

STOCHASTIC BEAMS AND WHERE TO FIND THEM

ARTICLE WOUTER KOOL, HERKE VAN HOOFF, MAX WELLING

PRESENTATION DAYANA SAVOSTIANOVA

NRU HSE

26 09 2019

Applications

- Image captioning
- Neural machine translation
- Speech recognition

MATH

THE CATEGORICAL DISTRIBUTION

$$I \sim \text{Categorical}(p_1, \dots, p_n)$$

$$P(I = i) = p_i \quad \forall i \in N$$

$$\phi_i, i \in N, \exp \phi_i \propto p_i$$

$$I \sim \text{Categorical} \left(\frac{\exp \phi_i}{\sum_{j \in N} \exp \phi_j}, i \in N \right). \quad (1)$$

THE GUMBEL DISTRIBUTION

$$U \sim \text{Uniform}(0, 1)$$

$$G = \phi - \log(-\log U) : G \sim \text{Gumbel}(\phi)$$

$$G' = G + \phi' \sim \text{Gumbel}(\phi + \phi')$$

Properties: we can *shift* Gumbel variables.

THE GUMBEL MAX-TRICK

The Gumbel-Max trick allows to sample from the categorical distribution by independently *perturbing* the log-probabilities ϕ_i with Gumbel noise and finding the largest element.

$G_i \sim \text{Gumbel}(0), i \in N, I^* = \text{argmax}_i \{\phi_i + G_i\}$

$I^* \sim \text{Categorical}(p_i, i \in N)$ with $p_i \propto \exp \phi_i$.

$G_{\phi_i} = G_i + \phi_i \sim \text{Gumbel}(\phi_i)$

$\forall B \subseteq N$:

$$\max_{i \in B} G_{\phi_i} \sim \text{Gumbel} \left(\log \sum_{j \in B} \exp \phi_j \right), \quad (2)$$

$$\text{argmax}_{i \in B} G_{\phi_i} \sim \text{Categorical} \left(\frac{\exp \phi_i}{\sum_{j \in B} \exp \phi_j}, i \in B \right). \quad (3)$$

Max and argmax are independent.

THE GUMBEL TOP-K-TRICK

Difference: *ordered sample of size k without replacement.*

Theorem

For $k \leq n$, let $l_1^*, \dots, l_k^* = \text{argtop}_k G_{\phi_l}$. Then l_1^*, \dots, l_k^* is an (ordered) sample without replacement from the Categorical $\left(\frac{\exp \phi_i}{\sum_{j \in N} \exp \phi_j}, i \in N \right)$ distribution, e.g. for a realization i_1^*, \dots, i_k^* it holds that

$$P(l_1^* = i_1^*, \dots, l_k^* = i_k^*) = \prod_{j=1}^k \frac{\exp \phi_{i_j^*}}{\sum_{\ell \in N_j^*} \exp \phi_\ell} \quad (4)$$

where $N_j^* = N \setminus \{i_1^*, \dots, i_{j-1}^*\}$ is the domain (without replacement) for the j -th sampled element.

SEQUENCE MODELS

Typically $p_{\theta}(y_t|\mathbf{y}_{1:t-1})$ is defined as a softmax normalization of unnormalized log-probabilities $\phi_{\theta}(y_t|\mathbf{y}_{1:t-1})$ with optional temperature T (default $T = 1$):

$$p_{\theta}(y_t|\mathbf{y}_{1:l-1}) = \frac{\exp(\phi_{\theta}(y_t|\mathbf{y}_{1:t-1})/T)}{\sum_{y'} \exp(\phi_{\theta}(y'|\mathbf{y}_{1:t-1})/T)}. \quad (5)$$

The normalization is w.r.t. a single token, so the model is *locally normalized*. The total probability of a (partial) sequence $\mathbf{y}_{1:l}$ follows from the chain rule of probability:

$$p_{\theta}(\mathbf{y}_{1:t}) = p_{\theta}(y_t|\mathbf{y}_{1:t-1}) \cdot p_{\theta}(\mathbf{y}_{1:t-1}) \quad (6)$$

$$= \prod_{t'=1}^t p_{\theta}(y_{t'}|\mathbf{y}_{1:t'-1}). \quad (7)$$

1. Limited-width breadth first search.
2. Often used as an approximation to finding the (single) sequence \mathbf{y} that maximizes, or as a way to obtain a set of high-probability sequences from the model.
3. Expands at every step $t = 0, 1, 2, \dots$ at most k partial sequences (those with highest probability) to compute the probabilities of sequences with length $t + 1$.
4. Terminates with a beam of k complete sequences.

