

# A simple polynomial-time approximation algorithm for the total variation distance between two product distributions

Weiming Feng<sup>1</sup>, Heng Guo<sup>1</sup>, Mark Jerrum<sup>2</sup>, Jiaheng Wang<sup>1</sup>

(1) University of Edinburgh (2) Queen Mary, University of London

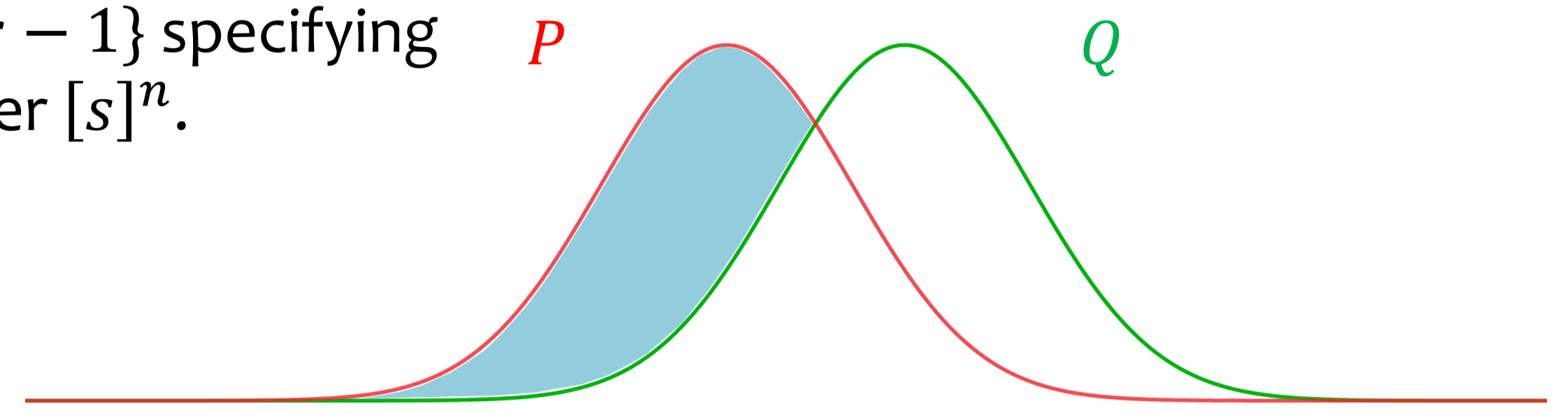
## The total variation distance between two product distributions

**Input:** distributions  $P_1, P_2, \dots, P_n$  and  $Q_1, Q_2, \dots, Q_n$  over finite domain  $[s] = \{0, 1, \dots, s-1\}$  specifying two **product distributions**  $P = P_1 \times P_2 \times \dots \times P_n$  and  $Q = Q_1 \times Q_2 \times \dots \times Q_n$  over  $[s]^n$ .

**Output:** the **total variation distance (TV distance)** between  $P, Q$

$$d_{TV}(P, Q) = \frac{1}{2} \sum_{x \in [s]^n} |P(x) - Q(x)| = \max_{A \subseteq [s]^n} |P(A) - Q(A)|.$$

**Challenge:** the size of sample space of  $P, Q$  is **exponentially large**  $s^n$ .



## Background & previous results [Bhattacharyya, Gayen, Meel, Myrasiotis, Pavan, Vinodchandran, 2022]

**Hardness:** the exact computation of the TV distance between two product distributions is **#P-complete**.

**Approximate the TV distance:** Given  $P, Q$  and an error bound  $0 < \epsilon < 1$ , an **FPRAS** outputs a random  $\hat{d}$  in time  $\text{poly}(n, 1/\epsilon)$  such that

$$\Pr[(1 - \epsilon)d_{TV}(P, Q) \leq \hat{d} \leq (1 + \epsilon)d_{TV}(P, Q)] \geq \frac{2}{3}.$$

**Algorithm:** there is an FPRAS if  $s = 2$  (Boolean domain) and  $\frac{1}{2} \leq P_i(1) < 1$  and  $0 < Q_i(1) \leq P_i(1)$  for all  $1 \leq i \leq n$ .

## Our results [Feng, Guo, Jerrum, Wang, SOSA 2023]

There is an **FPRAS** for the TV distance between two product distributions

- $O(n^2/\epsilon^2)$ -time assuming the cost of each arithmetic operation is  $O(1)$ ;
- each operation acts on  $\text{poly}(n)$ -bit numbers if input par. has  $\text{poly}(n)$  bits.

## Open problems

- **Deterministic** approximate algorithm (FPTAS)?
- FPTAS/FPRAS **beyond** the product distribution?

## TV distance and coupling

A **Coupling** of two distribution  $P, Q$  over  $\Omega$  is a pair of **joint random variables**

$$(X, Y) \in \Omega \times \Omega \text{ such that } X \sim P \text{ and } Y \sim Q$$

