

1. 准备数据

首先，确保你有一个包含文本数据的文件，这些数据需要被导入到 Hive 中进行处理。假设我们有一个文本文件 `input.txt` 包含需要处理的文本数据。

2. 创建外部表

在 Hive 中，我们可以创建一个外部表来指向我们的文本文件，然后使用该表进行数据处理。这里我们假设表的字段为 `word` 和 `count`。

```
sql
CREATE EXTERNAL TABLE IF NOT EXISTS word_count (
    word STRING,
    count INT
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
LOCATION '/path/to/your/input';
```

- `word`: 单词列，存储单词字符串。
- `count`: 整数列，存储单词的词频数。

`LOCATION` 指定了文本文件所在的路径。确保 Hive 能够访问该路径。

3. 加载数据到表中

一旦表创建好了，我们需要将数据加载到这个表中。假设数据文件 `input.txt` 已经在指定的路径 `/path/to/your/input` 下。

```
sql
LOAD DATA INPATH '/path/to/your/input/input.txt' OVERWRITE INTO TABLE word_count;
```

4. 执行统计查询

现在，我们可以执行 Hive 查询来统计每个单词的词频总和。

```
sql
SELECT word, SUM(count) AS total_count
FROM word_count
GROUP BY word
ORDER BY total_count DESC;
```

这条查询将对 `word_count` 表中的数据按单词进行分组，并计算每个单词的词频总和，然后按词频总和降序排序输出结果。