

# Cauchy-Schwarz Divergence Information Bottleneck

Shujian Yu

Vrije Universiteit Amsterdam

# Outline

- Information Bottleneck
- Deep Information Bottleneck
  - Cauchy-Schwarz divergence Information Bottleneck

# Information Bottleneck

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable



$X$

“cat” laying on a “laptop”

$T$

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable



“cat” laying on a “laptop”

$X$

$T$

- Tasks
  - Is there a cat? **relevant:** “cat”; **irrelevant:** “laptop”
  - How many pixels are there in the image? **irrelevant:** “cat” and “laptop”

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable
    - Related to task  $Y \rightarrow$  Useful for predicting  $Y$



“cat” laying on a “laptop”

$X$

$T$

- Tasks
  - Is there a cat? relevant: “cat”; irrelevant: “laptop”
  - How many pixels are there in the image? irrelevant: “cat” and “laptop”

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable
    - Related to task  $Y \rightarrow$  Useful for predicting  $Y$
  - How to define the optimal representation  $T$ ?

## Sufficient Statistics $S(X)$

$$I(S(X); Y) = I(X; Y)$$

A representation  $T$  of  $X$  is sufficient for  $Y$  if and only if  $I(X; Y) = I(T; Y)$ ;  $T$  contains **all** information regarding  $Y$  that can be obtained also from  $X$



# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable
    - Related to task  $Y \rightarrow$  Useful for predicting  $Y$
- How to define the optimal representation  $T$ ?

Sufficient Statistics  $S(X)$

$$I(S(X); Y) = I(X; Y)$$

Minimal Sufficient Statistics  $T(X)$

$$T(X) = \arg \min_{S(X): I(S(X); Y) = I(X; Y)} I(S(X); X)$$

$T$  contains **only** relevant information regarding  $Y$ ,  
but **least** information from  $X$ .

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable
    - Related to task  $Y \rightarrow$  Useful for predicting  $Y$
- How to define the optimal representation  $T$ ?

Sufficient Statistics  $S(X)$

$$I(S(X); Y) = I(X; Y)$$

Minimal Sufficient Statistics  $T(X)$

$$T(X) = \arg \min_{S(X): I(S(X); Y) = I(X; Y)} I(S(X); X)$$

## Sufficiency vs Minimality

# Information Bottleneck

- Let  $T$  be a representation of  $X$ 
  - Which  $T$  is useful?
    - Disentangled
    - Interpretable
    - Related to task  $Y \rightarrow$  Useful for predicting  $Y$
  - How to define the optimal representation  $T$ ?

Sufficient Statistics  $S(X)$

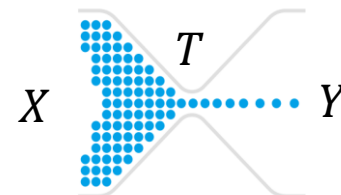
$$I(S(X); Y) = I(X; Y)$$

Minimal Sufficient Statistics  $T(X)$

$$T(X) = \arg \min_{S(X): I(S(X); Y) = I(X; Y)} I(S(X); X)$$

Information Bottleneck as an Approximation

$$\max_{p(t|x)} I(T; Y) - \beta I(T; X)$$

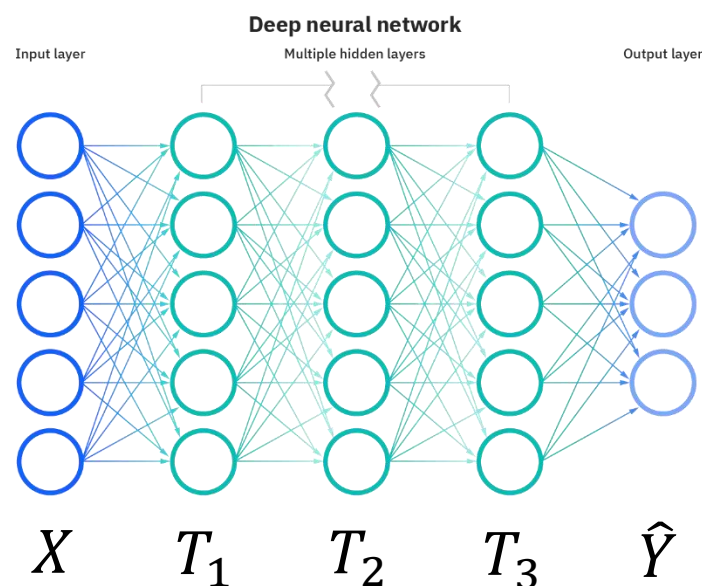


# Information Bottleneck Approach in Deep Neural Networks

# Information Bottleneck in DNNs

- Deep Information Bottleneck
  - Neural Network parameterization of IB

$$\max_{\theta} I(T; Y) - \beta I(T; X)$$



# Information Bottleneck in DNNs

- Deep Information Bottleneck
  - Neural Network parameterization of IB
    - $\max_{\theta} I(T; Y) - \beta I(T; X)$
    - How to estimate  $I(T; Y)$  and  $I(T; X)$
    - How to optimize?

