

# fn make\_base\_exponential\_candidates\_gumbel

Michael Shoemate

March 8, 2023

## 1 Introduction

The discrete exponential mechanism is passed a vector of candidate scores and approximately releases the index of the greatest score. This release can be later be used to index into the public candidate set (via postprocessing).

The naive implementation samples some index  $k$  from a categorical distribution, with probabilities assigned to each candidate relative to their score. We may use inverse transform sampling to select the smallest index  $k$  for which the cumulative probability is greater than some  $U \sim Uniform(0, 1)$ .

$$M(s) = argmin_k \sum_i^k p_i >= U \quad (1)$$

The probability of index  $k$  being selected is the normalization of its likelihood  $e^{s_k/\tau}$ . As a candidate's score  $s_k$  increases, the candidate becomes exponentially more likely to be selected.

$$p_k = \frac{e^{s_k/\tau}}{\sum_i e^{s_i/\tau}} \quad (2)$$

This equation introduces a new temperature parameter,  $\tau$ , which calibrates how distinguishable scores are from each other. As temperature increases, the categorical output distribution tends towards entropy/uniformity and becomes more privacy preserving. As temperature decreases, the categorical distribution tends towards a one-hot vector, becoming less private. Temperature is related to  $\epsilon$  and the sensitivity ( $\Delta$ ) of the scoring function as follows:

$$\tau = \Delta/\epsilon \quad (3)$$

When  $\epsilon$  increases, temperature decreases, and candidates become more distinguishable from each other. We also divide scores by their global sensitivity to normalize the sensitivity to one. In the differential privacy literature for the exponential mechanism, the sensitivity is often multiplied by two. In OpenDP this factor is bundled into the  $\Delta$  term, which is expressed in terms of a metric that captures monotonicity.

## 2 Stable Reparameterization

In practice, computing  $e^{s_i/\tau}$  is prone to zero underflow and overflow. Specifically, a scaled score of just  $-709$  underflows to zero and  $+710$  overflows to infinity when stored in a 64-bit float. A simple improvement is to shift the scores by subtracting the greatest score from all scores. In idealized arithmetic, the resulting probabilities are not affected by shifts in the underlying scores. On finite data types, this shift prevents a catastrophic overflow, but makes underflow more likely, causing tail values of the distribution to round to zero.

The inverse transform sampling is also subject to accumulated rounding errors from the arithmetic and sum, which influence the likelihood of being chosen.

The Gumbel-max trick may instead be used to privately select an index. Let  $K = \operatorname{argmax}_k G_k$ , a random variable representing the selected index. Denote the  $k^{\text{th}}$  noisy score as  $G_k \sim \text{Gumbel}(\mu = s_k/\tau)$ .  $K$  can be sampled via an inverse transform, where  $u_k$  is sampled iid uniformly from  $(0, 1)$ :

$$M(s) = \operatorname{argmax}_k (s_k/\tau - \log(-\log(u_k))) \quad (4)$$

**Theorem 2.1.** Sampling from  $K$  is equivalent to sampling from the softmax, because  $P(K = k) = p_k$ . [1]

$$\begin{aligned} P(K = k) &= P(G_k = \max_i G_i) && \text{by definition of } K \\ &= P(-\log(Z_k/N) = \max_i -\log(Z_k/N)) && \text{by 2.2} \\ &= P(\log(Z_k/N) = \min_i \log(Z_k/N)) && \text{since } \max - a_i = -\min_i a_i \\ &= P(Z_k = \min_i Z_i) && \text{simplify monotonic terms} \\ &= P(Z_k \leq \min_{i \neq k} Z_i) \\ &= P(Z_k \leq Q) && \text{by 2.3 where } Q \sim \text{Exp}(\sum_{i \neq k} p_i) \\ &= \frac{p_k}{p_k + \sum_{i \neq k} p_i} && \text{by 2.4} \\ &= p_k && \text{since } p_k + \sum_{i \neq k} p_i = 1 \end{aligned}$$

