

NoiseThresholdPrivacyMap<L02I<AbsoluteDistance<UBig>, Approximate<ZeroConcentratedDivergence>> for ZExpFamily<2>

Michael Shoemate

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

Proves soundness of the implementation of `NoiseThresholdPrivacyMap` for `ZExpFamily<2>` in `mod.rs` at commit `f5bb719` (outdated¹).

1 Hoare Triple

Precondition

Compiler-Verified

`NoiseThresholdPrivacyMap` is parameterized as follows:

- MI, the input metric, is of type `L02I<AbsoluteDistance<UBig>`
- MO, the output measure, is of type `Approximate<ZeroConcentratedDivergence>`

User-Verified

None

Pseudocode

```
1 # analogous to impl NoiseThresholdPrivacyMap<L02I<AbsoluteDistance<RBig>>, Approximate<
  ZeroConcentratedDivergence>> for ZExpFamily<2> in Rust
2 class ZExpFamily2:
3     def noise_threshold_privacy_map(
4         self,
5         _input_metric: L02InfDistance[AbsoluteDistance[RBig]],
6         output_measure: Approximate[ZeroConcentratedDivergence],
7         threshold: UBig,
8     ) -> PrivacyMap[L02InfDistance[AbsoluteDistance[RBig]], Approximate[
  ZeroConcentratedDivergence]]:
9         #
10        noise_privacy_map = self.noise_privacy_map(L2Distance.default(), output_measure[0])
11        scale = self.scale
12
13        def privacy_map(d_in: tuple[u32, RBig, RBig]):
14            l0, l2, l1 = d_in
```

¹See new changes with `git diff f5bb719..cde9e54 rust/src/measurements/noise_threshold/distribution/gaussian/mod.rs`

```

15         li_sign, li = li.floor().into_parts() #
16         if li_sign != Sign.Positive: #
17             raise f"l-infinity sensitivity ({li}) must be non-negative"
18
19
20         l2 = l2.min(li * RBig.try_from(f64.from_(10).inf_sqrt())) #
21         li = li.min(l2.floor()) #
22
23         if l2.is_zero(): #
24             return 0.0, 0.0
25
26         if scale.is_zero(): #
27             return f64.INFINITY, 1.0
28
29         rho = noise_privacy_map.eval(l2) #
30
31         if li > threshold: #
32             raise f"threshold must not be smaller than {li}"
33
34         d_instability = threshold - li #
35
36         delta_single = conservative_discrete_gaussian_tail_to_alpha(
37             scale,
38             d_instability,
39         )
40
41         delta_joint: f64 = (1.0).inf_sub(
42             (1.0).neg_inf_sub(delta_single).neg_inf_powi(IBig.from_(10)),
43         )
44
45         # delta is only sensibly at most 1
46         return rho, delta_joint.min(1.0)
47
48     return PrivacyMap.new_fallible(privacy_map)

```

Postcondition

Theorem 1.1. Given a distribution `self`, returns `Err(e)` if `self` is not a valid distribution. Otherwise the output is `Ok(privacy_map)` where `privacy_map` observes the following:

Define `function(x)` as a function that updates each pair $(k_i, v_i + Z_i)$, where Z_i are iid samples from `self`, and discards pairs where $v_i + Z_i$ has smaller magnitude than $|\text{threshold}|$. The ordering of returned pairs is independent from the input ordering.

For every pair of elements x, x' in `VectorDomain<AtomDomain<IBig>>`, and for every pair $(d_{\text{in}}, d_{\text{out}})$, where d_{in} has the associated type for `input_metric` and d_{out} has the associated type for `output_measure`, if x, x' are d_{in} -close under `input_metric`, `privacy_map(d_in)` does not raise an exception, and `privacy_map(d_in) ≤ d_out`, then `function(x), function(x')` are d_{out} -close under `output_measure`.

Proof. Line 9 rejects `self` if `self` does not represent a valid distribution, satisfying the error conditions of the postcondition.

