

# Sliced Mutual Information: A Scalable Measure of Statistical Dependence

Ziv Goldfeld

Cornell University

Joint work with Kristjan Greenewald, Theshani Nuradha, and Galen Reeves

2022 ITA Workshop

May 24th, 2022

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

Fundamental measure of statistical dependence:

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties
- Emerges as solution to operational problems

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties
- Emerges as solution to operational problems
- Applications in information theory, statistics, machine learning.

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties
- Emerges as solution to operational problems
- Applications in information theory, **statistics, machine learning**.

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties
- Emerges as solution to operational problems
- Applications in information theory, **statistics, machine learning**.

**Challenge:** MI estimation in high dim. sample complexity  $n^*(\epsilon, d) \asymp \epsilon^{-d}$

# Mutual Information: Virtues and Challenges

## Definition (Shannon'48)

The mutual information (MI) between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$I(X;Y) := \int_{\mathbb{R}^{d_x}} \int_{\mathbb{R}^{d_y}} \log \left( \frac{dP_{XY}}{dP_X \otimes P_Y} \right) dP_{X,Y} = D_{KL}(P_{XY} \| P_X \otimes P_Y)$$

## Fundamental measure of statistical dependence:

- Meaningful units & structural properties
- Emerges as solution to operational problems
- Applications in information theory, **statistics, machine learning**.

**Challenge:** MI estimation in high dim. sample complexity  $n^*(\epsilon, d) \asymp \epsilon^{-d}$

⊗ **Goal:** Scalable MI surrogate that preserves its structure

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

# Sliced Mutual Information (SMI)

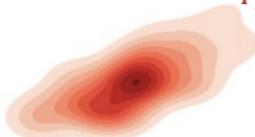
## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

## Illustration:

$$P_X \in \mathcal{P}(\mathbb{R}^{d_x})$$



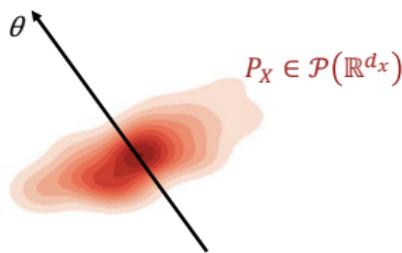
# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

### Illustration:



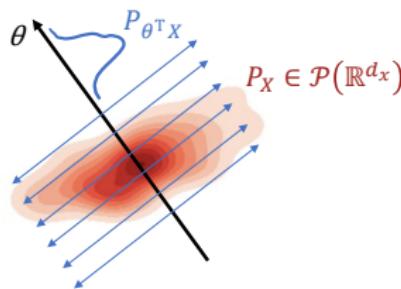
# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

## Illustration:



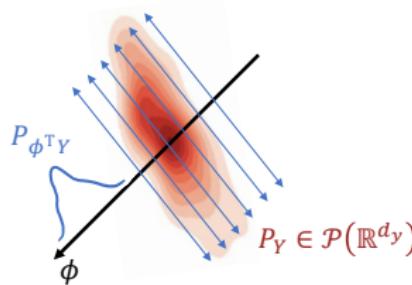
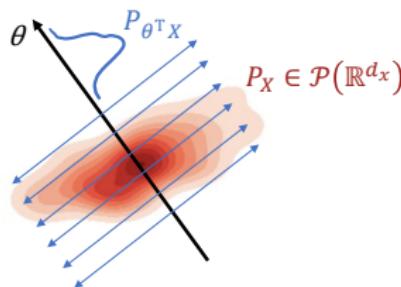
# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

### Illustration:



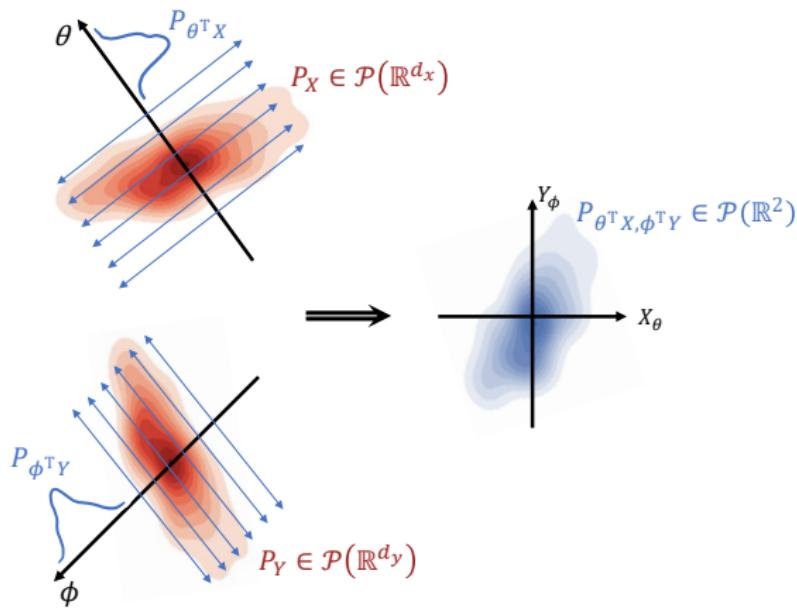
# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

### Illustration:



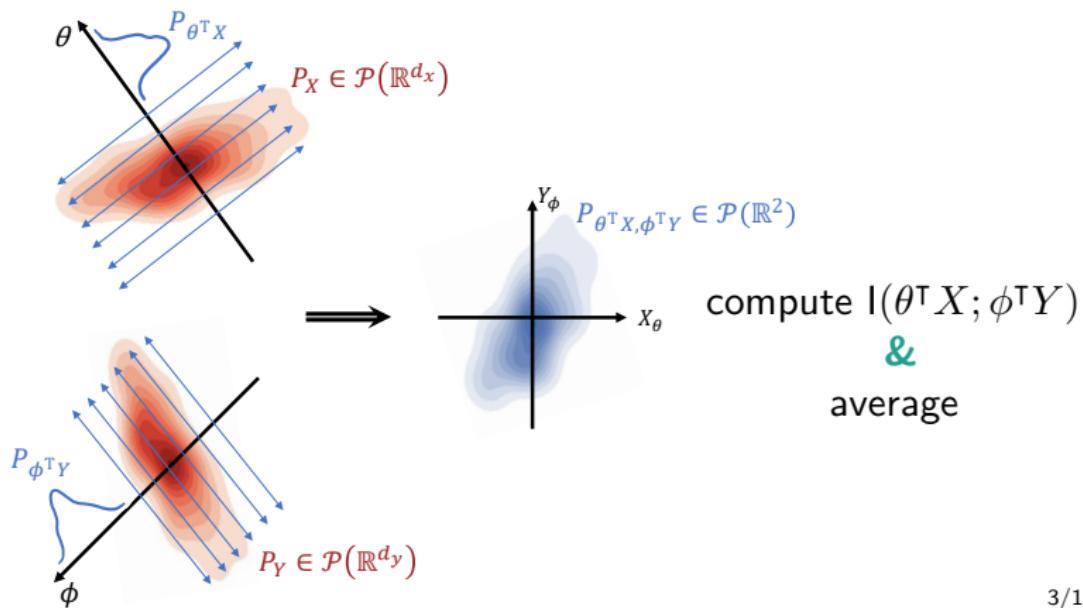
# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI btw.  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is (here  $\sigma_d = \text{Unif}(\mathbb{S}^{d-1})$ )

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi).$$

### Illustration:



# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

*SMI preserves many of the properties of MI, including:*

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

*SMI preserves many of the properties of MI, including:*

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

SMI preserves many of the properties of MI, including:

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$

$$\theta^\top X \perp \phi^\top Y, \quad \forall (\theta, \phi) \in \mathbb{S}^{d_x-1} \times \mathbb{S}^{d_y-1} \iff X \perp Y$$