**Lemma 2.2.**  $G_k = -\log(Z_k/N)$  where  $Z_k \sim \text{Exp}(p_k)$  and normalization term  $N = \sum_i e^{s_i/\tau}$ .

$$\begin{aligned} G_k &= s_k/\tau - \log(-\log(U_k)) && \text{Gumbel PDF centered at } s_k/\tau \\ &= \log(e^{s_k/\tau}) - \log(-\log(U_k)) \\ &= \log(p_k N) - \log(-\log(U_k)) && \text{since } p_k = e^{s_k/\tau}/N \\ &= \log(p_k N / (-\log(U_k))) \\ &= -\log(-\log(U_k) / (p_k N)) \\ &= -\log(Z_k/N) && \text{substitute } Z_k = -\log(U_k)/p_k \end{aligned}$$

**Lemma 2.3.** If  $X_1 \sim \text{Exp}(\lambda_1)$ ,  $X_2 \sim \text{Exp}(\lambda_2)$  and  $Z \sim \text{Exp}(\lambda_1 + \lambda_2)$ , then  $\min(X_1, X_2) \sim Z$ .

$$\begin{aligned} P(\min(X_1, X_2) \geq x) &= P(X_1 \geq x)P(X_2 \geq x) && \text{by independence} \\ &= e^{-\lambda_1 x} e^{-\lambda_2 x} && \text{substitute exponential density} \\ &= e^{-(\lambda_1 + \lambda_2)x} \\ &= P(Z \geq x) && \text{substitute exponential density} \end{aligned}$$

**Lemma 2.4.** If  $X_1 \sim \text{Exp}(\lambda_1)$ ,  $X_2 \sim \text{Exp}(\lambda_2)$ , then  $P(X_1 \leq X_2) = \frac{\lambda_1}{\lambda_1 + \lambda_2}$ .

$$\begin{aligned} P(X_1 \leq X_2) &= \int_0^\infty \int_{x_1}^\infty \lambda_1 \lambda_2 e^{-\lambda_1 x_1} e^{-\lambda_2 x_2} dx_1 dx_2 \\ &= \int_0^\infty -\lambda e^{-(\lambda_1 + \lambda_2)x_1} dx_1 \\ &= \frac{\lambda_1}{\lambda_1 + \lambda_2} \end{aligned}$$

### 3 Function

```

1 # type: ignore
2 def make_base_exponential_candidates_gumbel(temperature: TI):
3     def function(scores: List[TI]):
4         noisy_score = lambda i: scores[i] * temperature - ln( -ln(sample_uniform()))
5         return max(range(len(scores)), key=noisy_score)
6
7     def privacy_relation(d_in: u32, d_out: T0):
8         return d_out * temperature >= d_in
9
10    return Transformation(
11        input_domain=VectorDomain(AllDomain(TI)),
12        output_domain=AllDomain(T0),
13        function=function,
14        input_metric=InfDifferenceDistance(TI),
15        output_metric=MaxDivergence(TI),
16        privacy_relation=privacy_relation,
17    )

```

#### 3.1 Metric

We need a metric that captures the distance between score vectors  $s$  and  $s'$  respectively on neighboring datasets. The  $i^{th}$  element of each score vector is the score for the  $i^{th}$  candidate. Traditionally, the sensitivity of the scoring function is measured in terms of **InfDistance**, the greatest that any one score may change (or more formally, the  $L_\infty$  norm on distances). The sensitivity is defined as follows:

$$\Delta_\infty = \max_{s \sim s'} \max_i |s_i - s'_i| \quad (5)$$

Unfortunately, this choice of metric always results in a loosening by a factor of 2 when evaluating the privacy guarantee of the exponential mechanism. This is because both the  $i^{th}$  likelihood and normalization term may vary in opposite directions, resulting in a more distinguishing event. However, this loosening is not necessary if we can prove that the scoring function is monotonic, because the  $i^{th}$  likelihood and normalization term will always vary in the same direction.

