

# Chebyshev polynomials, moment matching and optimal estimation of the unseen

Yihong Wu

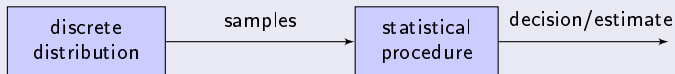
Department of ECE  
University of Illinois at Urbana-Champaign  
`yihongwu@illinois.edu`

Joint work with Pengkun Yang (Illinois)

Mar 17, 2014

## Task

Given samples from a discrete distribution, how to make statistical inference on certain **property** of the distribution?



- Support size:

$$S(P) = \sum_i \mathbf{1}_{\{p_i > 0\}}$$

- Example:

$$S \left( \text{---} \left( \begin{array}{cccccccccc} \circ & \circ & \circ & \circ & \circ & \circ & \circ & \circ & \circ & \circ \\ | & & & | & | & & & & & | \\ \circ & & & \circ & \circ & & & & & \circ \end{array} \right) \text{---} \right) = 5$$

- $\Leftrightarrow$  estimating the number of unseens (SEEN + UNSEEN =  $S(P)$ )



- maybe the Egyptians have studied it...

- maybe the Egyptians have studied it...
- Ecology:

THE RELATION BETWEEN THE NUMBER OF SPECIES AND  
THE NUMBER OF INDIVIDUALS IN A RANDOM SAMPLE  
OF AN ANIMAL POPULATION

BY R. A. FISHER (*Galton Laboratory*), A. STEVEN CORBET (*British Museum, Natural History*)  
AND C. B. WILLIAMS (*Rothamsted Experimental Station*)

- maybe the Egyptians have studied it...
- Ecology:

THE RELATION BETWEEN THE NUMBER OF SPECIES AND  
THE NUMBER OF INDIVIDUALS IN A RANDOM SAMPLE  
OF AN ANIMAL POPULATION

BY R. A. FISHER (*Galton Laboratory*), A. STEVEN CORBET (*British Museum, Natural History*)  
AND C. B. WILLIAMS (*Rothamsted Experimental Station*)

- Linguistics, numismatics, etc:

**Estimating the number of unseen species: How many  
words did Shakespeare know?**

BY BRADLEY EFRON AND RONALD THISTED  
*Department of Statistics, Stanford University, California*

- maybe the Egyptians have studied it...
- Ecology:

THE RELATION BETWEEN THE NUMBER OF SPECIES AND  
THE NUMBER OF INDIVIDUALS IN A RANDOM SAMPLE  
OF AN ANIMAL POPULATION

BY R. A. FISHER (*Galton Laboratory*), A. STEVEN CORBET (*British Museum, Natural History*)  
AND C. B. WILLIAMS (*Rothamsted Experimental Station*)

- Linguistics, numismatics, etc:

**Estimating the number of unseen species: How many  
words did Shakespeare know?**

BY BRADLEY EFRON AND RONALD THISTED  
*Department of Statistics, Stanford University, California*

- Will not discuss probability estimation  
[Good-Turing, Orlitsky et al., ...]



- Data:  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$

- Data:  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$
- Estimate:  $\hat{S} = \hat{S}(X_1, \dots, X_n)$  close to  $S(P)$  in prob or expectation

- Data:  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$
- Estimate:  $\hat{S} = \hat{S}(X_1, \dots, X_n)$  close to  $S(P)$  in prob or expectation
- Goal: find minimal sample size & fast algorithms

- Data:  $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P$
- Estimate:  $\hat{S} = \hat{S}(X_1, \dots, X_n)$  close to  $S(P)$  in prob or expectation
- Goal: find minimal sample size & fast algorithms
- Need to assume minimum non-zero mass

## Space of distributions

$\mathcal{D}_k \triangleq \{\text{prob distributions whose non-zero mass is at least } 1/k\}$

## Space of distributions

$\mathcal{D}_k \triangleq \{\text{prob distributions whose non-zero mass is at least } 1/k\}$

## Sample complexity

$n^*(k, \epsilon) \triangleq \min\{n : \exists \hat{S}, \text{s.t. } \mathbb{P}[|\hat{S} - S(P)| \leq \epsilon k] \geq 0.5, \forall P \in \mathcal{D}_k\}$

## Space of distributions

$\mathcal{D}_k \triangleq \{\text{prob distributions whose non-zero mass is at least } 1/k\}$

## Sample complexity

$n^*(k, \epsilon) \triangleq \min\{n : \exists \hat{S}, \text{ s.t. } \mathbb{P}[|\hat{S} - S(P)| \leq \epsilon k] \geq 0.5, \forall P \in \mathcal{D}_k\}$

## Remarks

- Upgrade the confidence:  $n \rightarrow n \log \frac{1}{\delta} \Rightarrow 0.5 \rightarrow 1 - \delta$  (subsample + median + Hoeffding)

## Space of distributions

$\mathcal{D}_k \triangleq \{\text{prob distributions whose non-zero mass is at least } 1/k\}$

## Sample complexity

$n^*(k, \epsilon) \triangleq \min\{n : \exists \hat{S}, \text{s.t. } \mathbb{P}[|\hat{S} - S(P)| \leq \epsilon k] \geq 0.5, \forall P \in \mathcal{D}_k\}$

## Remarks

- Upgrade the confidence:  $n \rightarrow n \log \frac{1}{\delta} \Rightarrow 0.5 \rightarrow 1 - \delta$  (subsample + median + Hoeffding)
- Zero error ( $\epsilon = 0$ ):  $n^*(k, 0) \asymp k \log k$  (coupon collector)



- WYSIWYE:

$$\hat{S}_{\text{seen}} = \text{number of seen symbols}$$

- WYSIWYE:

$$\hat{S}_{\text{seen}} = \text{number of seen symbols}$$

- underestimate:

$$\hat{S}_{\text{seen}} \leq S(P), \quad P\text{-a.s.}$$

- severely underbiased in the **sublinear**-sampling regime:  $n \ll k$

# Do we have to estimate the distribution itself?

From a statistical perspective

- **high-dimensional** problem
  - ▶ estimating  $P$  provably requires  $n = \Theta(k)$  samples
  - ▶ empirical distribution is optimal up to constants
- **functional estimation**
  - ▶ **scalar** functional (support size)  $\stackrel{?}{\Rightarrow} n = o(k)$  suffices
  - ▶ plug-in is frequently suboptimal

