Unified Maximum Likelihood Estimation of Symmetric Distribution Properties

Jayadev Acharya Hirakendu Das Alon Orlitsky Ananda Suresh

Cornell Yahoo UCSD Google

Frontiers in Distribution Testing Workshop FOCS 2017

Special Thanks to...



I don't smile



I only smile

Symmetric properties

 $\mathcal{P} = \{ (p_1, \dots p_k) \}$ - collection of distributions over $\{1, \dots, k\}$

Distribution property: $f: \mathcal{P} \rightarrow \mathbb{R}$

f is symmetric if unchanged under input permutations

Entropy
$$H(p) \triangleq \sum p_i \log \frac{1}{p_i}$$

Support size $S(p) \triangleq \sum_{i} \mathbb{I}_{\{p_i > 0\}}$

Rényi entropy, support coverage, distance to uniformity, ...

Determined by the probability multiset $\{p_1, p_2, ..., p_k\}$

Symmetric properties

 $p = (p_1, ..., p_k)$, (*k* finite or infinite) - discrete distribution

f(p), a property of p

f(p) is symmetric if unchanged under input permutations

Entropy $H(p) \triangleq \sum p_i \log \frac{1}{p_i}$

Support size $S(p) \triangleq \sum_{i} \mathbb{I}_{\{p_i > 0\}}$

Rényi entropy, support coverage, distance to uniformity, ...

Determined by the probability multiset $\{p_1, p_2, ..., p_k\}$

Property estimation

 $oldsymbol{p}$ unknown distribution in ${\mathcal P}$

Given independent samples $X_1, X_2, ..., X_n \sim p$

Estimate **f**(**p**)

Sample complexity $S(f, k, \varepsilon, \delta)$

Minimum *n* necessary to

Estimate $f(p) \pm \varepsilon$

With error probability $< \delta$

Plug-in estimation

Use $X_1, X_2, \dots, X_n \sim p$ to find an estimate \hat{p} of p

Estimate f(p) by $f(\hat{p})$

How to estimate **p**?

Sequence Maximum Likelihood (SML)

$$p^{\text{sml}} = \arg \max_{p} p(x^{n}) = \arg \max_{p(x)} \Pi_{i} p(x_{i})$$

$$x^{3} = h, h, t$$

$$p^{\text{sml}}_{h,h,t} = \arg \max p^{2}(h) \cdot p(t)$$

$$p^{\text{sml}}_{h,h,t}(h) = 2/3 \qquad p^{\text{sml}}_{h,h,t}(t) = 1/3$$

Same as empirical-frequency distribution

Multiplicity N_x - # times x appears in x^n

$$p^{\mathrm{sml}}(x) = \frac{N_x}{n}$$

Prior Work

For several important properties

Empirical-frequency plugin requires $\Theta(\mathbf{k})$ samples New complex (non-plugin) estimators need $\Theta\left(\frac{k}{\log k}\right)$ samples

Property	Notation	SML	Optimal	References
Entropy	H(p)	$\frac{k}{\varepsilon}$	$\frac{k}{\log k} \frac{1}{\varepsilon}$	(Valiant & Valiant, 2011a; Wu &
				Yang, 2016; Jiao et al., 2015)
Support size	$\frac{S(p)}{k}$	$k \log \frac{1}{\varepsilon}$	$\frac{k}{\log k}\log^2\frac{1}{\varepsilon}$	(Wu & Yang, 2015)
Support coverage	$\frac{S_m(p)}{m}$	m	$\frac{m}{\log m}\log \frac{1}{\varepsilon}$	(Orlitsky et al., 2016)
Distance to <i>u</i>	$ p - u _1$	$\frac{k}{\varepsilon^2}$	$\frac{k}{\log k} \frac{1}{\varepsilon^2}$	(Valiant & Valiant, 2011b; Jiao
			5	et al., 2016)

Different estimator for each property

Use sophisticated approximation theory results

Entropy estimation

- SML estimate of entropy = $\sum_{x} \frac{N_x}{n} \log \frac{n}{N_x}$
- Sample complexity: $\Theta(k/\varepsilon)$
- Various corrections proposed: Miller-Maddow, Jackknifed estimator, Coverage adjusted, ...
- Sample complexity: $\Omega(k)$ for all the above estimators

Entropy estimation

[Paninski'03]: o(k) sample complexity (existential)

[ValiantValiant'11a]: Constructive LP based methods: $\Theta_{\varepsilon}\left(\frac{k}{\log k}\right)$