# Information Bottleneck in DNNs

- Estimate  $I(T; Y)$ 
  - $\max_{\theta} I(T; Y) \Leftrightarrow \min_{\theta} D_{\text{KL}}(p(y|\mathbf{x}); q_{\theta}(\hat{y}|\mathbf{x}))$
  - If  $q_{\theta}(\hat{y}|\mathbf{x}) \sim \mathcal{N}(h_{\theta}(\mathbf{x}), \sigma^2 I)$ , we obtain  $\max_{\theta} \mathbb{E}(\|y - h_{\theta}(\mathbf{x})\|_2^2)$
  - How to evaluate  $I(T; Y)$  without any parametric assumption?

# Information Bottleneck in DNNs

- Estimate  $I(T; X)$ 
  - An upper bound of  $I(T; X)$

Variational upper bound (Variational IB)

$$\begin{aligned} I(T; X) &= \mathbb{E}_{p(x,t)} \log p(t|x) - \mathbb{E}_{p(t)} \log p(t) \\ &\leq \mathbb{E}_{p(x,t)} \log p(t|x) - \mathbb{E}_{p(t)} \log v(t) = D_{\text{KL}}(p(t|x); v(t)) \end{aligned}$$



# Information Bottleneck in DNNs

- Estimate  $I(T; X)$ 
  - An upper bound of  $I(T; X)$

Non-parametric upper bound (Nonlinear IB)

$$I(T; X) \leq -\frac{1}{N} \sum_{i=1}^N \log \frac{1}{N} \sum_{j=1}^N \exp \left( -D_{\text{KL}} \left( p(t|x_i); p(t|x_j) \right) \right)$$

# Information Bottleneck in DNNs

- Estimate  $I(T; X)$ 
  - An upper bound of  $I(T; X)$
  - How to evaluate  $I(T; X)$  without using an upper bound or variational approximation?

How to carry out nonlinear Information Bottleneck *without* variational approximation or parametric distributional assumption?

original IB  
Lagrangian

$$\max_{p(t|x)} I(T; Y) - \beta I(T; X)$$

deep IB objective

$$\min_{p(t|x)} D_{\text{KL}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I(T; X)$$

Cauchy-Schwarz  
divergence IB

$$\min_{p(t|x)} D_{\text{CS}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I_{\text{CS}}(T; X)$$

How to carry out nonlinear Information Bottleneck ***without*** variational approximation or parametric distributional assumption?

original IB  
Lagrangian

$$\max_{p(t|x)} I(T; Y) - \beta I(T; X)$$

deep IB objective

$$\min_{p(t|x)} D_{\text{KL}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I(T; X)$$

Cauchy-Schwarz  
divergence IB

$$\min_{p(t|x)} D_{\text{CS}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I_{\text{CS}}(T; X)$$

Conditional CS Divergence  
between  $p(y|x)$  and  $q_{\theta}(\hat{y}|x)$

How to carry out nonlinear Information Bottleneck *without* variational approximation or parametric distributional assumption?

original IB  
Lagrangian

$$\max_{p(t|x)} I(T; Y) - \beta I(T; X)$$

deep IB objective

$$\min_{p(t|x)} D_{\text{KL}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I(T; X)$$