- Histogram:

$$N_j = \sum_i \mathbf{1}_{\{X_i=j\}} : \# \text{ of occurrences of } j^{\text{th}} \text{ symbol}$$

- Histogram:

$$N_j = \sum_i \mathbf{1}_{\{X_i=j\}} : \# \text{ of occurrences of } j^{\text{th}} \text{ symbol}$$

- Histogram<sup>2</sup>/fingerprints/profiles:

$$h_i = \sum_j \mathbf{1}_{\{N_j=i\}} : \# \text{ of symbols that occurred exactly } i \text{ times}$$

- Histogram:

$$N_j = \sum_i \mathbf{1}_{\{X_i=j\}} : \# \text{ of occurrences of } j^{\text{th}} \text{ symbol}$$

- Histogram<sup>2</sup>/fingerprints/profiles:

$$h_i = \sum_j \mathbf{1}_{\{N_j=i\}} : \# \text{ of symbols that occurred exactly } i \text{ times}$$

- $h_0$ : # of unseens

Estimators that are **linear in the fingerprints**:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Estimators that are **linear in the fingerprints**:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Classical procedures:

- **Plug-in:**

$$\hat{S}_{\text{seen}} = h_1 + h_2 + h_3 + \dots$$



Estimators that are **linear in the fingerprints**:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Classical procedures:

- Plug-in:

$$\hat{S}_{\text{seen}} = h_1 + h_2 + h_3 + \dots$$

- Good-Toulmin '56: empirical Bayes

$$\hat{S}_{\text{GT}} = th_1 - t^2h_2 + t^3h_3 - t^4h_4 + \dots$$

Estimators that are **linear in the fingerprints**:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Classical procedures:

- Plug-in:

$$\hat{S}_{\text{seen}} = h_1 + h_2 + h_3 + \dots$$

- Good-Toulmin '56: empirical Bayes

$$\hat{S}_{\text{GT}} = th_1 - t^2h_2 + t^3h_3 - t^4h_4 + \dots$$

- Efron-Thisted '76: Bayesian

$$\hat{S}_{\text{ET}} = \sum_{j=1}^J (-1)^{j+1} t^j b_j h_j$$

where  $b_j = \mathbb{P}[\text{Binomial}(J, 1/(t+1)) \geq j]$

- $\hat{\mathcal{S}}_{\text{seen}}$ :  $n^*(k, \epsilon) \leq k \log \frac{1}{\epsilon}$

- $\hat{S}_{\text{seen}}$ :  $n^*(k, \epsilon) \leq k \log \frac{1}{\epsilon}$
- Valiant '08, Raskhodnikova et al. '09, Valiant-Valiant '11-'13: sublinear is possible.
  - ▶ Upper bound:  $n^*(k, \epsilon) \lesssim \frac{k}{\log k} \frac{1}{\epsilon^2}$  by LP [Efron-Thisted '76]
  - ▶ Lower bound:  $n^*(k, \epsilon) \gtrsim \frac{k}{\log k}$

- $\hat{S}_{\text{seen}}$ :  $n^*(k, \epsilon) \leq k \log \frac{1}{\epsilon}$
- Valiant '08, Raskhodnikova et al. '09, Valiant-Valiant '11-'13: sublinear is possible.
  - ▶ Upper bound:  $n^*(k, \epsilon) \lesssim \frac{k}{\log k} \frac{1}{\epsilon^2}$  by LP [Efron-Thisted '76]
  - ▶ Lower bound:  $n^*(k, \epsilon) \gtrsim \frac{k}{\log k}$

## Theorem (W.-Yang '14)

$$n^*(k, \epsilon) \asymp \frac{k}{\log k} \log^2 \frac{1}{\epsilon}$$

## Theorem (W.-Yang '14)

$$\inf_{\hat{S}} \sup_{P \in \mathcal{D}_k} \mathbb{E}[(\hat{S} - S(P))^2] \asymp k^2 \exp\left(-\sqrt{\frac{n \log k}{k}} \vee \frac{n}{k}\right)$$

## Objectives

- a principled way to obtain **rate-optimal linear** estimator
- a natural **lower bound** to establish optimality via **duality**

Best polynomial approximation



- $\mathcal{P}_L = \{\text{polynomials of degree at most } L\}$ .
- $I = [a, b]$ : a finite interval.
- Optimal approximation error

$$E_L(f, I) \triangleq \inf_{p \in \mathcal{P}_L} \sup_{x \in I} |f(x) - p(x)|$$

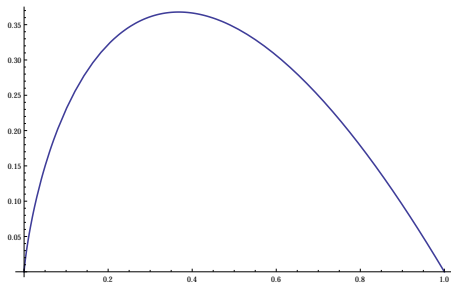
# Best polynomial approximation

- $\mathcal{P}_L = \{\text{polynomials of degree at most } L\}$ .
- $I = [a, b]$ : a finite interval.
- Optimal approximation error

$$E_L(f, I) \triangleq \inf_{p \in \mathcal{P}_L} \sup_{x \in I} |f(x) - p(x)|$$

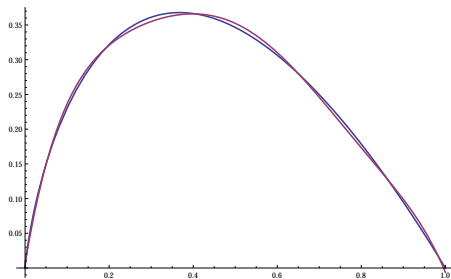
- Stone-Weierstrass theorem:  $f$  continuous  $\Rightarrow E_L(f, I) \xrightarrow{L \rightarrow \infty} 0$
- Speed of convergence related to modulus of continuity.
- Finite-dim convex optimization/Infinite-dim LP
- Many fast algorithms (e.g., Remez)

