

Sample Complexity of Robust Hypothesis Testing

Shankar Vallinayagam Alston Xu

Summer Research Festival

Supervised by Professor Varun Jog

Simple hypothesis testing

Simple hypothesis testing

- Let p and q be known distributions over $\{1, 2, \dots, k\}$

Simple hypothesis testing

- Let p and q be known distributions over $\{1, 2, \dots, k\}$
- Observe n i.i.d. samples from one of them, equally likely

Simple hypothesis testing

- Let p and q be known distributions over $\{1, 2, \dots, k\}$
- Observe n i.i.d. samples from one of them, equally likely
- **Goal:** Decide whether p or q using samples

Simple hypothesis testing

- Let p and q be known distributions over $\{1, 2, \dots, k\}$
- Observe n i.i.d. samples from one of them, equally likely
- **Goal:** Decide whether p or q using samples
- **Sample complexity n^* :** Smallest n such that $\mathbb{P}(\text{error}) \leq 0.1$ for optimal test

Characterising n^* in terms of TV and Hellinger distance

Characterising n^* in terms of TV and Hellinger distance

- TV-distance defined as

$$d_{TV}(p, q) = \sup_S p(S) - q(S) = \frac{1}{2} \sum |p(i) - q(i)|$$

Characterising n^* in terms of TV and Hellinger distance

- **TV-distance** defined as

$$d_{TV}(p, q) = \sup_S p(S) - q(S) = \frac{1}{2} \sum |p(i) - q(i)|$$

Sample complexity using d_{TV}

$$\frac{1}{d_{TV}(p, q)} \lesssim n^* \lesssim \frac{1}{d_{TV}(p, q)^2}$$

Characterising n^* in terms of TV and Hellinger distance

- **TV-distance** defined as

$$d_{TV}(p, q) = \sup_S p(S) - q(S) = \frac{1}{2} \sum |p(i) - q(i)|$$

Sample complexity using d_{TV}

$$\frac{1}{d_{TV}(p, q)} \lesssim n^* \lesssim \frac{1}{d_{TV}(p, q)^2}$$

- **Hellinger distance** defined as

$$d_h^2(p, q) = \sum_x (\sqrt{p(x)} - \sqrt{q(x)})^2 = 2(1 - \langle \sqrt{p}, \sqrt{q} \rangle)$$

Characterising n^* in terms of TV and Hellinger distance

- **TV-distance** defined as

$$d_{TV}(p, q) = \sup_S p(S) - q(S) = \frac{1}{2} \sum |p(i) - q(i)|$$

Sample complexity using d_{TV}

$$\frac{1}{d_{TV}(p, q)} \lesssim n^* \lesssim \frac{1}{d_{TV}(p, q)^2}$$

- **Hellinger distance** defined as

$$d_h^2(p, q) = \sum_x (\sqrt{p(x)} - \sqrt{q(x)})^2 = 2(1 - \langle \sqrt{p}, \sqrt{q} \rangle)$$

Sample complexity is characterized by Hellinger

$$n^* = \Theta\left(\frac{1}{d_h^2(p, q)}\right)$$

Robust statistics

- **Robust statistical procedures** are those that demonstrate insensitivity to small deviations from the assumptions made.

- **Robust statistical procedures** are those that demonstrate insensitivity to small deviations from the assumptions made.
- The optimal solution to this hypothesis testing problem in the non-robust setting is the **likelihood ratio test**:

$$\prod_{i=1}^k \frac{p(i)}{q(i)}$$

- **Robust statistical procedures** are those that demonstrate insensitivity to small deviations from the assumptions made.
- The optimal solution to this hypothesis testing problem in the non-robust setting is the **likelihood ratio test**:

$$\prod_{i=1}^k \frac{p(i)}{q(i)}$$

- Our adversaries are
 - (Oblivious) Huber contamination
 - TV-contamination
 - Adaptive Huber contamination

Adversaries

Huber Adversary

$$B_{\text{Huber}}(p, \epsilon) = \{p^* : p^* = (1 - \epsilon)p + \epsilon h\}$$

Given the underlying distribution is p , adversary selects distribution in this ball for data to be drawn from

Adversaries

Huber Adversary

$$B_{\text{Huber}}(p, \epsilon) = \{p^* : p^* = (1 - \epsilon)p + \epsilon h\}$$

Given the underlying distribution is p , adversary selects distribution in this ball for data to be drawn from

TV Adversary

$$B_{\text{TV}}(p, \epsilon) = \{p^* : d_{\text{TV}}(p, p^*) \leq \epsilon\}$$

Given the underlying distribution is p , adversary selects distribution in this ball for data to be drawn from

Adversaries

Huber Adversary

$$B_{\text{Huber}}(p, \epsilon) = \{p^* : p^* = (1 - \epsilon)p + \epsilon h\}$$

Given the underlying distribution is p , adversary selects distribution in this ball for data to be drawn from

TV Adversary

$$B_{\text{TV}}(p, \epsilon) = \{p^* : d_{\text{TV}}(p, p^*) \leq \epsilon\}$$

Given the underlying distribution is p , adversary selects distribution in this ball for data to be drawn from

