Frequency Distribution testing

Michel de Rougemont University Paris II & IRIF CNRS

(Joint work with Claire Mathieu)

Motivations

- 1. Decide very quickly Statistical constraints
- 2. SAT panorama (n variables with densities, m clauses)



Heuristics adapted for each area, associated with a statistical constraint.

- 3. Non worst-case complexity
- 4. L.L.M. How to represent $P(t_k | t_1, t_2, \dots, t_{k-1})$

How is the distribution given?

- 1. Distribution as an Oracle
 - Sample D : output (a,b) with probability p(a,b)
 - Sample a with probability p(a | b)

- 2. Streaming data
 - Stream of objects: a,b,a,c,a,a,b.....
 - Stream of texts: all the text on the web.....

Distribution D defined by the frequencies of letters, words,.... Reservoir sampling gives samples on D.

Standard problems

Domain of size n

- 1. Is the distribution D close to:
 - A fixed distribution (power law, for example)
 - Uniform distribution

2. Is
$$D = D_1 * D_2$$
?
 $D_1 = D_2$?

Frequency distributions

Stream of objects: a,c,a,c,a,a,b,c c,a,b,b,a,b,b,a.....





Frequency distribution: Domain={1,2,...n} Decreasing function



Plan

- 1. Distribution Testing
 - Total variation distance
 - Is D close to the Uniform distribution
- 2. Frequency Distribution Testing
 - Relative Frechet distance
 - Is D close to the Uniform distribution: $\Omega(n)$ space
 - Smooth Decreasing distributions: O(log² n) space

Property Testing framework: algorithm A s.t.

- If D=Uniform distribution, A(D)=1
- If D is ϵ -far from the Uniform distribution, Prob [A(D)=0]> 1- δ

1. Uniformity testing (Cannone 2023)

	Sample complexity	Notes	References
Collision-based	$rac{k^{1/2}}{arepsilon^2}$	"Natural"	Goldreich and Ron, 2000; Diakonikolas <i>et al.</i> , 2019b
Unique elements	$rac{k^{1/2}}{arepsilon^2}$	Low sensitivity $\varepsilon \gg 1/k^{1/4}$	Paninski, 2008
Modified χ^2	$rac{k^{1/2}}{arepsilon^2}$	(None)	Valiant and Valiant, 2017; Acharya <i>et al.</i> , 2015; Diakonikolas <i>et al.</i> , 2015
Empirical distance to uniform	$rac{k^{1/2}}{\varepsilon^2}$	Low sensitivity	Diakonikolas et al., 2018
Random binary hashing	$\frac{k}{\varepsilon^2}$	Suboptimal, but fast	Acharya et al., 2020d
Bipartite collisions	$\frac{k^{1/2}}{\varepsilon^2}$	Tradeoff possible	Diakonikolas et al., 2019a
Empirical subset weighting	$\frac{k^{1/2}}{\varepsilon^2}$	Tradeoff possible $\varepsilon \gg 1/k^{1/4}$	Acharya et al., 2022

Collision-based tester

Sample X_1 , X_2 , ... X_m CountCollisions

Compute:
$$Z = \frac{\sum_{1 \le s \le t \le n} 1\{X_s = X_t\}}{n(n-1)/2}$$
 $E(Z) = ||p||_2^2$

General: $||p - q||_2 \le \frac{1}{2} ||p - q||_1 \le \frac{\sqrt{n}}{2} ||p - q||_2$ If q is uniform: $||p - q||_2^2 = \sum_i (p(i) - 1/n)^2 = ||p||_2^2 - 1/n$

 $||p-q||_1 \ge \varepsilon \rightarrow ||p||_2^2 \ge (1+4\varepsilon^2)/n$

Study Var(Z), apply Chebyshev inequality

2. Frequency distributions

Stream of objects: a,c,a,c,a,a,b,c (a,b),(a,c),(b,c).....

Streaming algorithm in space **O(poly(log n))** to decide if the frequency distribution g of the stream is close to a given distribution f?

Frequency distribution (absolute values):



Distances between distributions

- 1. L_1 , L_2 ... L_p
- 2. EMD
 - Earth-moving distance
- 3. Fréchet
 - Absolute/relative
- 4. Kullback-Liebler Divergence

Fréchet relative distance y*(1+ ε₂) x*(1+ ε₁) 2 3 1 n Degree distributions: f(i)=#v: degree(v)=i

Are And I close ?

Results (n items)

Is g close to the uniform distribution f ?
Ω(n) space

Communication Complexity (reduction from Index)

- 2. Smooth and decreasing distributions f
 - Step compatible (Smooth): every point belongs to a complete Frechet rectangle

•
$$\gamma - decreasing: f(\gamma, t) \leq \frac{f(t)}{2}$$

Main result: Streaming tester with O(log² n) space





Frequent items

- 1. Deterministic algorithms
 - Misra-Gries
 - Spacesaving (K=3, additive error)



- 2. Randomized algorithms
 - Count-Min-Sketch

Simplified Tester: log n substreams with Spacesaving $z_{i} = \left| \left(1 + \varepsilon_{1}^{2}\right)^{i} \right| \qquad a_{i} = \left| \frac{z_{i} \cdot \varepsilon_{1}^{2}}{\log \log n} \right|$



Tester

- 1. Details
 - Ignore z_i close to the discontinuities of f
 - Spacesaving table of size O(log n . log log n)
 - Exact counting if $\varepsilon_2 f(z_i) \le f(n)$
- 2. Key points
 - If f is γ -decreasing, Spacesaving guarantees relative errors.
 - If f,g are smooth and $\neg (f \sim_{3\varepsilon_1, 3\varepsilon_2} g)$
 - \rightarrow \exists ($\varepsilon_1, \varepsilon_2$) separating rectangle
 - Some z_i will witness the separating rectangle
 - The tester will reject with h.p.



Extensions

- 1. Tuples in dimension d: (a_1, a_2, \dots, a_d)
 - Domain size: n^d
 - Frequency of the Marginals: hash on the projections
- 2. Questions
 - Independence of subsets
 - Marginals are close

$$D = D_1 * D_2 ?$$
$$D_1 = D_2$$

Conclusion

1. Relative Frechet distance

- 2. Frequency distribution testing
 - Uniform testing is hard
 - Smooth and decreasing frequency distributions are « easy »

(Zipf, Power law,)

3. Extensions to other properties