deg-6 approximation



Chebyshev alternation theorem

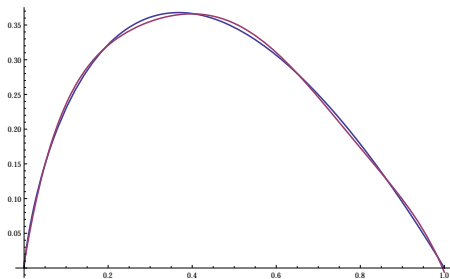
deg-6 approximation



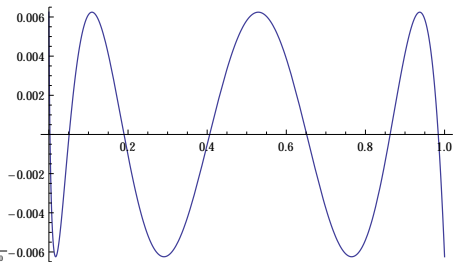
Chebyshev alternation theorem

# Example

deg-6 approximation



Chebyshev alternation theorem



$$\begin{aligned}\mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E} [f(U)] - \mathbb{E} [f(U')] \\ &\text{s.t. } \mathbb{E} [U^j] = \mathbb{E} [U'^j], \quad j = 1, \dots, L \\ &\quad U, U' \in I\end{aligned}$$

$$\begin{aligned} \mathcal{E}_L(f, I) &\triangleq \sup \int f d\mu - \int f d\mu' \\ \text{s.t.} \quad &\int f d\mu = \int f d\mu', \quad j = 1, \dots, L, \\ &\mu, \mu' \text{ supported on } I \end{aligned}$$

$$\begin{aligned}\mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E} [f(U)] - \mathbb{E} [f(U')] \\ &\text{s.t. } \mathbb{E} [U^j] = \mathbb{E} [U'^j], \quad j = 1, \dots, L, \\ &\quad U, U' \in I\end{aligned}$$



$$\begin{aligned} \mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E} [f(U)] - \mathbb{E} [f(U')] \\ &\text{s.t. } \mathbb{E} [U^j] = \mathbb{E} [U'^j], \quad j = 1, \dots, L, \quad \lambda_j \in \mathbb{R} \\ &\quad U, U' \in I \end{aligned}$$

Infinite-dim linear programming. **Dual:**

$$\inf_{\lambda_1^L} \sup_{U, U' \in I} \mathbb{E} [f(U)] - \mathbb{E} [f(U')] + \sum_j \lambda_j (\mathbb{E} [U^j] - \mathbb{E} [U'^j])$$

$$\begin{aligned} \mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E} [f(U)] - \mathbb{E} [f(U')] \\ &\text{s.t. } \mathbb{E} [U^j] = \mathbb{E} [U'^j], \quad j = 1, \dots, L, \quad \lambda_j \in \mathbb{R} \\ &\quad U, U' \in I \end{aligned}$$

Infinite-dim linear programming. **Dual:**

$$\inf_{\lambda_1^L} \sup_{U, U' \in I} \mathbb{E} [f(U)] - \mathbb{E} [f(U')] + \sum_j \lambda_j (\mathbb{E} [U^j] - \mathbb{E} [U'^j])$$

$$\begin{aligned}\mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E}[f(U)] - \mathbb{E}[f(U')] \\ &\text{s.t. } \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j = 1, \dots, L, \quad \lambda_j \in \mathbb{R} \\ &\quad U, U' \in I\end{aligned}$$

Infinite-dim linear programming. **Dual:**

$$\begin{aligned}&\inf_{\lambda_1^L} \sup_{U, U' \in I} \mathbb{E}[f(U)] - \mathbb{E}[f(U')] + \sum_j \lambda_j (\mathbb{E}[U^j] - \mathbb{E}[U'^j]) \\ &= \inf_{\lambda_1^L} \sup_{U \in I} \mathbb{E}\left[f(U) - \sum_j \lambda_j U^j\right] - \inf_{U' \in I} \mathbb{E}\left[f(U') - \sum_j \lambda_j U'^j\right]\end{aligned}$$

$$\begin{aligned}\mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E}[f(U)] - \mathbb{E}[f(U')] \\ &\text{s.t. } \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j = 1, \dots, L, \quad \lambda_j \in \mathbb{R} \\ &\quad U, U' \in I\end{aligned}$$

Infinite-dim linear programming. **Dual:**

$$\begin{aligned}&\inf_{\lambda_1^L} \sup_{U, U' \in I} \mathbb{E}[f(U)] - \mathbb{E}[f(U')] + \sum_j \lambda_j (\mathbb{E}[U^j] - \mathbb{E}[U'^j]) \\ &= \inf_{\lambda_1^L} \sup_{U \in I} \mathbb{E}\left[f(U) - \sum_j \lambda_j U^j\right] - \inf_{U' \in I} \mathbb{E}\left[f(U') - \sum_j \lambda_j U'^j\right] \\ &= \inf_{\lambda_0^L} \left( \sup_{u \in I} f(u) - \sum_j \lambda_j u^j \right) - \left( \inf_{u \in I} f(u) - \sum_j \lambda_j u^j \right)\end{aligned}$$

$$\begin{aligned}\mathcal{E}_L(f, I) &\triangleq \sup \mathbb{E}[f(U)] - \mathbb{E}[f(U')] \\ &\text{s.t. } \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j = 1, \dots, L, \quad \lambda_j \in \mathbb{R} \\ &\quad U, U' \in I\end{aligned}$$