THE GUMBEL-TOP-K TRICK ON A TREE

We represent the sequence model as a tree.

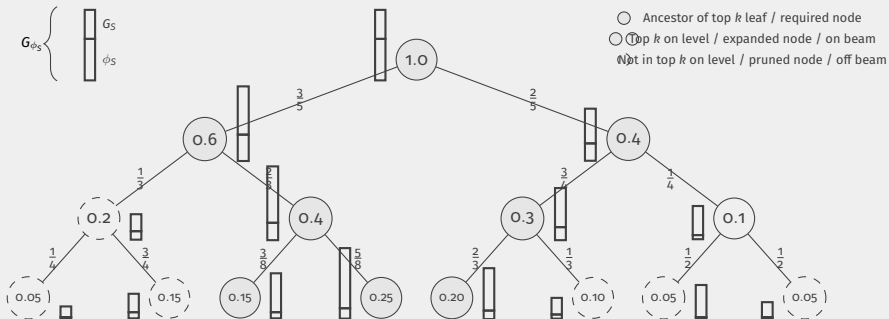
Internal nodes at level t represent partial sequences $\mathbf{y}_{1:t}$

Leaf nodes represent completed sequences.

Normalized log-probability $\phi_i = \log p_\theta(\mathbf{y}^i)$.

- Compute $\phi_i = \log p_\theta(\mathbf{y}^i)$ for all sequences $\mathbf{y}^i, i \in N$. (Complete probability tree is instantiated)
- Sample $G_{\phi_i} \sim \text{Gumbel}(\phi_i)$, so G_{ϕ_i} can be seen as the perturbed log-probability of sequence \mathbf{y}^i .
- Let $i_1^*, \dots, i_k^* = \arg \text{top } k G_{\phi_i}$, then $\mathbf{y}^{i_1^*}, \dots, \mathbf{y}^{i_k^*}$ is a sample of sequences without replacement.

EXAMPLE ON TREE



PERTURBING LOG PROBABILITY

Internal or leaf – the set S of leaves in the corresponding subtree, and \mathbf{y}^S – the corresponding sequence.

$$\phi_S = \log p_\theta(\mathbf{y}^S) = \log \sum_{i \in S} \exp \phi_i. \quad (8)$$

$\forall S$ G_{ϕ_S} – maximum of the perturbed log-probabilities G_{ϕ_i} in the subtree leaves S .

$$G_{\phi_S} = \max_{i \in S} G_{\phi_i} \sim \text{Gumbel}(\phi_S) \quad (9)$$

Gumbel noise $G_S \sim \text{Gumbel}(0)$

$$G_{\phi_S} = \phi_S + G_S.$$

Children(S) – set of direct children of the node S Rule:

$$G_{\phi_S} = \max_{S' \in \text{Children}(S)} G_{\phi_{S'}}. \quad (10)$$

If we want to sample G_{ϕ_S} for all nodes, we can use the *bottom-up* sampling procedure: sample the leaves $G_{\phi_{\{i\}}} = G_{\phi_i}$, $i \in N$ and recursively compute G_{ϕ_S} . This is effectively sampling from the degenerate (constant) distribution resulting from conditioning on the children.

TOP DOWN

1. Initialize from the root
2. Rule: (from bottom up)
3. Sample the children conditionally on the parent variable (preserving max).
4. Sampling a set of Gumbels conditionally on their maximum being equal to a certain value can be done by first sampling the $\arg \max$ and then sampling the individual Gumbels conditionally on both the \max and $\arg \max$.
5. For all leaves $(G_{\phi_{\{i\}}} =) G_{\phi_i} \sim \text{Gumbel}(\phi_i)$ is *independent*
6. The benefit of using top-down sampling is that if we are interested only in obtaining the top k leaves, we do *not* need to instantiate the complete tree.