Adaptive Huber adversary

Take a sample of size n , then adversary selects ϵn of these samples and changes them

Optimal test for oblivious adversaries

Optimal test for oblivious adversaries

- We want to find a test that minimises

$$\max_{p^* \in B(p, \epsilon), q^* \in B(q, \epsilon)} n^*(p^*, q^*)$$

Optimal test for oblivious adversaries

- We want to find a test that minimises

$$\max_{p^* \in B(p, \epsilon), q^* \in B(q, \epsilon)} n^*(p^*, q^*)$$

- The way Huber, 1964 does this is by finding **Least Favourable Distributions** that maximise error for all LRTs

Optimal test for oblivious adversaries

- We want to find a test that minimises

$$\max_{p^* \in B(p, \epsilon), q^* \in B(q, \epsilon)} n^*(p^*, q^*)$$

- The way Huber, 1964 does this is by finding **Least Favourable Distributions** that maximise error for all LRTs
- In both the oblivious cases this leads to a **Censored Likelihood Ratio Test**; we set thresholds c', c'' such that for any observation if the likelihood ratio is above c'' or below c' we censor down to one of these values

Optimal test for oblivious adversaries

- We want to find a test that minimises

$$\max_{p^* \in B(p, \epsilon), q^* \in B(q, \epsilon)} n^*(p^*, q^*)$$

- The way Huber, 1964 does this is by finding **Least Favourable Distributions** that maximise error for all LRTs
- In both the oblivious cases this leads to a **Censored Likelihood Ratio Test**; we set thresholds c', c'' such that for any observation if the likelihood ratio is above c'' or below c' we censor down to one of these values
- The problem reduces to understanding the Hellinger distance between the LFDs

Small ϵ Regime

Small ϵ Regime

- For small epsilon $\epsilon \leq \frac{d_H(p,q)^2}{9}$

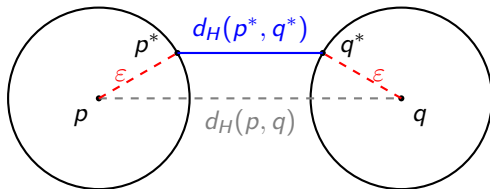
$$n_{TV}^*(\epsilon) \asymp n_{Huber}^*(\epsilon) \asymp n^*(p, q)$$

Small ϵ Regime

- For small epsilon $\epsilon \leq \frac{d_H(p,q)^2}{9}$

$$n_{TV}^*(\epsilon) \asymp n_{Huber}^*(\epsilon) \asymp n^*(p, q)$$

- This follows trivially from the triangle inequality



Conjectures

- We conjectured the sample complexity for Huber and TV stays within constant for more general ϵ : $n_{TV}^*(\epsilon) \asymp n_{Huber}^*(\epsilon) \quad \forall \epsilon$

- We conjectured the sample complexity for Huber and TV stays within constant for more general ϵ : $n_{TV}^*(\epsilon) \asymp n_{Huber}^*(\epsilon) \quad \forall \epsilon$
- We also conjectured the sample complexity stays within constants if we scale ϵ to $\epsilon/2$:

$$\begin{aligned} n_{TV}^*(\epsilon) &\asymp n_{TV}^*(\epsilon/2) \\ n_{Huber}^*(\epsilon) &\asymp n_{Huber}^*(\epsilon/2) \end{aligned}$$

Counterexample

Counterexample

- We disproved the above statements by coming up with the following pair of 7-point distribution counterexample.

i	p	q	$\frac{p_i}{q_i}$
3	$2kt^{1+\delta}$	$kt^{1+\delta}$	2
2	$k(1+t^{1-\delta})t^{2\delta}$	$kt^{2\delta}$	$1+t^{1-\delta}$
1	$14k(1+t^{1+\delta})$	$14k$	$1+t^{1+\delta}$
0	$1 - (\dots)$	$1 - (\dots)$	1
-1	$14k$	$14k(1+t^{1+\delta})$	$\frac{1}{1+t^{1+\delta}}$
-2	$kt^{2\delta}$	$k(1+t^{1-\delta})t^{2\delta}$	$\frac{1}{1+t^{1-\delta}}$
-3	$kt^{1+\delta}$	$2kt^{1+\delta}$	$\frac{1}{2}$

