

fn match_num_groups_predicate

Michael Shoemate

April 22, 2025

This proof resides in “**contrib**” because it has not completed the vetting process.

Proves soundness of `match_num_groups_predicate` in `mod.rs` at commit `f5bb719` (outdated¹).

1 Hoare Triple

Precondition

None

Function

```
1 def match_per_group_predicate(  
2     enumeration: Expr,  
3     partition_by: Vec[Expr],  
4     identifier: Expr,  
5     threshold: u32,  
6 ) -> Optional[Bound]:  
7     # reorderings of an enumeration are still enumerations  
8     if isinstance(enumeration, Expr.Function) and isinstance( #  
9         enumeration.options.collect_groups, ApplyOptions.GroupWise  
10    ):  
11         input = enumeration.input  
12         function = enumeration.function  
13  
14         # FunctionExprs that may reorder data  
15         if function == FunctionExpr.Reverse:  
16             is_reorder = True  
17         elif isinstance(function, FunctionExpr.Random):  
18             method = str(function.method)  
19             is_reorder = method == "shuffle"  
20         else:  
21             is_reorder = False  
22  
23         if is_reorder:  
24             enumeration = input[0]  
25  
26     elif isinstance(enumeration, Expr.SortBy):  
27         for key in enumeration.by:  
28             check_infallible(key, Resize.Ban)  
29         enumeration = enumeration.expr#  
30  
31     # in Rust, the != results in a boolean comparison, not a "ne" expression
```

¹See new changes with `git diff f5bb719..12e5590c rust/src/transformations/make_stable_lazyframe/truncate/matching/mod.rs`

```

32     if enumeration != int_range(lit(0), len(partition_by), 1, DataType.Int64): #
33         return None
34
35     # we now know this is a per group predicate,
36     # and can raise more informative error messages
37
38     # check if the function is limiting partition contributions
39     ids, by = partition(lambda expr: expr == identifier, partition_by) #
40
41     if not ids:
42         raise "failed to find identifier column in per_group predicate condition"
43
44     return Bound(by=by, per_group=threshold, num_groups=None) #

```

Postcondition

Theorem 1.1 (Postcondition). If `ranks` is a dense ranking of grouping columns, and `partition_by` is a singleton of `identifier`, then returns the bound on per-identifier contributions, or an error if the truncation is mis-specified.

Proof. Reorderings of an enumeration, like reversal, shuffling and sorting are also valid enumerations, so lines 8-29 ignore these reorderings and re-assign the enumeration to the input of the reordering.

Due to the ambiguity between matching predicates that bound `num_groups` or `per_group`, an error is only raised if the predicate is unambiguously a `num_groups` truncation predicate. This check happens on line 32.

Line 39 splits the partition by expressions into two sets, one containing the `identifier` and the other containing the grouping columns.

This predicate corresponds to a `per_group` truncation predicate, because the over expression groups by both the `identifier` column and any expressions in `by`, and within each group, an enumeration is applied to count the number of rows. If the only rows kept are those whose enumeration is less than `threshold`, then each user identifier will have at most `threshold` rows when their records are grouped by `by`. Reorderings of the enumeration, like reversal, shuffling and sorting can be used to choose which records from each individual are kept.

Therefore the bound on user contributions constructed on line 44 is valid. □