## How Much is Unseen Depends Chiefly On Information About the Seen

### Seongmin Lee and Marcel Böhme

### **Abstract & Contribution**



### **Missing Mass Problem.**

Given the sample  $X^n$  from the unknown multinomial distribution  $\mathcal{D}$ , what is the probability  $M_0$  that the next sample  $X_{n+1}$  has never been seen before?

• Missing mass represents the <u>representativeness of the training data</u>; if the missing mass regarding the training data is <u>high</u> 😨, the model trained by the data will <u>often face unseen labels</u> in the test data 😥.

### In this work,

- We show exactly to what extent  $\mathbb{E}[M_0]$  can be estimated from the sample  $X^n$ and how much remains  $\Rightarrow$  minimal bias estimator  $\hat{M}_0^B$  of.
- We define a class of nearly unbiased estimators of  $M_0$  is from different representations of  $\mathbb{E}[M_0]$ .
- We cast the distribution-free estimation of  $M_0$  as a search problem from  $\mathbf{m}$ whose goal is to find a distribution-specific estimator with a minimized MSE  $\Rightarrow$  minimal MSE estimator  $M_0^{Evo}$   $\textcircled{\bullet}$ .

### **Problem Definition**

• Frequency: Example. Sample  $X^n$  of size 10  $N_x(X^n) = \sum \mathbf{1}(X_i = x)$  $X_2 \mid X_3 \mid X_4 \mid$  $X_5$  $X_1$ • Frequency of frequencies:  $\Phi_k(X^n) = \sum \mathbf{1}(N_x(X^n) = k)$ •  $N_1 = 3, N_2 = 2, • \Phi_1 = 3,$  $N_3 = 2, N_4 = 1, \Phi_2 = 2,$ • Expected  $\Phi$ :  $f_k(n) = \mathbb{E}_{X^n \sim \mathcal{D}} \left[ \Phi_k(X^n) \right]$  $N_5 = 1, N_6 = 1$  $\Phi_{3} = 1$  $= \sum \mathbb{E}_{X^n \sim \mathcal{D}} \left[ \mathbf{1}(N_x(X^n) = k) \right]$ • Missing Mass:  $M_0$  = - - - - - - - - - -Good-Turing estim





K <sub>7</sub>	$X_8$	<i>X</i> <sub>9</sub>	<i>X</i> <sub>10</sub>	
5	6	3	1	

 $X_6$ 

**Missing Mass:**  $Pr(X_{11} \notin [1,6])$ ?

$$= \sum_{x=1}^{S} p_{x} \mathbf{1}(N_{x}(X^{n}) = 0)$$
  
hator [1]:  $\hat{M}_{0}^{G} = \frac{\Phi_{1}}{m}$ 

ľι

### (1) Expected Missing Mass

$$\Xi[M_0] = \sum_{x=1}^{S} p_x (1 - p_x)^n = g_1(n + 1)$$



$$\hat{M}_{0}^{B} = -\sum_{i=1}^{n} (-1)^{i} \frac{\Phi_{i}}{\binom{n}{i}} =$$



"The average missing mass  $\mathbb{E}[M_0]$  depends chiefly on the average frequencies of frequencies  $\mathbb{E}[\Phi_k]!''$ 

[1] Irving J Good. The population frequencies of species and the estimation of population parameters. Biometrika, 40(3-4):237–264, 1953. 2] Alon Orlitsky and Ananda Theertha Suresh. Competitive distribution estimation: Why is good-turing good. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015

# MAX PLANCK INSTITUTE FOR SECURITY AND PRIVACY



"The MSE of  $\hat{M}^B_0$  is larger than  $\hat{M}^G_0$  due to the variance terms,  $Var(\Phi_k)$ ,  $Cov(\Phi_k, \Phi_l)$ ."

7.97e-03 0.92 72

0.89 7

8.02e-03 0.87 81

.10e-02

9.90e-03

diri-0.5

Avg.

		n = S				n = 2S		
0	$MSE(\hat{M}_0^G)$	$MSE(M_0^{\rm Evo})$	$\hat{A}_{12}$	Ratio	$MSE(\hat{M}_0^G)$	$MSE(M_0^{\rm Evo})$	$\hat{A}_{12}$	Ratio
$\sim$	6.05e-03	4.29e-03	0.97	70%	1.93e-03	1.73e-03	0.96	89%
%	5.46e-03	4.07e-03	0.98	74%	1.57e-03	1.42e-03	0.93	90%
%	3.42e-03	3.04e-03	0.89	88%	1.26e-03	1.08e-03	0.94	85%
%	5.23e-03	4.16e-03	0.96	79%	1.73e-03	1.54e-03	0.97	88%
%	4.36e-03	3.47e-03	0.92	79%	1.23e-03	1.05e-03	0.91	85%
%	3.47e-03	2.86e-03	0.88	82%	9.41e-04	8.08e-04	0.86	85%
%			0.93	79%			0.93	87%

the Good-Turing estimator  $\hat{M}_0^G$ ."