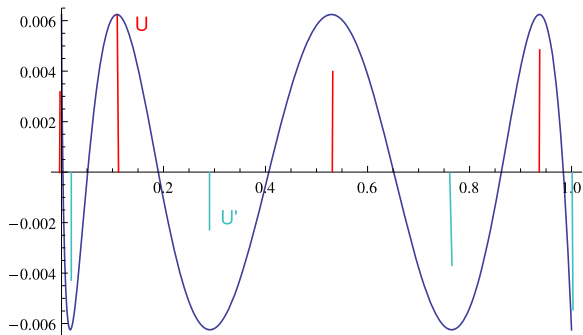
Infinite-dim linear programming. **Dual:**

$$\begin{aligned}&\inf_{\lambda_1^L} \sup_{U, U' \in I} \mathbb{E}[f(U)] - \mathbb{E}[f(U')] + \sum_j \lambda_j (\mathbb{E}[U^j] - \mathbb{E}[U'^j]) \\ &= \inf_{\lambda_1^L} \sup_{U \in I} \mathbb{E}\left[f(U) - \sum_j \lambda_j U^j\right] - \inf_{U' \in I} \mathbb{E}\left[f(U') - \sum_j \lambda_j U'^j\right] \\ &= \inf_{\lambda_0^L} \left( \sup_{u \in I} f(u) - \sum_j \lambda_j u^j \right) - \left( \inf_{u \in I} f(u) - \sum_j \lambda_j u^j \right) \\ &= 2 \inf_{p \in \mathcal{P}_L} \sup_{u \in I} |f(u) - p(u)|\end{aligned}$$

$$\mathcal{E}_L(f, I) = 2E_L(f, I)$$

# Moment matching $\Leftrightarrow$ best polynomial approximation

$$\mathcal{E}_L(f, I) = 2E_L(f, I)$$



Optimal estimator



- Poisson sampling model
  - ▶ draw sample size  $n' \sim \text{Poi}(n)$
  - ▶ draw  $n'$  i.i.d. samples from  $P$ .
- Histograms are **independent**:  $N_i \stackrel{\text{ind}}{\sim} \text{Poi}(np_i)$
- sample complexity/minimax risks remain unchanged within constant factors

Recall

$$\text{MSE} = \text{bias}^2 + \text{variance}$$

Main problem of  $\hat{S}_{\text{seen}}$ : huge bias.

# Unbiased estimators?

Unbiased estimator for  $f(P)$  from  $n$  samples:

- Independent sampling:  $f(P)$  is **polynomial** of degree  $\leq n$
- Poissonized sampling:  $f(P)$  is **real analytic**.

Unbiased estimator for  $f(P)$  from  $n$  samples:

- Independent sampling:  $f(P)$  is **polynomial** of degree  $\leq n$
- Poissonized sampling:  $f(P)$  is **real analytic**.

## Example

- Flip a coin with bias  $p$  for  $n$  times and estimate  $f(p)$
- Sufficient stat:  $Y \sim \text{Binomial}(n, p)$ .
- Unbiased estimator exists  $\Leftrightarrow f(p)$  is a **polynomial of degree  $\leq n$**

$$\mathbb{E}[\hat{f}(Y)] = \sum_{k=0}^n \hat{f}(k) \binom{n}{k} p^k (1-p)^{n-k}.$$

$$S(P) = \sum_i \mathbf{1}_{\{p_i > 0\}}$$

$$S(P) = \sum_i \mathbf{1}_{\{p_i > 0\}}$$

- Approximate  $\mathbf{1}_{\{x > 0\}}$  by  $q(x) = \sum_{m=0}^L a_m x^m$
- Find an unbiased estimator for the proxy

$$\tilde{S}(P) = \sum_i q(p_i)$$

- $|\text{bias}| \leq \text{uniform approx error}$

$$S(P) = \sum_i \mathbf{1}_{\{p_i > 0\}}$$

- Approximate  $\mathbf{1}_{\{x > 0\}}$  by  $q(x) = \sum_{m=0}^L a_m x^m$
- Find an unbiased estimator for the proxy

$$\tilde{S}(P) = \sum_i q(p_i)$$

- $|\text{bias}| \leq \text{uniform approx error}$
- But the function is discontinuous...



Consider estimators that are linear in the fingerprints:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Guidelines:

- $f(0) = 0$

Consider estimators that are linear in the fingerprints:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Guidelines:

- $f(0) = 0$
- $f(j) = 1$  for sufficiently large  $j > L$

Consider estimators that are linear in the fingerprints:

$$\hat{S} = \sum_i f(N_i) = \sum_{j \geq 1} f(j)h_j$$

Guidelines:

- $f(0) = 0$
- $f(j) = 1$  for sufficiently large  $j > L$
- How to choose  $f(1), \dots, f(L)$ ?

Choose

- $L = c_0 \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j), \quad N_j \sim \text{Poi}(np_j)$

Bias:

$$\mathbb{E}[\hat{S} - S] = \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]\mathbf{1}_{\{p_i > 1/k\}}$$

Choose

- $L = c_0 \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned} \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}] \mathbf{1}_{\{p_i > 1/k\}} \\ &\approx \sum \underbrace{\mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]}_{\approx 0} \mathbf{1}_{\{2L/n > p_i > 1/k\}} \end{aligned}$$

Choose

- $L = c_0 \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned} \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}] \mathbf{1}_{\{p_i > 1/k\}} \\ &\approx \sum \underbrace{\mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]}_{e^{-np_i} \times \text{poly of deg } L!} \mathbf{1}_{\{2L/n > p_i > 1/k\}} \end{aligned}$$

- Observe

$$\mathbb{E}[(f(N) - 1)\mathbf{1}_{\{N \leq L\}}] = e^{-\lambda} \underbrace{\sum_{j \geq 0} \frac{f(j) - 1}{j!} \lambda^j}_{q(\lambda)}$$

- Then

$$|\text{bias}| \leq k \sup_{n/k \leq \lambda \leq c \log k} |q(\lambda)|$$

- Choose the best deg- $L$  polynomial  $q$  s.t.  $q(0) = -1$

- Observe

$$\mathbb{E}[(f(N) - 1)\mathbf{1}_{\{N \leq L\}}] = e^{-\lambda} \underbrace{\sum_{j \geq 0} \frac{f(j) - 1}{j!} \lambda^j}_{q(\lambda)}$$

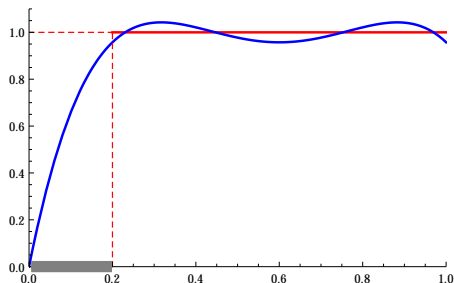
- Then

$$|\text{bias}| \leq k \sup_{n/k \leq \lambda \leq c \log k} |q(\lambda)|$$

- Choose the best deg- $L$  polynomial  $q$  s.t.  $q(0) = -1$
- Solution: [Chebyshev polynomial](#)



