

fn make_sized_quantile_score_candidates

Michael Shoemate

April 12, 2023

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

Vetting History

- [Pull Request #456](#)

Proves soundness of `make_sized_quantile_score_candidates` in `mod.rs` at commit `f5bb719` (outdated¹). `make_sized_quantile_score_candidates` returns a Transformation that takes a numeric vector database and a vector of numeric quantile candidates, and returns a vector of scores, where higher scores correspond to more accurate candidates.

It is the same as the unsized version, `make_quantile_score_candidates`, but with a sized input domain, and different stability map. The unsized version has a more thorough explanation of the utility function.

1 Hoare Triple

Precondition

- TIA (input atom type) is a type with trait `Number`.
- TOA (output atom type) is a type with trait `Float`.

Function

```
1 def make_sized_quantile_score_candidates(size: usize, candidates: List[TIA], alpha: TOA):
2     for i in range(len(candidates) - 1):
3         assert candidates[i] < candidates[i + 1]
4
5     alpha_num, alpha_den = alpha.into_frac(size=size)
6     if alpha_num > alpha_den or alpha_den == 0:
7         raise ValueError("alpha must be within [0, 1]")
8
9     # ensures that the function will not overflow
10    size * alpha_den
11
12    def function(arg: List[TIA]):
13        return score(arg, candidates, alpha_num, alpha_den, size)
14
15    def stability_map(d_in: IntDistance):
16        return TOA.inf_cast(d_in // 2).inf_mul(4).inf_mul(alpha_den)
17
18    return Transformation(
```

¹See new changes with `git diff f5bb719..5fe96270 rust/src/transformations/quantile_score_candidates/mod.rs`

```

19     input_domain=SizedDomain(VectorDomain(AtomDomain(TIA)), size),
20     output_domain=VectorDomain(AtomDomain(TOA)),
21     function=function,
22     input_metric=SymmetricDistance(),
23     output_metric=InfDifferenceDistance(),
24     stability_map=stability_map,
25 )

```

Postcondition

For every setting of the input parameters (`size`, `candidates`, `alpha`) to `make_sized_quantile_score_candidates` such that the given preconditions hold, `make_sized_quantile_score_candidates` 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 v in `input_domain`, `function(v)` is in `output_domain` or raises a data-independent runtime exception.
2. (Domain-metric compatibility). The domain `input_domain` matches one of the possible domains listed in the definition of `input_metric`, and likewise `output_domain` matches one of the possible domains listed in the definition of `output_metric`.
3. (Stability guarantee). For every pair of elements u, v in `input_domain` and for every pair (d_{in}, d_{out}) , where d_{in} has the associated type for `input_metric` and d_{out} has the associated type for `output_metric`, if u, v are d_{in} -close under `input_metric` and `stability_map(d_in) ≤ d_out`, then `function(u), function(v)` are d_{out} -close under `output_metric`.

2 Proof

2.1 Appropriate Output Domain

The raw type and domain are equivalent, save for potential nullity in the atomic type. The scalar scorer structurally cannot emit null, because the input argument is non-null.

2.2 Domain-metric compatibility

On the input side, `SymmetricDistance` is well-defined on `VectorDomain`. On the output side, `InfDifferenceDistance` is well-defined on `VectorDomain` consisting of numeric elements.

2.3 Stability Guarantee

Lemma 2.1. If $d_{CO}(X, X') \leq 1$, then $d_{\infty}(\text{function}(X), \text{function}(X')) \leq 2$.