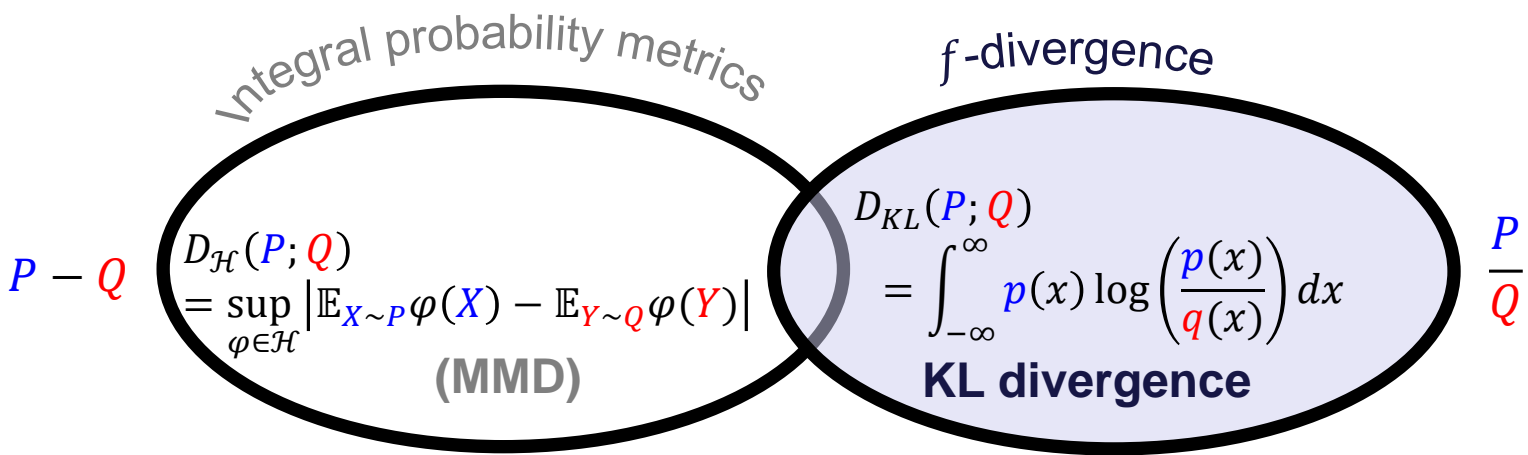
Cauchy-Schwarz  
divergence IB

$$\min_{p(t|x)} D_{\text{CS}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I_{\text{CS}}(T; X)$$

Conditional CS Divergence  
between  $p(y|x)$  and  $q_{\theta}(\hat{y}|x)$

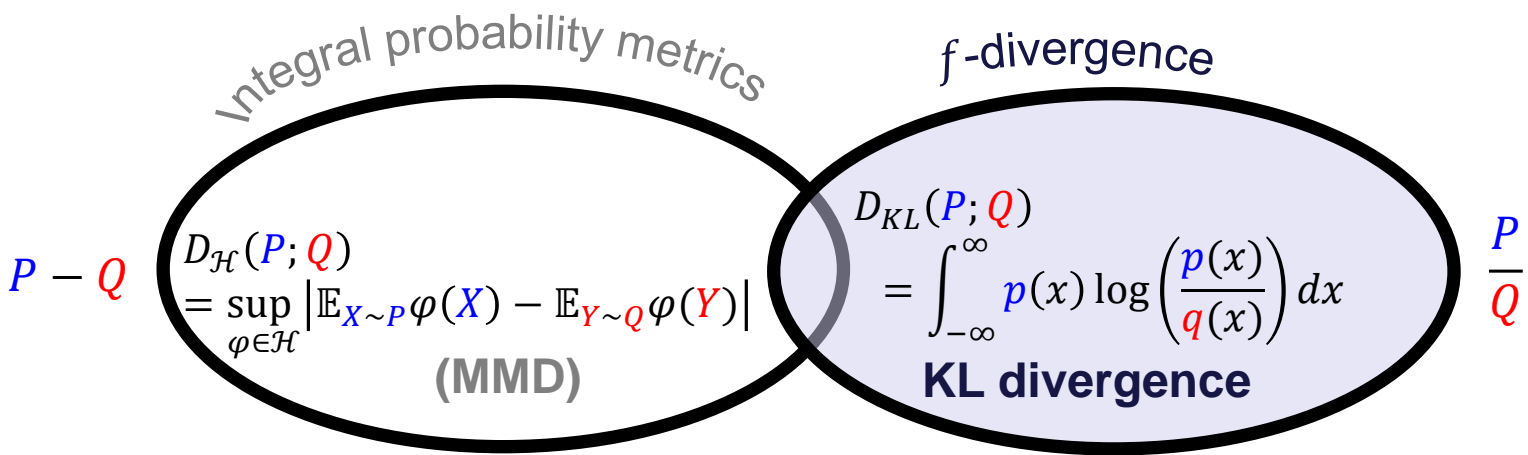
$I_{\text{CS}}(T; X) = D_{\text{CS}}(p(x, t); p(x)p(t))$   
CS Quadratic Mutual Information

# Cauchy-Schwarz Divergence



$$D_f(P; Q) = \int_{-\infty}^{\infty} q(x) f \left( \frac{p(x)}{q(x)} \right) dx$$