[ValiantValiant11b, WuYang'14, HanJiaoVenkatWeissman'14]:

Simplified algorithms, and growth rate: $\Theta\left(\frac{k}{\epsilon \log k}\right)$

New (as of August) Results

Unified, simple, sample-optimal approach for all above problems

Plug-in estimator, replace **sequence** maximum likelihood with **profile** maximum likelihood

Profiles

h, *h*, *t* or *h*, *t*, *h* or *t*, *h*, *t* \rightarrow same estimate One element appeared once, on appeared twice

Profile: Multi-set of multiplicities: $\Phi(X_1^n) = \{N_x : x \in X_1^n\}$

$$\Phi(\boldsymbol{h},\boldsymbol{h},\boldsymbol{t}) = \Phi(\boldsymbol{t},\boldsymbol{h},\boldsymbol{t}) = \{1,2\}$$

$$\Phi(\alpha, \gamma, \beta, \gamma) = \{1, 1, 2\}$$

Sufficient statistic for symmetric properties

Profile maximum likelihood [+svz'04]

Profile probability

$$p(\Phi) = \sum_{\Phi(x_1^n) = \Phi} p(x_1^n)$$

Maximize the profile probability

$$p_{\Phi}^{pml} = \arg\max_{p} p(\Phi(X_1^n))$$

See "On estimating the probability multiset", Orlitsky, Santhanam, Viswanathan, Zhang for a detailed treatment, and an argument for competitive distribution estimation.

Profile maximum likelihood (PML) [+SVZ '04]

Profile probability

$$p(\Phi) = \sum_{x^n: \Phi(x_1^n) = \Phi} p(x^n)$$

Distribution Maximizing the profile probability $p_{\Phi}^{pml} = \arg\max_{p} p(\Phi)$

PML competitive for distribution estimation

PML example

 $X^{3} = h, h, t$

 $p_{h,h,t}^{sml}(h)=2/3$ $p_{h,h,t}^{sml}(t)=1/3$ $\Phi(h, h, t) = \{1, 2\}$ $p(\Phi = \{1,2\}) = p(s,s,d) + p(s,d,s) + p(d,s,s)$ = 3p(s,s,d) $= \binom{3}{1} \left(\sum_{x \neq y} p^2(x) p(y) \right)$ $p^{sml}(\{1,2\}) = {\binom{3}{1}} \left({\binom{2}{3}}^2 {\binom{1}{3}} + {\binom{1}{3}}^2 {\binom{2}{3}} \right) = \frac{18}{27} = \frac{2}{3}$

PML of {1,2}

 $P({1,2}) = p(s,s,d)+p(s,d,s)+p(d,s,s) = 3 p(s,s,d)$

$$(1/2, 1/2) \rightarrow p(s, s, d) = 1/8 + 1/8 = 1/4$$

$$p(s, s, d) = \Sigma_{x \neq y} p^{2}(x) p(y)$$

= $\Sigma_{x} p^{2}(x) (1 - p(x))$
 $\leq \frac{1}{4} \Sigma_{x} p(x) = \frac{1}{4}$

 p^{pml} (s,s,d) = $\frac{1}{4}$

 p^{pml} ({1,2}) = ³/₄ Recall: p^{sml} ({1,2})=2/3

PML({1,1,2})

$$\Phi(\alpha, \gamma, \beta, \gamma) = \{1, 1, 2\}$$
$$p^{pml}(\{1, 1, 2\}) = U[5]$$

PML can predict existence of new symbols

Profile maximum likelihood

PML of {1,2} is {½, ½}

$$p^{pml}(\{1,2\}) = {3 \choose 1} \left(\left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right) \right) = \frac{3}{4} > \frac{18}{27}$$

$$\Sigma_{x \neq y} p^2(x) p(y) = \Sigma_x p^2(x) \left(1 - p(x)\right) \le \frac{1}{4} \Sigma_x p(x) = \frac{1}{4}$$

$$X^n = \alpha, \gamma, \beta, \gamma, \quad \Phi(X^n) = \{1,1,2\}$$

$$p^{pml}(\{1,1,2\}) = U[5]$$

PML can predict existence of new symbols

PML Plug-in

To estimate a symmetric property fFind $p^{pml}(\Phi(X^n))$ Output $f(p^{pml})$

Simple

Unified

No tuning parameters

Some experimental results (c. 2009)