STOCHASTIC BEAM SEARCH ALGORITHM

1. Apply the top-down sampling procedure
2. At each level we only need to expand the k nodes with the highest perturbed log-probabilities G_{ϕ_S}
3. By the Gumbel-Top- k trick the result is a sample without replacement from the sequence model

We can also think of G_{ϕ_S} as the *stochastic score* of the partial sequence \mathbf{y}^S .

WHY BETTER?

1. A sampling procedure
2. Principled way to randomize a beam search.

Randomized beam search problems:

- Low-probability partial sequence need to be re-chosen, independently, again with low probability at each step. The result is a much lower probability to sample this sequence than assigned by the model. Intuitively, we should somehow *commit* to a sampling ‘decision’ made at step t .
- Stochastic Beam Search is better suited because it makes a *soft* commitment to a partial sequence by propagating the Gumbel perturbation of the log-probability consistently down the subtree.

EXPERIMENTS

Different algorithms

- Beam Search
- Sampling
- Stochastic Beam Search
- Diverse Beam Search with G groups

Metrics

- Diversity – the fraction of unique n -grams in the k translations
- Quality of the sample – maximum BLEU score
- Diversity – mean BLEU score

N-gram diversity

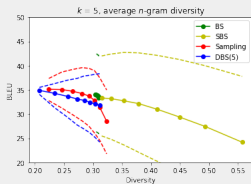
$$d_n = \frac{\text{\# of unique } n\text{-grams in } k \text{ translations}}{\text{total \# of } n\text{-grams in } k \text{ translations}}.$$

BLEU

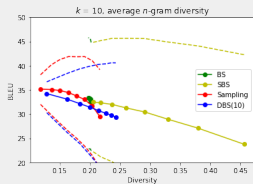
Bilingual evaluation understudy

"The closer a machine translation is to a professional human translation, the better it is"

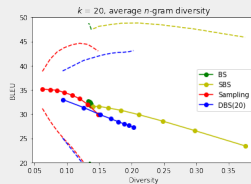
The wmt14.v2.en-fr.newstest2014 test set consisting of 3003 sentences



(a) $k = 5$



(b) $k = 10$



(c) $k = 20$

Figure: Minimum, mean and maximum BLEU score vs. diversity for different sample sizes k . Points indicate different temperatures/diversity strengths, from 0.1 (low diversity, left in graph) to 0.8 (high diversity, right in graph).

BLUE: $f(\mathbf{y}) = \text{BLEU}(\mathbf{y}, \mathbf{x})$

Entropy: $f(\mathbf{y}) = -\log p_{\theta}(\mathbf{y}|\mathbf{x})$

Evaluation methods

■ Monte Carlo

$$\mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y}|\mathbf{x})} [f(\mathbf{y})] \approx \frac{1}{k} \sum_{i \in S} f(\mathbf{y}^i). \quad (11)$$

■ StochasticBeamSearch (Normalized)

$$\mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y}|\mathbf{x})} [f(\mathbf{y})] \approx \sum_{i \in S} \frac{p_{\theta}(\mathbf{y}^i|\mathbf{x})}{q_{\theta, \kappa}(\mathbf{y}^i|\mathbf{x})} f(\mathbf{y}^i) \text{ if norm } \left(\frac{1}{W(S)} \right) \quad (12)$$

■ Beam search (Norm)

$$\sum_{i \in S} p_{\theta}(\mathbf{y}^i|\mathbf{x}) f(\mathbf{y}^i) \text{ if norm } \left(\frac{1}{\sum_{i \in S} p_{\theta}(\mathbf{y}^i|\mathbf{x})} \right) \quad (13)$$

