

fn make_quantile_score_candidates

Michael Shoemate

Christian Covington

Ira Globus-Harrus

April 12, 2023

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

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

Vetting History

- [Pull Request #456](#)

1 Intuition

The quantile score function scores each c in a set of candidates C .

$$s_i = -|(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \quad (1)$$

Where $\#(X < C_i) = |\{x \in X | x < C_i\}|$ is the number of values in X less than C_i , and similarly for other variations of inequalities. The scalar score function can be equivalently stated:

$$s_i = -|(1 - \alpha) \cdot \#(X < c) - \alpha \cdot \#(X > c)| \quad (2)$$

$$= -|(1 - \alpha) \cdot \#(X < c) - \alpha \cdot (|X| - \#(X < c) - \#(X = c))| \quad (3)$$

$$= -|\#(X < c) - \alpha \cdot (|X| - \#(X = c))| \quad (4)$$

It has an intuitive interpretation as $-|candidate_rank - ideal_rank|$, where the absolute distance between the candidate and ideal rank penalizes the score. The ideal rank does not include values in the dataset equal to the candidate. This scoring function considers higher scores better, and the score is maximized at zero when the candidate rank is equivalent to the rank at the ideal α -quantile.

The scalar scorer is almost equivalent to Smith’s^[1], but adjusts for a source of bias when there are values in the dataset equal to the candidate. For comparison, we can equivalently write the OpenDP scorer as if there were some α -discount on dataset entries equal to the candidate.

$$\begin{array}{ll} OpenDP & -|\#(X < c) + \alpha \cdot \#(X = c) - \alpha \cdot |X|| \\ Smith & -|\#(X < c) + 1 \cdot \#(X = c) - \alpha \cdot |X|| \end{array}$$

Observing that $\#(X \leq c) = \#(X < c) + 1 \cdot \#(X = c)$.

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

1.1 Examples

Let $X = \{0, 1, 2, 3, 4\}$ and $\alpha = 0.5$ (median):

$$\text{score}(X, 0, \alpha) = -|0 - .5 \cdot (5 - 1)| = -2$$

$$\text{score}(X, 1, \alpha) = -|1 - .5 \cdot (5 - 1)| = -1$$

$$\text{score}(X, 2, \alpha) = -|2 - .5 \cdot (5 - 1)| = -0$$

$$\text{score}(X, 3, \alpha) = -|3 - .5 \cdot (5 - 1)| = -1$$

$$\text{score}(X, 4, \alpha) = -|4 - .5 \cdot (5 - 1)| = -2$$

The score is maximized by the candidate at the true median.

Let $X = \{0, 1, 2, 3, 4, 5\}$ and $\alpha = 0.5$ (median):

$$\text{score}(X, 0, \alpha) = -|0 - .5 \cdot (6 - 1)| = -2.5$$

$$\text{score}(X, 1, \alpha) = -|1 - .5 \cdot (6 - 1)| = -1.5$$

$$\text{score}(X, 2, \alpha) = -|2 - .5 \cdot (6 - 1)| = -0.5$$

$$\text{score}(X, 3, \alpha) = -|3 - .5 \cdot (6 - 1)| = -0.5$$

$$\text{score}(X, 4, \alpha) = -|4 - .5 \cdot (6 - 1)| = -1.5$$

$$\text{score}(X, 5, \alpha) = -|5 - .5 \cdot (6 - 1)| = -2.5$$

The two candidates nearest the median are scored equally and highest.

Let $X = \{0, 1, 2, 3, 4\}$ and $\alpha = 0.25$ (first quartile):

$$\text{score}(X, 0, \alpha) = -|0 - .25 \cdot (5 - 1)| = -1$$

$$\text{score}(X, 1, \alpha) = -|1 - .25 \cdot (5 - 1)| = -0$$

$$\text{score}(X, 2, \alpha) = -|2 - .25 \cdot (5 - 1)| = -1$$

$$\text{score}(X, 3, \alpha) = -|3 - .25 \cdot (5 - 1)| = -2$$

$$\text{score}(X, 4, \alpha) = -|4 - .25 \cdot (5 - 1)| = -3$$

As expected, the score is maximized when $c = 1$.

