

The Broad Optimality of Profile Maximum Likelihood

Yi Hao and Alon Orlitsky

University of California San Diego {yih179, alon}@ucsd.edu

Discrete Distributions

- Discrete support set \mathcal{X}
 $\{\text{heads, tails}\} = \{\text{h, t}\} \quad \{\dots, -1, 0, 1, \dots\} = \mathbb{Z}$
- Distribution p over \mathcal{X} , probability p_x for $x \in \mathcal{X}$
 $p_x \geq 0 \quad \sum_{x \in \mathcal{X}} p_x = 1$
 $p = (p_{\text{h}}, p_{\text{t}}) \quad p_{\text{h}} = .6, p_{\text{t}} = .4$
- \mathcal{P} collection of distributions
- $\mathcal{P}_{\mathcal{X}}$ all distributions over \mathcal{X}
 $\mathcal{P}_{\{\text{h, t}\}} = \{(p_{\text{h}}, p_{\text{t}})\} = \{(.6, .4), (.4, .6), (.5, .5), (0, 1), \dots\}$

Distribution Property (Functional)

- $f : \mathcal{P}_{\mathcal{X}} \rightarrow \mathbb{R}$
- Maps distribution to real value

Shannon entropy	$H(p)$	$\sum_x p_x \log \frac{1}{p_x}$
Rényi entropy	$H_{\alpha}(p)$	$\frac{1}{1-\alpha} \log (\sum_x p_x^{\alpha})$
Support size	$S(p)$	$\sum_x \mathbb{1}_{p_x > 0}$
Support coverage	$S_m(p)$	$\sum_x (1 - (1 - p_x)^m)$
Expected # distinct symbols in m samples		
Distance to uniformity	$L_{\text{uni}}(p)$	$\sum_x p_x - \frac{1}{ \mathcal{X} } $

Property Estimation

- Given:** support set \mathcal{X} , property f
- Unknown:** $p \in \mathcal{P}_{\mathcal{X}}$
- Estimate:** $f(p)$
Entropy of English words
Given: $\mathcal{X} = \{\text{English words}\}$, unknown: p , estimate: $H(p)$
species in habitat
Given: $\mathcal{X} = \{\text{bird species}\}$, unknown: p , estimate: $S(p)$
- How to estimate $f(p)$ when p is unknown?**
- Many applications:
vocabulary and population estimation, database similarity, graphical model learning, neural spike trains, property testing ...

Learn from Examples

- Observe n independent samples
 $X^n = X_1, \dots, X_n \sim p$
- Estimate $f(p)$
- Estimator:** $f^{\text{est}} : \mathcal{X}^n \rightarrow \mathbb{R}$
- Estimate for $f(p)$: $f^{\text{est}}(X^n)$

Empirical Estimator

- n samples – X^n
- $p_x^{\text{emp}}(X^n) := (N_x : \# \text{ times } x \text{ appears in } X^n) / n$
 $\mathcal{X} = \{a, b, c\} \quad p = (p_a, p_b, p_c) = (.5, .3, .2)$
 $X^{10} = c, a, b, a, b, a, b, a, b, c$
 $p_a^{\text{emp}} = \frac{4}{10}, p_b^{\text{emp}} = \frac{4}{10}, p_c^{\text{emp}} = \frac{2}{10}$
- Maximum likelihood estimator (**MLE**): $p^{\text{ml}} = p^{\text{emp}}$

Empirical Plug-In Estimator

- $f^{\text{emp}}(X^n) = f(p^{\text{emp}}(X^n))$
- Advantages:
Plug-and-play: simple two steps
Universal and Intuitive: applies to all properties
- Best-known, most-used property estimator

Sample Complexity

- Probably Approximately Correct (PAC)
- Allowed additive approximation error $\epsilon > 0$
- Allowed error probability $\delta > 0$
- $n_f(f^{\text{est}}, \epsilon, \delta)$: number of samples f^{est} needs to approximate property f well:
 $|f^{\text{est}}(X^n) - f(p)| \leq \epsilon$
with probability $\geq 1 - \delta$, for all $p \in \mathcal{P}$

