



BEATING NATURAL DISTRIBUTION ESTIMATION WITH SIDE INFORMATION*

EE Research Symposium

Haricharan Balasundaram, Prof. Andrew Thangaraj

*Partly presented at Allerton 2025 and partly submitted to ISIT 2026

DICKENS' SIMULATOR



- Machine that simulates Charles Dickens' writing
- Can be trained on previous Dickens books (data)

How to get a 'Good' Simulator?

NEXT WORD PREDICTION PROBLEM

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of

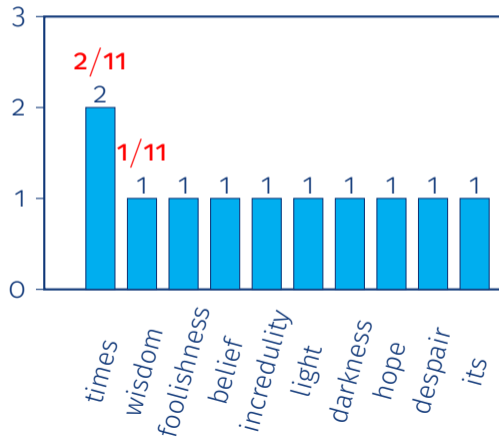
Given the previous word, what is the Probability of the next word being something?

Words occurring after 'of'	Probabilities
Surprise	0.3
Incredulity	0.2
Amazement	0.05

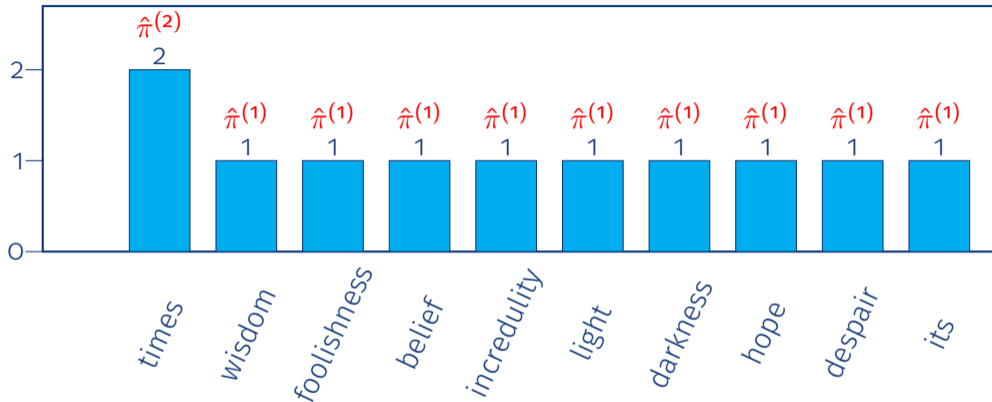
Table 1: Example estimate for the Word occurring after 'of': $\hat{\pi}(\cdot|of)$

THE EMPIRICAL ESTIMATOR FOR $\pi(\cdot|‘OF’)$

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair ...that some of its noisiest authorities insisted on its being received, for good or for evil ...



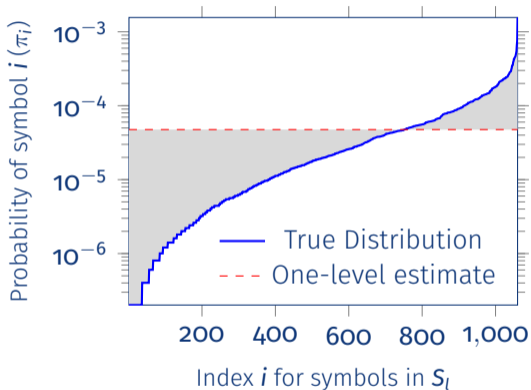
WHAT IS A NATURAL ESTIMATOR?



Natural Estimators (Empirical, smoothed, profile-based, etc.)

Letters occurring l times (S_l) assigned the same estimator probability mass $\hat{\pi}^{(l)}$

WHY DO NATURAL ESTIMATORS HAVE ERROR?