U[500], 350x, 12 experiments





U[500], 700x, 12 experiments



Staircase

15K elements, 5 steps, ~3x 30K samples Observe 8,882 elts 6,118 missing



Zipf



1990 Census - Last names

SMITH	1.006	1.006	1
JOHNSON	0.810	1.816	2
WILLIAMS	0.699	2.515	3
JONES	0.621	3.136	4
BROWN	0.621	3.757	5
DAVIS	0.480	4.237	6
MILLER	0.424	4.660	7
WILSON	0.339	5.000	8
MOORE	0.312	5.312	9
TAYLOR	0.311	5.623	10
AMEND	0.001	77.478	18835
ALPHIN	0.001	77.478	18836
ALLBRIGHT	0.001	77.479	18837
AIKIN	0.001	77.479	18838

ACRES

ZUPAN

ZEOLLA

ZUCHOWSKI

0.001 77.480 18839

0.000 77.480 18840

0.000 77.481 18841

18842

0.000 77.481

18,83	anames
77.48	% population
~230	million

1990 Census - Last names

18,839 last names based on ~230 million35,000 samples, observed 9,813 names



Coverage (# new symbols)

Zipf distribution over 15K elements

Sample 30K times

Estimate: # new symbols in sample of size λ * 30K

 $\begin{array}{l} \lambda < 1\\ \text{Good-Toulmin:}\\ \lambda > 1\\ \text{Estimate PML & predict}\\ \text{Extends to } \lambda > 1\\ \text{Applies to other properties} \end{array}$





Finding the PML distribution

EM algorithm [+ Pan, Sajama, Santhanam, Viswanathan, Zhang '05 - '08]

Approximate PML via Bethe Permanents [Vontobel]

Extensions of Markov Chains [Vatedka, Vontobel]

No provable algorithms known

Motivated Valiant & Valiant

Maximum Likelihood Estimation Plugin

General property estimation technique

 ${\mathcal P}$ - collection of distributions over domain ${\mathcal Z}$

 $f: \mathcal{P} \to \mathbb{R}$ any property (say entropy)

MLE estimator

Given
$$z \in \mathcal{Z}$$

Determine $p_z^{\text{MLE}} \triangleq \arg \max_{p \in \mathcal{P}} p(z)$
Output $f(p_z^{\text{MLE}})$

How good is MLE?

Competitiveness of MLE plugin

 ${\mathcal P}$ - collection of distributions over domain ${\mathcal Z}$

 $\hat{f}: \mathbb{Z} \to \mathbb{R}$ any estimator such that $\forall p \in \mathcal{P}, \ \mathbb{Z} \sim p$ $\Pr(|f(p) - \hat{f}(\mathbb{Z})| > \varepsilon) < \delta$

MLE plugin error bounded by

$$\Pr(\left|f(p) - f(p_z^{\text{MLE}})\right| > 2 \cdot \varepsilon) < \delta \cdot |\mathcal{Z}|$$

Simple, universal, competitive with any \hat{f}

Quiz: Probability of unlikely outcomes

6-sided die, $p=(p_1, p_2, ..., p_6)$

 $p_i \ge 0$, and $\Sigma p_i = 1$, otherwise arbitrary

 $Z \sim p$

 $P_r(p_z \le 1/6) \quad \text{Can be anything}$ (1,0,...,0) $\rightarrow \Pr(p_z \le 1/6) = 0$ (1/6, ..., 1/6) $\rightarrow \Pr(p_z \le 1/6) = 1$ $P_r(p_z \le 0.01) \le 0.06$ $P_r(p_z \le 0.01) = \Sigma_{i:p_i \le 0.01} p_i \le 6 \cdot 0.01 = 0.06$ Competitiveness of MLE plugin - proof $\hat{f}: \mathcal{Z} \to \mathbb{R}: \forall p \in \mathcal{P}, Z \sim p \to \Pr(|f(p) - \hat{f}(Z)| > \varepsilon) < \delta$ then $\Pr(|f(p) - f(p_Z^{\text{MLE}})| > 2 \cdot \varepsilon) < \delta \cdot |\mathcal{Z}|$

For all z such that $p(z) \ge \delta$: 1) $|f(p) - \hat{f}(z)| \le \varepsilon$

2) $p_z^{\text{MLE}}(z) \ge p(z) > \delta$, hence $|f(p_z^{\text{MLE}}) - \hat{f}(z)| \le \varepsilon$