We instead use a slight adjustment to this metric, **InfDifferenceDistance**, with sensitivity as follows:

$$\Delta_{\infty'} = \max_{s \sim s'} \max_{ij} |(s_i - s'_i) - (s_j - s'_j)| \quad (6)$$

Consider when the scoring function is not monotonic. The sensitivity is maximized when  $s_i - s'_i$  and  $s_j - s'_j$  vary maximally in opposite directions, resulting in the same loosening factor of 2. On the other hand, when the scoring function is monotonic, the sign of the  $s_i - s'_i$  term matches the sign of the  $s_j - s'_j$  term, and their magnitudes cancel. Therefore, when the scorer is monotonic, the sensitivity is maximized when one term is zero. It is shown in 4.2 that a tighter analysis of the exponential mechanism is compatible with a score vector whose sensitivity is expressed in terms of this metric.

## 4 Proof

### 4.1 Domain-metric compatibility

On the input side, **SymmetricDistance** is well-defined on **VectorDomains**. On the output side, **InfDifferenceDistance** is well-defined on **VectorDomains**.

### 4.2 Privacy Guarantee

A mechanism  $M$  satisfies  $\epsilon$ -differential privacy if, for any adjacent score vectors  $s$  and  $s'$  and for all  $K \subseteq \text{Range}(M)$

$$P(M(s) \in K) \leq e^\epsilon P(M(s') \in K) \quad (7)$$

which implies

$$\ln \left( \frac{P(M(s) \in K)}{P(M(s') \in K)} \right) \leq \epsilon \quad (8)$$

The following proves that  $M$  as defined in 1 satisfies differential privacy.

$$\begin{aligned}
\ln \left( \frac{P(M(s) = k)}{P(M(s') = k)} \right) &= \ln \left( \frac{\exp\left(\frac{\epsilon s_k}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \bigg/ \frac{\exp\left(\frac{\epsilon s'_k}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)} \right) && \text{substitute 2, assuming } \tau \geq \Delta/\epsilon \\
&= \ln \left( \frac{\exp\left(\frac{\epsilon s_k}{\Delta}\right)}{\exp\left(\frac{\epsilon s'_k}{\Delta}\right)} \frac{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&= \ln \left( \frac{\exp\left(\frac{\epsilon s_k}{\Delta}\right)}{\exp\left(\frac{\epsilon s'_k}{\Delta}\right)} \right) + \ln \left( \frac{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&= \frac{\epsilon(s_k - s'_k)}{\Delta} + \ln \left( \frac{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&\leq \frac{\epsilon(s_k - s'_k)}{\Delta} + \frac{-\epsilon \max_j (s_j - s'_j)}{\Delta} && \text{by 4.1} \\
&\leq \frac{\epsilon \max_{ij} |(s_i - s'_i) - (s_j - s'_j)|}{\Delta} \\
&\leq \epsilon
\end{aligned}$$

**Lemma 4.1.**  $\ln \left( \frac{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \geq \frac{\epsilon \max_j (s_j - s'_j)}{\Delta}.$

$$\begin{aligned}
\ln \left( \frac{\sum_i \exp\left(\frac{\epsilon s'_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) &= \ln \left( \frac{\sum_i \exp\left(\frac{\epsilon(s'_i - s_i + s_i)}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&= \ln \left( \frac{\sum_i \exp\left(\frac{\epsilon(s'_i - s_i)}{\Delta}\right) \exp\left(\frac{\epsilon s_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&\geq \ln \left( \frac{\exp\left(\frac{\epsilon \min_j (s'_j - s_j)}{\Delta}\right) \sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)}{\sum_i \exp\left(\frac{\epsilon s_i}{\Delta}\right)} \right) \\
&= \frac{\epsilon \min_j (s'_j - s_j)}{\Delta} \\
&= -\frac{\epsilon \max_j (s_j - s'_j)}{\Delta}
\end{aligned}$$

## References

- [1] Andrés Muñoz Medina and Jennifer Gillenwater. Duff: A dataset-distance-based utility function family for the exponential mechanism. *ArXiv*, abs/2010.04235, 2020.