

# Learning PBD Powers

---

Kontonis Vasilis

November 13, 2017

Corelab, NTUA

1. Introduction
2. Binomial Powers
3. Learning the Parameters of a PBD

# Introduction

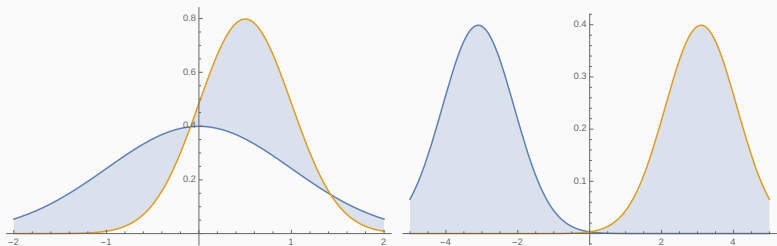
---

# Distribution Learning

- Draw samples from an unknown distribution  $D$ .
- Output an approximation of the density function of  $D$ .

## Distances of Distributions

$$d_{\text{tv}}(Q, P) = \frac{1}{2} \int |p(x) - q(x)| dx = \sup_{A \in \mathcal{A}} |P(A) - Q(A)|.$$



# Dvoretzky Kiefer Wolfowitz Inequality

## Kolmogorov Distance

$$d_{\text{kol}}(Q, P) = \sup_{x \in \mathbb{R}} |F_Q(x) - F_P(x)|$$

## Empirical CDF and DKW Inequality

$$\hat{F}_N(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{[X_i \leq x]}, \quad x \in \mathbb{R}$$

For every  $\varepsilon > 0$

$$\mathbb{P} \left[ d_{\text{kol}}(F_P, \hat{F}_N) > \varepsilon \right] \leq 2e^{-2N\varepsilon^2}$$

- With  $N = O(1/\varepsilon^2)$  samples we approximate the unknown CDF.
- **Question:** Is there anything left to do?

## Poisson Binomial Distributions

- $X_i$ 's are 0/1 Bernoulli with  $\mathbf{E}[X_i] = p_i$ .
- $X = \sum_{i=1}^n X_i$  is a  $n$ -PBD with probability vector  $\mathbf{p} = (p_1, \dots, p_n)$ .
- $\mathbf{E}[X] = \sum p_i$ ,  $\mathbf{V}[X] = \sum p_i(1 - p_i)$ .
- If  $p_i$ 's are small then  $X$  is close in TVD to a  $\text{Pois}(\sum p_i)$ .
- If  $\mathbf{V}[X]$  is large, then  $X$  is close in TVD to a (Discretized) Normal.

# Learning PBDs

Simply running Birge's Algorithm [1] is *not* good enough,  $O(\log n/\varepsilon^3)$  samples are needed.

- Learn Sparse:
  - Truncate the support: Draw  $O(1/\varepsilon^2)$  samples, sort them and let  $a, b$  be such that  $X(a \leq i \leq b) = 1 - \varepsilon$ .
  - If  $b - a > 1/\varepsilon^3$  then output fail, else run Birge's unimodal algorithm on  $X_{[a,b]}$ .
- Learn Heavy:
  - Estimate the variance  $\hat{\sigma}^2$  and the mean  $\hat{\mu}$  of  $X$  using  $O(1/\varepsilon^2)$  samples.
  - The discretized normal  $DN(\hat{\mu}, \hat{\sigma}^2)$  is  $\varepsilon$ -close, so just output  $\hat{\mu}$  and  $\hat{\sigma}^2$ .
- Choose between Sparse and Heavy.

# Tough Curriculum

- set of  $m$  items, e.g. set of students.
- a set of  $n$  items, e.g. set of courses.
- Each students has passed course  $i$  with probability  $p_i$  independently from other students.

**Question:** What is the distribution of the number of different courses that a set of  $k$  students will have passed?

**Answer:** For course  $i$  not to be in the set we need to exclude from all students. This happens with probability  $(1 - p_i)^k$ . Therefore  $i$  is included in the union of courses with probability  $1 - (1 - p_i)^k$ .



## PBD Powers

- Let  $P$  be a  $n$ -PBD defined by  $\mathbf{p}$ .
- $P^i$  is the  $i$ -th PBD power of  $P$  defined by  $\mathbf{p}^i$ .