# Chebyshev polynomial



best approximation to one by  
polynomial passing through origin is  
Chebyshev polynomial

$$p_L(x) = 1 - \frac{\cos L \arccos x}{\cos L \arccos a}$$

- Chebyshev polynomial:  $r \triangleq c_1 \log k$  and  $l \triangleq \frac{n}{k}$ ,

$$-\frac{\cos L \arccos\left(\frac{2}{r-l}x - \frac{r+l}{r-l}\right)}{\cos L \arccos\left(-\frac{r+l}{r-l}\right)} \triangleq \sum_{j=0}^L a_m x^m.$$

- Chebyshev polynomial:  $r \triangleq c_1 \log k$  and  $l \triangleq \frac{n}{k}$ ,

$$-\frac{\cos L \arccos\left(\frac{2}{r-l}x - \frac{r+l}{r-l}\right)}{\cos L \arccos\left(-\frac{r+l}{r-l}\right)} \triangleq \sum_{j=0}^L a_j x^j.$$

- Choose

$$f(j) = \begin{cases} 0 & j = 0 \\ 1 + a_j j! & j = 1, \dots, L \\ 1 & j > L. \end{cases}$$

- Chebyshev polynomial:  $r \triangleq c_1 \log k$  and  $l \triangleq \frac{n}{k}$ ,

$$-\frac{\cos L \arccos\left(\frac{2}{r-l}x - \frac{r+l}{r-l}\right)}{\cos L \arccos\left(-\frac{r+l}{r-l}\right)} \triangleq \sum_{j=0}^L a_m x^m.$$

- Choose

$$f(j) = \begin{cases} 0 & j = 0 \\ 1 + a_j j! & j = 1, \dots, L \\ 1 & j > L. \end{cases}$$

- Linear estimator (precomputable coefficients): **no sample splitting!!**

$$\hat{S} = \sum_{j=1}^L f(j)h_j + \sum_{j>L} h_j$$

- Chebyshev polynomial:  $r \triangleq c_1 \log k$  and  $l \triangleq \frac{n}{k}$ ,

$$-\frac{\cos L \arccos\left(\frac{2}{r-l}x - \frac{r+l}{r-l}\right)}{\cos L \arccos\left(-\frac{r+l}{r-l}\right)} \triangleq \sum_{j=0}^L a_m x^m.$$

- Choose

$$f(j) = \begin{cases} 0 & j = 0 \\ 1 + a_j j! & j = 1, \dots, L \\ 1 & j > L. \end{cases}$$

- Linear estimator (precomputable coefficients): **no sample splitting!!**

$$\hat{S} = \sum_{j=1}^L f(j) h_j + \sum_{j>L} h_j$$

- Significantly faster than LP [Efron-Thisted '76, Valiant-Valiant '11]

- ① bias  $\leq$  approximation error of Chebyshev polynomial:

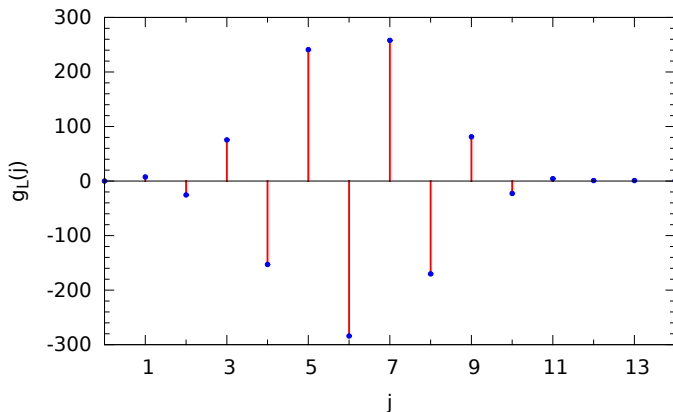
$$\frac{1}{|\cos M \arccos(-\frac{r+l}{r-l})|} \asymp \exp\left(-c\sqrt{\frac{n \log k}{k}}\right),$$

- ② variance  $\approx \text{poly}(k)$ .

# Optimal estimator

Plot of coefficients ( $k = 10^6$  and  $n = 2 \times 10^5$ ):

$$\hat{S} = \sum_{j \geq 1} f(j)h_j$$



# Why oscillatory and alternating?

$$\hat{S} = \sum_{j \geq 1} f(j)h_j$$

The same oscillation also happens in:

- Good-Toulmin '56: empirical Bayes

$$\hat{S}_{\text{GT}} = th_1 - t^2h_2 + t^3h_3 - t^4h_4 + \dots$$

- Efron-Thistle '76: Bayesian

$$\hat{S}_{\text{ET}} = \sum_{j=1}^J (-1)^{j+1} t^j b_j h_j$$

I HAVE NO EXPLANATION!



Impossibility results

$$n^*(k, \epsilon) \gtrsim \frac{k}{\log k} \log^2 \frac{1}{\epsilon}$$

- $\text{TV}(P_0, P_1) = \frac{1}{2} \int |dP_0 - dP_1|$
- optimal error probability for testing  $P_0$  vs  $P_1$

$$1 - \text{TV}(P_0, P_1) = \min_{\psi} P_0[\psi = 1] + P_1[\psi = 0]$$

Given  $U \sim \mu$ ,

$$\mathbb{E}[\text{Poi}(U)] = \int_{\mathbb{R}_+} \text{Poi}(\lambda) \mu(d\lambda)$$

Two-prior argument (composite HT):

- draw random distribution  $P \xrightarrow{\text{Poisson}} N_i \stackrel{\text{ind}}{\sim} \text{Poi}(np_i)$
- draw random distribution  $P' \xrightarrow{\text{Poisson}} N'_i \stackrel{\text{ind}}{\sim} \text{Poi}(np'_i)$

Two-prior argument (composite HT):

- draw random distribution  $P \xrightarrow{\text{Poisson}} N_i \stackrel{\text{ind}}{\sim} \text{Poi}(np_i)$
- draw random distribution  $P' \xrightarrow{\text{Poisson}} N'_i \stackrel{\text{ind}}{\sim} \text{Poi}(np'_i)$