Support constraint due to ratio

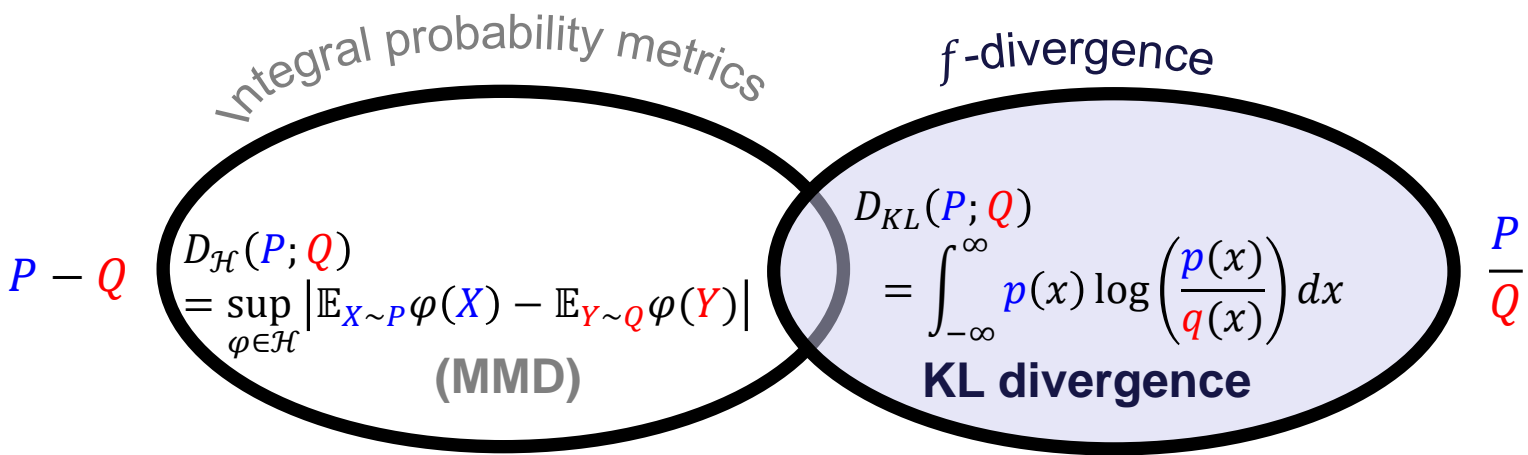


Cauchy-Schwarz inequality  
for square-integral functions

$$\left( \int p(x) q(x) dx \right)^2 \leq \int p^2(x) dx \int q^2(x) dx$$

$$\frac{\int p^2(x) dx \int q^2(x) dx}{\left( \int p(x) q(x) dx \right)^2}$$

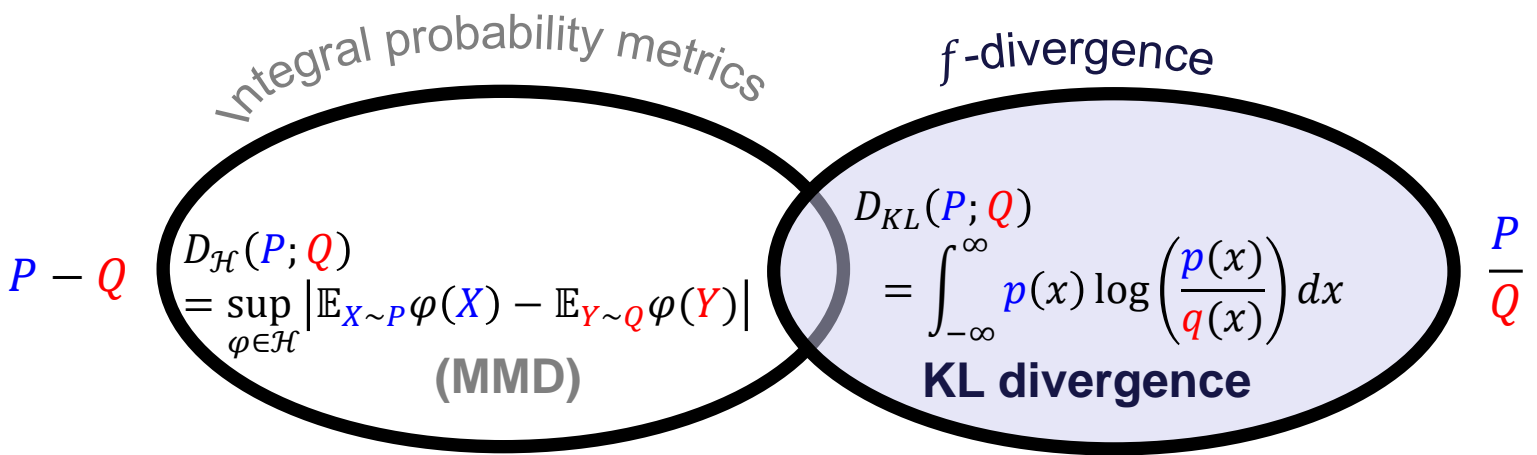




Cauchy-Schwarz inequality  
for square-integral functions

$$\left( \int p(x) q(x) dx \right)^2 \leq \int p^2(x) dx \int q^2(x) dx$$

$$\log \left( \frac{\int p^2(x) dx \int q^2(x) dx}{\left( \int p(x) q(x) dx \right)^2} \right)$$



The equation for  $D_{CS}(P; Q)$  is enclosed in a light green ellipse. To the right of the ellipse is the text  $\int PQ$  with  $P$  in blue and  $Q$  in red.

