# Cheat sheets

Monday, 13 February, 2023    4:39 PM

| No | Objective | Language | Script |
|----|-----------|----------|--------|
| 1 | *Get last month end date* | SQL | select left(dateadd(day, -1, to_date(left(to_char(GETDATE(),'YYYY-MM-DD'),8)+'01', 'YYYY-MM-DD')),10) as lst_mth_date<br><br>select left(dateadd(day, -1, to_date(left(to_char(DATEADD(Month,-1,GETDATE()),'YYYY-MM-DD'),8)+'01', 'YYYY-MM-DD')),10) as lst_2mth_date |
| 2 | *Create new columns based on the values in the column & assign 1* | Pyspark | for p_id in pack_id:<br>    df = pp_recharges.withColumn(p_id, when(col('l9_package_id')==p_id, 1).otherwise(0))<br>    pp_recharges = df |
| 3 | *Group by without assigning the columns name (need to have the list of columns)* | Pyspark | exprs = {x: "max" for x in pack_id}<br>df1 = df.groupBy("account_no").agg(exprs) |
| 4 | *Find the possible combinations from the list (need to have the list of elements)* | Pyspark | from itertools import combinations<br>sample_list = ['a','b','c']<br>list_combinations = list()<br>for n in range(len(pack_id)+1):<br>    list_combinations += list(combinations(pack_id,2)) |
| 5 | *Get the group by and the aggregation based on the columns of the dataframe* | | def view_dur_pivot(prefix):<br>    df_ivpevo_sum = df_ivpevo.filter(col('content_pillar_type') == prefix).groupby('account_no') \\<br>    .pivot('month_key') \\<br>    .agg(round(sum('duration_in_sec')/60, 2)) \\<br>    .na.fill(0)<br>    df_ivpevo_sum = df_ivpevo_sum.select([f.col(c).alias(prefix + "_" + "view_dur_" + c) \\<br>            if c not in {'account_no'} else 'account_no' for c in df_ivpevo_sum.columns])<br>    return df_ivpevo_sum<br>def view_tag_pivot(prefix):<br>    df_ivpevo_sum = df_ivpevo.filter(col('content_pillar_type') == prefix).groupby('account_no') \\<br>    .pivot('month_key') \\<br>    .agg(max('view_cnt')) \\<br>    .na.fill(0)<br>    df_ivpevo_sum = df_ivpevo_sum.select([f.col(c).alias(prefix + "_" + "view_tag_" + c) \\<br>            if c not in {'account_no'} else 'account_no' for c in df_ivpevo_sum.columns])<br>    return df_ivpevo_sum<br>def view_sum_pivot(prefix):<br>    df_ivpevo_sum = df_ivpevo.filter(col('content_pillar_type') == prefix).groupby('account_no') \\<br>    .pivot('month_key') \\<br>    .agg(sum('view_cnt')) \\<br>    .na.fill(0)<br>    df_ivpevo_sum = df_ivpevo_sum.select([f.col(c).alias(prefix + "_" + "view_cnt_" + c) \\<br>            if c not in {'account_no'} else 'account_no' for c in df_ivpevo_sum.columns])<br>    return df_ivpevo_sum |
| 6 | *Get a list of date range* | Pyspark | months = pd.date_range(start='07-01-2022', periods=6, freq='M').strftime('%Y%m').tolist() |
| 7 | *Save the file into s3* | Pyspark | df.repartition(1).write.csv(s3_path + file_name, header="true", mode="overwrite")<br>df.write.mode("overwrite").partitionBy("month_key").format('orc').save(s3_path + file_name, header=True)<br>df.repartition(1).write.parquet(s3_path+'.parquet') |
| 8 | *Get the week from date* | SQL | concat ('WK', lpad (date_part (w, cast ('2022-01-22' as date)), 2, '0')) as week |
| 9 | *Pivot* | Pyspark | cols = engage.drop(col('viewer_id')).columns<br>pivot = engage.groupBy(cols).agg(count('viewer_id').alias('cust_cnt')) |
| 10 | *Assign date to variable* | Pyspark | month_ref = '202204'<br>last_3 = (datetime.strptime(month_ref, '%Y%m') - relativedelta(months=6)).strftime('%Y-%m-%d')<br>end = (datetime.strptime(month_ref, '%Y%m') - relativedelta(days=1)).strftime('%Y-%m-%d') |
| 11 | *Append new column from another df* | Pyspark | def append_dfs(df1,df2):<br>        list1 = df1.columns<br>        list2 = df2.columns<br>        For col in list2:<br>            If(col not in list1):<br>                Df1 = df1.withColumn(col, f.lit(None))<br>        For col in list1:<br>            If(col not in list2): |

| | | | | Df2 = df2.withColumn(col, f.lit(None)) Return df1.unionByName(df2) |
| --- | --- | --- | --- | --- |
| | | | | |