**Question:** Given samples from a subset of the powers of  $P$  can we learn the other powers? Can we do better than learning each one of them separately?

# Binomial Powers

---

# Approximating Binomials

## Binomial TVD [3]

We want to approximate all the distributions  $B(n, p^i), i \in \mathbb{N}$ .

$$|p - q| \leq \varepsilon \sqrt{\frac{p(1-p)}{n}} = \text{err}(n, p, \varepsilon) \implies d_{\text{tv}}(B(n, p), B(n, q)) \leq \varepsilon$$

## Approximating $p$

- Estimator:

$$\hat{p} = \frac{\sum_{i=1}^N X_i}{Nn}$$

- Sample Complexity:

Choosing  $N = O(\ln(1/\delta)/\varepsilon^2)$  from Chernoff's bound we have:

$$\mathbb{P}[|\hat{p} - p| > \text{err}(n, p, \varepsilon)] < \delta.$$

# Real Powers of $p$

Assume that  $p \approx 1 - \frac{1}{n}$  or  $p \approx \frac{1}{n}$ .

$$p = 0.\underbrace{99\dots9}_{\# \log n} \underbrace{458382}_{\text{"constant" part}}$$

$$p = 0.\underbrace{00\dots0}_{\# \log n} \underbrace{235711}_{\text{"constant" part}}$$

- Sampling from the first power reveals the first part of  $p$ , since  $\sqrt{p(1-p)/n} \approx 1/n$ .
- Is this good enough to approximate all binomial powers ?
- $0.9995^{1000} \approx 0.6064, 0.9997^{1000} \approx 0.7407$

## A blast from the past

$$|x - y| = \frac{|x^{2^i} - y^{2^i}|}{\prod_{j=0}^{i-1} (x^{2^j} + y^{2^j})}$$

## Finding the Sweet Spot

By Mean Value Theorem applied to mapping  $x \mapsto x^l$  we obtain

$$p^l - \hat{q}_1^l \leq l p^{l-1} (p - \hat{q}_1), l \in (1, +\infty)$$

and

$$\hat{q}_2^l - p^l \leq l p^{l-1} (\hat{q}_2 - p), l \in (0, 1)$$

We need to find a function  $u(p)$  such that for all  $l > 0$ :

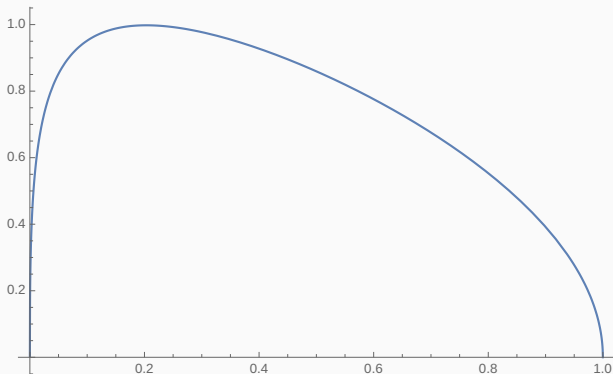
$$\begin{aligned} u(p) l p^{l-1} \text{err}(n, p, \varepsilon) &\leq \text{err}(n, p^l, \varepsilon) & (1) \\ u(p) l p^{l-1} \sqrt{\frac{p(1-p)}{n}} &\leq \sqrt{\frac{p^l(1-p^l)}{n}} \\ u^2(p) &\leq \frac{p}{1-p} \frac{p^{-l} - 1}{l^2} \end{aligned}$$

## Finding the Sweet Spot

$f(l) = \frac{p^{-l}-1}{l^2}$  is convex attaining its minimum at  $\bar{l} = -\frac{C}{\ln p}$ ,  
 $f(\bar{l}) = C \ln^2(1/p)$ .

Now we can choose:

$$u(p) = D \sqrt{\frac{p}{1-p}} \ln(1/p), \quad D \approx 1.24$$



## Getting in Range

Magic Power:  $a = -\frac{1}{\ln p}$ .

**Question:** Can we guess the "magic" power using samples from the first one for all values of  $p$ ?

**Answer:** No, if  $p$  is very close to 1 then we cannot hope to learn the number of 9's in its decimal representation.

We approximate  $a$  with  $\hat{a} = -\frac{1}{\ln \hat{p}}$ .

