法和2级数法仍然保持了较高的F-measure值，此时二分法的F-measure值很小，也就是说，特殊级数法和2级数法更好地确保了聚类结果的可用性。



本次报告在K-means算法的基础上实现了DPK-means算法，并采用了多种隐私预算的分配方式对隐私预算进行分配。同时，我们采用了多个数据集对算法进行实验测试，并用噪声评估因子和F-measure值对实验结果进行分析。最终我们得出一个结论：特殊级数法和2级数法（p级数法）的隐私预算分配方式较为优秀。



南京邮电大学
Nanjing University of Posts and Telecommunications

感谢老师 请老师批评指正！

报告人：陈瑾瑜