

### 7.1.3 HBase 列式存储的数据模型

HBase 根据列族来存储数据，一个列族对应物理存储上的一个 HFile，列族包含多列，列族在创建表的时候被指定。图 7-4 简单对比了关系型数据库 RDBMS 表和 HBase 表的关系。

RDBMS表				VS	HBase表		
Primary Key	列1	列2	列3		RowKey	列族1	列族2
数据1	XX	XX	XX		数据1	列1, 列2.....列n	列1
数据2	XX	XX	XX		数据2	列1, 列2	列1, 列2
数据3	XX	XX	XX		数据3	列1, 列2, 列3	列2

图 7-4 RDBMS 和 HBase 的数据模型对比

#### 1. Column Family

Column Family 即列族，HBase 基于列族划分数据的物理存储，一个列族可以包含任意多列。HBase 在创建表的时候就必须指定列族。HBase 的列族不是越多越好，官方推荐一个表的列族数量最好小于或者等于 3，过多的列族不利于 HBase 数据的管理和索引。

HBase 的列族在表结构上与关系型数据库中的列类似，但是两者是完全不同的概念。HBase 基于列族来完成列式数据的存储，而关系型数据库基于行来完成数据的存储，列只是一种数据结构上的表示。

#### 2. RowKey

RowKey 的概念与关系型数据库中的主键相似，HBase 使用 RowKey 来唯一标识某行的数据。HBase 只支持 3 种查询方式：基于 RowKey 的单行查询、基于 RowKey 的范围查询和全表查询。

#### 3. Region

Region 的概念与关系型数据库表的横向分区相似（比如 MySQL 根据 id 的一致性 Hash 值将数据存储在不同数据库中）。HBase 将表中的数据基于 RowKey 的不同范围划分到不同 Region 上，每个 Region 都负责一定范围的数据存储和访问。

HBase 数据分区的过程与其他数据库的 Shard 类似，这样即使有一个包括上百亿条数据的表，由于数据被划分到不同的 Region 上，每个 Region 都可以独立地进行写入和查