We now construct the privacy map. The `li` sensitivity can be floored, as neighboring integer datasets always differ in whole integer increments. This doesn't affect the check on line 17 to ensure `li` is non-negative.

There are two ways to bound the `l2` sensitivity given the three bounds:

1. The `l2` bound directly.
2. Define $x \sim x'$ as $\|x - x'\|_0 \leq 10$ and $\|x - x'\|_\infty \leq li$. Then

$$\max_{x \sim x'} \|x - x'\|_2 = \max_{x \sim x'} \sqrt{\sum_i^n (x_i - x'_i)^2} \leq \sqrt{10 \cdot li^2} = \sqrt{10} \cdot li \quad (1)$$

Line 20 updates 12 to the tighter of these bounds.

1i similarly has multiple bounds:

1. The 1i bound directly.
2. Define $x \sim x'$ as $\|x - x'\|_2 \leq 12$. Then

$$\max_{x \sim x'} \|x - x'\|_\infty = \max_{x \sim x'} \max_i |x_i - x'_i| \leq \max_i [0, \dots, \lfloor 12 \rfloor, \dots, 0] = \lfloor 12 \rfloor \quad (2)$$

Line 21 updates 1i to the tighter of these bounds.

Now the 12 sensitivity is zero if any of the three bounds are zero, due to the tightening of the 12 bound on line 20. Line 23 then checks for zero sensitivity to return a privacy loss of zero rho, zero delta, because all neighboring datasets have identical pairs, and the ordering is randomized.

Otherwise sensitivity is nonzero, so if the scale is zero, the privacy loss is unbounded, shown on line 26.

Now that the edge cases are handled, all sensitivities and scales are finite and strictly positive.

The mechanism that adds noise to the values in the hashmap, as described in the conditions of the postcondition of **NoiseThresholdPrivacyMap**, matches the mechanism that adds noise to a vector, as described in the conditions of the postcondition of **NoisePrivacyMap**. Therefore releasing the values incurs the privacy loss computed by **noise_privacy_map** as defined on line 29, The privacy map also ensures that 12 is non-negative.

We now focus on computing the remaining privacy loss parameter delta. To ensure that the distance to instability remains positive, the threshold must be at least as large as the sensitivity, as checked on line 31.

Adapting the proof from [Rog23] (Theorem 7). Consider S to be the set of labels that are common between x and x' . Define event E to be any potential outcome of the mechanism for which all labels are in S (where only stable partitions are released). We then lower bound the probability of the mechanism returning an event E . In the following, c_j denotes the exact count for partition j , and Z_j is a random variable distributed according to **self**.

$$\begin{aligned} \Pr[E] &= \prod_{j \in x \setminus x'} \Pr[c_j + Z_j \leq T] \\ &\geq \prod_{j \in x \setminus x'} \Pr[\Delta_\infty + Z_j \leq T] \\ &\geq \Pr[\Delta_\infty + Z_j \leq T]^{\Delta_0} \end{aligned}$$

The probability of returning a set of stable partitions ($\Pr[E]$) is the probability of not returning any of the unstable partitions. We now solve for the choice of threshold T such that $\Pr[E] \geq 1 - \delta$.

$$\begin{aligned} \Pr[\Delta_\infty + Z_j \leq T]^{\Delta_0} &= \Pr[Z_j \leq T - \Delta_\infty]^{\Delta_0} \\ &= (1 - \Pr[Z_j > T - \Delta_\infty])^{\Delta_0} \end{aligned}$$

Let **d_instability** denote the distance to instability of $T - \Delta_\infty$. By the postcondition of **conservative_discrete_gaussian_tail_to_alpha**, the probability that a random noise sample exceeds **d_instability** is at most **delta_single**. Therefore $\delta = 1 - (1 - \text{delta_single})^{\Delta_0}$.

The privacy map returns **(rho, delta.min(1))**, as δ is bounded by 1. It is shown that **function(x)**, **function(x')** are **d_out**-close under **output_measure**. \square

References

- [Rog23] Ryan Rogers. A unifying privacy analysis framework for unknown domain algorithms in differential privacy, 2023.