- $p^a = \hat{p}^{\hat{a}} = 1/e$ .
- If  $|p - \hat{p}| \leq \text{err}(n, p, \epsilon)$  then  $\frac{1}{e^2} \leq p^{\hat{a}} \leq \frac{1}{e^{3/2}}$ .
- Works for  $p \in [\epsilon^2/n, 1 - \epsilon^2/n]$ .

**Question:** What if  $p$  is closer to 1 or 0?

**Answer:** For  $p \in [\epsilon^2/n^d, 1 - \epsilon^2/n^d]$  we need  $O(\log(d)/\epsilon^2)$  samples.

---

## Algorithm 1 Binomial Powers

---

**Input :**  $O(\ln(1/\delta)^2/\varepsilon^2)$  samples from the powers of  $B(n, p)$ .

**Output :**  $\hat{a}, \hat{q}_1, \hat{q}_2$ .

- 1: Draw  $O(\ln(1/\delta)/\varepsilon^2)$  samples from  $B(n, p)$  to obtain the approximation  $\hat{p}$ .
  - 2: Let  $\hat{a} \leftarrow -1/\ln(\hat{p})$ .
  - 3: Draw  $O(\ln(1/\delta)^2/(\varepsilon^2\psi(p^{\hat{a}})^2))$  samples from  $B(n, p^{\hat{a}})$  to get estimations  $\hat{q}_1, \hat{q}_2$  of  $p$ ,  $\hat{q}_1 \leq p \leq \hat{q}_2$ .
  - 4: **return**  $\hat{a}, \hat{q}_1, \hat{q}_2$
- 

**Question:** How do we obtain  $\hat{q}_1, \hat{q}_2$  ?



# Learning the Parameters of a PBD

---

# Learning the Powers vs Learning the Parameters

**Question:** PBD powers  $\Leftrightarrow$  Parameter Estimation ?

- Assuming that all  $p_i$ 's are well separated we can learn them by sampling from the powers of a PBD.
- Is there a PBDs where learning its **powers** is **easy** but learning its **parameters** is **hard**?

# Learning the Parameters

- $P(x) = \prod (x - p_i) = x^n + c_{n-1}x^{n-1} + \dots + c_0$ .
- $\mu_j = \mathbf{E}P_j = \sum p_i^j$ .

## Newton Identities

$$\begin{pmatrix} 1 & & & & & \\ \mu_1 & 2 & & & & \\ \mu_2 & \mu_1 & 3 & & & \\ \vdots & \vdots & \ddots & \ddots & & \\ \mu_{n-1} & \mu_{n-2} & \dots & \mu_1 & n & \end{pmatrix} \begin{pmatrix} c_{n-1} \\ c_{n-2} \\ c_{n-3} \\ \vdots \\ c_0 \end{pmatrix} = \begin{pmatrix} -\mu_1 \\ -\mu_2 \\ -\mu_3 \\ \vdots \\ -\mu_n \end{pmatrix} \Leftrightarrow \mathbf{Ac} = \mathbf{b}$$

We know the  $\mu_i$ 's only approximately from sampling the PBD powers.

**Question:** How to measure the impact of noise to the solution of the system?

$$\|\mathbf{c} - \hat{\mathbf{c}}\|_{\infty} \leq u O\left(n^{3/2}2^n\right)$$

## Pan's Algorithm

- $P(x) = \sum_{i=0}^n c_i x^i = c_n \prod_{i=1}^n (x - p_i)$ ,  $c_n \neq 0$ .
- $|p_j| \leq 1$  for all  $j$ .
- Computes roots such that  $|\hat{p}_j - p_j| < \varepsilon$  for  $j = 1, \dots, n$ .
- Precision needed:

$$\|\mathbf{c} - \hat{\mathbf{c}}\|_{\infty} = 2^{O(-n \max(\log(1/\varepsilon), \log(n)))}. \quad (2)$$

- Overall we need  $2^{O(n \max(\log(1/\varepsilon), \log(n)))}$  samples.

# Le Cam's Inequality

**Question:** How do we show a sampling complexity Lower Bound ?

## Hypothesis Testing

- We know that samples can come from either  $P_1$  or  $P_2$ .
- How many samples do we need in order to decide whether they come from  $P$  or  $Q$ ?
- $\Psi$  is a testing function,  $\Psi : \mathcal{X} \rightarrow \{1, 2\}$
- Let  $V$  be the random variable of the choice of the unknown distribution.