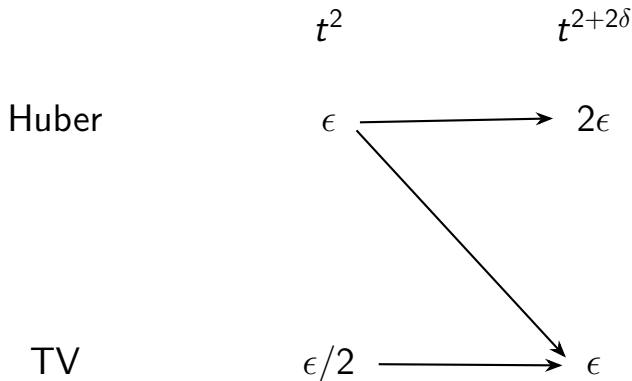
Counterexample

Counterexample

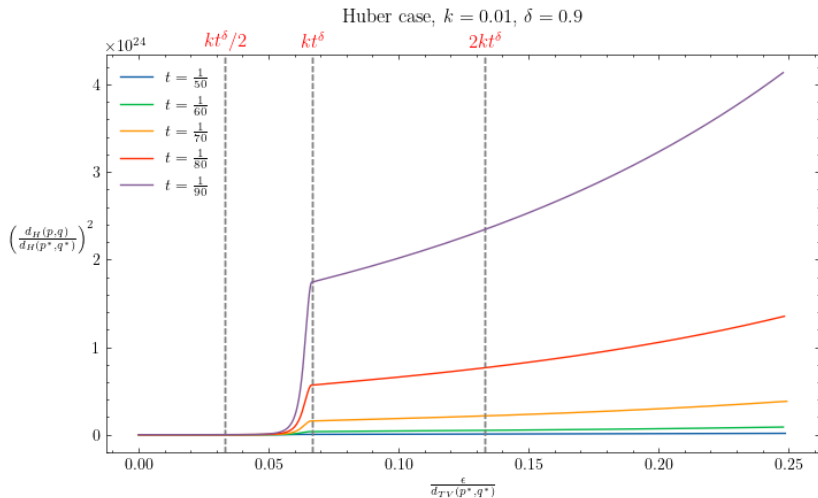
- **Key idea:** Clipping at $1 + t$ and $1 + t^{1-\delta}$ and take $t \rightarrow 0$

Counterexample

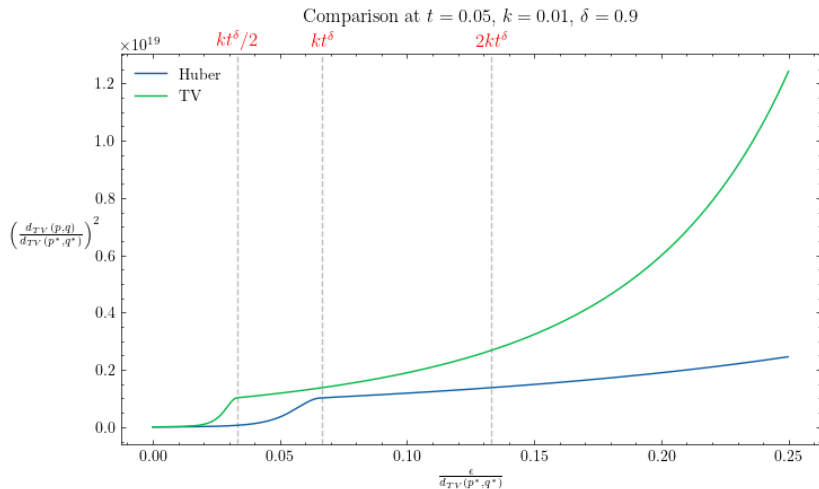
- **Key idea:** Clipping at $1 + t$ and $1 + t^{1-\delta}$ and take $t \rightarrow 0$
- Here, ϵ is up to a negligible error term kt^δ



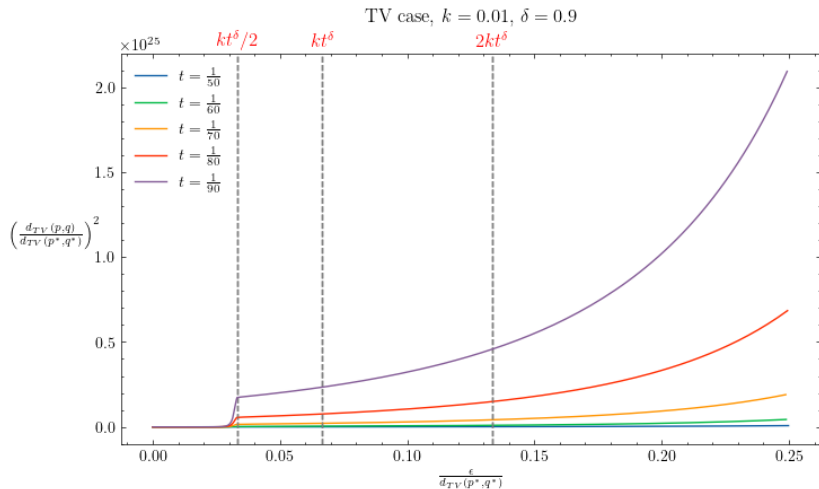
Huber contamination



Huber vs TV



TV contamination



Additional/Ongoing results

- Motivated by our counter example, we are working on proving:

$$n_{TV}^*(\epsilon/2) \lesssim n_{Huber}^*(\epsilon)$$

- Motivated by our counter example, we are working on proving:

$$n_{TV}^*(\epsilon/2) \lesssim n_{Huber}^*(\epsilon)$$

- We are also working on the following result:

$$n_{TV}^*(\epsilon/2) \lesssim n_{Adv}^*(\epsilon) \lesssim n_{TV}^*(\epsilon)$$

Thank you!