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

*SMI preserves many of the properties of MI, including:*

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$
- ② **Sliced entropy:**  $\text{SI}(X; Y) = \text{sh}(X) - \text{sh}(X|Y) = \dots$   
where  $\text{sh}(X) := h(\Theta^\top X | \Theta)$

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

SMI preserves many of the properties of MI, including:

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$
- ② **Sliced entropy:**  $\text{SI}(X; Y) = \text{sh}(X) - \text{sh}(X|Y) = \dots$   
where  $\text{sh}(X) := h(\Theta^\top X | \Theta)$
- ③ **Max sliced entropy:** ▶  $\Sigma_P \preccurlyeq \Sigma \implies P^* = \mathcal{N}(0, \Sigma)$   
▶  $\text{spt}(P) \subseteq \mathbb{B}_d(r) \implies P^* = \text{Unif}(r\mathbb{S}^{d-1})$

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

SMI preserves many of the properties of MI, including:

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$
- ② **Sliced entropy:**  $\text{SI}(X; Y) = \text{sh}(X) - \text{sh}(X|Y) = \dots$   
where  $\text{sh}(X) := h(\Theta^\top X | \Theta)$
- ③ **Max sliced entropy:** ▶  $\Sigma_P \preccurlyeq \Sigma \implies P^* = \mathcal{N}(0, \Sigma)$   
▶  $\text{spt}(P) \subseteq \mathbb{B}_d(r) \implies P^* = \text{Unif}(r\mathbb{S}^{d-1})$
- ④ **Tensorization:**  $(X_1, Y_1), \dots, (X_n, Y_n)$  are mutually independent  
 $\implies \text{SI}(X_1, \dots, X_n; Y_1, \dots, Y_n) = \sum_{i=1}^n \text{SI}(X_i; Y_i)$

# Properties of Sliced Mutual Information

## Theorem (ZG-Greenewald'21)

SMI preserves many of the properties of MI, including:

- ① **Independence:**  $\text{SI}(X; Y) = 0 \iff X \perp Y$
- ② **Sliced entropy:**  $\text{SI}(X; Y) = \text{sh}(X) - \text{sh}(X|Y) = \dots$   
where  $\text{sh}(X) := h(\Theta^\top X | \Theta)$
- ③ **Max sliced entropy:** ▶  $\Sigma_P \preccurlyeq \Sigma \implies P^* = \mathcal{N}(0, \Sigma)$   
▶  $\text{spt}(P) \subseteq \mathbb{B}_d(r) \implies P^* = \text{Unif}(r\mathbb{S}^{d-1})$
- ④ **Tensorization:**  $(X_1, Y_1), \dots, (X_n, Y_n)$  are mutually independent  
 $\implies \text{SI}(X_1, \dots, X_n; Y_1, \dots, Y_n) = \sum_{i=1}^n \text{SI}(X_i; Y_i)$
- ⑤ **DV:**  $\text{SI}(X; Y) = \sup_f \mathbb{E}[f(\Theta, \Phi, \Theta^\top X, \Phi^\top Y)] - \log \left( \mathbb{E} \left[ e^{f(\Theta, \Phi, \Theta^\top X, \Phi^\top Y)} \right] \right)$

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI    $I(X; Y) \geq I(f(X); Y)$

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $SI(X; Y) < SI(f(X); Y)$

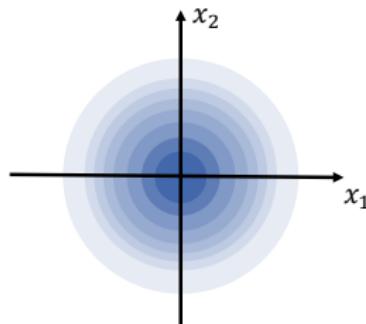
# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, \mathbf{I}_2)$



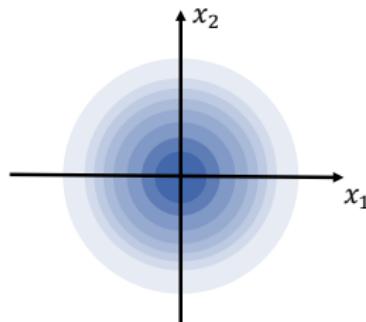
# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- **Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, \mathbf{I}_2)$ ,  $Y = X_1$



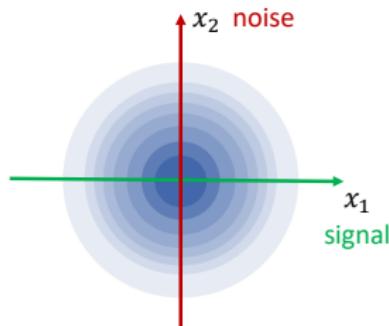
# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = X_1$



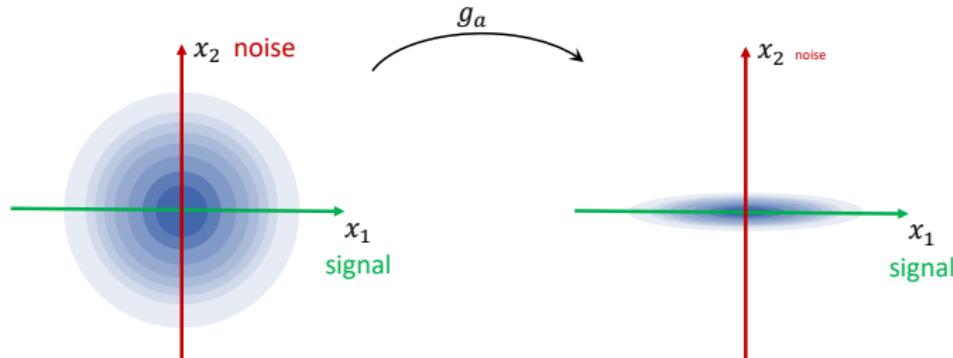
# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = X_1$ ,  $g_a(x_1, x_2) = (x_1 \ ax_2)^\top$



# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

$\implies$  MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $SI(X; Y) < SI(f(X); Y)$
- Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = X_1$ ,  $g_a(x_1, x_2) = (x_1 \ ax_2)^\top$   
 $\implies SI(X; Y) < SI(g_a(X); Y), \quad \forall a < 1$

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- **Example:**  $X = (\textcolor{green}{X}_1 \textcolor{red}{X}_2)^\top \sim \mathcal{N}(0, \text{I}_2)$ ,  $\textcolor{green}{Y} = \textcolor{green}{X}_1$ ,  $g_a(x_1, x_2) = (x_1 \ ax_2)^\top$   
⇒  $\text{SI}(X; Y) < \text{SI}(g_a(X); Y), \quad \forall a < 1$

⇒ SMI **can increase** via processing (violates DPI)

# Sliced Mutual Information & Processing

Mutual information: Satisfies DPI  $I(X; Y) \geq I(f(X); Y)$

⇒ MI **cannot grow** via processing

Sliced mutual information: Follows the DPI?

- Only considers projections of RVs
- If  $f(X)$  has more informative projections:  $\text{SI}(X; Y) < \text{SI}(f(X); Y)$
- **Example:**  $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = X_1$ ,  $g_a(x_1, x_2) = (x_1 \ ax_2)^\top$   
⇒  $\text{SI}(X; Y) < \text{SI}(g_a(X); Y), \quad \forall a < 1$

⇒ SMI **can increase** via processing (violates DPI)

Can be used for feature extraction via SMI maximization

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

- ① Take scalar MI estimator  $\hat{I}(A^n, B^n)$  w/ error  $\delta(n)$  over some class

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

- ① Take scalar MI estimator  $\hat{I}(A^n, B^n)$  w/ error  $\delta(n)$  over some class
- ② Sample  $m$  random directions  $\Theta^m$  and  $\Phi^m$  from corresponding spheres

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

- ① Take scalar MI estimator  $\hat{I}(A^n, B^n)$  w/ error  $\delta(n)$  over some class
- ② Sample  $m$  random directions  $\Theta^m$  and  $\Phi^m$  from corresponding spheres
- ③ Compute  $(\Theta_i^\top X)^n := (\Theta_i^\top X_1, \dots, \Theta_i^\top X_n)$ ,  $i \in [m]$ , and  $(\Phi_i^\top Y)^n$ .