$D_{CS}(P; Q) =$   
 $\log \left( \frac{\int p^2(x) dx \int q^2(x) dx}{(\int p(x) q(x) dx)^2} \right)$

$\int PQ$

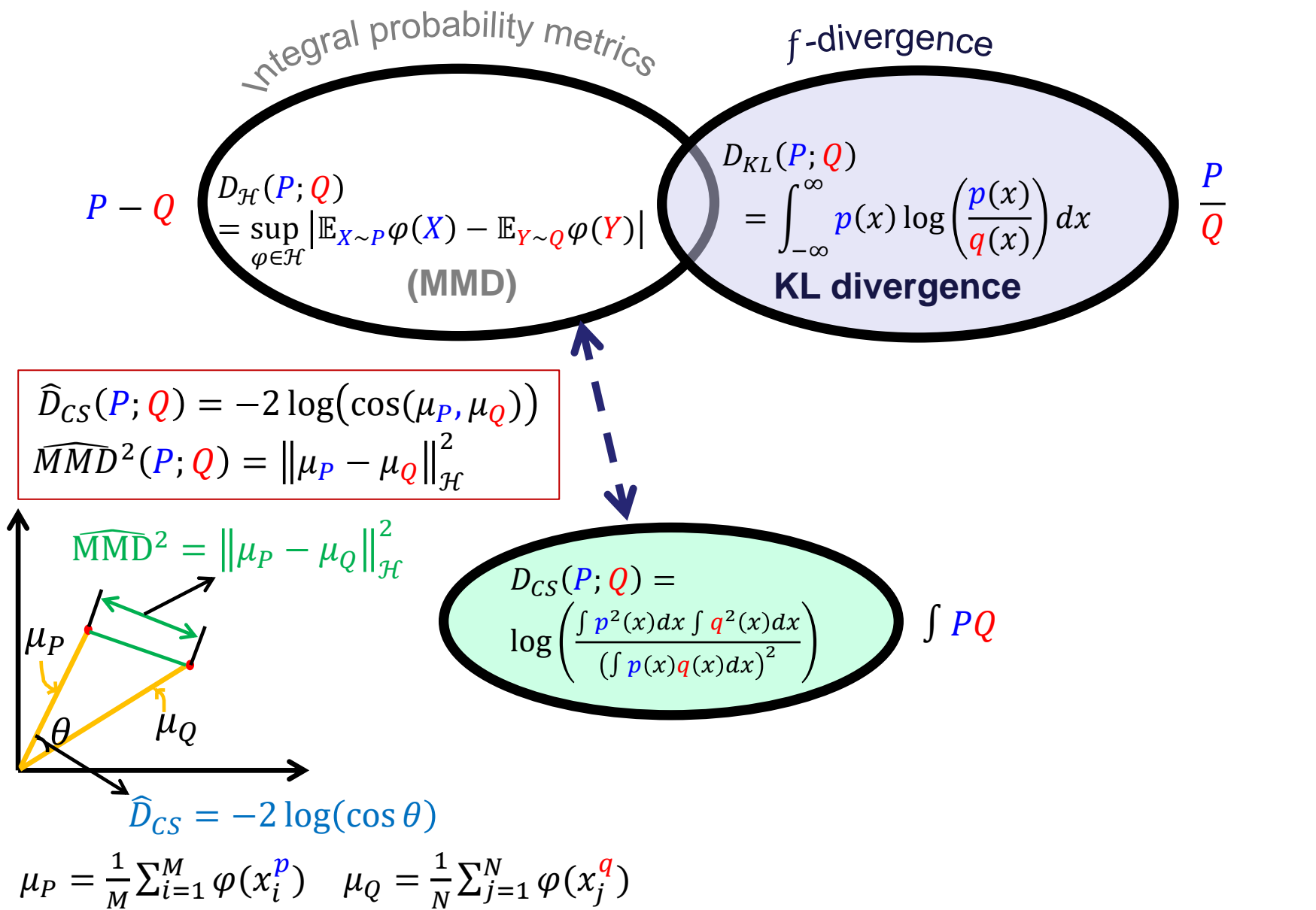
$$(\int p(x) q(x) dx)^2 \leq \int p^2(x) dx \int q^2(x) dx$$

