# Testing Closeness of distributions

- Poissonization method
- Flattening Technique
- Estimating $L2$ distance
- $L1$ closeness tester

Problem: sample access to $P$ and $q$

Test whether $\begin{cases} P = q \\ \quad vs \\ \|P - q\|_1 \geq \varepsilon \end{cases}$

$\boxed{\text{Poissonization method}} \leftarrow$ A general method that facilitates the analysis of distribution testing algorithm by making the numbers of instances of different elements independent.

Sample set: $S = \{ s_1, \ldots, s_m \}$     $X_i := \#$instances of $i$ in $S$

- Example    $S = \{ 2, 5, 3, 2, 3 \}$     $X_2 = X_3 = 2$, $X_5 = 1$

Main Difficulty:

For a fixed $m$,   $X_i$'s are $\boxed{not}$ independent

e.g. if $X_4 = \frac{m}{2}$, then $X_3 \leq \frac{m}{2}$.

- Can we make $X_i$'s independent? Yes, via Poissonization...

(2) For $i = 1, \ldots, n$

$\qquad X_i \sim Poi(s \cdot p_i)$

independent $X_i$'s

$=$

\* These two processes result in the same distribution on $X_i$'s

(1) $\hat{m} \leftarrow Poi(m)$

For $i = 1, \ldots, \hat{m}$

$\qquad s_i \leftarrow$ Draw a sample from $p$.

$S = \{ s_1, s_2, \ldots, s_{\hat{m}} \}$

Compute $X_i$'s from $S$.

Proof of **\***    For any $c = 0, 1, 2, \ldots$

Recall:



$X \sim Poi(\lambda)$

$Pr[X = k] = \dfrac{e^{-\lambda} \lambda^k}{k!}$

$E[X] = Var[X] = \lambda$

$Pr[X_i = c \text{ according to } (1)]$

$$= \sum_{k=c}^{\infty} Pr[\hat{m} = k] \cdot \binom{k}{c} \cdot p_i^{c} \cdot (1-p_i)^{k-c}$$

$$= \sum_{k=c}^{\infty} \frac{e^{-m} m^k}{k!} \cdot \frac{k!}{(k-c)! \, c!} \cdot p_i^{c} \cdot (1-p_i)^{k-c}$$

$$= \frac{e^{-m} m^c p_i^{c}}{c!} \underbrace{\sum_{k=c}^{\infty} \frac{m^{k-c} (1-p_i)^{k-c}}{(k-c)!}}_{} = \sum_{k'=0}^{\infty} \frac{(m(1-p_i))^{k'}}{k'!} = e^{m(1-p_i)}$$

$$= e^{-m + m(1-p_i)} \frac{(mp_i)^c}{c!} = e^{mp_i} \frac{(mp_i)^c}{c!}$$

$$= Pr[X_i = c \text{ according to } (2)]$$

↳ Probability of observing $X_i = c$ when $X_i \sim Poi(mp_i)$

Done!

$\boxed{L_p \text{ norm}}$

$$\| q \|_p = \left( \sum_i (q^{(i)})^p \right)^{1/p}$$

$\boxed{L_p \text{ distance}}$

$$\| q_1 - q_2 \|_p = \left( \sum_i (q_1^{(i)} - q_2^{(i)})^p \right)^{1/p}$$

# Reducing the L2 norm of distributions

Goal: transform a distribution, $p$, to another distribution, $p'$, such that $\|p'\|_2^2$ is low.

We use the following randomized process:

$S \leftarrow$ Draw $Poi(k)$ samples from $p$

$b_x \leftarrow$ the number of instances of $x \in S$ $\qquad \forall x = 1, \ldots, n$

For each element $x$ in the domain of, we assign $b_x + 1$ elements in the domain of $p'$ to $x$

To generate a sample from $p'$:

1) Draw $x \sim P$

2) Pick $y$ uniformly at random from $[b_i + 1]$

3) Output $(x, y)$

Example: $\qquad S = \{3, 3, 1\}$

$P = $

| 0.2 | 0.05 | 0.75 |
|---|---|---|
| 1 | 2 | 3 |

$p' = $

| 0.1 | 0.1 | 0.05 | 0.25 | 0.25 | 0.25 |
|---|---|---|---|---|---|
| (1,1) | (1,2) | 2 | (3,1) | (3,2) | (3,3) |

Facts about $p'$

- Domain $= \{(x,y) \mid x \in [n] \land y \in [b_x + 1]\}$

- Domain size $= n + k$.