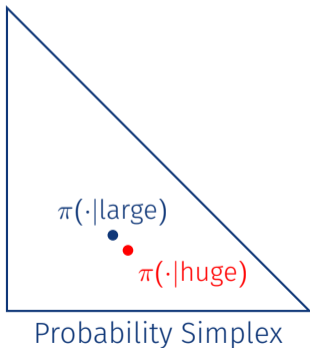


- Over species occurring l times, we assign a single probability mass S_l , while the real probabilities could be anything!
- This error is a random variable that depends on S_l
- We estimate this error!

Error in estimation of unigram probabilities for 10000 contiguous words taken from TOTC (Total 100,000 words!)

AN EXAMPLE OF SIDE INFORMATION

TOTC: 'huge' occurs once, 'large' occurs 44 times
Huge and Large are synonyms!



- Probability distribution of words after 'Huge' ($\pi(\cdot|huge)$) and words after 'Large' ($\pi(\cdot|large)$)
- The closeness can be used to estimate $\pi(\cdot|huge)$ accurately
- Estimator:

$$\alpha \cdot \hat{\pi}(\cdot|huge) + (1 - \alpha) \cdot \hat{\pi}(\cdot|large)$$

OUR RESULTS

THE SETTING

Discrete distribution $\pi = (\pi_1, \pi_2, \dots, \pi_d)$

X^n : n i.i.d. samples drawn from π

Distribution Estimation

Devise an estimate $\hat{\pi}(X^n)$ to minimize the squared ℓ_2 -norm error between $\hat{\pi}(X^n)$ and π .

$$R(\pi, \hat{\pi}) = E[(\hat{\pi}(X^n) - \pi)^2]$$

OUR RESULTS (INFORMAL)¹

$$\underbrace{\sum_{i \in S_l} (\pi_i - \hat{\pi}^{(l)})^2}_{\text{Error over } S_l} = \underbrace{\sum_{i \in S_l} (\pi_i - \pi^{(\text{or})})^2}_{\text{Unavoidable (Oracle) Error}} + \underbrace{\frac{(M_l - \hat{M}_l)^2}{\phi_l}}_{\text{Estimator Error}}$$

$$\phi_l = |S_l| \quad \pi^{(\text{or})} = \frac{M_l}{\phi_l}$$

We provide an estimator T_l for the unavoidable oracle error Q_l based on the **Good-Turing estimator**:

- Show low bias
- Show error probability at most δ with (i) $n \geq O(d^2/\delta)^{\frac{1}{3}}$ samples in case of $l = o$, and (ii) $n \geq O(1/\delta)$ samples for $l > o$

¹Estimating Error in Natural Distribution Estimation. H. Balasundaram, A. Thangaraj. Allerton 2025.

ERROR IN THE ESTIMATOR FOR VARIOUS DISTRIBUTIONS

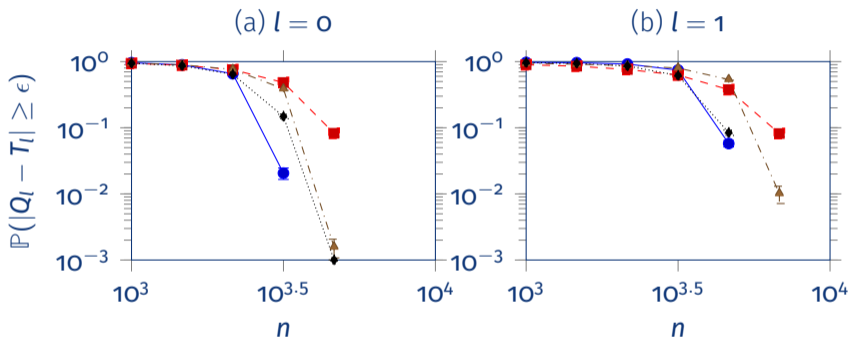
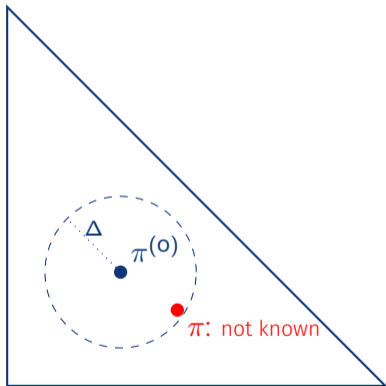


Figure 1: $\mathbb{P}(|Q_l - T_l| \geq \epsilon)$ vs n for various distributions, with $d = 1000$ and $\epsilon = 0.001$. uniform, zipf, uniform mixture, and random distributions.