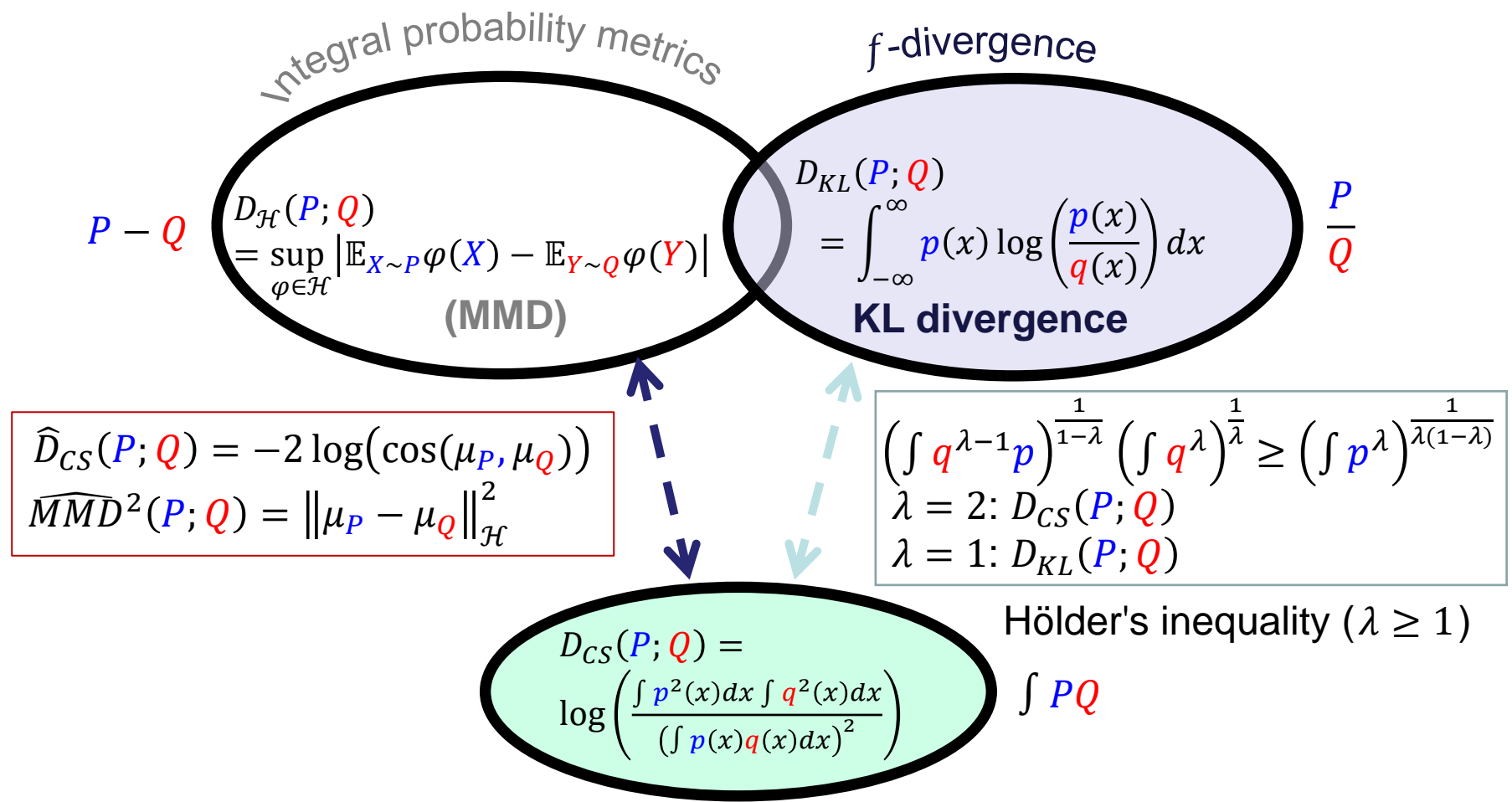
# Sample Estimator of the CS Divergence

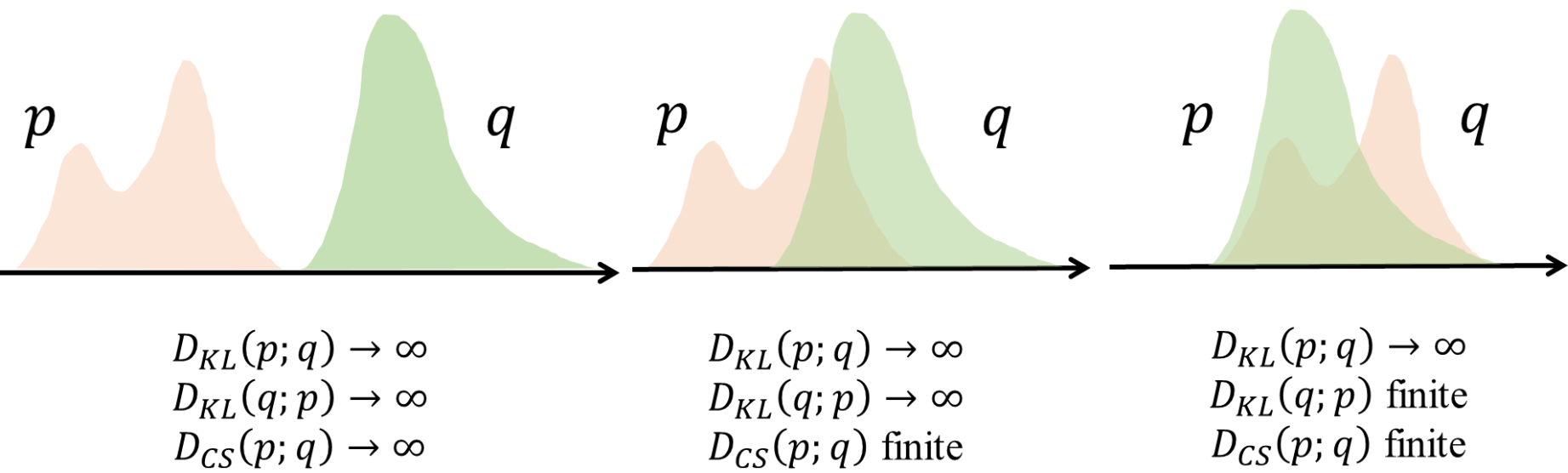
- Cauchy-Schwarz divergence between  $p(x)$  and  $q(x)$ ,  $x \in \mathbb{R}^d$ 
  - $\{x_i^p\}_{i=1}^M \sim p$ ,  $\{x_i^q\}_{i=1}^N \sim q$