Le Cam's lemma applies if

- $S(P)$  and  $S(P')$  differ with high probability
- Distributions of  $N$  and  $N'$  are indistinguishable ( $k$ -dim Poisson mixtures)

Two-prior argument (composite HT):

- draw random distribution  $P \xrightarrow{\text{Poisson}} N_i \stackrel{\text{ind}}{\sim} \text{Poi}(np_i)$
- draw random distribution  $P' \xrightarrow{\text{Poisson}} N'_i \stackrel{\text{ind}}{\sim} \text{Poi}(np'_i)$

Le Cam's lemma applies if

- $S(P)$  and  $S(P')$  differ with high probability
- Distributions of  $N$  and  $N'$  are indistinguishable ( $k$ -dim Poisson mixtures)

Main hurdle: difficult to work with distributions on high-dimensional probability simplex.

# Key construction: reduction to one dimension

- Given  $U, U'$  with **unit mean**:

$$P = \frac{1}{k} \underbrace{(U_1, \dots, U_k)}_{\underset{\sim}{\text{i.i.d.}}_U}, \quad P' = \frac{1}{k} \underbrace{(U'_1, \dots, U'_k)}_{\underset{\sim}{\text{i.i.d.}}_{U'}}$$

- By LLN,
  - ▶  $P$  and  $P'$  are not, but **close to**, probability distributions.



# Key construction: reduction to one dimension

- Given  $U, U'$  with **unit mean**:

$$P = \frac{1}{k} \underbrace{(U_1, \dots, U_k)}_{\underset{\sim}{\text{i.i.d.}}_U}, \quad P' = \frac{1}{k} \underbrace{(U'_1, \dots, U'_k)}_{\underset{\sim}{\text{i.i.d.}}_{U'}}$$

- By LLN,
  - ▶  $P$  and  $P'$  are not, but **close to**, probability distributions.
  - ▶ **support size concentrates** on the mean:

$$\mathbb{E}[S(P)] - \mathbb{E}[S(P')] = k(\mathbb{P}\{U > 0\} - \mathbb{P}\{U' > 0\})$$

# Key construction: reduction to one dimension

- Given  $U, U'$  with **unit mean**:

$$P = \frac{1}{k} \underbrace{(U_1, \dots, U_k)}_{\text{i.i.d. } U}, \quad P' = \frac{1}{k} \underbrace{(U'_1, \dots, U'_k)}_{\text{i.i.d. } U'}$$

- By LLN,
  - ▶  $P$  and  $P'$  are not, but **close to**, probability distributions.
  - ▶ **support size concentrates** on the mean:

$$\mathbb{E}[S(P)] - \mathbb{E}[S(P')] = k(\mathbb{P}\{U > 0\} - \mathbb{P}\{U' > 0\})$$

- Sufficient statistic are **iid**:

$$N_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{E}[\text{Poi}(nU/k)], \quad N'_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{E}[\text{Poi}(nU'/k)].$$

- Suffice to show  $\text{TV}(\underbrace{\mathbb{E}[\text{Poi}(nU/k)], \mathbb{E}[\text{Poi}(nU'/k)]}_{\text{one-dimensional Poisson mixtures}}) = o(1/k)$ .

## Lemma

- $U, U' \in [0, \frac{k \log k}{n}]$
- $\mathbb{E}[U^j] = \mathbb{E}[U'^j], j = 1, \dots, L = C \log k$
- *Then*

$$\text{TV}(\mathbb{E}[\text{Poi}(nU/k)], \mathbb{E}[\text{Poi}(nU'/k)]) = o(1/k)$$

# Optimize the lower bound

Let  $\lambda = k \log k/n$ .

Choose the best  $U, U'$ :

$$\begin{aligned} & \sup \mathbb{P} \{U = 0\} - \mathbb{P} \{U' = 0\} \\ & \text{s.t. } \mathbb{E} [U] = \mathbb{E} [U'] = 1 \\ & \quad \mathbb{E} [U^j] = \mathbb{E} [U'^j], \quad j \in [L] \\ & \quad U, U' \in \{0\} \cup [1, \lambda] \end{aligned}$$

# Optimize the lower bound

Let  $\lambda = k \log k/n$ .

Choose the best  $U, U'$ :

$$\begin{aligned} & \sup \mathbb{P}\{U = 0\} - \mathbb{P}\{U' = 0\} \\ & \text{s.t. } \mathbb{E}[U] = \mathbb{E}[U'] = 1 \\ & \quad \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j \in [L] \\ & \quad U, U' \in \{0\} \cup [1, \lambda] \end{aligned}$$

# Optimize the lower bound

Let  $\lambda = k \log k/n$ .

Choose the best  $U, U'$ :

$$\begin{aligned} & \sup \mathbb{P}\{U = 0\} - \mathbb{P}\{U' = 0\} \\ & \text{s.t. } \mathbb{E}[U] = \mathbb{E}[U'] = 1 \\ & \quad \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j \in [L] \\ & \quad U, U' \in \{0\} \cup [1, \lambda] \\ = & \sup \mathbb{E}[1/X] - \mathbb{E}[1/X'] \\ & \text{s.t. } \mathbb{E}[X^j] = \mathbb{E}[X'^j], \quad j \in [L] \\ & \quad X, X' \in [1, \lambda], \end{aligned}$$

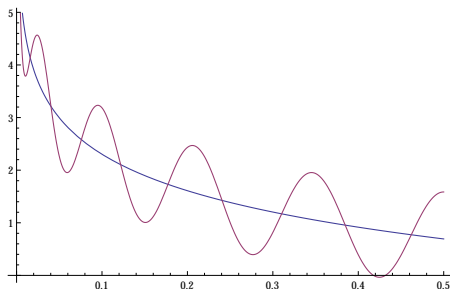
$$P_U(du) = (1 - \mathbb{E}[\frac{1}{X}]) \delta_0(du) + \frac{1}{u} P_X(du)$$

# Optimize the lower bound

Let  $\lambda = k \log k/n$ .

