

## Bucketing Table

- Create bucketed table -
  - **create table cases\_bucket(Case\_Id int, Hospital\_Code int, Patient\_Id int, Ward\_Type char(1), Ward\_Facility\_Code char(1), Bed\_Grade int, Type\_of\_Admission string, Visitors\_with\_Patient int, Age\_Group string, Admission\_Deposit float, Stay string )**  
**partitioned by ( Department string, Severity\_of\_Illness string )**  
**clustered by (patient\_id) into 10 buckets**  
**row format delimited**  
**fields terminated by '\$'**  
**TBLPROPERTIES ('serialization.null.format'='');**
  - We are creating 10 buckets based on the patient\_id within each partition.
- Put file to HDFS -
  - **hdfs dfs -put /home/itv180149/Hive/Partitions/Hospital\_Data/case.txt /user/itv180149**
- Load data into Hive table -
  - **load data inpath "/user/itv180149/case.txt" into table cases\_bucket;**
  - Notice that since we are bucketing the data into 10 buckets, Hive is using 10 reducers to accomplish this job.
- Let's check the table description -
  - **describe formatted cases\_bucket;**
  - Notice the partitions and buckets here.

- We have 10 buckets or files within each partition -
  - **hdfs dfs -ls**  
**/user/itv180149/warehouse/patient\_db.db/cases\_bucket/department=anesthesia/severity\_of\_illness=Extreme**
- So if we run queries on the partitioned and bucketed columns, then those queries will run much faster.
  - Check the records where the patient id is 306180 and department is gynecology and severity of illness is Minor -
    - **select hospital\_code, department, severity\_of\_illness from cases\_bucket where patient\_id=306180 and department='gynecology' and severity\_of\_illness='Minor';**