Let $X = \{0, 1, 2, 3, 4, 5\}$ and $\alpha = 0.25$ (first quartile):

$$\text{score}(X, 0, \alpha) = -|0 - .25 \cdot (6 - 1)| = -1.25$$

$$\text{score}(X, 1, \alpha) = -|1 - .25 \cdot (6 - 1)| = -0.25$$

$$\text{score}(X, 2, \alpha) = -|2 - .25 \cdot (6 - 1)| = -0.75$$

$$\text{score}(X, 3, \alpha) = -|3 - .25 \cdot (6 - 1)| = -1.75$$

$$\text{score}(X, 4, \alpha) = -|4 - .25 \cdot (6 - 1)| = -2.75$$

$$\text{score}(X, 5, \alpha) = -|5 - .25 \cdot (6 - 1)| = -3.75$$

The ideal rank is 1.25. The nearest candidate, 1, has the greatest score, followed by 2, and then 0.

2 Finite Data Types

The previous equation assumes the existence of real numbers to represent α . We instead assume α is rational, such that $\alpha = \frac{\alpha_{num}}{\alpha_{den}}$. Multiply the equation through by α_{den} to get the following, which only uses integers:

$$\text{score}(X, c, \alpha_{\text{num}}, \alpha_{\text{den}}) = -|\alpha_{\text{den}} \cdot \#(X < c) - \alpha_{\text{num}} \cdot (|X| - \#(X = c))| \quad (5)$$

This adjustment also increases the sensitivity by a factor α_{den} , but does not affect the utility. We now make the scoring strictly non-negative.

- Drop the negation and instead configure the exponential mechanism to minimize the score.
- Compute the absolute difference in a function that swaps the order of arguments to keep the sign positive.

$$\text{score}(X, c, \alpha_{\text{num}}, \alpha_{\text{den}}) = \text{abs_diff}(\alpha_{\text{den}} \cdot \#(X < c), \alpha_{\text{num}} \cdot (|X| - \#(X = c))) \quad (6)$$

To prevent a numerical overflow when computing the arguments, first add a limit on the dataset size, l . α_{den} is chosen sufficiently small such that $\alpha_{\text{den}} \cdot l$ is representable. Instead of sampling the dataset size down, limit both counts to be at most l and replace $|X|$ with l .

$$\text{score}(X, c, \alpha_{\text{num}}, \alpha_{\text{den}}, l) = \text{abs_diff}(\alpha_{\text{den}} \cdot \min(\#(X < c), l), \alpha_{\text{num}} \cdot \min(|X| - \#(X = c), l)) \quad (7)$$

Should we compute counts with a 64-bit integer, we might choose α_{den} to be 10,000. This would allow for a fine fractional approximation of alpha, while still leaving enough room for datasets on the order of 10^{15} elements.

3 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_quantile_score_candidates(candidates: List[TIA], alpha: TOA):
2
3     for i in range(len(candidates) - 1):
4         assert candidates[i] < candidates[i + 1]
5
6     alpha_number, alpha_denom = alpha.into_frac(size=None)
7     if alpha_number > alpha_denom or alpha_denom == 0:
8         raise ValueError("alpha must be within [0, 1]")
9
10    # ensures that the function will not overflow
11    size_limit = size * alpha_denom
12
13    abs_dist_const = max(alpha_number, alpha_denom.inf_sub(alpha_number))
14    sup_dist_const = abs_dist_const.inf_mul(2)
15
16    def function(arg: List[TIA]):
17        return score(arg, candidates, alpha_number, alpha_denom, size_limit)
18
19    def stability_map(d_in: u32):
20        return TOA.inf_cast(d_in).alerting_mul(sup_dist_const)
21
22    return Transformation(
23        input_domain=VectorDomain(AtomDomain(TIA)),

```