Triangle inequality: $|f(p_z^{\text{MLE}}) - f(p)| \le 2\varepsilon$

If $|f(p_z^{\text{MLE}}) - f(p)| > 2\varepsilon$ then $p(z) < \delta$,

$$\Pr\left(\left|f\left(p_{z}^{\text{MLE}}\right) - f(p)\right| > 2\varepsilon\right) \le \Pr(p(Z) < \delta) \le \sum_{p(Z) < \delta} p(Z) \le \delta \cdot |\mathcal{Z}|$$

PML performance bound

If $n = S(f, k, \varepsilon, \delta)$, then $S^{pml}(f, k, 2 \cdot \varepsilon, |\Phi^n| \cdot \delta) \le n$

 $|\Phi^n|$: number of profiles of length n

Profile of length n: partition of n

 $\{3\}, \{1,2\}, \{1,1,1\} \rightarrow 3, 2+1, 1+1+1$

 $|\Phi^n| = partition \# of n$

Hardy-Ramanujan: $|\Phi^n| < e^{3\sqrt{n}}$

Easy: $e^{\log n \cdot \sqrt{n}}$

If $n = S(f, k, \varepsilon, e^{-4\sqrt{n}})$, then $S^{pml}(f, k, 2\varepsilon, e^{-\sqrt{n}}) \le n$

Summary

Symmetric property estimation PML plug-in approach Simple Universal Sample optimal for known sublinear properties

Future directions

Provably efficient algorithms

Independent proof technique

Thank You!

PML for symmetric f

For any symmetric property f,

if $n = S(f, k, \varepsilon, 0.1)$, then $S^{pml}(f, k, 2 \cdot \varepsilon, 0.1) = O(n^2)$.

Proof. By median trick, $S(f, k, \varepsilon, e^{-m}) = O(n \cdot m)$. Therefore, $S^{pml}(f, k, 2 \cdot \varepsilon, e^{3\sqrt{n \cdot m} - m}) = O(n \cdot m)$,

Plugging in $m = C \cdot n$, gives the desired result.

Better error probabilities – warm up

Estimating a distribution p over [k] to L1 distance ε w.p. > 0.9 requires $\Theta(k/\varepsilon^2)$ samples.

Proof: Exercise.

Estimating a distribution p over [k] to L1 distance ε w.p. $1 - e^{-k}$ requires $\Theta(k/\varepsilon^2)$ samples.

Proof:

- Empirical estimator \hat{p}
- $|\hat{p} p|$ has bounded difference constant (b.d.c.) 2/n
- Apply McDiarmid's inequality

Better error probabilities

Recall

$$S(H,\varepsilon,k,2/3) = \Theta\left(\frac{k}{\varepsilon \cdot \log k}\right)$$

- Existing optimal estimators: high b.d.c.
- Modify them to have small b.d.c., and still be optimal
- In particular, can get b.d.c. = $n^{-0.95}$ (exponent close to 1)
- With twice the samples error drops **super-fast**

$$S(H,\varepsilon,k,e^{-n^{0.9}}) = \Theta\left(\frac{k}{\varepsilon \cdot \log k}\right)$$

Similar results for other properties

Even approximate PML works!

- Perhaps finding exact PML is hard.
- Even approximate PML works.

Find a distribution q such that

$$q(\Phi(X_1^n)) \ge e^{-n^{\cdot 8}} \cdot p^{pml}(\Phi(X_1^n))$$

Even this is optimal (for large k)

In Fisher's words ...

Of course nobody has been able to prove that MLE is best under all circumstances. MLE computed with all the information available may turn out to be inconsistent. Throwing away a substantial part of the information may render them consistent.

R. A. Fisher

Proof of PML performance

If $n = S(f, k, \varepsilon, \delta)$, then $S^{pml}(f, k, 2 \cdot \varepsilon, |\Phi^n| \cdot \delta) \le n$ $S(f, k, \varepsilon, \delta)$, achieved by an estimator $\hat{f}(\Phi(X^n))$

• Profiles $\Phi(X^n)$ such that $p(\Phi(X^n)) > \delta$,

 $p^{PML}(\Phi) \ge p(\Phi) > \delta$

 $\left|f\left(p_{\Phi}^{PML}\right) - f(p)\right| \leq \left|f\left(p_{\Phi}^{PML}\right) - \hat{f}(\Phi)\right| + \left|\hat{f}(\Phi) - f(p)\right| < 2\varepsilon$

• Profiles with $p(\Phi(X^n)) < \delta$,

 $p(p(\Phi(X^n) < \delta) < \delta \cdot |\Phi^n|$