Empirical and Optimal Complexity

- \mathcal{P}_k all k -symbol distributions, $\epsilon \gtrsim n^{-0.1}, \delta = 1/3$

Property	$n_f(f^{\text{emp}}, \epsilon)$	$n_f(f^{\text{opt}}, \epsilon)$
Entropy	$k \cdot \frac{1}{\epsilon}$	$\frac{k}{\log k} \cdot \frac{1}{\epsilon}$
Support coverage	m	$\frac{m}{\log m} \cdot \log \frac{1}{\epsilon}$
Distance to uniform	$k \cdot \frac{1}{\epsilon^2}$	$\frac{k}{\log k} \cdot \frac{1}{\epsilon^2}$
Support size	$k \cdot \log \frac{1}{\epsilon}$	$\frac{k}{\log k} \cdot \log^2 \frac{1}{\epsilon}$

- For support size, $\mathcal{P}_{\geq 1/k} := \{p \mid p_x \geq 1/k, \forall x \in \mathcal{X}\}$
- Support size and coverage normalized by k and m

Profiles

- iid:** order doesn't matter
- Symmetric properties:** labels don't matter
- Profile: # elements appearing any given # times
- (h,h,t), (t,t,h), (h,t,h), (t,h,t), (t,h,h) same entropy
1 element appeared once, 1 element twice
Profile: $\varphi = \{1, 2\}$
- $\varphi(x^n)$: multiset of symbol frequencies in x^n

Profile Probability

- Probability of observing φ when sampling from p

$$p(\varphi) := \sum_{y^n: \varphi(y^n) = \varphi} p(y^n) = \sum_{y^n: \varphi(y^n) = \varphi} \prod_{i=1}^n p(y_i)$$

- $p, q \quad p + q = 1$
 $\Pr(\varphi = \{1, 2\}) = 3(p^2q + q^2p)$

Profile Maximum Likelihood (PML)

- Distribution maximizing profile probability
- Maps x^n to $p_{\varphi(x^n)}^{\text{pml}} := \operatorname{argmax}_{p \in \mathcal{P}} p(\varphi(x^n))$

Broad optimality of PML

- PML – unified, time- and sample-optimal for
- Additive property estimation
 - Non-additive property – Rényi entropy estimation
 - Sorted distribution estimation
 - Identity/Uniformity testing

Additive Property Estimation

- Additive property: $f(p) = \sum_x f_x(p_x)$
entropy, support size, distance to uniformity ...
- For **all** symmetric, additive, properly Lipschitz, properties, for $n \geq n_f(|\mathcal{X}|, \epsilon, 1/3)$ and $\epsilon \gtrsim n^{-0.1}$,
 $\Pr\left(\left|f\left(p_{\varphi(X^n)}^{\text{pml}}\right) - f(p)\right| > 5\epsilon\right) \leq \exp(-\sqrt{n})$
- With **four** times the **optimal** # samples for error probability 1/3, PML plug-in achieves **much lower** error probability
- Near **linear-time** (A)PML approximation [CSS19]

Rényi, Distribution, Testing

- Rényi entropy**
Integer $\alpha > 1$, PML has optimal $k^{1-1/\alpha}$ complexity
Non-integer $\alpha > 3/4$, PML improves best-known
- Sorted distribution estimation**
(A)PML yields optimal $\Theta(k/(\epsilon^2 \log k))$ sample complexity under ℓ_1 distance
- Uniformity testing:** $p = p_u$ v.s. $|p - p_u| \geq \epsilon$;
optimal $\Theta(\sqrt{k}/\epsilon^2)$ up to logarithmic factors of k :

Input: params k, ϵ , and a sample with profile φ
If $\exists N_x \geq 3 \max\{1, n/k\} \log k$, return 1
If $\|p_{\varphi}^{\text{pml}} - p_u\|_2 \geq 3\epsilon/(4\sqrt{k})$, return 1; else, 0