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

- ① Take scalar MI estimator  $\hat{I}(A^n, B^n)$  w/ error  $\delta(n)$  over some class
- ② Sample  $m$  random directions  $\Theta^m$  and  $\Phi^m$  from corresponding spheres
- ③ Compute  $(\Theta_i^\top X)^n := (\Theta_i^\top X_1, \dots, \Theta_i^\top X_n)$ ,  $i \in [m]$ , and  $(\Phi_i^\top Y)^n$ .
- ④ **Estimate SMI via MC:** 
$$\widehat{\text{SI}}_{m,n} := \frac{1}{m} \sum_{i=1}^m \hat{I}((\Theta_i^\top X)^n, (\Phi_i^\top Y)^n)$$

# Scalable Estimation from Samples

Estimator: Given samples  $(X^n, Y^n)$  i.i.d. from  $P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$

- ① Take scalar MI estimator  $\hat{I}(A^n, B^n)$  w/ error  $\delta(n)$  over some class
- ② Sample  $m$  random directions  $\Theta^m$  and  $\Phi^m$  from corresponding spheres
- ③ Compute  $(\Theta_i^\top X)^n := (\Theta_i^\top X_1, \dots, \Theta_i^\top X_n)$ ,  $i \in [m]$ , and  $(\Phi_i^\top Y)^n$ .
- ④ Estimate SMI via MC: 
$$\widehat{\text{SI}}_{m,n} := \frac{1}{m} \sum_{i=1}^m \hat{I}((\Theta_i^\top X)^n, (\Phi_i^\top Y)^n)$$

## Theorem (ZG-Greenewald-Reeves'22)

If  $P_{XY}$  has finite 2nd moments & Fisher information  $J(P_{XY}) < \infty$ , then

$$\mathbb{E} \left[ \left| \text{SI}(X; Y) - \widehat{\text{SI}}_{m,n} \right| \right] \leq C(P_{XY}) \sqrt{\frac{d_x + d_y}{d_x d_y} m^{-\frac{1}{2}}} + \delta(n),$$

where  $C(P_{XY}) = 21 \sqrt{\|J_F(P_{XY})\|_{\text{op}} (\|\Sigma_X\|_{\text{op}} \vee \|\Sigma_Y\|_{\text{op}})}$ .

## Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

## Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

Compare to [Polyanskiy-Wu'16]

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

Compare to [Polyanskiy-Wu'16]  
• Relax  $(c_1, c_2)$ -regularity

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

Compare to [Polyanskiy-Wu'16]

- Relax  $(c_1, c_2)$ -regularity
- Sharp constant (optimal)

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

- ①  $\mathsf{I}_{XY}$  is Lipschitz on  $\mathbb{S}^{d_x-1} \times \mathbb{S}^{d_y-1}$

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|\right] \leq \underbrace{\mathbb{E}\left[|\widehat{\mathsf{SI}}_{m,n} - \circledast|\right]}_{(I)} + \underbrace{\mathbb{E}\left[|\circledast - \mathsf{SI}(X; Y)|\right]}_{(II)}$

Bound: For each term

$$(I) \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[ |\widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i)| \right] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$(II) \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

- ①  $\mathsf{I}_{XY}$  is Lipschitz on  $\mathbb{S}^{d_x-1} \times \mathbb{S}^{d_y-1}$
- ② Concentration of Lip. functions on  $\mathbb{S}^{d-1}$  & Efron-Stein-Steele ineq.

# Estimation Error Bound: Proof Outline

Define:  $\mathsf{I}_{XY}(\theta, \phi) := \mathsf{I}(\theta^\top X, \phi^\top Y)$  &  $\circledast := \frac{1}{m} \sum_{i=1}^m \mathsf{I}_{XY}(\Theta_i, \Phi_i)$

Decompose:  $\mathbb{E}[|\widehat{\mathsf{SI}}_{m,n} - \mathsf{SI}(X; Y)|] \leq \underbrace{\mathbb{E}[|\widehat{\mathsf{SI}}_{m,n} - \circledast|]}_{\text{(I)}} + \underbrace{\mathbb{E}[|\circledast - \mathsf{SI}(X; Y)|]}_{\text{(II)}}$

Bound: For each term

$$\text{(I)} \leq \frac{1}{m} \sum_{i=1}^m \mathbb{E} [\left| \widehat{\mathsf{I}}_{XY}(\Theta_i, \Phi_i) - \mathsf{I}_{XY}(\Theta_i, \Phi_i) \right|] \leq \delta(n) \quad (\text{assumption on } \widehat{\mathsf{I}})$$

$$\text{(II)} \leq m^{-1/2} \sqrt{\text{Var}(\mathsf{I}_{XY}(\Theta, \Phi))} \quad (\mathbb{E}[\circledast] = \mathsf{SI}(X; Y) \text{ & } L^1(P) \leq L^2(P))$$

Variance: New continuity result  $\mathsf{h}(P) - \mathsf{h}(Q) \leq \sqrt{\mathsf{J}(Q)} \mathsf{W}_2(P, Q)$  (via HWI)

- ①  $\mathsf{I}_{XY}$  is Lipschitz on  $\mathbb{S}^{d_x-1} \times \mathbb{S}^{d_y-1}$
- ② Concentration of Lip. functions on  $\mathbb{S}^{d-1}$  & Efron-Stein-Steele ineq.

$$\implies \text{Var}(\mathsf{I}_{XY}(\Theta, \Phi)) \leq C(P_{XY})^2 (d_x^{-1} + d_y^{-1})$$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

Neural est.: Minimax optimal  $D_{KL}$  est. over Barron class [Sreekumar-ZG'22]

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

Neural est.: Minimax optimal  $D_{KL}$  est. over Barron class [Sreekumar-ZG'22]

- $I(X; Y) = \sup_f \mathbb{E}[f(X, Y)] - \log(\mathbb{E}[e^{f(\tilde{X}, \tilde{Y})}])$

# SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

Neural est.: Minimax optimal  $D_{KL}$  est. over Barron class [Sreekumar-ZG'22]

- $I(X; Y) = \sup_f \mathbb{E}[f(X, Y)] - \log(\mathbb{E}[e^{f(\tilde{X}, \tilde{Y})}])$   
 $\approx \sup_{g \in \mathcal{F}_{nn}^k} \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i) - \log\left(\frac{1}{n} \sum_{i=1}^n e^{g(X_{\sigma(i)}, Y_i)}\right)$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

Neural est.: Minimax optimal  $D_{KL}$  est. over Barron class [Sreekumar-ZG'22]

- $I(X; Y) = \sup_f \mathbb{E}[f(X, Y)] - \log(\mathbb{E}[e^{f(\tilde{X}, \tilde{Y})}])$   
 $\approx \sup_{g \in \mathcal{F}_{nn}^k} \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i) - \log\left(\frac{1}{n} \sum_{i=1}^n e^{g(X_{\sigma(i)}, Y_i)}\right)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{k,m,n}^{(\text{NE})}\right|\right] \lesssim m^{-1/2} + k^{-1/2} + n^{-1/2}$$

## SMI Estimation: Examples

Upshot: The SMI estimator attains error  $\lesssim m^{-1/2} + \delta(n)$

KDE & BPA: Minimax optimal  $h(X)$  est. over Lipschitz balls [Han et al.'20]

