# fn make_stable_group_by

Michael Shoemate

April 22, 2025

> This proof resides in **"contrib"** because it has not completed the vetting process.

Proves soundness of `make_stable_group_by` in `mod.rs` at commit f5bb719 (outdated[1]).

# 1 Hoare Triple

## Precondition

### Caller Verified

None

## Function

```python
def make_stable_group_by(
    input_domain: DslPlanDomain, input_metric: FrameDistance[M], plan: DslPlan
) -> Transformation[DslPlanDomain, DslPlanDomain, FrameDistance[M], FrameDistance[M]]:
    match plan:
        case DslPlan.GroupBy(input, keys, aggs, apply, maintain_order, options):
            pass
        case _:
            raise "Expected group by in logical plan"

    if apply is not None:
        raise "apply is not currently supported"

    if maintain_order:
        raise "maintain_order is wasted compute because row ordering is protected
    information"

    if options != GroupbyOptions.default():
        raise "options is not currently supported"

    t_prior = input.make_stable(input_domain, input_metric)
    middle_domain, middle_metric = t_prior.output_space()

    expr_domain = WildExprDomain(
        columns=middle_domain.series_domains,
        context=Context.RowByRow,
    )

    for key in keys:  #
        key.make_stable(expr_domain, PartitionDistance(middle_metric[0]))
```

---

[1]See new changes with `git diff f5bb719..12e5590c rust/src/transformations/make_stable_lazyframe/group_by/mod.rs`

```
29
30    for agg in aggs:   #
31        check_infallible(agg, Resize.Allow)
32
33    if middle_metric[0].identifier().is_some():   #
34        raise "identifier is not currently supported"
35
36    # prepare output domain series
37    output_schema = middle_domain.simulate_schema(   #
38        lambda lf: lf.group_by(keys).agg(aggs)
39    )
40    series_domains = [SeriesDomain.new_from_field(f) for f in output_schema]
41
42    # prepare output domain margins
43    h_keys = list(keys)
44
45    def without_invariant(m: Margin) -> Margin:
46        m.invariant = None
47        return m
48
49    margins = [
50        without_invariant(m) for m in middle_domain.margins if m.by.is_subset(h_keys)
51    ]
52
53    output_domain = FrameDomain.new_with_margins(series_domains, margins)
54
55    def stability_map(d_in: Bounds) -> Bounds:
56        #
57        contributed_rows = d_in.get_bound(set()).per_group
58        #
59        contributed_groups = d_in.get_bound(h_keys).num_groups
60
61        influenced_groups = option_min(contributed_rows, contributed_groups)
62        if influenced_groups.is_none():
63            return "an upper bound on the number of contributed rows or groups is required"
64
65        if per_group is not None: #
66            per_group = per_group.inf_mul(2)
67
68        bound = Bound(by=set(), per_group=per_group, num_groups=None)
69        return Bounds([bound]) #
70
71    t_group_agg = Transformation.new(
72        middle_domain,
73        output_domain,
74        lambda plan: DslPlan.GroupBy(
75            input=plan,
76            keys=keys,
77            aggs=aggs,
78            apply=None,
79            maintain_order=False,
80            options=options,
81        ),
82        middle_metric,
83        middle_metric,
84        StabilityMap.new_fallible(stability_map),
85    )
86
87    return t_prior >> t_group_agg
```

## Postcondition

**Theorem 1.1.** For every setting of the input parameters (`input_domain`, `input_metric`, `plan`) to `make_stable_group_by` such that the given preconditions hold,

`make_stable_group_by` raises an exception (at compile time or run time) or returns a valid transformation. A valid transformation has the following properties:

1. (Appropriate output domain). For every element $x$ in `input_domain`, `function`$(x)$ is in `output_domain` or raises a data-independent runtime exception.

2. (Stability guarantee). For every pair of elements $x, x'$ in `input_domain` and for every pair $(\texttt{d\_in}, \texttt{d\_out})$, where `d_in` has the associated type for `input_metric` and `d_out` has the associated type for `output_metric`, if $x, x'$ are `d_in`-close under `input_metric`, `stability_map(d_in)` does not raise an exception, and `stability_map(d_in)` $\leq$ `d_out`, then `function`$(x)$, `function`$(x')$ are `d_out`-close under `output_metric`.

*Appropriate Output Domain.* By line 27 the grouping keys are stable row-by-row transformations of the data. By line 30 the aggregates are infallible. Therefore the function does not raise data-dependent errors.

By the postcondition of `DslPlan.simulate_schema`, `series_domains` follows the expected schema of members in the output domain. Notice that this is a very conservative output domain, as no data descriptors are preserved except for the schema itself. On the other hand, this comes with the benefit that aggregations are black-boxes, allowing for any infallible group-wise data processing.

For the same reason, the only margins preserved are those that are a subset of the grouping keys. Among these margins, invariants are discarded. A more careful proof may be able to preserve invariants, but this is not necessary for the soundness of the transformation.

It has been shown that for every element $x$ in `input_domain`, `function`$(x)$ is in `output_domain` or raises a data-independent runtime exception. $\square$

The stability guarantee doesn't attempt to claim the broadest set of possible bounds on output distances, rather it only claims a simple bound that might be useful for the user. This can be extended in the future, but is sufficient for select queries.

*Stability guarantee.* We first simplify the problem on line 33 by only considering datasets that differ by rows, not by identifiers.

The only bound derived is when there are no grouping keys.

We first retrieve optional upper bounds on the number of rows and groups an individual may contribute on lines 56 and 58. Both of these bounds directly correspond to the number of rows an individual may influence in the resulting dataset.

$$\max_{x \sim x'} d_{\text{FrameDistance<M>}}(\texttt{function}(x), \texttt{function}(x')) \tag{1}$$

$$\leq \text{option\_min}(\texttt{contributed\_rows}, \texttt{contributed\_groups}) \cdot 2 \tag{2}$$

Since each influenced row accounts for one addition and one removal, the distance is twice the number of influenced rows. This is reflected on line 65. The resulting bound is returned on line 69, satisfying the postcondition. $\square$