```

24     output_domain=VectorDomain(AtomDomain(TOA)),
25     function=function,
26     input_metric=SymmetricDistance(),
27     output_metric=InfDifferenceDistance(TOA),
28     stability_map=stability_map,
29 )

```

Postcondition

For every setting of the input parameters (`candidates`, `alpha`) to `make_quantile_score_candidates` such that the given preconditions hold, `make_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`.

4 Proof

4.1 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.

4.2 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.

4.3 Stability Guarantee

Lemma 4.1. If $d_{Sym}(X, X') = 1$, then $d_{\infty}(\text{function}(X), \text{function}(X')) \leq \max(\alpha_{num}, \alpha_{den} - \alpha_{num})$.

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

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

Consider each of the three cases of adding or removing an element in X .

Case 1. Assume X' is equal to X , but with some $X_j < C_i$ added or removed.

$$\begin{aligned}
&= \alpha_{den} \cdot \max_i |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - |(1 - \alpha) \cdot (\#(X < C_i) \pm 1) - \alpha \cdot \#(X > C_i)| \\
&\leq \alpha_{den} \cdot \max_i |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - (|(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| + |\pm (1 - \alpha)|) \\
&= \alpha_{den} \cdot \max_i |1 - \alpha| \quad \text{by triangle inequality} \\
&\quad \text{scores cancel} \\
&= \alpha_{den} - \alpha_{num} \quad \text{since } \alpha \leq 1
\end{aligned}$$

Case 2. Assume X' is equal to X , but with some $X_j > C_i$ added or removed.

$$\begin{aligned}
&= \alpha_{den} \cdot \max_i |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot (\#(X > C_i) \pm 1)| \\
&\leq \alpha_{den} \cdot \max_i |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - (|(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| + |\pm \alpha|) \\
&= \alpha_{den} \cdot \max_i |\alpha| \quad \text{by triangle inequality} \\
&\quad \text{scores cancel} \\
&= \alpha_{num} \quad \text{since } \alpha \geq 0
\end{aligned}$$

Case 3. Assume X' is equal to X , but with some $X_j = C_i$ added or removed.

$$\begin{aligned}
&= \alpha_{den} \cdot \max_i |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&\quad - |(1 - \alpha) \cdot \#(X < C_i) - \alpha \cdot \#(X > C_i)| \\
&= 0 \quad \text{no change in score}
\end{aligned}$$

Take the union bound over all cases.

$$\max_{X \sim X'} d_\infty(s_i, s'_i) \leq \max(\alpha_{num}, \alpha_{den} - \alpha_{num})$$

□

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) ≤ d_out`.

$$\begin{aligned}
d_{\text{out}} &= \max_{X \sim X'} d_{IDD}(s, s') && \text{where } s = \text{function}(X) \\
&= \max_{X \sim X'} \max_{ij} |(s_i - s'_i) - (s_j - s'_j)| && \text{by definition of } \text{InfDifferenceDistance} \\
&\leq 2 \max_{X \sim X'} \max_i |s_i - s'_i| && s \text{ is not monotonic; take looser bound} \\
&\leq 2 \sum_j^{d_{\text{in}}} \max_{Z_j \sim Z_{j+1}} \max_i |s_{i,j} - s_{i,j+1}| && \text{by path property where } d_{\text{Sym}}(Z_i, Z_{i+1}) = 1, X = Z_0 \text{ and } Z_{d_{\text{in}}} = X' \\
&\leq 2 \sum_j^{d_{\text{in}}} \max(\alpha_{\text{num}}, \alpha_{\text{den}} - \alpha_{\text{num}}) && \text{by 4.1} \\
&\leq 2 \cdot d_{\text{in}} \cdot \max(\alpha_{\text{num}}, \alpha_{\text{den}} - \alpha_{\text{num}})
\end{aligned}$$

It is shown that $\text{function}(X)$, $\text{function}(X')$ are d_{out} -close under output_metric .

References

- [1] Adam Smith. Privacy-preserving statistical estimation with optimal convergence rates. In *Proceedings of the Forty-Third Annual ACM Symposium on Theory of Computing*, STOC '11, page 813–822, New York, NY, USA, 2011. Association for Computing Machinery.