- Plug-in KDE of  $f_X$  into best polynomial approximation of  $h(X)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{m,n}^{(\text{Lip})}\right|\right] \lesssim m^{-1/2} + n^{-1/2}(1 + (\log n)^{1/4})$$

Neural est.: Minimax optimal  $D_{KL}$  est. over Barron class [Sreekumar-ZG'22]

- $I(X; Y) = \sup_f \mathbb{E}[f(X, Y)] - \log(\mathbb{E}[e^{f(\tilde{X}, \tilde{Y})}])$   
 $\approx \sup_{g \in \mathcal{F}_{nn}^k} \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i) - \log\left(\frac{1}{n} \sum_{i=1}^n e^{g(X_{\sigma(i)}, Y_i)}\right)$

$$\implies \mathbb{E}\left[\left|\text{SI}(X; Y) - \widehat{\text{SI}}_{k,m,n}^{(\text{NE})}\right|\right] \lesssim m^{-1/2} + k^{-1/2} + n^{-1/2}$$

- **No curse of dimensionality:** Compare to classic MI rate  $n^{-1/(d_x+d_y)}$

## Experiments: Independence Testing

Recall:  $\text{SI}(X; Y) = 0 \iff (X, Y) \text{ independent}$

## Experiments: Independence Testing

Recall:  $\text{SI}(X; Y) = 0 \iff (X, Y)$  independent

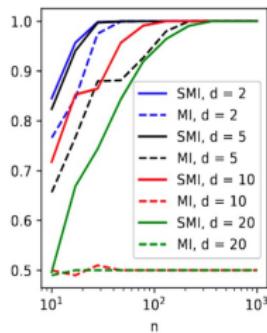
$\implies$  Compute SMI & threshold for independence testing

# Experiments: Independence Testing

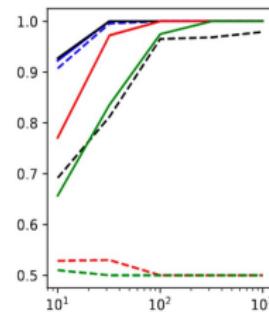
Recall:  $\text{SI}(X; Y) = 0 \iff (X, Y) \text{ independent}$

$\implies$  Compute SMI & threshold for independence testing

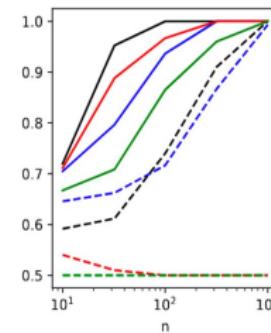
**Figure: Area under the ROC curve**



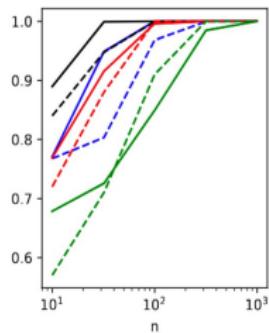
(a)  $Y$  encodes a single feature.



(b)  $Y$  encodes two features.



(c) Low rank common signal.



(d) Independent coordinates.

$$Y = \frac{1}{\sqrt{d}}(\mathbf{1}^\top X)\mathbf{1} + Z$$

$$Y_i = \begin{cases} \frac{1}{d}(\mathbf{1}_{[d/2]} 0 \dots 0)^\top X + Z_i, & i \leq \frac{d}{2} \\ \frac{1}{d}(0 \dots 0 \mathbf{1}_{[d/2]})^\top X + Z_i, & i > \frac{d}{2} \end{cases}$$

$$\begin{aligned} X &= P_1 V + Z_1 \\ Y &= P_2 V + Z_2 \end{aligned}$$

$$Y = X + Z$$

## Experiments: Feature Extraction

Goal: Maximize  $\text{SI}(AX; AY)$ ,  $X, Y$  i.i.d. from same MNIST class (0 or 1)

## Experiments: Feature Extraction

Goal: Maximize  $\text{SI}(AX; AY)$ ,  $X, Y$  i.i.d. from same MNIST class (0 or 1)

- **Upper bound:**  $\text{SI}(AX; AY) \leq I(X; Y) = 1$  bit

## Experiments: Feature Extraction

Goal: Maximize  $\text{SI}(AX; AY)$ ,  $X, Y$  i.i.d. from same MNIST class (0 or 1)

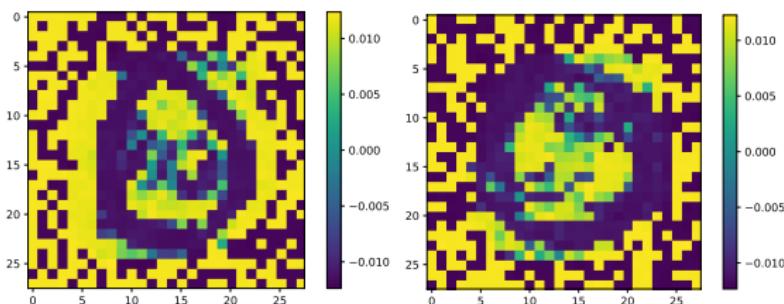
- **Upper bound:**  $\text{SI}(AX; AY) \leq I(X; Y) = 1$  bit
- **Interpretation:**  $A$  is linear features most useful for label classification

# Experiments: Feature Extraction

Goal: Maximize  $\text{SI}(AX; AY)$ ,  $X, Y$  i.i.d. from same MNIST class (0 or 1)

- **Upper bound:**  $\text{SI}(AX; AY) \leq I(X; Y) = 1$  bit
- **Interpretation:**  $A$  is linear features most useful for label classification

Figure: Feature extraction for MNIST



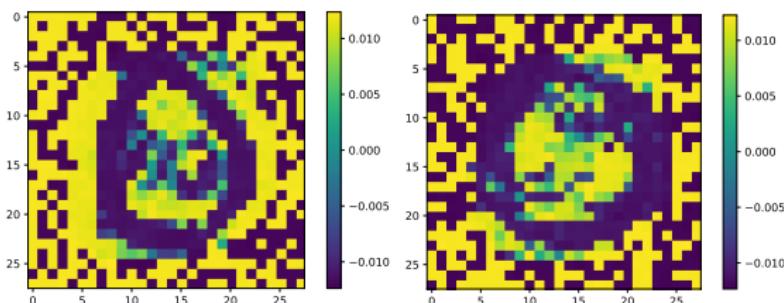
Rows 0 and 1 of optimized  $A$  (rearranged as MNIST image)

# Experiments: Feature Extraction

Goal: Maximize  $\text{SI}(AX; AY)$ ,  $X, Y$  i.i.d. from same MNIST class (0 or 1)

- **Upper bound:**  $\text{SI}(AX; AY) \leq I(X; Y) = 1$  bit
- **Interpretation:**  $A$  is linear features most useful for label classification

Figure: Feature extraction for MNIST



Rows 0 and 1 of optimized  $A$  (rearranged as MNIST image)

$$\implies \text{SI}(A^*X; A^*Y) = 0.68 \text{ (compare to 0.0752 for random } A\text{)}$$

## Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

## Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI

## Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)

## Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate

## Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

# Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

Future directions: Theoretical and applied

# Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

Future directions: Theoretical and applied

- Extensions to  $k$ -dimensional projections

# Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

Future directions: Theoretical and applied

- Extensions to  $k$ -dimensional projections
- Decomposition to Gaussian SMI plus residual (negligible?)

# Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

Future directions: Theoretical and applied

- Extensions to  $k$ -dimensional projections
- Decomposition to Gaussian SMI plus residual (negligible?)
- Applications to more complex learning tasks (InfoGAN, InfoMAX, ...)

