fn truncate_id_bound

Michael Shoemate

June 4, 2025

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

Proves soundness of truncate_id_bound in mod.rs at commit f5bb719 (outdated¹). truncate_id_bound returns an uppper bound on dataset distances in terms of the symmetric distance metric.

1 Hoare Triple

Precondition

Caller Verified

• id_bound.by = truncate.by

Function

```
def truncate_id_bound(
      id_bound: Bound,
      truncation: Bound,
      total_ids: Optional[int],
  ) -> Bound:
      # Once truncated, max contributions when grouped by "over" are bounded
      row_bound = Bound.by(truncation.by)
      # In each group, the worst-case row contributions is the
      # the number of ids contributed (known from id_bound)
10
      # times the number of rows contributed under each id (known from truncation),
11
      num_ids, num_rows = id_bound.per_group, truncation.per_group
12
      if num_ids is not None and num_rows is not None:
13
14
          row_bound = row_bound.with_per_group(num_ids.inf_mul(num_rows)) #
15
      # Worst case number of groups contributed is the
      # total number of ids contributed (total_ids)
17
      # times the number of groups contributed under each id (known from truncation).
18
      num_groups_via_truncation = None #
19
      if total_ids is not None and truncation.num_groups is not None:
20
21
          num_groups_via_truncation = total_ids.inf_mul(truncation.num_groups)
22
      # Alternatively, the number of groups contributed may be known outright from id_bound.
      # Use the smaller of the two if both are known.
24
      num_groups = option_min(num_groups_via_truncation, id_bound.num_groups)
26
      if num_groups is not None:
          row_bound = row_bound.with_num_groups(num_groups) #
27
      return row_bound
```

 $^{^1\}mathrm{See}$ new changes with git diff f5bb719..d584b91 rust/src/transformations/make_stable_lazyframe/truncate/mod.rs

Postcondition

Theorem 1.1 (Postcondition). Let g vary over groups when partitioned by id_bound.by. If for any two datasets x, x' we have that

```
\begin{split} ||d_{\operatorname{SymId}}(\operatorname{function}(x)_g, \operatorname{function}(x')_g)||_{\infty} & \leq \operatorname{id\_bound.per\_group}, \\ ||d_{\operatorname{SymId}}(\operatorname{function}(x)_g, \operatorname{function}(x')_g)||_0 & \leq \operatorname{id\_bound.num\_groups}, \\ d_{\operatorname{SymId}}(\operatorname{function}(x), \operatorname{function}(x')) & \leq \operatorname{total\_ids}, \end{split}
```

and function truncates a dataset such that,

$$\max_{id} ||d_{\mathrm{Sym}}(\mathtt{function}(x)_{id,g},\mathtt{function}(x')_{id,g})||_{\infty} \leq \mathtt{truncation.per_group},$$

$$\max_{id} ||d_{\mathrm{Sym}}(\mathtt{function}(x)_{id,g},\mathtt{function}(x')_{id,g})||_{0} \leq \mathtt{truncation.num_groups},$$

then we have that

$$||d_{\mathrm{Sym}}(\mathtt{function}(x)_g,\mathtt{function}(x')_g)||_{\infty} \leq \mathtt{row_bound.per_group},$$

 $||d_{\mathrm{Sym}}(\mathtt{function}(x)_g,\mathtt{function}(x')_g)||_{0} \leq \mathtt{row_bound.num_groups},$

where row_bound denotes the return value.

Proof. Assume the preconditions are met, as well as the conditions of the postcondition. We first prove the per-group bound.

```
\begin{aligned} &||d_{\operatorname{Sym}}(\operatorname{function}(x)_g,\operatorname{function}(x')_g)||_{\infty} \\ &\leq ||d_{\operatorname{SymId}}(\operatorname{function}(x)_g,\operatorname{function}(x')_g)||_{\infty}\cdot \max_{id}||d_{\operatorname{Sym}}(\operatorname{function}(x)_{id,g},\operatorname{function}(x')_{id,g})||_{\infty} \\ &\leq \operatorname{id\_bound.per\_group} \cdot \operatorname{truncation.per\_group} \\ &= \operatorname{row\_bound.per\_group} \qquad \qquad \operatorname{by\ line\ 14} \end{aligned}
```

We now prove the number of groups bound. There are two ways to bound the number of contributed groups. We first reason by the total number of identifiers.

```
||d_{\operatorname{Sym}}(\operatorname{function}(x)_g, \operatorname{function}(x')_g)||_0
\leq d_{\operatorname{SymId}}(\operatorname{function}(x), \operatorname{function}(x')) \cdot \max_{id} ||d_{\operatorname{Sym}}(\operatorname{function}(x)_{id,g}, \operatorname{function}(x')_{id,g})||_0
\leq \operatorname{total\_ids} \cdot \operatorname{truncation.num\_groups}
= \operatorname{num\_groups\_via\_truncation} by line 19
```

We can also directly bound the number of contributed groups by id_bound.num_groups:

```
\begin{split} &||d_{\mathrm{Sym}}(\mathtt{function}(x)_g,\mathtt{function}(x')_g)||_0\\ &=||d_{\mathrm{SymId}}(\mathtt{function}(x)_g,\mathtt{function}(x')_g)||_0\\ &\leq \mathtt{id\_bound.num\_groups} \end{split}
```

This is a valid upper bound on the number of contributed groups, because truncation is applied independently to each group.

Therefore, we take the minimum of the two upper bounds.

```
||d_{\mathrm{Sym}}(\mathtt{function}(x)_g,\mathtt{function}(x')_g)||_0

\min(\mathtt{num\_groups\_via\_truncation},\mathtt{id\_bound.num\_groups})

=\mathtt{row\_bound.num\_groups} by line 27
```

In the implementation, any of the input bounds could be missing. When a bound is missing, any output bounds that require it are not claimed. \Box