Choose the best  $U, U'$ :

$$\begin{aligned} & \sup \mathbb{P}\{U = 0\} - \mathbb{P}\{U' = 0\} \\ & \text{s.t. } \mathbb{E}[U] = \mathbb{E}[U'] = 1 \\ & \quad \mathbb{E}[U^j] = \mathbb{E}[U'^j], \quad j \in [L] \\ & \quad U, U' \in \{0\} \cup [1, \lambda] \\ = & \sup \mathbb{E}[1/X] - \mathbb{E}[1/X'] \\ & \text{s.t. } \mathbb{E}[X^j] = \mathbb{E}[X'^j], \quad j \in [L] \\ & \quad X, X' \in [1, \lambda], \\ = & 2E_L(1/x, [1, \lambda]) \gtrsim e^{-c\sqrt{\frac{n \log k}{k}}} \end{aligned}$$



Our inspiration: earlier work on Gaussian models

- Ibragimov-Nemirovskii-Khas'minskii '87: smooth functions
- Lepski-Nemirovski-Spokoiny '99:  $L_q$  norm of Gaussian regression function
- Cai-Low '11:  $L_1$  norm of normal mean



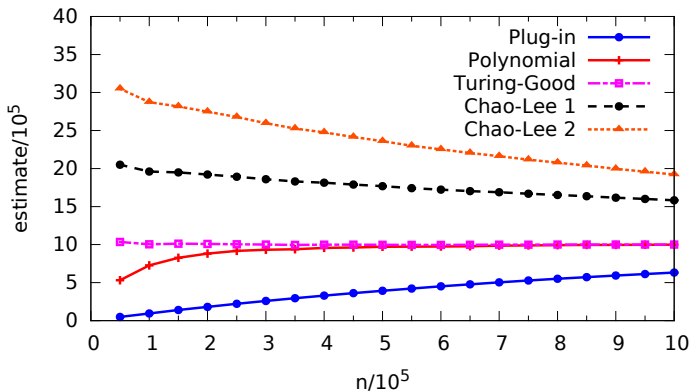
Lower bound in [Valiant-Valiant '11]

- Deal with fingerprints – **high-dim** distribution with dependent components
- Approximate distribution by quantized Gaussian
- Bound distance between mean and covariance matrices

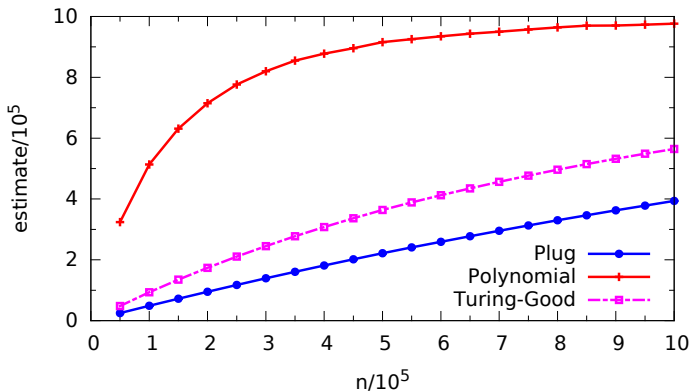
Lower bound here: reduce to **one dimension**

Experiments

# Uniform over 1 million elements

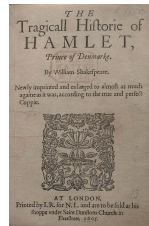
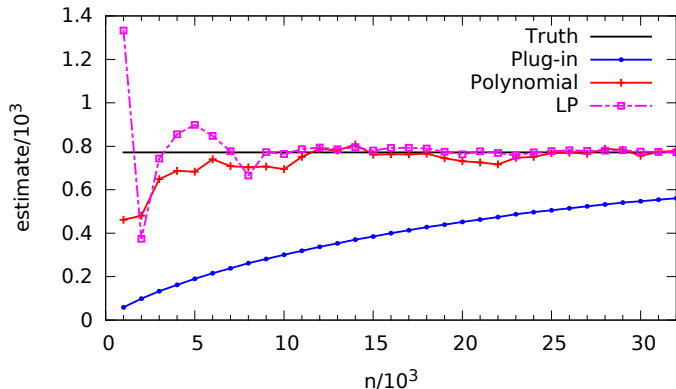


# Uniform mixed with point mass

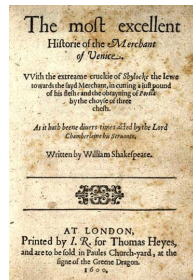
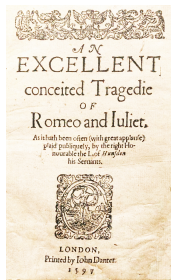
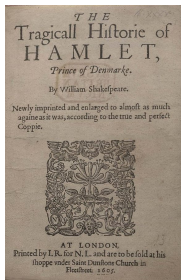


# How many words did Shakespeare know?

- Hamlet: total words 32000, total distinct words  $\sim 7700$ ,
- deg-10 Chebyshev polynomial
- sampling with replacement
- compare with LP [Efron-Thisted '76, Valiant-Valiant '13]



# How many words did Shakespeare know?



Feed the entire Shakespearean canon into the estimator:

- $\hat{S} = 68944 \sim 73257$
- Efron-Thisted '76: 66534

## Formulation

Given an urn containing  $k$  balls, estimate the number of distinct colors  $S$  by sampling (e.g. with replacement).

- Special case of support size estimation:  $p_i \in \{0, \frac{1}{k}, \frac{2}{k}, \dots\}$ .
- Same sample complexity as DISTINCT-ELEMENT problem in TCS.

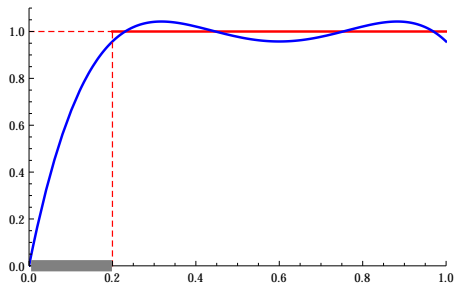
## Formulation

Given an urn containing  $k$  balls, estimate the number of distinct colors  $S$  by sampling (e.g. with replacement).