# Summary

Sliced mutual information: Avg. scalar MI terms btw. 1D projections

- **Structure:** Preserves many properties of classic MI
- **Processing:** Violates DPI and can increase from processing (useful?)
- **Estimation:** Efficiently computable & fast to estimate
- **Applications:** Independence testing & feature extraction

Future directions: Theoretical and applied

- Extensions to  $k$ -dimensional projections
- Decomposition to Gaussian SMI plus residual (negligible?)
- Applications to more complex learning tasks (InfoGAN, InfoMAX, ...)

Thank you!

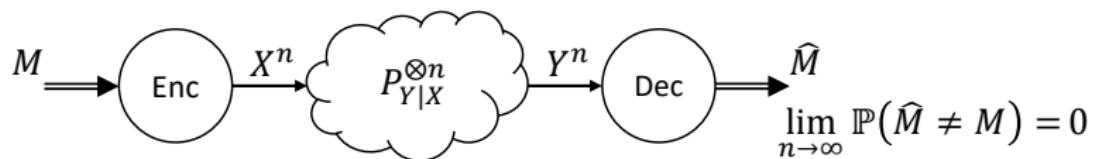
# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

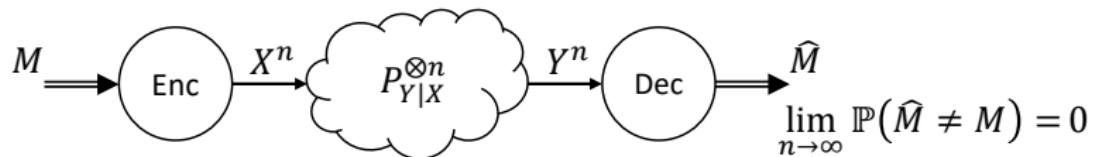
- **Noisy communication:** Channel capacity  $C(P_{Y|X}) = \max_{P_X} I(X; Y)$



# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

- **Noisy communication:** Channel capacity  $C(P_{Y|X}) = \max_{P_X} I(X; Y)$

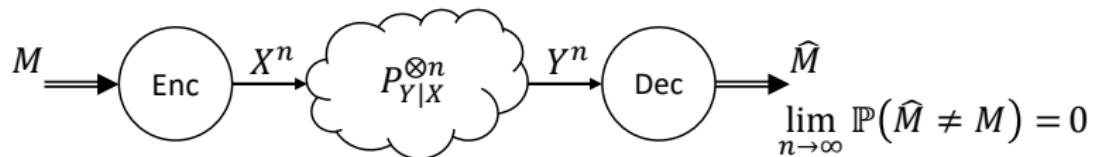


- **Compression:** Rate-distortion  $R(D, P_X) = \min_{P_{\hat{X}|X}: \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X})$

# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

- **Noisy communication:** Channel capacity  $C(P_{Y|X}) = \max_{P_X} I(X; Y)$

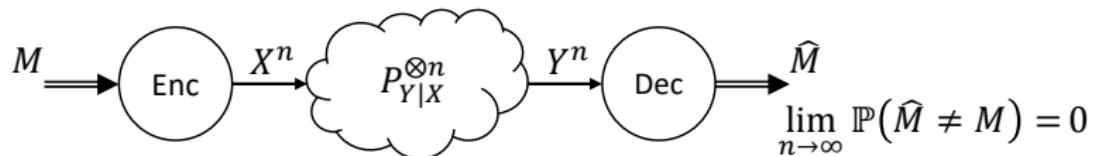


- **Compression:** Rate-distortion  $R(D, P_X) = \min_{P_{\hat{X}|X}: \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X})$
- **More:** Distribution simulation, privacy & security, common info...

# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

- **Noisy communication:** Channel capacity  $C(P_{Y|X}) = \max_{P_X} I(X; Y)$



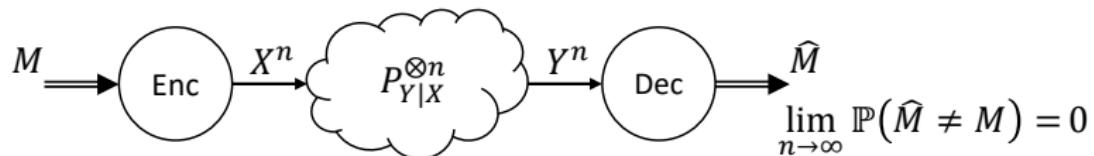
- **Compression:** Rate-distortion  $R(D, P_X) = \min_{P_{\hat{X}|X}: \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X})$
- **More:** Distribution simulation, privacy & security, common info...

Statistics: Ind. testing, impossibility results, dependence quantification...

# Mutual Information: Applications

Information Theory: Emerges as solution to operation problems

- **Noisy communication:** Channel capacity  $C(P_{Y|X}) = \max_{P_X} I(X; Y)$



- **Compression:** Rate-distortion  $R(D, P_X) = \min_{P_{\hat{X}|X}: \mathbb{E}[d(X, \hat{X})] \leq D} I(X; \hat{X})$
- **More:** Distribution simulation, privacy & security, common info...

Statistics: Ind. testing, impossibility results, dependence quantification...

Machine Learning: Hosts of modern applications

# Mutual Information in Machine Learning

Representation learning: InfoMax principle

# Mutual Information in Machine Learning

Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m, m \ll d$

# Mutual Information in Machine Learning

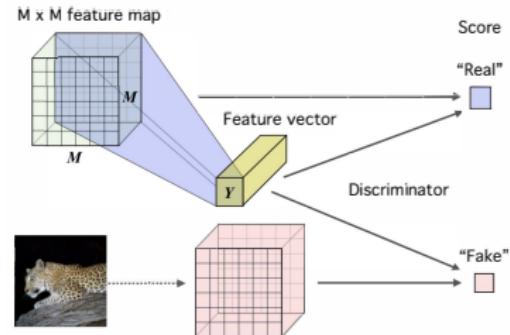
## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m, m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m, m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$

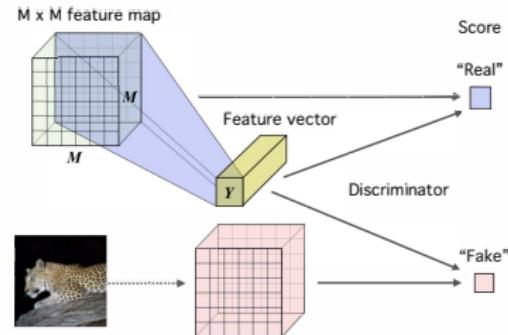


- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$

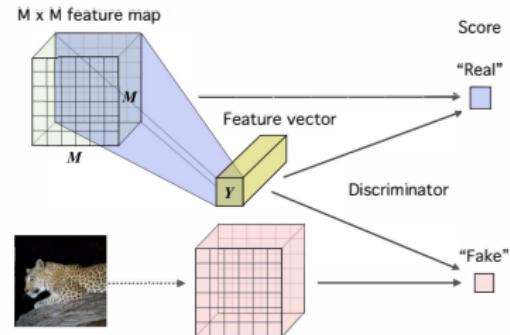


- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound
- ➋ **CPC [Oord et al.'19]:** Contrastive  $\mathcal{L}_{\text{InfoNCE}}$  lower bound

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$



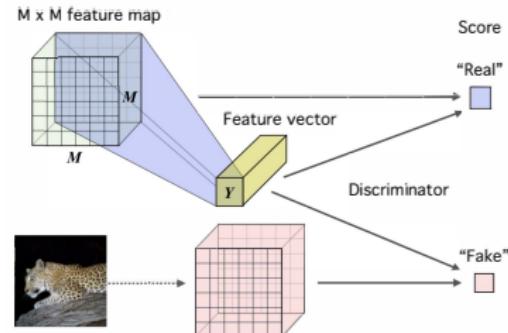
- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound
- ➋ **CPC [Oord et al.'19]:** Contrastive  $\mathcal{L}_{\text{InfoNCE}}$  lower bound