LOCAL MIN-MAX PROBLEM FORMULATION



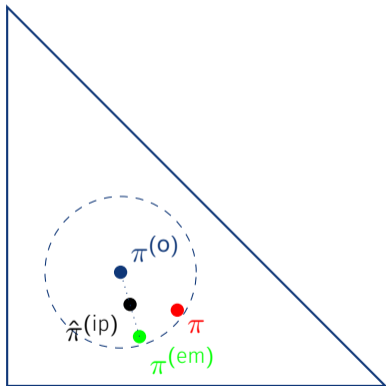
- We know $\pi^{(0)}$ such that

$$\|\pi - \pi^{(0)}\|_2 \leq \Delta.$$

- Estimator $\hat{\pi}$ is a function of $(X^n, \pi^{(0)}, \Delta)$

How to minimize the squared ℓ_2 -error?

INTERPOLATION ESTIMATOR



$$\hat{\pi}^{(ip)} = \alpha \pi^{(em)} + (1 - \alpha) \pi^{(o)}$$

Theorem: $\hat{\pi}^{(ip)}$ Error

Worst-case Risk:

$$\leq \min \left(\Delta^2, \frac{1 - (\|\pi^{(o)}\| - \Delta)^2}{n} \right).$$

SUMMARY OF RISK LOWER BOUNDS²

Setting	Upper Bound	Lower Bound
General $\pi^{(0)}$	$O\left(\min\left(\Delta^2, \frac{1 - (\ \pi^{(0)}\ - \Delta)^2}{n}\right)\right)$	$\Omega\left(\frac{(1 - \ \pi^{(0)}\ ^2)}{d} \cdot \min\left(\Delta^2, \frac{1}{n}\right)\right)$ LeCam's Method
Deterministic $\pi^{(0)} = (\mathbf{1}, \mathbf{0}, \dots, \mathbf{0})$	$O(\min(\Delta^2, \frac{\Delta}{n}))$	$\Omega(\min(\Delta^2, \frac{\Delta}{n}))$ LeCam's Method
Uniform $\pi^{(0)} = (\frac{1}{d}, \frac{1}{d}, \dots, \frac{1}{d})$	$O(\min(\Delta^2, \frac{1}{n}))$	$\Omega(\min(\Delta^2, \frac{1}{n}))$ Assouad's Lemma

²Distribution Estimation with Side Information. H. Balasundaram, A. Thangaraj. Submitted to ISIT.

SIMULATIONS

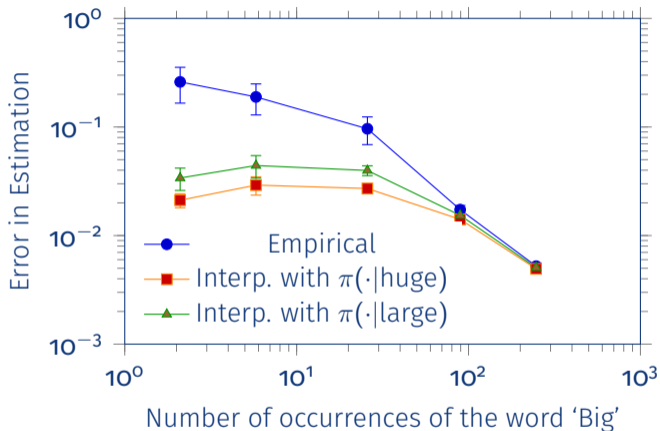


Figure 2: Estimation errors vs. number of samples for the Empirical and Interpolation Estimators for $\pi^{(0)}$ from the dataset.

QUESTIONS?
