Predictive Querying for Autoregressive Neural Sequence Models Padhraic Smyth1,2 Sam Showalter^{2*} Stephan Mandt^{1,2} Alex Boyd^{1*} Selected as Oral

¹ Department of Statistics, ² Department of Computer Science - University of California, Irvine * Denotes equal contribution {alexjb, showalte}@uci.edu

Abstract

- Autoregressive models capture detailed information about future events.
- To-date, investigation is centered solely on next-event prediction.
- Goal: Formalize complex queries such as "event A occurs before B"
- Goal: Leverage our formalism to produce methods for estimating complex queries that exist in exponentially large sequence spaces



We seek to approximate queries of interest with the following

$$p_{\theta}^*(X_{1:K} \in \mathcal{Q}) = \sum_{x_{1:K} \in \mathcal{Q}} p_{\theta}^*(X_{1:K} = x_{1:K}) = \sum_{x_{1:K} \in \mathcal{Q}} \prod_{k=1}^{K} p_{\theta}^*(X_k = x_k | X_k)$$

• With this query formalism, we can decompose this process:

$$p_{\theta}^*(X_{1:K} \in \mathcal{Q}) = \sum_i p_{\theta}^*(X_{1:K} \in \mathcal{Q}^{(i)}) = \sum_i p_{\theta}^*(\bigcap_{k=1}^K X_k \in \mathcal{D})$$

• We seek other estimation methods beyond Monte-Carlo sampling: $p_{\theta}^*(X_{1:K} \in \mathcal{Q}) = \mathbb{E}_{x_{1:K} \sim p_{\theta}^*} \left[\mathbf{1}(x_{1:K} \in \mathcal{Q}) \right]$

Probabilistic Queries

#	Question	Probabilistic Query	$\operatorname{Cost}\left(K\cdot \mathcal{Q} \right)$
Q 1	Next event?	$p_{ heta}^*(X_1)$	$\mathcal{O}(1)$
Q 2	Event K steps from now?	$p_{\theta}^*(X_K)$	$\mathcal{O}(V^{K-1})$
Q3	Next instance of a?	$p_{\theta}^*(\tau(a) = K)$	$\mathcal{O}((V-1)^{K-1})$
Q 4	Will a happen before b?	$p_{\theta}^*(\tau(a) < \tau(b))$	$\mathcal{O}((V-2)^K)^\dagger$
Q 5	How many instances of a in K steps?	$p_{\theta}^*(N_a(K) = n)$	$\mathcal{O}\left(\binom{K}{n}(V-1)^{K-n}\right)$

• There are many queries of interest in sequence modeling, but all require marginalization over exponentially large path spaces



• **Beam search** can *guarantee coverage* by finding beams $\mathscr{B} \subset \mathscr{Q}$

$$p_{\theta}^*(X_{1:K} \in \mathcal{Q}) = \sum_{x_{1:K} \in \mathcal{Q}} p_{\theta}^*(X_{1:K} = x_{1:K}) \ge \sum_{x_{1:K} \in \mathcal{B}} p_{\theta}^*(X_{1:K} = x_{1:K})$$

Our novel hybrid method merges search and sampling by noting that:

$$p_{\theta}^*(X_{1:K} \in \mathcal{Q}) = \sum_{x_{1:K} \in \mathcal{B}_K} p_{\theta}^*(X_{1:K} = x_{1:K}) + \sum_{x_{1:K} \in \mathcal{Q} \setminus \mathcal{B}_K} p_{\theta}^*(X_{1:K} = x_{1:K})$$

- Where \mathscr{B}_{K} is the sequence set from beam search. We first search for likely sequences and then sample from the remaining space with $q(X_{1:K} = x_{1:K} | X_{1:K} \notin \mathscr{B}_K)$
- Ground-truth is often *intractable* to compute. We evaluate our methods against pseudo-ground truth (PGT) estimates generated with importance sampling
- PGT leverages a high compute budget and the convergence guarantees of the CLT

$$= x_{k} | X_{
$$= x_{k} | X_{$$$$

$$K = x_{1:K} \Big] = \mathbb{E}_{x_{1:K} \sim q} \left[\frac{p_{\theta}^*(X_{1:K} = x_{1:K})}{q(X_{1:K} = x_{1:K})} \right]$$

for $x_{1:K}^{(1)}, \dots, x_{1:K}^{(M)} \stackrel{iid}{\sim} q$



 $p_{\theta}^*(\tau(a) < \tau(b))$







Median relative absolute error (RAE) for for hitting-time query estimates over 4 datasets and 3 estimation methods. Hybrid method outperforms others including naive sampling.

Relative Efficiency ¹⁶ Shakespeare Mobile Apps MOOCs

Median rel. efficiency of importance sampling for all datasets. For the regime of K where computing ground truth is intractable, we see a significant boost over naive sampling.

Query Estimation and Entropy



(left) Median relative absolute error (RAE) vs restricted entropy per query (with best linear fits), (right) Median RAE v. model temperature T for Mobile Apps data. Beam search error highly correlates with entropy, and increasing T induces 100% error.

Composability and Saving Computation

 More complex queries (e.g. Q4, Q5) can be decomposed into operations over hitting-time queries (Q3)

$$p) = \sum_{k=1}^{\infty} p_{\theta}^* (\tau(a) = k, \tau(b) > k) = \sum_{k=1}^{\infty} p_{\theta}^* (X_k = a, X_{< k} \in (\mathbb{V} \setminus \{a, b\})^{k-1})$$

• For this reason, experiments focus on hitting-time queries • One can save computation between queries Q and Q' by re-using sub-sequences of $p_{A}^{*}(X_{1:K} \in \mathbb{Q})$ on \mathbb{Q}' if $\mathcal{V}_{i} = \mathcal{V}'_{i}$ for i = 1, ..., K-1