## Generative modeling : Disentangled latent space

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$



- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound
- ➋ **CPC [Oord et al.'19]:** Contrastive  $\mathcal{L}_{\text{InfoNCE}}$  lower bound

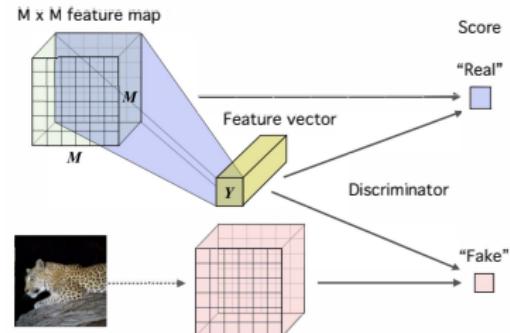
## Generative modeling : Disentangled latent space

- **Goal:** Learn  $g_\theta(Z) \sim Q_\theta$  that mimics  $X \sim P_X$

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$



- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound
- ➋ **CPC [Oord et al.'19]:** Contrastive  $\mathcal{L}_{\text{InfoNCE}}$  lower bound

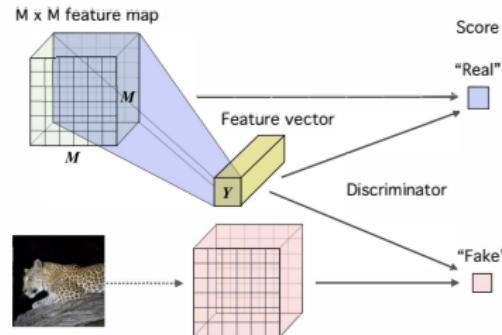
## Generative modeling : Disentangled latent space

- **Goal:** Learn  $g_\theta(Z) \sim Q_\theta$  that mimics  $X \sim P_X$
- **GAN:**  $\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{L}_{\text{GAN}}(g_\theta, d_\phi)$

# Mutual Information in Machine Learning

## Representation learning: InfoMax principle

- **Data:**  $X \sim P_X \in \mathcal{P}(\mathbb{R}^d)$
- **Encoder:**  $e_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^m$ ,  $m \ll d$
- **Objective:**  $\sup_{\theta \in \Theta} I(e_\theta(X); X)$

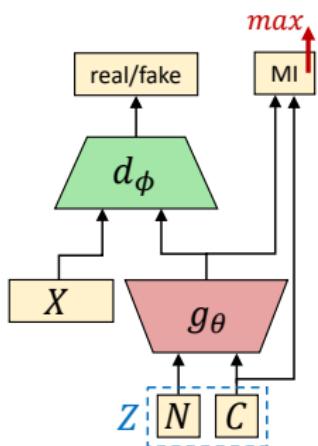


- ➊ **DIM [Hjelm et al.'18]:** Parameterized DV lower bound
- ➋ **CPC [Oord et al.'19]:** Contrastive  $\mathcal{L}_{\text{InfoNCE}}$  lower bound

## Generative modeling : Disentangled latent space

- **Goal:** Learn  $g_\theta(Z) \sim Q_\theta$  that mimics  $X \sim P_X$
- **GAN:**  $\min_{\theta \in \Theta} \max_{\phi \in \Phi} \mathcal{L}_{\text{GAN}}(g_\theta, d_\phi)$
- **InfoGAN [Chen et al.'16], [Belghazi et al.'18]:**

- ▶ Model  $Z = (N \ C)$  & regularize  $\mathcal{L}_{\text{GAN}}(g_\theta, d_\phi)$  by  $\beta I(g_\theta(N, C); C)$



# Sliced InfoGAN: MNIST Results

Regular InfoGAN (MI)



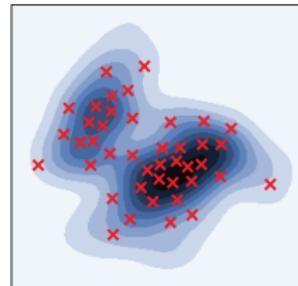
Sliced InfoGAN (SMI)



Codes:  $C_1 \in [0 : 9]$  (digits),  $C_2 \in [-2, 2]$  (rotation),  $C_3 \in [-2, 2]$  (width)

# Mutual Information Estimation

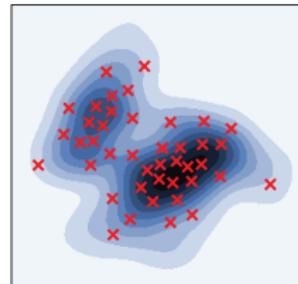
In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$



# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

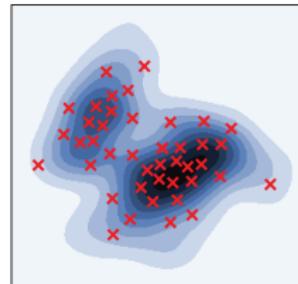


# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?

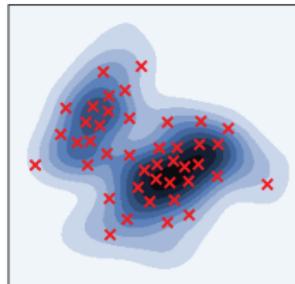


# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?



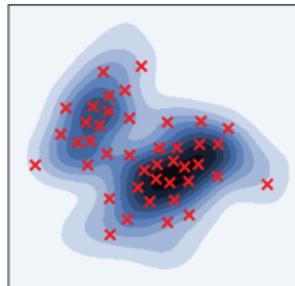
**Main challenge:** Estimation in high dim. is infeasible

# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?



**Main challenge:** Estimation in high dim. is infeasible

- **Sample complexity:**  $n^*(\epsilon, d) \asymp \epsilon^{-d}$  (under regularity)

- ▶ Hölder smooth [Jiao-Gao-Han'18]
- ▶ (Gen.) Lipschitz smooth [Han-Jiao-Weissman-Wu'20]

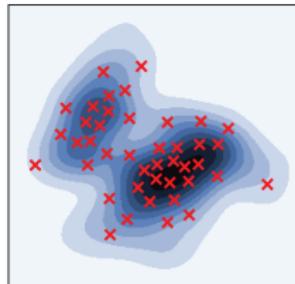
⋮

# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?



**Main challenge:** Estimation in high dim. is infeasible

- **Sample complexity:**  $n^*(\epsilon, d) \asymp \epsilon^{-d}$  (under regularity)

- ▶ Hölder smooth [Jiao-Gao-Han'18]
- ▶ (Gen.) Lipschitz smooth [Han-Jiao-Weissman-Wu'20]

⋮

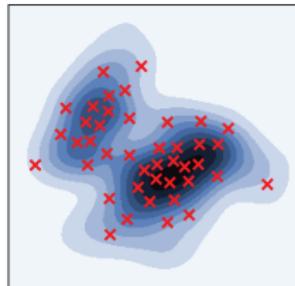


# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?



**Main challenge:** Estimation in high dim. is infeasible

- **Sample complexity:**  $n^*(\epsilon, d) \asymp \epsilon^{-d}$  (under regularity)

- ▶ Hölder smooth [Jiao-Gao-Han'18]
- ▶ (Gen.) Lipschitz smooth [Han-Jiao-Weissman-Wu'20]

⋮

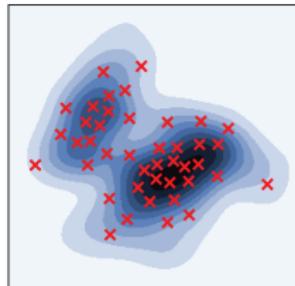
- **Formal limitations:**  $n^*(\epsilon, d) \gtrsim \exp(I(X; Y))$  [McAllester-Stratos'20]

# Mutual Information Estimation

In practice: Don't have  $P_{XY}$  but samples  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} P_{XY}$

- **Estimation:**  $\hat{I}(X^n, Y^n)$  via  $k$ -NN, KDE, etc.