## Le Cam's Inequality

$$\inf_{\Psi} \mathbb{P} [\Psi(X^N) \neq V] = 1 - d_{\text{tv}}(P_1^N, P_2^N)$$

# Minimax Risk, Sample Complexity

- $\mathfrak{P}$  is a family of distributions.
- $P \in \mathfrak{P}$  is a distribution.
- $\Theta$  is the space of the parameter we want to estimate.
- $\theta : \mathfrak{P} \rightarrow \Theta$ ,  $\theta(P)$  is the parameter of  $P$  we want to estimate.
- $\hat{\theta} : \mathcal{X}^N \rightarrow \Theta$  is the estimator.
- $X^N = (X_1, \dots, X_N) \sim P^N$  is the sample vector,  $N$  is the number of samples.
- $\rho$  is a semimetric on  $\Theta$ .

$$\mathfrak{M}_N(\theta(\mathfrak{P}), \rho) := \inf_{\hat{\theta}} \sup_{P \in \mathfrak{P}} \mathbf{E}_{P^N} \left[ \rho \left( \hat{\theta}(X^N), \theta(P) \right) \right].$$

$$n(\varepsilon, \theta(\mathfrak{P})) = \inf \{ N : \mathfrak{M}_N(\theta(\mathfrak{P}), \rho) \leq \varepsilon \}$$

# Generalization of Minimax Risk

- $\mathfrak{P}$  is a family of *sequences* of distributions.
- $\mathcal{P} \in \mathfrak{P}$  is a *sequence* of distributions.
- $\Theta$  is the space of the parameter we want to estimate.
- $\theta : \mathfrak{P} \rightarrow \Theta$ ,  $\theta(P)$  is the parameter of  $P$  we want to estimate.
- $\hat{\theta} : \mathcal{X}^m \rightarrow \Theta$  is the estimator.
- $X^m = (X_1, \dots, X_{m_1}, \dots, X_{m_k}) \sim P^m$  is the sample vector,  $N$  is the number of samples.
- $\rho$  is a semimetric on  $\Theta$ .

## Definition

$$\mathfrak{M}_N(\theta(\mathfrak{P}), \rho) := \inf_{\hat{\theta}} \inf_{|m|=N} \sup_{\mathcal{P} \in \mathfrak{P}} \mathbf{E}_{P^m} \left[ \rho \left( \hat{\theta}(X^m), \theta(\mathcal{P}) \right) \right]. \quad (3)$$

# From Estimation to Testing

## Canonical Hypothesis Testing

- "nature" chooses  $V$  uniformly from  $\mathcal{V}$ .
- Conditioned on  $V = v$ , we draw the sample  $X^m$  from the  $N$ -fold product distribution  $P_v^m$ .

Given  $X^m$  our goal is to determine  $V$ .

### Lower Bound

Let  $\mathfrak{F}_{\mathcal{V}} \subseteq \mathfrak{P}$  be a family of sequences of distributions indexed by  $v \in \mathcal{V}$  such that  $\rho(\theta(\mathcal{P}_v, \mathcal{P}_u)) \geq 2\delta$  for all  $\mathcal{P}_v, \mathcal{P}_u \in \mathfrak{F}_{\mathcal{V}}$ , where,  $v \neq u \in \mathcal{V}$  and  $\delta > 0$ . Then

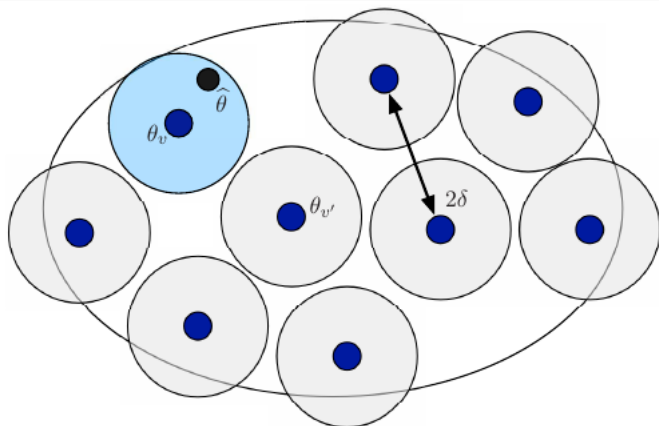
$$\mathfrak{M}_N(\theta(\mathfrak{P}), \rho) \geq \delta \inf_{m=|N|} \inf_{\Psi} \mathfrak{v}^m(\Psi(X^m) \neq V).$$