- $p'(x,y) = \dfrac{p(x)}{b_i + 1}$

3

- $E[\|P'\|_2^2] \le \dfrac{1}{k}$  (over the randomness of $S$)

$$E[\|P'\|_2^2] = E\left[\sum_{x=1}^{n}\sum_{y=1}^{b_x+1} P'(x,y)^2\right] = E\left[\sum_{x}\sum_{y}\frac{P(x)^2}{(b_x+1)^2}\right]$$

$$= E\left[\sum_{x}\frac{P(x)^2}{b_x+1}\right] \overset{*}{\le} \sum_{x}\frac{P(x)^{\prime}}{k\cdot P(x)} = \frac{1}{k}$$

Proof of $*$ in [DK '16]

Let $Z \sim Poi(\lambda)$, then

$$E[a^Z] = \sum_{z=0}^{\infty}\frac{e^{-\lambda}(\lambda a)^z}{z!} = e^{\lambda(a-1)}\sum_{z=0}^{\infty}\frac{e^{-\lambda a}(\lambda a)^z}{z!} = e^{\lambda(a-1)}$$

$$E\left[\frac{1}{Z+1}\right] = E\left[\int_0^1 a^Z\, da\right] = \int_0^1 E[a^Z]\, da = \int_0^1 e^{\lambda(a-1)}\, da$$

$$= \frac{1}{\lambda}e^{\lambda(a-1)}\Big|_{a=0}^{1} = \frac{1}{\lambda}(1-e^{-\lambda}) \le \frac{1}{\lambda}$$

Alternative proof of $*$

$$E\left[\frac{1}{Z+1}\right] = \sum_{z=0}^{\infty}\frac{e^{-\lambda}\lambda^z}{(Z+1)!} = \frac{1}{\lambda}\sum_{z=0}^{\infty}\frac{e^{-\lambda}\lambda^{z+1}}{(Z+1)!} = \frac{1}{\lambda}\sum_{z'=1}^{\infty}\frac{e^{-\lambda}\lambda^{z'}}{z'!}$$

$$= \frac{1}{\lambda}\cdot\left(\sum_{z'=0}^{\infty}\frac{e^{-\lambda}\lambda^{z'}}{z'!}\right) - \frac{e^{-\lambda}}{\lambda} = \frac{1-e^{-\lambda}}{\lambda}$$

— Let $q'$ be the transformed version of $q$ with samples in $S$.

Then $\|P-q\|_1 = \|p'-q'\|_1$

$$\|P'-q'\|_1 = \sum_{x=1}^{n}\sum_{y=1}^{b_x+1}|p'(x)-q'(x)| = \sum_{x}\sum_{y}\frac{|P(x)-q(x)|}{b_x+1}$$

$$= \sum_{x}|p(x)-p(y)| = \|P-q\|_1$$

4

- For a known distribution $q$:

$$b_x = \lfloor \ln q(x) \rfloor$$

Then, we have

$$\|q'\|_2^2 = \sum_{x=1}^{\hat{n}} \sum_{y=1}^{b_x+1} q'(x)^2 = \sum_x \frac{q(x)^2}{\lfloor \ln q_i \rfloor + 1}$$

$$\leq \sum_{\substack{x \\ q(x) \neq 0}} \frac{q(x)^2}{n\, q(x)} \leq \frac{1}{n}$$

new domain size $\sum_{x=1}^{n} b_x + 1 = n + \sum_x \lfloor \ln q_i \rfloor$

$$\leq n + n \cdot \sum_x p(x) \leq 2n$$



General Framework of [DK'16]

$p', q' \leftarrow$ Flatten $p$ and $q$ using $Poi(k)$ from $p$

Estimate $\|P\|_2^2$ and $\|q\|_2^2$ within constant factor error

if the estimation of $\|P\|_2^2$ and $\|q\|_2^2$ are more than a constant apart from each other:

$\qquad$ infer $p \neq q$ and reject.

Else

$\qquad$ Test $p' = q'$ Given that $\underbrace{\|P\|_2^2 \pm \|q\|_2^2} \leq \theta(\frac{1}{k})$

Note that $\|P\|_2^2$ is low due to flattening

and $\|q\|_2^2$ is low because it is within a constant factor of $\|P\|_2^2$

## $L_2$ distance estimator.

[Chan, Diakonikolas, Valiant, Valiant '14]

Let $p, q$ be two distribution that we have sample access to
They provide a statistic $Z$ where

$$Pr\left[ \left| Z/_{m^2} - \|p-q\|_2^2 \right| \geq \varepsilon^2 \right] \leq \frac{2}{3}$$

- How to compute the estimation of $\|p-q\|_2^2$ ?