⇒ Can we approximate  $I(X; Y) \approx \hat{I}(X^n, Y^n)$ ?



**Main challenge:** Estimation in high dim. is infeasible

- **Sample complexity:**  $n^*(\epsilon, d) \asymp \epsilon^{-d}$  (under regularity)

- ▶ Hölder smooth [Jiao-Gao-Han'18]
- ▶ (Gen.) Lipschitz smooth [Han-Jiao-Weissman-Wu'20]

⋮

- **Formal limitations:**  $n^*(\epsilon, d) \gtrsim \exp(I(X; Y))$  [McAllester-Stratos'20]



**Goal:** Scalable MI surrogate that preserves its structure

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

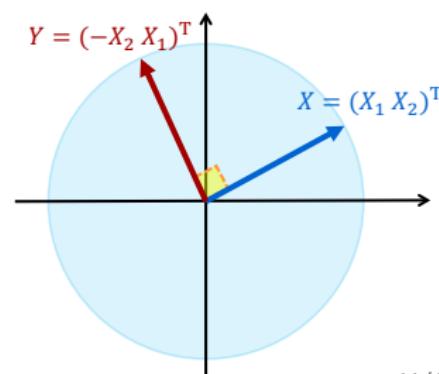
$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$
- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = (-X_2 \ X_1)^\top$



# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

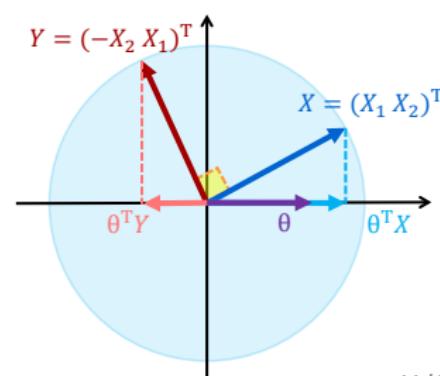
Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$

- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, \mathbf{I}_2)$ ,  $Y = (-X_2 \ X_1)^\top$



# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

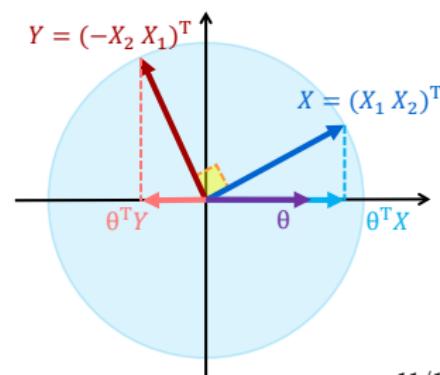
- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$

- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = (-X_2 \ X_1)^\top$

►  $\text{cov}(\theta^\top X, \theta^\top Y) = 0 \implies \tilde{\text{SI}}(X; Y) = 0$



# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$

# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

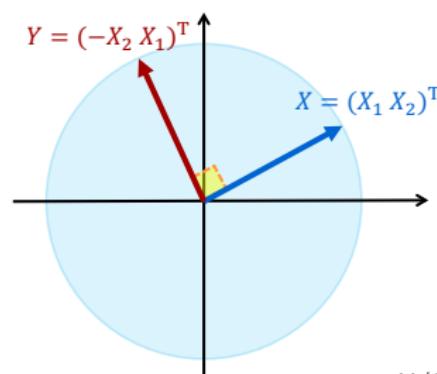
$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$
- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = (-X_2 \ X_1)^\top$



# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

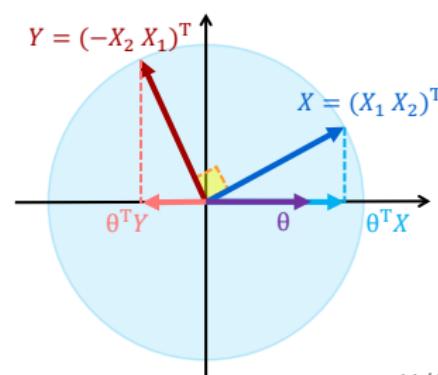
$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$
- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = (-X_2 \ X_1)^\top$



# Sliced Mutual Information (SMI)

## Definition (ZG-Greenewald'21)

The SMI between  $(X, Y) \sim P_{XY} \in \mathcal{P}(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$  is

$$\text{SI}(X; Y) := \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_x(\theta) d\sigma_y(\phi),$$

Sliced divergences:  $\delta(P, Q)$  btw  $P, Q \in \mathcal{P}(\mathbb{R}^d)$  (Wasserstein,  $f$ -div., IPM)

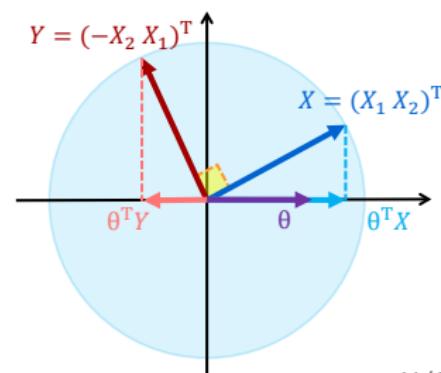
- **Slicing:**  $\bar{\delta}(P, Q) := \int_{\mathbb{S}^{d-1}} \delta(\mathfrak{p}_\sharp^\theta P, \mathfrak{p}_\sharp^\theta Q) d\sigma(\theta)$ , where  $\mathfrak{p}^\theta(x) = \theta^\top x$   
⇒ Scalable & preserves properties

Difference: 2 proj. for  $\text{SI}(X; Y)$  vs. 1 for  $\bar{\delta}(P, Q)$

- 1-proj. SMI can nullify btw dependent  $(X, Y)$

- $X = (X_1 \ X_2)^\top \sim \mathcal{N}(0, I_2)$ ,  $Y = (-X_2 \ X_1)^\top$

►  $\text{cov}(\theta^\top X, \theta^\top Y) = 0 \implies \tilde{\text{SI}}(X; Y) = 0$



# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t,s) = \varphi_X(t)\varphi_Y(s), \quad \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t,s) = \varphi_X(t)\varphi_Y(s), \quad \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

Consider:  $\text{SI}(X;Y) = \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} \mathsf{I}(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi) = 0$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t,s) = \varphi_X(t)\varphi_Y(s), \quad \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

Consider:  $\text{SI}(X;Y) = \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi) = 0$

$$\iff$$

$$I(\theta^\top X; \phi^\top Y) = 0, \quad \forall \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t,s) = \varphi_X(t)\varphi_Y(s), \quad \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

Consider:  $\text{SI}(X;Y) = \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi) = 0$

$$\iff$$

$$I(\theta^\top X; \phi^\top Y) = 0, \quad \forall \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{\theta^\top X, \phi^\top Y}(u, v) = \varphi_{\theta^\top X}(u)\varphi_{\phi^\top Y}(v), \quad \forall u, v \in \mathbb{R}, \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t,s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t,0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t,s) = \varphi_X(t)\varphi_Y(s), \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

Consider:  $\text{SI}(X;Y) = \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi) = 0$

$$\iff$$

$$I(\theta^\top X; \phi^\top Y) = 0, \quad \forall \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{\theta^\top X, \phi^\top Y}(u, v) = \varphi_{\theta^\top X}(u)\varphi_{\phi^\top Y}(v), \quad \forall u, v \in \mathbb{R}, \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{X,Y}(u\theta, v\phi) = \varphi_X(u\theta)\varphi_Y(v\phi), \quad \forall u, v \in \mathbb{R}, \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