$$\begin{aligned} \widehat{D}_{\text{CS}}(p; q) = & \underbrace{\log \left( \frac{1}{M^2} \sum_{i,j=1}^M \kappa_{\sigma} \left( x_i^p - x_j^p \right) \right)}_{\text{within distr. similarity}} + \underbrace{\log \left( \frac{1}{N^2} \sum_{i,j=1}^N \kappa_{\sigma} \left( x_i^q - x_j^q \right) \right)}_{\text{within distr. similarity}} \\ & - 2 \underbrace{\log \left( \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N \kappa_{\sigma} \left( x_i^p - x_j^q \right) \right)}_{\text{cross distr. similarity}} \end{aligned}$$

$$\begin{aligned} \widehat{MMD}^2(p; q) = & \frac{1}{M^2} \sum_{i,j=1}^M \kappa_{\sigma} \left( x_i^p - x_j^p \right) + \frac{1}{N^2} \sum_{i,j=1}^N \kappa_{\sigma} \left( x_i^q - x_j^q \right) \\ & - \frac{2}{MN} \sum_{i=1}^M \sum_{j=1}^N \kappa_{\sigma} \left( x_i^p - x_j^q \right) \end{aligned}$$







How to carry out nonlinear Information Bottleneck *without* variational approximation or parametric distributional assumption?

original IB  
Lagrangian

$$\max_{p(t|x)} I(T; Y) - \beta I(T; X)$$

deep IB objective

$$\min_{p(t|x)} D_{\text{KL}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I(T; X)$$

Cauchy-Schwarz  
divergence IB

$$\min_{p(t|x)} D_{\text{CS}}(p(y|x); q_{\theta}(\hat{y}|x)) + \beta I_{\text{CS}}(T; X)$$

Conditional CS Divergence  
between  $p(y|x)$  and  $q_{\theta}(\hat{y}|x)$

$I_{\text{CS}}(T; X) = D_{\text{CS}}(p(x, t); p(x)p(t))$   
CS Quadratic Mutual Information

- Estimate  $D_{\text{CS}}(p(y|\mathbf{x}); q_{\theta}(\hat{y}|\mathbf{x}))$

Proposition 1. Given  $\{\mathbf{x}_i, y_i, \hat{y}_i\}_{i=1}^N$ , let  $K$ ,  $L^1$  and  $L^2$  denote respectively the kernel matrices for variables  $\mathbf{x}$ ,  $y$  and  $\hat{y}$ . Further, let  $L^{21}$  denote the cross-kernel matrix for  $y$  and  $\hat{y}$  (i.e.,  $L_{ij}^{21} = \kappa(\hat{y}_i, y_j)$ ).  $D_{\text{CS}}(p(y|\mathbf{x}); q_{\theta}(\hat{y}|\mathbf{x}))$  is estimated by:

$$D_{\text{CS}}(p(y|\mathbf{x}); q_{\theta}(\hat{y}|\mathbf{x})) = \log \left( \sum_{j=1}^N \left( \frac{\sum_{i=1}^N K_{ji} L_{ji}^1}{(\sum_{i=1}^N K_{ji})^2} \right) \right) + \log \left( \sum_{j=1}^N \left( \frac{\sum_{i=1}^N K_{ji} L_{ji}^2}{(\sum_{i=1}^N K_{ji})^2} \right) \right) - 2 \log \left( \sum_{j=1}^N \left( \frac{\sum_{i=1}^N K_{ji} L_{ji}^{21}}{(\sum_{i=1}^N K_{ji})^2} \right) \right)$$



- Estimate  $I_{\text{CS}}(T; X)$

$$\begin{aligned} I_{\text{CS}}(T; X) &= D_{\text{CS}}(\mathbf{p}(\mathbf{x}, \mathbf{t}); \mathbf{p}(\mathbf{x})\mathbf{p}(\mathbf{t})) \\ &= \log \left( \frac{\int \mathbf{p}^2(\mathbf{x}, \mathbf{t}) d\mathbf{x} d\mathbf{t} \int \mathbf{p}^2(\mathbf{x}) \mathbf{p}^2(\mathbf{t}) d\mathbf{x} d\mathbf{t}}{(\int \mathbf{p}(\mathbf{x}, \mathbf{t}) \mathbf{p}(\mathbf{x}) \mathbf{p}(\mathbf{t}) d\mathbf{x} d\mathbf{t})^2} \right) \end{aligned}$$