  - Draw $Poi(m)$ sample from $p$ and $q$
  - Let $\begin{cases} X_i & \text{be \# instances of } i \text{ in the samples from } p. \\ Y_i & \text{\textquotedblright} \quad \text{\textquotedblright} \quad \text{\textquotedblright} \quad q. \end{cases}$

  - $Z = \sum_{i=1}^{\hat{}} (X_i - Y_i)^2 - X_i - Y_i$

  - output $Z/m^2$.

- steps of the analysis

  $b$ is $\max\left( \|p\|_2^2, \|q\|_2^2 \right)$

  ↗ next page.

  1) $E[Z] = m^2 \|p-q\|_2^2$ , $Var[Z] \leq 8m^3 \|p-q\|_4^2 \sqrt{b} + 8m^2 b$

  2) Chebyshev's inequality:

  $$Pr\left[ \left| \frac{Z}{m^2} - \|p-q\|_2^2 \right| \geq \varepsilon^2 \right] = Pr\left[ |Z - E[Z]| \geq m^2 \varepsilon^2 \right]$$

  $$\leq \frac{Var[Z]}{m^4 \varepsilon^4} \leq \frac{1}{10}$$

  $\Rightarrow$ for $m \geq \theta\left( \frac{\sqrt{b} \|p-q\|_4^2}{\varepsilon^4} + \frac{\sqrt{b}}{\varepsilon^2} \right)$

Let's analyze $Z$: $\begin{cases} X_i \sim Poi(mp_i) \\ Y_i \sim Poi(mq_i) \end{cases}$ $\qquad E[2X_iY_i] \overset{by\ indep.}{=}$

$$E[Z] = \sum_{i=1}^{n} E[(X_i - Y_i)^2 - X_i - Y_i] \qquad \overset{\nearrow}{} \quad 2E[X_i]\cdot E[Y_i]$$
$$= 2(m\,p_i)\cdot(m\,q_i)$$

$$= \sum_{i=1}^{n} E[X_i^2 - X_i] \underbrace{-2E[X_i]\cdot E[Y_i]}_{} + E[Y_i^2 - Y_i]$$

$$= \sum_{i=1}^{n} m^2 p_i^2 - 2m^2 p_i q_i + m^2 q_i^2$$

$$= \sum_{i=1}^{n} m^2 (p_i - q_i)^2 = m^2 \| p - q \|_2^2$$

if $X \sim Poi(\lambda)$
$$E[X^2 - X] = E[X^2] - \lambda$$
$$= Var[X] + E[X]^2 - \lambda$$
$$= \lambda + \lambda^2 - \lambda = \lambda^2$$

So, $E[Z/m^2]$ is exactly the distance that we want!

$$Var[Z] \overset{\text{by independence}}{=} \sum_{i=1}^{n} Var[(X_i - Y_i)^2 - X_i - Y_i]$$

$$= \sum_{i=1}^{n} E[((X_i - Y_i)^2 - X_i - Y_i)^2] - E[(X_i - Y_i)^2 - X_i - Y_i]^2$$

$$\vdots$$
$$= \sum_{i=1}^{n} \text{bunch of term lik } E[X_i^\ell] \text{ where } \ell = 1, 2, 3, 4$$
$$\vdots \qquad \to \text{moments of } Poi(\lambda) \to \begin{array}{c} \text{we can look} \\ \text{them up!} \end{array}$$

$$= \sum_{i=1}^{n} 4m^3 (p_i - q_i)^2 (p_i + q_i) + 2m^2 (p_i + q_i)^2$$

$$\leq \sum_{i=1}^{n} 4m^3 \sqrt{\sum_{i=1}^{n} (p_i - q_i)^4 \sum_{i} (p_i + q_i)^2} + 2m^2 \underbrace{(p_i + q_i)^2}_{}$$

Cauchy-Schwarz

$$\leq 8m^3 \| p - q \|_4^2 \sqrt{b} + 8m^2 b$$

$$\leq \sum_{i} 2(p_i^2 + q_i^2)$$
$$\leq 4\max(\|p\|_2^2 ; \|q\|_2^2)$$
$$\underbrace{\qquad}_{b}$$
Let's call this $b$

