



## 建模锯齿是什么意思

建模锯齿是指在建模过程中的一种常见的效应，其中模型的输出可能会产生锯齿状的波动。这种效应通常是由于模型的不稳定性或过度拟合导致的。

在建模过程中，锯齿可能会使模型的表现变差，并且在预测新数据时也可能出现不一致的结果。因此，在建模时需要注意避免出现锯齿状的波动。

一种常用的方法是使用正则化来限制模型的复杂度，以避免过度拟合。此外，还可以使用交叉验证来评估模型的稳定性并选择最优的模型。



## 数据粒度的锯齿操作

数据粒度的锯齿操作，也称为数据锯齿化操作，是将数据从一种粒度转换到另一种粒度的过程。例如，将日数据转换为月数据或将小时数据转换为天数据。

这种操作常常用于数据分析和可视化，以便获得不同粒度下数据的视图。例如，将数据从每小时转换为每天可能会更容易看出数据的趋势或周期性。

在实现数据锯齿化操作时，通常需要指定转换的粒度，以及如何计算转换后数据的值。例如，将日数据转换为月数据时，可能需要计算每月的平均值或总和。

有关数据锯齿化操作的更多信息，可以查看相关的文献或使用数据分析工具（如 Excel 或 R）来实现此操作。



## 上游表删除，表结构尽失，物理删除，下游表依赖上游表，这时候怎么保证下游任务正常产出。

如果你的下游表依赖于上游表的结构，那么在上游表被物理删除之后，下游表就可能无法正常工作了。在这种情况下，你需要找出导致下游表依赖上游表的原因，并考虑重新设计表结构，使得下游表不再依赖于上游表。

具体来说，你可以考虑以下几种方法：

在下游表中缓存上游表的数据：如果你的下游表依赖于上游表的数据，你可以考虑在下游表中缓存这些数据。这样，即使上游表被删除，下游表也可以继续使用缓存的数据。

建立视图：如果你的下游表依赖于上游表的结构，你可以考虑在下游表中建立一个视图，视图可以看作是一张虚表，它可以让你像查询一张表一样查询视图。这样，即使上游表被删除，下游表也可以通过视图访问上游表的数据。

使用数据库触发器：如果你的下游表依赖于上游表的数据，你可以考虑在上游表中使用数据库触发器。



数仓中没有指标可不可以建模，没有指标可不可以建宽表。

在数仓中建模和建立宽表时，可以使用没有指标的原始数据。实际上，您可以使用任何可用的数据来建模或建立宽表，只要该数据有意义并且可以用于您的分析目的即可。

然而，指标通常是非常有用的，因为它们可以汇总和概括原始数据，使其更易于理解和分析。如果可能的话，建议您使用合适的指标来帮助更好地理解数据。



数仓数据域划分方式

在数据仓库中，数据域划分是指将数据仓库中的数据分割成若干个不同的部分，以便于数据管理和分析。常见的数据域划分方式包括：

说法1

- 1、按业务线划分：将数据按照业务线来划分，比如销售业务、金融业务等。
- 2、按数据源划分：将数据按照数据源来划分，比如来自不同的系统或数据库。
- 3、按数据类型划分：将数据按照数据类型来划分，比如事实数据、指标数据等。

4、按数据时间划分：将数据按照数据的时间来划分，比如近期数据、历史数据等。

## 说法2

1、按照业务或业务过程划分:比如一个靠销售广告位置的门户网站主题域可能会有广告域客户域等，而广告域可能就会有广告的库存，销售分析、内部投放分析等主题

2、根据需求方划分:比如需求方为财务部，就可以设定对应的财务主题域，而财务主题域里面可能就会有员工工资分析，投资回报比分析等主题

3、按照功能或应用划分:比如微信中的朋友圈数据域、群聊数据域等，而朋友圈数据域可能就会有用户动态信息主题、广告主题等:

4、按照部门划分:比如可能会有运营域、技术域等，运营域中可能会有工资支出分析、活动宣传效果分析等主题



## 数仓一致性是如何保证的

数仓一致性是指多个数仓系统之间的数据保持一致性。这通常是通过使用两阶段提交（2PC）或多版本并发控制（MVCC）来实现的。

两阶段提交是一种分布式事务处理协议，用于在多个数据库系统之间保证事务的一致性。它包括两个阶段：第一阶段是预提交阶段，所有参与者都会把事务的更改提交到本地数据库中，但是这些更改并不会对其他系统可见。在第二阶段，所有参与者都会把事务提交到其他系统中，从而使得所有系统的数据保持一致。

多版本并发控制是另一种用于维护数仓一致性的方法。它通过在每个事务中保存一个版本号来实现。当一个事务开始时，它会读取当前的版本号，然后在更新数据时将版本号加 1。这样，如果另一个事务试图在相同的数据上进行更新，它会检测到版本号不匹配，从而防止冲突。



## 数仓优化需要考虑哪些方面？趋向于重构，应该考虑哪些方面

数仓优化可以从以下几个方面考虑：

- 1、数据模型：对于数仓的数据模型，可以考虑进行重构，以提高数据的质量和可用性。
- 2、数据流：可以考虑优化数据的加载、处理和存储流程，以提高数仓的性能。
- 3、计算和存储：可以考虑使用更高效的计算和存储技术，比如使用内存数据库或图数据库。
- 4、数据安全：应该考虑数据安全性，确保数仓中的数据不会被未经授权的访问。
- 5、可维护性：应该考虑如何使数仓的维护和管理变得更加容易。

在进行数仓重构时，需要考虑以上所有方面，以确保重构能够真正提升数仓的性能和可用性。



数据的全生命周期指的是数据从出生到死亡的整个过程。这个过程包括数据的收集、存储、处理、分析、可视化和清理等环节。

在数据的收集阶段，可以通过各种方式收集数据，如通过传感器、日志文件、Web 表单、社交媒体等获取数据。在存储阶段，数据通常会被存储到数据库、数据仓库或云端存储中。

在处理阶段，数据可能需要进行清洗、格式化、转换等操作，以便进行分析和可视化。在分析阶段，可以使用各种工具和技术，如机器学习、统计学、数据挖掘等，来对数据进行分析。

在可视化阶段，可以使用各种图表、地图和其他工具来展示数据，使其更加直观和易于理解。

最后，在清理阶段，可以将无效或过期的数据从数据库中删除。

数据的全生命周期是一个持续的过程，并且在整个过程中都需要进行合理的管理和控制，以保证数据的准确性和可靠性。

面试整理PDF可以关注球球，加我微信获取。