**Coupling lemma:** for any coupling  $(X, Y)$  of  $P, Q$ ,

$$d_{TV}(P, Q) \leq \Pr_{\text{coupling}}[X \neq Y],$$

and there exists an **optimal coupling** of  $P, Q$  such that

$$d_{TV}(P, Q) = \Pr_{\text{opt}}[X \neq Y].$$

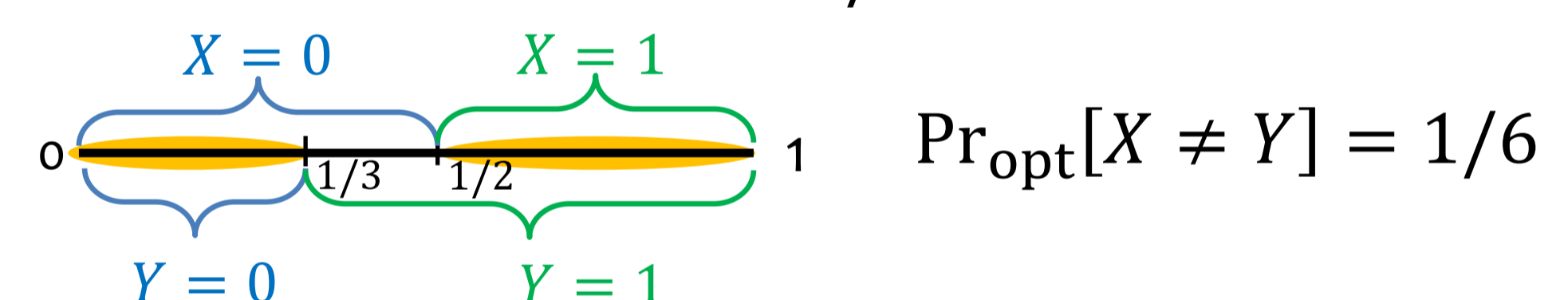
## Examples of coupling

**Distributions:**  $P(0) = P(1) = \frac{1}{2}$  and  $Q(0) = \frac{1}{3}, Q(1) = \frac{2}{3}$ .

**Ind. Coupling:** sample  $X \sim P, Y \sim Q$  independently

$$\Pr_{\text{ind}}[X \neq Y] = 1/2.$$

**Opt. Coupling:** sample real  $r \in [0, 1]$  uniformly at random  
 $X = 0$  iff  $r \leq 1/2$  and  $Y = 0$  iff  $r \leq 1/3$ .



## Greedy coupling for product distributions

**Product distributions:**  $P = P_1 \times P_2 \times \dots \times P_n$  and  $Q = Q_1 \times Q_2 \times \dots \times Q_n$ .

**Greedy coupling:** couple each  $(P_i, Q_i)$  **optimally** and **independently**.

**Example:** Boolean distributions  $P, Q \in \{0, 1\}^n$

- sample  $r_i \in [0, 1]$  uniformly and ind. for all  $i$ ;
- $X_i = 0$  iff  $r_i \leq P_i(0)$  and  $Y_i = 0$  iff  $r_i \leq Q_i(0)$ ;
- $X = (X_1, X_2, \dots, X_n)$  and  $Y = (Y_1, Y_2, \dots, Y_n)$ .

## Properties of the greedy coupling

**Non-optimal:** Ex  $P = P_1 \times P_2$  and  $Q = Q_1 \times Q_2$  s.t.  $\forall i$

$$P_i(0) = P_i(1) = \frac{1}{2} \text{ and } Q_i(0) = \frac{1}{2} - \delta, Q_i(1) = \frac{1}{2} + \delta.$$

$$\Pr[X \neq Y] = \delta(2 - \delta) > \delta(1 + \delta) = d_{TV}(P, Q).$$

**Poly(n)-approximation:** for any product distributions,

$$d_{TV}(P, Q) \leq \Pr_{\text{greedy}}[X \neq Y] \leq n \cdot d_{TV}(P, Q).$$

## Our idea: estimate the ratio

**Fact:** Prob. of  $X \neq Y$  in greedy coupling is **easy to compute**

$$\Pr_{\text{greedy}}[X \neq Y] = 1 - \Pr_{\text{greedy}}[X = Y] = 1 - \prod_{i=1}^n (1 - d_{TV}(P_i, Q_i)).$$

**Lemma:** There is FPRAS for the ratio

$$R = \frac{d_{TV}(P, Q)}{\Pr_{\text{greedy}}[X \neq Y]} \geq \frac{1}{n}.$$

**Estimator** for the ratio  $R$ :

- $\pi$ : distribution of  $X$  in the **greedy coupling conditional on  $X \neq Y$**

$$\forall \sigma \in [s]^n, \quad \pi(\sigma) = \frac{\Pr_{\text{greedy}}[X = \sigma \mid X \neq Y]}{\Pr_{\text{greedy}}[X \neq Y]}.$$

- $f$ : a function  $[s]^n \rightarrow \mathbb{R}_{>0}$  such that  $\forall \sigma \in [s]^n$

$$f(\sigma) = \frac{\Pr_{\text{opt}}[X = \sigma \mid X \neq Y]}{\Pr_{\text{greedy}}[X = \sigma \mid X \neq Y]} = \frac{\max\{0, P(\sigma) - Q(\sigma)\}}{\Pr_{\text{greedy}}[X = \sigma \mid X \neq Y]}.$$

- Estimator:  $W = f(\sigma)$  where  $\sigma \sim \pi$ .

## Properties of the estimator

- **Correct expectation**

$$\mathbb{E}_{\sigma \sim \pi}[f(\sigma)] = \frac{d_{TV}(P, Q)}{\Pr_{\text{greedy}}[X \neq Y]} = R \geq \frac{1}{n}.$$

- **Low variance**

$$\text{Var}_{\sigma \sim \pi}[f(\sigma)] \leq 1.$$

- **Efficient computation**

- sample  $\sigma \sim \pi$  in time  $O(n)$ ;
- given any  $\sigma \in \{0, 1\}^n$ , compute  $f(\sigma)$  in time  $O(n)$ .

## Algorithm

- draw  $\sigma_1, \sigma_2, \dots, \sigma_m \sim \pi$  for  $m = (n/\epsilon^2)$ ;
- return the average  $\hat{R} = \frac{1}{m} \sum_{i=1}^m f(\sigma_i)$ .

**Correctness:** Chebyshev ineq.  $\oplus \text{Var}[f] \leq 1 \oplus \mathbb{E}[f] \geq \frac{1}{n}$

**Efficiency:** efficient computation property.