$\boxed{L_1 \text{ closeness tester base on } \|P-q\|_2^2 \text{ estimation}}$
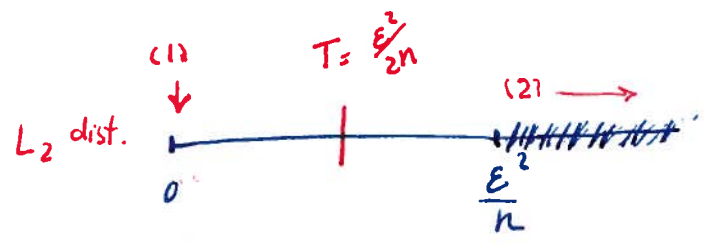
How $L_1$ distance is related to $L_2$ distance ?

$$\|P-q\|_2^2 \le \|p-q\|_1^2 \le n \cdot \|P-q\|_2^2$$

$L_p$ norm inequality ↲          Cauchy-schwarz ↓

$L_1$ distance closeness tester distinguishes :

$$\begin{cases} (1) \|P-q\|_1 = 0 & \Rightarrow \|P-q\|_2^2 = 0 \\[2em] (2) \|P-q\|_1 \ge \varepsilon & \Rightarrow \|P-q\|_2^2 \ge \dfrac{\varepsilon^2}{n} \end{cases}$$

$(1)$   $T = \dfrac{\varepsilon^2}{2n}$    $(2) \longrightarrow$    with probability $0.9$ for large $m$

$L_2$ dist.

$\overset{(1)\downarrow}{0} \longmapsto \overset{|}{\underset{\frac{\varepsilon^2}{n}}{\;}} \overset{(2)\longrightarrow}{\text{HHHHHHH}}$

$1) \Rightarrow \frac{Z}{m^2} \le T$

$2) \Rightarrow \frac{Z}{m^2} \ge T$

$(1)$
$$\Pr\left[ \left| \frac{Z}{m^2} - 0 \right| \ge \frac{\varepsilon^2}{2n} \right] \le \frac{n^2 \, Var[Z]}{m^4 \varepsilon^4}$$

$$\le \frac{n^2}{m^4 \varepsilon^4} \cdot \left( 8m^3 \underbrace{\|P-q\|_4^2}_{=0} \sqrt{b} + 8m^2 b \right)$$

$$\le \frac{n^2 b}{\varepsilon^4 m^2} \le \frac{1}{10} \quad \Leftarrow \quad m = \theta\left( \frac{n\sqrt{b}}{\varepsilon^2} \right)$$

$(2)$
$$\Pr\left[ \frac{Z}{m^2} \ge \frac{\varepsilon^2}{2n} \right] \le \Pr\left[ \left| \frac{Z}{m^2} - \|P-q\|_2^2 \right| \ge \frac{\|P-q\|_2^2}{2} \right]$$

$$\le \frac{4 \cdot Var[Z]}{m^4 \|P-q\|_2^4} \le \frac{32 \|P-q\|_4^2 \sqrt{b}}{m \|P-q\|_2^{*2}} + \frac{32\, b}{m^2 \|P-q\|_2^4}$$

$$\le \frac{32 n \sqrt{b}}{m\, \varepsilon^2} + \frac{32\, b\, n^2}{m^2\, \varepsilon^4} \le \frac{1}{10} \quad \Leftarrow \quad m = \theta\left( \frac{n\sqrt{b}}{\varepsilon^2} \right)$$

✱ Putting everything together ....

by Markov $\Pr\left[\|P\|_2^2 \geq \frac{10}{k}\right] \leq \frac{1}{10}$.

Let say with probability $0.9$ we estimate $\|P\|_2^2$ and $\|Q\|_2^2$ correctly

with probability $0.9$ $\ell_7$ tester works!

$\Rightarrow$ with probability $0.7 \geq \frac{2}{3}$ the tester works

(union bound) ↙

Sample complexity

$$\Theta\left(k + \frac{n\sqrt{b}}{\varepsilon^2}\right) = \Theta\left(k + \frac{n}{\sqrt{k}\,\varepsilon^2}\right) \quad \overset{optimize}{\leadsto} \quad \begin{array}{l} \text{for } k \\ k \leq n \end{array}$$

$$\Rightarrow \boxed{\Theta\left(\frac{n^{2/3}}{\varepsilon^{4/3}} + \frac{\sqrt{n}}{\varepsilon^2}\right)}$$

optimal sample complexity
for testing closeness