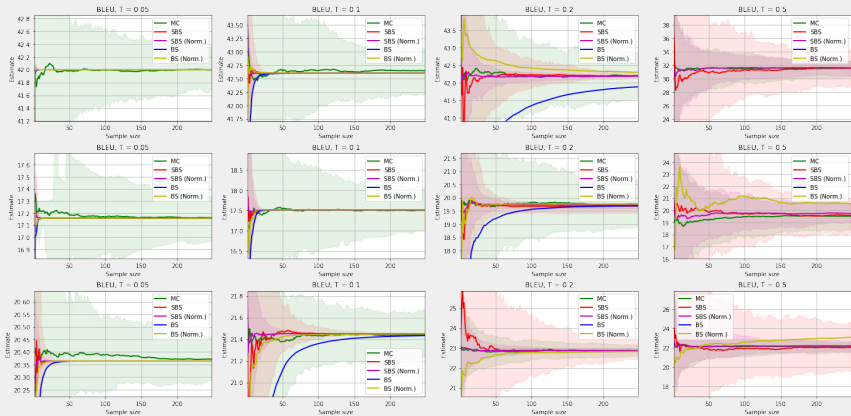


Figure: BLEU score estimates for three sentences sampled/decoded by different estimators for different temperatures.

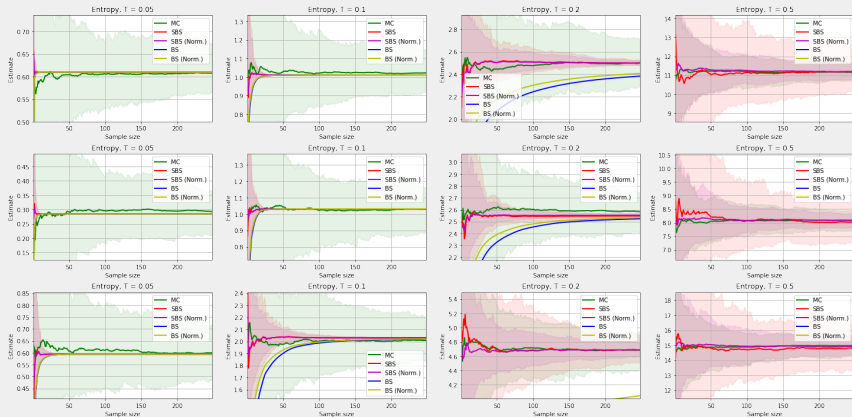


Figure: Entropy score estimates for three sentences sampled/decoded by different estimators for different temperatures.

Stochastic Beam Search shares some of the benefits of these heuristic variants, such as the ability to control diversity or produce randomized output.

Benefits

- Linear towards k and the length of the sequence
- Ability to control diversity
- Produce randomized output
- Easy to implement on top of a beam search as a way to sample sequences without replacement

Thank you for your **attention!**

Algorithm 1 StochasticBeamSearch(p_{θ}, k)

```
1: Input: one-step probability distribution  $p_{\theta}$ , beam/sample size  $k$ 
2: Initialize BEAM empty
3: add ( $\mathbf{y}^N = \emptyset, \phi_N = 0, G_{\phi_N} = 0$ ) to BEAM
4: for  $t = 1, \dots$ , steps do
5:   Initialize EXPANSIONS empty
6:   for ( $\mathbf{y}^S, \phi_S, G_{\phi_S}$ )  $\in$  BEAM do
7:      $Z \leftarrow -\infty$ 
8:     for  $S' \in \text{Children}(S)$  do
9:        $\phi_{S'} \leftarrow \phi_S + \log p_{\theta}(\mathbf{y}^{S'} | \mathbf{y}^S)$ 
10:       $G_{\phi_{S'}} \sim \text{Gumbel}(\phi_{S'})$ 
11:       $Z \leftarrow \max(Z, G_{\phi_{S'}})$ 
12:    end for
13:    for  $S' \in \text{Children}(S)$  do
14:       $\tilde{G}_{\phi_{S'}} \leftarrow -\log(\exp(-G_{\phi_S}) - \exp(-Z) + \exp(-G_{\phi_{S'}}))$ 
15:      add ( $\mathbf{y}^{S'}, \phi_{S'}, \tilde{G}_{\phi_{S'}}$ ) to EXPANSIONS
16:    end for
17:  end for
18:  BEAM  $\leftarrow$  take top  $k$  of EXPANSIONS according to  $\tilde{G}$ 
19: end for
20: Return BEAM
```