- Special case of support size estimation:  $p_i \in \{0, \frac{1}{k}, \frac{2}{k}, \dots\}$ .
- Same sample complexity as DISTINCT-ELEMENT problem in TCS.
- Use Chebyshev:  $\frac{k}{\log k}$  samples can achieve achieve  $0.1k$
- Converse:  $\frac{k}{\log k}$  samples are necessary to achieve  $0.1k$  [Valiant '12]



# Can we do better?

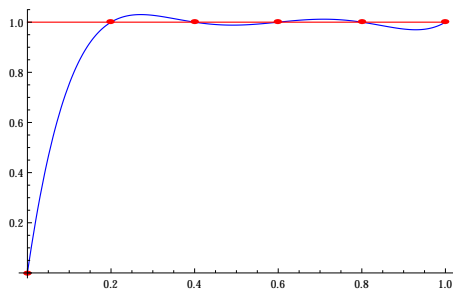


Use Lagrange interpolation polynomial to achieve **zero** bias

- **Uniform approximation:**  
 $\epsilon \lesssim \exp(-c\sqrt{\log k})$
- **Interpolation:**  
 $\epsilon \lesssim \exp(-c \log k)$ .

$$q_L(x) = 1 - \frac{\prod_{j=1}^L (j - x)}{L!}$$

# Can we do better?



Use Lagrange interpolation polynomial to achieve **zero** bias

- **Uniform approximation:**  
 $\epsilon \lesssim \exp(-c\sqrt{\log k})$
- **Interpolation:**  
 $\epsilon \lesssim \exp(-c \log k)$ .

$$q_L(x) = 1 - \frac{\prod_{j=1}^L (j - x)}{L!}$$

$$\text{minimax risk} \gtrsim k^2 \exp\left(-c \frac{n \log k}{k}\right)$$

- Tight when  $n = 0.1k$
- Compare to general support size:

$$\text{minimax risk} \asymp k^2 \exp\left(-c \sqrt{\frac{n \log k}{k}}\right)$$

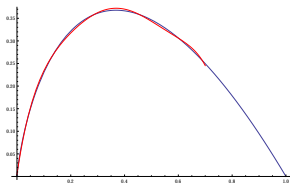
$$H(P) = \sum p_i \log \frac{1}{p_i}$$

## Theorem (W.-Yang '14)

Sample complexity to estimate within  $\epsilon$  bits:  $n \asymp \max \left\{ \frac{k}{\epsilon \log k}, \frac{\log^2 k}{\epsilon^2} \right\}$   
(upper bound also in Jiao et al. '14)

## Strategy

- degree:  $L \sim \log k$
- small masses: polynomial approximation
- large masses: plug-in with bias correction
- coeff's bounded by Chebyshev



- Estimating  $H_\alpha(P) = \frac{1}{1-\alpha} \log \sum p_i^\alpha$  [Jiao et al. '14, Acharya et al. '14]

To estimate

$$F(P) = \sum f(p_i)$$

Sample complexity is roughly governed by the following convex optimization problem (over logarithmic variables):

$$\begin{aligned} \mathcal{F}(\lambda) &\triangleq \sup \quad \mathbb{E}[f(U)] - \mathbb{E}[f(U')] \\ & \quad s.t. \quad \mathbb{E}[U^j] = \mathbb{E}[U'^j] \quad j = 1, \dots, \log k, \\ & \quad \mathbb{E}[U] \leq 1/k, \\ & \quad U, U' \in [0, \log k/n], \end{aligned}$$

- Lower bound: primal program (inapproximability result)
- Upper bound: dual program (approximability result)

- Many open problems and directions
  - ▶ Confidence intervals
  - ▶ Adaptive estimation
  - ▶ How to go beyond iid sampling
  - ▶ How to incorporate structures

## References

- W. & P. Yang (2014). *Minimax rates of entropy estimation on large alphabets via best polynomial approximation*. [arXiv:1407.0381](#)
- W. & P. Yang (2015). *Chebyshev polynomials, moment matching, and optimal estimation of the unseen*. [arXiv:1503.xxxx](#)

Choose

- $M = c \log k.$
- $\hat{S} = \sum_{j \geq 1} f(N_j), \quad N_j \sim \text{Poi}(np_j)$

Bias:

$$\mathbb{E}[\hat{S} - S] = \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}}$$



Choose

- $M = c \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j), \quad N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned} \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}} \\ &\stackrel{f(0)=0}{=} \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 0\}}] \end{aligned}$$

Choose

- $M = c \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned} \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}} \\ &\stackrel{f(0)=0}{=} \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 0\}}] \\ &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 1/k\}}] \end{aligned}$$

Choose

- $M = c \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned}
 \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}} \\
 &\stackrel{f(0)=0}{=} \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 0\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 1/k\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]\mathbf{1}_{\{p_i > 1/k\}}
 \end{aligned}$$

Choose

- $M = c \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned}
 \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}} \\
 &\stackrel{f(0)=0}{=} \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 0\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 1/k\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]\mathbf{1}_{\{p_i > 1/k\}} \\
 &\stackrel{whp}{=} \underbrace{\sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]}_{\mathbf{1}_{\{L/2n > p_i > 1/k\}}}
 \end{aligned}$$

Choose

- $M = c \log k$ .
- $\hat{S} = \sum_{j \geq 1} f(N_j)$ ,  $N_j \sim \text{Poi}(np_j)$

Bias:

$$\begin{aligned}
 \mathbb{E}[\hat{S} - S] &= \sum \mathbb{E}[f(N_i)] - \mathbf{1}_{\{p_i > 0\}} \\
 &\stackrel{f(0)=0}{=} \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 0\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{p_i > 1/k\}}] \\
 &= \sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]\mathbf{1}_{\{p_i > 1/k\}} \\
 &\stackrel{whp}{=} \underbrace{\sum \mathbb{E}[(f(N_i) - 1)\mathbf{1}_{\{N_i \leq L\}}]}_{\text{poly of deg } L} \mathbf{1}_{\{L/2n > p_i > 1/k\}}
 \end{aligned}$$

**Observe:**  $g(\lambda) \triangleq \mathbb{E}[(f(N) - 1)\mathbf{1}_{\{N \leq L\}}] = e^{-\lambda} \times \text{poly of deg } L$