Proof. Assume $d_{CO}(X, X') \leq 1$. For convenience, let $s = \text{function}(X)$.

$$\begin{aligned}
d_{\infty}(s_i, s'_i) &= \max_i |s_i - s'_i| && \text{by definition of } d_{\infty} \\
&= \max_i |\text{abs_diff}(\alpha_{den} \cdot \min(\#(X < C_i), l), \alpha_{num} \cdot \min(|X| - \#(X = C_i), l))| && \text{by definition of function} \\
&\quad |\text{abs_diff}(\alpha_{den} \cdot \min(\#(X' < C_i), l), \alpha_{num} \cdot \min(|X'| - \#(X' = C_i), l))| \\
&= \alpha_{den} \cdot \max_i |\min(\#(X < C_i), l) - \alpha \cdot \min(|X| - \#(X = C_i), l)| \\
&\quad |\min(\#(X' < C_i), l) - \alpha \cdot \min(|X'| - \#(X' = C_i), l)| \\
&= \alpha_{den} \cdot \max_i |\#(X < C_i) - \alpha \cdot (|X| - \#(X = C_i))| \\
&\quad |\#(X' < C_i) - \alpha \cdot (|X'| - \#(X' = C_i))|
\end{aligned}$$

Consider each of the four cases of changing a row in X .

Case 1. Assume X' is equal to X , but with some $X_j < C_i$ replaced with $X'_j > C_i$.

$$\begin{aligned}
&= 2 \cdot \alpha_{den} \cdot \max_i ||(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - (1 - \alpha) \cdot (\#(X < C_i) - 1) - \alpha \cdot (\#(X > C_i) + 1)|| \\
&\leq 2 \cdot \alpha_{den} \cdot \max_i ||(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - (|(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| + |1|)| \\
&= 2 \cdot \alpha_{den} \cdot \max_i |1| \\
&= 2 \cdot \alpha_{den}
\end{aligned}$$

by definition of **function**

by triangle inequality
scores cancel

Case 2. Assume X' is equal to X , but with some $X_j > C_i$ replaced with $X'_j < C_i$.

$$= 2 \cdot \alpha_{den}$$

by symmetry, follows from Case 1.

Case 3. Assume X' is equal to X , but with some $X_j \neq C_i$ replaced with C_i .

$$\leq 2 \cdot \max(\alpha_{num}, \alpha_{den} - \alpha_{num})$$

equivalent to one removal (see **make**)

Case 4. Assume X' is equal to X , but with some $X_j = C_i$ replaced with $X'_j \neq C_i$.

$$\leq 2 \cdot \max(\alpha_{num}, \alpha_{den} - \alpha_{num})$$

equivalent to one addition (see **make**)

Take the union bound over all cases.

$$d_\infty(s_i, s'_i) \leq \max(2 \cdot \alpha_{den}, 2 \cdot \max(\alpha_{num}, \alpha_{den} - \alpha_{num})) = 2 \cdot \alpha_{den}$$

since $\max(\alpha, 1 - \alpha) \leq 1$

□

Take any two elements X, X' in the **input_domain** and any pair (d_in, d_out) , where d_in has the associated type for **input_metric** and d_out has the associated type for **output_metric**. Assume X, X' are d_in -close under **input_metric** and that **stability_map**(d_in) $\leq d_out$.

$$\begin{aligned}
d_out &= \max_{X \sim X'} d_{IDD}(s, s') \\
&= \max_{X \sim X'} \max_{ij} |(s_i - s'_i) - (s_j - s'_j)| && \text{by definition of } \mathbf{InfDifferenceDistance} \\
&\leq 2 \max_{X \sim X'} \max_i |s_i - s'_i| && s \text{ is not monotonic; take looser bound} \\
&\leq 2 \sum_j^{d_in/2} \max_{Z_j \sim Z_{j+1}} \max_i |s_{i,j} - s_{i,j+1}| && \text{by path property where } d_{CO}(Z_i, Z_{i+1}) = 1, X = Z_0 \text{ and } Z_{d_in} = X' \\
&\leq 2 \sum_j^{d_in/2} 2 \cdot \alpha_{den} && \text{by 2.1} \\
&\leq 4(d_in/2) \cdot \alpha_{den}
\end{aligned}$$

It is shown that **function**(X), **function**(X') are d_out -close under **output_metric**.