# From Estimation to Testing

**Testing Function:**

$$\Psi(X^m) := \operatorname{argmin}_{v \in \mathcal{V}} \{\rho(\hat{\theta}, \theta_v)\}, \rho(\hat{\theta}, \theta_v) \leq \delta \Leftrightarrow \Psi(\hat{\theta}) = v.$$



**Main Idea:** Find **two sequences** of distributions that are **close** in **Total Variation** but their **parameters** are **far**.

- $\mathcal{P}, \mathcal{Q} \in \mathfrak{F}$  and  $\delta > 0$ .
- $\rho(\theta(\mathcal{P}), \theta(\mathcal{Q})) \geq 2\delta$  then

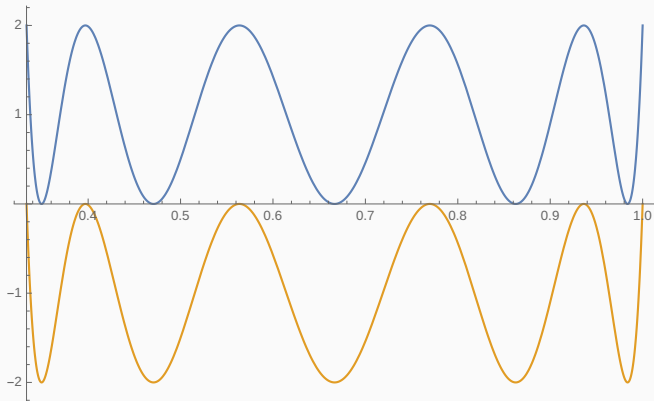
After  $N$  observations (samples) the minimax risk has lower bound

$$\mathfrak{M}_N(\theta(\mathfrak{F}), \rho) \geq \frac{\delta}{2} (1 - \sqrt{2} \sqrt{1 - (1 - d_{\text{tv}}(\mathcal{P}, \mathcal{Q}))^N}).$$

# The Lower Bound

Construction of [2]:

- $n = \Theta(\log(N/\varepsilon))$
- $p_j := (1 + \cos(\frac{2\pi j}{n}))/8$ ,  $q_j := (1 + \cos(\frac{2\pi j + \pi}{n}))/8$ ,  $j \in [n]$ .
- $j = n/4 + O(1)$ , we have that  $|p_i - q_j| = \Omega(1/\log(N/\varepsilon))$



# The Lower Bound

- $p_i$ 's are the roots of  $T_n(8x - 1) - 1$ .
- $q_j$ 's are the roots of  $T_n(8x - 1) + 1$ .
- for all  $l \in \{1, 2, \dots, n - 1\}$ ,  $\sum_{i=1}^n p_i^l = \sum_{i=1}^n q_i^l$
- for  $l \geq n$ ,  $3^l (\sum_{i=1}^n (p_i^l - q_i^l)) \leq n(3/4)^n = \log(N/\epsilon)(3/4)^{\log(N/\epsilon)}$ .

All PBD powers of  $\mathbf{p}, \mathbf{q}$  are very close in TVD.

$$d_{\text{tv}}(P_i, Q_i) \leq c/N$$

. Therefore the number of samples  $N$  should be  $\Omega(2^{1/\epsilon})$ .

**Questions?**



L. Birgé.

**Estimating a Density under Order Restrictions: Nonasymptotic Minimax Risk.**

*The Annals of Statistics*, 15(3):995–1012, Sept. 1987.



I. Diakonikolas, D. Kane, and A. Stewart.

**Properly learning poisson binomial distributions in almost polynomial time.**

In *Proceedings of the 29th Conference on Learning Theory, (COLT'16)*, pages 850–878, 2016.



B. Roos.

**Improvements in the Poisson approximation of mixed Poisson distributions.**

*Journal of Statistical Planning and Inference*, 113:467–483, 2003.