# Independence Identification: Proof Outline

Char. functions:  $\varphi_{X,Y}(t, s) := \mathbb{E}[e^{it^\top X + is^\top Y}]$  &  $\varphi_X(t) := \varphi_{X,Y}(t, 0)$

Recall:  $X \perp Y \iff \varphi_{X,Y}(t, s) = \varphi_X(t)\varphi_Y(s), \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$

Consider:  $\text{SI}(X; Y) = \int_{\mathbb{S}^{d_x-1}} \int_{\mathbb{S}^{d_y-1}} I(\theta^\top X; \phi^\top Y) d\sigma_{d_x}(\theta) d\sigma_{d_y}(\phi) = 0$

$$\iff$$

$$I(\theta^\top X; \phi^\top Y) = 0, \quad \forall \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{\theta^\top X, \phi^\top Y}(u, v) = \varphi_{\theta^\top X}(u)\varphi_{\phi^\top Y}(v), \quad \forall u, v \in \mathbb{R}, \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{X,Y}(u\theta, v\phi) = \varphi_X(u\theta)\varphi_Y(v\phi), \quad \forall u, v \in \mathbb{R}, \theta \in \mathbb{S}^{d_x-1}, \phi \in \mathbb{S}^{d_y-1}$$

$$\iff$$

$$\varphi_{X,Y}(t, s) = \varphi_X(t)\varphi_Y(s), \quad \forall t \in \mathbb{R}^{d_x}, s \in \mathbb{R}^{d_y}$$

# Sliced Mutual Information & Processing (Cont.)

## Proposition (ZG-Greenwald'21)

*Extracting maximum-SMI linear feature:*

$$\sup_{A_x, A_y, b_x, b_y} \text{SI}(A_x X + b_x; A_y Y + b_y) = \sup_{\theta, \phi} I(\theta^\top X; \phi^\top Y).$$

Also, if  $(\theta_\star, \phi_\star) \in \operatorname{argmax} I(\theta^\top X; \phi^\top Y)$ , then  $[A_x^\star]_{1:} = \theta_\star^\top$ ,  $[A_y^\star]_{1:} = \phi_\star^\top$ .

# Sliced Mutual Information & Processing (Cont.)

## Proposition (ZG-Greenwald'21)

*Extracting maximum-SMI linear feature:*

$$\sup_{A_x, A_y, b_x, b_y} \text{SI}(A_x X + b_x; A_y Y + b_y) = \sup_{\theta, \phi} I(\theta^\top X; \phi^\top Y).$$

Also, if  $(\theta_*, \phi_*) \in \operatorname{argmax} I(\theta^\top X; \phi^\top Y)$ , then  $[A_x^*]_{1:} = \theta_*^\top$ ,  $[A_y^*]_{1:} = \phi_*^\top$ .

Extensions: Similar results for

# Sliced Mutual Information & Processing (Cont.)

## Proposition (ZG-Greenwald'21)

Extracting maximum-SMI linear feature:

$$\sup_{A_x, A_y, b_x, b_y} \text{SI}(A_x X + b_x; A_y Y + b_y) = \sup_{\theta, \phi} I(\theta^\top X; \phi^\top Y).$$

Also, if  $(\theta_*, \phi_*) \in \operatorname{argmax} I(\theta^\top X; \phi^\top Y)$ , then  $[A_x^*]_{1:} = \theta_*^\top$ ,  $[A_y^*]_{1:} = \phi_*^\top$ .

Extensions: Similar results for

- Rank-constrained matrices

# Sliced Mutual Information & Processing (Cont.)

## Proposition (ZG-Greenewald'21)

Extracting maximum-SMI linear feature:

$$\sup_{A_x, A_y, b_x, b_y} \text{SI}(A_x X + b_x; A_y Y + b_y) = \sup_{\theta, \phi} I(\theta^\top X; \phi^\top Y).$$

Also, if  $(\theta_*, \phi_*) \in \operatorname{argmax} I(\theta^\top X; \phi^\top Y)$ , then  $[A_x^*]_{1:} = \theta_*^\top$ ,  $[A_y^*]_{1:} = \phi_*^\top$ .

Extensions: Similar results for

- Rank-constrained matrices
- Shallow NNs

# Sliced Mutual Information & Processing (Cont.)

## Proposition (ZG-Greenewald'21)

Extracting maximum-SMI linear feature:

$$\sup_{A_x, A_y, b_x, b_y} \text{SI}(A_x X + b_x; A_y Y + b_y) = \sup_{\theta, \phi} I(\theta^\top X; \phi^\top Y).$$

Also, if  $(\theta_*, \phi_*) \in \operatorname{argmax} I(\theta^\top X; \phi^\top Y)$ , then  $[A_x^*]_{1:} = \theta_*^\top$ ,  $[A_y^*]_{1:} = \phi_*^\top$ .

Extensions: Similar results for

- Rank-constrained matrices
- Shallow NNs
- Analysis extends to other non-linear settings

## Experiments: Empirical Convergence

Estimator:  $\widehat{\text{SI}}_{m,n}$  with  $k$ -NN MI estimator [Kozachenko-Leonenko'87]

## Experiments: Empirical Convergence

Estimator:  $\widehat{\text{SI}}_{m,n}$  with  $k$ -NN MI estimator [Kozachenko-Leonenko'87]

- Let  $Z \sim \mathcal{N}(0, I_{15})$  and define:

- **$d = 3$ :**  $X = Z_{[1:3]}$  &  $Y = Z_{[2:4]}$
- **$d = 10$ :**  $X = Z_{[1:10]}$  &  $Y = Z_{[5:15]}$

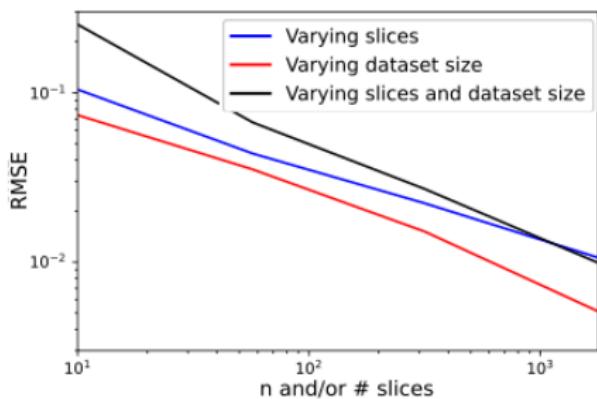
# Experiments: Empirical Convergence

Estimator:  $\widehat{\text{SI}}_{m,n}$  with  $k$ -NN MI estimator [Kozachenko-Leonenko'87]

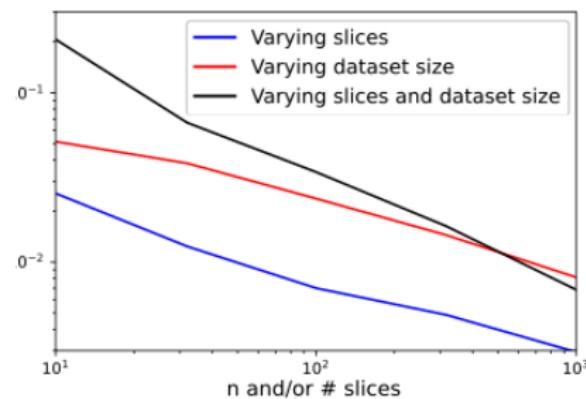
- Let  $Z \sim \mathcal{N}(0, I_{15})$  and define:

- $d = 3$ :  $X = Z_{[1:3]}$  &  $Y = Z_{[2:4]}$
- $d = 10$ :  $X = Z_{[1:10]}$  &  $Y = Z_{[5:15]}$

Figure: Empirical convergence rates



$$d = 3$$



$$d = 10$$