Estimation of Entropy in Constant Space

Maryam Aliakbarpour

Boston University and Northeastern University

Joint work Andrew McGregor, Jelani Nelson, and Erik Waingarten.



Image from: https://tilics.dmi.unibas.ch/the-turing-machine

General inference model



Example

Memory << size of data



Estimation with memory constraints





Prior work:

Problems: parity learning, learning PDFs, learning concept classes, robust estimation of statistics, distribution testing, estimating moments,

[Raz. FOCS16][Crouch, McGregor, Valiant, Woodruff, ESA 2016] [Guha, McGregor. AISTATS 2007], [Chien, Ligett, McGregor. ITS 2010] [Steinhardt, Valiant, Wagner. COLT 2016]
[Esfandiari, Hajiaghayi, Liaghat, Monemizadeh. SODA 2015] [Moshkovitz, Moshkovitz. COLT 2017] [Kol, Raz, Tal. STOC 2017] [Raz. FOCS 2017] [Garg, Kothari, Raz. STOC 2018] [Sharam, Sidford, Valiant. STOC 2019] [Diakonikolas, Gouleakis, Kane, Rao. COLT 2019] [Garg, Raz, Tal. Complexity 2019] [Acharya, Bhadane, Indyk, Sun, NeurIPS 2019] [Garg, Kothari, Raz. RANDOM 2020] [Garg et al. RANDOM 2021] [Brown, Bun, Smith. COLT 2022]...

This work: estimating entropy

Shannon's Entropy of $D = (p_1, p_2, ..., p_n)$:

$$H(D) \coloneqq \sum_{i=1}^{n} p_i \log_2 1/p_i$$

Entropy

Information theory

In information theory, the entropy of a random variable is the average level of "information", "surprise", or "uncertainty" inherent to the variable's possible outcomes. Wikipedia

Feedback

This work: estimating entropy Shannon's Entropy of $D = (p_1, p_2, ..., p_n)$:



Goal:

Memory constrains: O(1) words

Previous results

No memory constraint:

$$\Theta\left(\frac{n}{\epsilon \log n} + \frac{\log^2 n}{\epsilon^2}\right)$$

[Batu, Dasgupta, Kumar, Rubinfeld. STOC 2002][Valiant, Valiant. FOCS 2011] [Valiant, Valiant. JACM 2017] [Wu, Yang. IEEE Trans. IT 2016] [Jiao et al. IEEE Trans. IT 2015] (and many more)

Previous results

No memory constraint:

$$\Theta\left(\frac{n}{\epsilon \log n} + \frac{\log^2 n}{\epsilon^2}\right)$$

[Batu, Dasgupta, Kumar, Rubinfeld. STOC 2002][Valiant, Valiant. FOCS 2011] [Valiant, Valiant. JACM 2017] [Wu, Yang. IEEE Trans. IT 2016] [Jiao et al. IEEE Trans. IT 2015] (and many more)

O(1) words of memory:

$$0\left(\frac{n\log{(1/\epsilon)^3}}{\epsilon^3}\right)$$

[Acharya, Bhadane, Indyk, Sun, NeurIPS 2019]

Our results

This work O(1) words of memory:

$$0\left(\frac{n\log(1/\epsilon)^4}{\epsilon^2}\right)$$

O(1) words of memory:

$$0\left(\frac{n\log\left(1/\epsilon\right)^3}{\epsilon^3}\right)$$

[Acharya, Bhadane, Indyk, Sun, NeurIPS 2019]

A closely related problem

Estimating empirical entropy in the data streaming



Entropy estimation with no memory constraint

No memory constraint

Algorithm [Valiant and Valiant'11]:

1. Compute the fingerprint of the samples Count numbers of elements appeared *i* times

<u>List</u>

3 1 3 8 7 3 1 5

<u>Fingerprints</u> three elements appeared once.

One element appeared twice.

One element appeared three times.

No memory constraint

Algorithm [Valiant, Valiant'11]:

- 1. Compute the fingerprint of the samples
- 2. Come up with a histogram of a distribution that is likely to generate



Plots from [Valiant, Valiant'11]

No memory constraint

Algorithm [Valiant, Valiant'11]:

- 1. Compute the fingerprint of the samples
- 2. Come up with a histogram of a distribution that is likely to generate
- 3. Output a distribution that is compatible with the histogram

Works well ignoring the labels!

Entropy ✓ Support size ✓

Adding memory constraints

Computing fingerprint is hard when we cannot memorize

List 3 1 3 8 7 3 1 5

<u>Fingerprints</u> three elements appeared once.

One element appeared twice.

One element appeared three times.

Entropy estimation with no memory constraint

A simple approach

Simple algorithm

$$H(D) \coloneqq \sum_{i=1}^{n} p_i \cdot \log 1/p_i = \mathbb{E}_{i \sim D}[\log 1/p_i]$$

- 1. Repeat r times
 - 1. Draw $i \sim D$.
 - 2. $\hat{p}_i \leftarrow \text{Estimate } p_i$
- 2. Output average of log $1/\hat{p}_i$'s.

Simple algorithm

$$H(D) \coloneqq \sum_{i=1}^{n} p_i \cdot \log 1/p_i = \mathbb{E}_{i \sim D}[\log 1/p_i]$$

- 1. Repeat *r* times
 - 1. Draw $i \sim D$.
 - 2. $\hat{p}_i \leftarrow \text{Estimate } p_i$

Via negative binomial distribution Draw samples until *t* copies of *i* are observed. $X_i \leftarrow \frac{1}{t} \cdot (\# \text{ Observed samples})$ $E[X_i]$ is precisely $1/p_i$.

2. Output average of log $1/\hat{p}_i$'s.

Simple algorithm

$$H(D) \coloneqq \sum_{i=1}^{n} p_i \cdot \log 1/p_i = \mathbb{E}_{i \sim D}[\log 1/p_i]$$

- 1. Repeat *r* times
 - 1. Draw $i \sim D$.
 - 2. $\hat{p}_i \leftarrow \text{Estimate } p_i$

Via negative binomial distribution Draw samples until *t* copies of *i* are observed. $X_i \leftarrow \frac{1}{t} \cdot (\# \text{ Observed samples})$ $E[X_i]$ is precisely $1/p_i$.

2. Output average of log $1/\hat{p}_i$'s.





Entropy estimation with no memory constraint

A simple better approach

Remove bias

Idea: Estimate bias and decrease it from
$$\widehat{H}$$
.
Bias = $|E_{i\sim D}[\log 1/p_i] - E_{i\sim D}[\log X_i]| = |E_{i\sim D}[\log Y_i]|$

 $E_{i\sim D}[Y_i] = 1$. Taylor expansion around Y = 1:

Bias =
$$E_{i \sim D}[\log Y_i] = E\left[Y_i - 1 - \frac{(Y_i - 1)^2}{2} + \frac{(Y_i - 1)^3}{3} - \cdots\right]$$





 $\Pr[k \text{ samples are equal}] = p_i^k$

Remove $\log n$ factors



Buckets of large X_i can be computed with less accuracy.

Conclusion

This work O(1) words of memory:

$$0\left(\frac{n\log(1/\epsilon)^4}{\epsilon^2}\right)$$

Open question: can we improve the lower bound to $\Omega\left(\frac{n}{\epsilon^2}\right)$?

Thank you.