Proposition 2. Given  $\{\mathbf{x}_i, \mathbf{t}_i\}_{i=1}^N$ , let  $K$  and  $Q$  denote respectively the kernel matrices for variables  $\mathbf{x}$  and  $\mathbf{t}$ , and  $\mathbf{1}$  denote a  $N \times 1$  vector of ones. The empirical estimator of  $I_{\text{CS}}(T; X)$  is given by:

$$I_{\text{CS}}(T; X) = \log \left( \frac{1}{N^2} \text{tr}(KQ) \right) + \log \left( \frac{1}{N^4} \mathbf{1}^T K \mathbf{1} \mathbf{1}^T Q \mathbf{1} \right) - 2 \log \left( \frac{1}{N^3} \mathbf{1}^T K Q \mathbf{1} \right)$$

# Information Bottleneck in DNNs

- Prediction – Compression Tradeoff

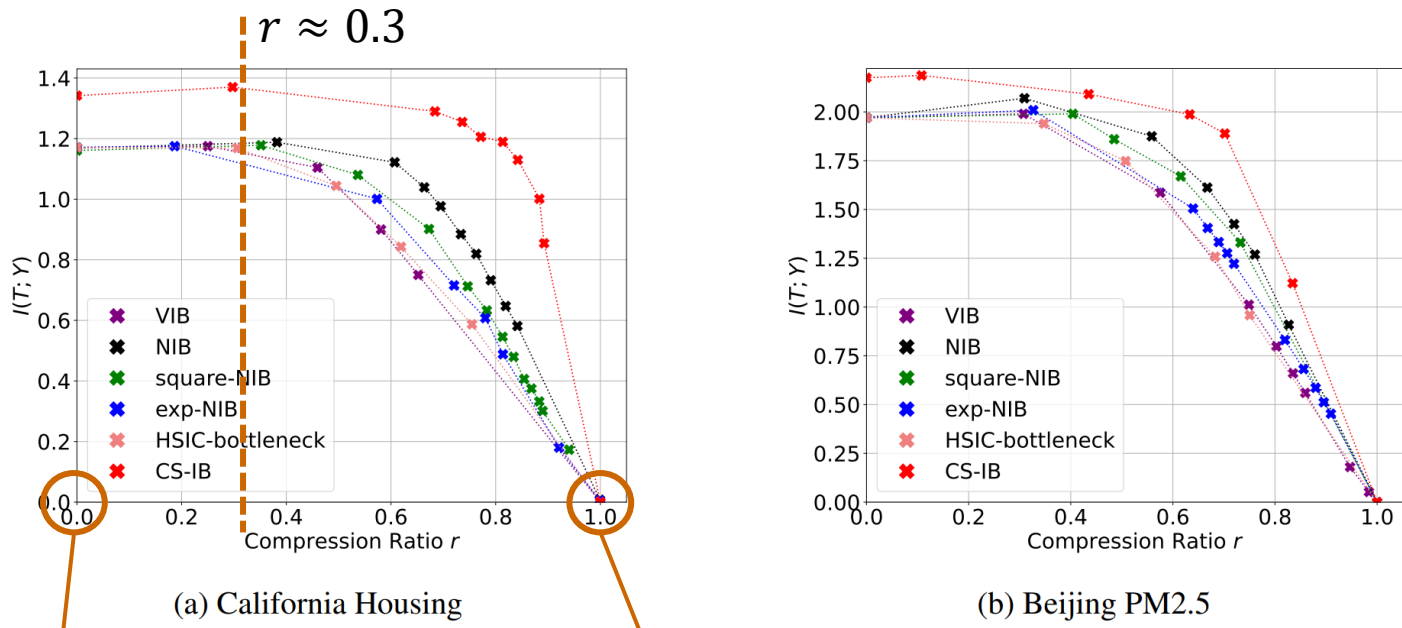


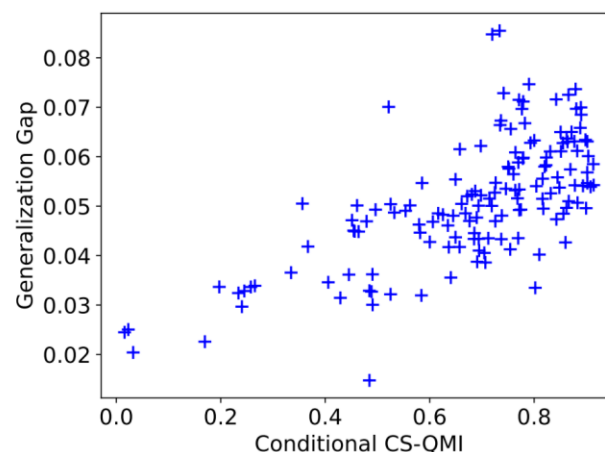
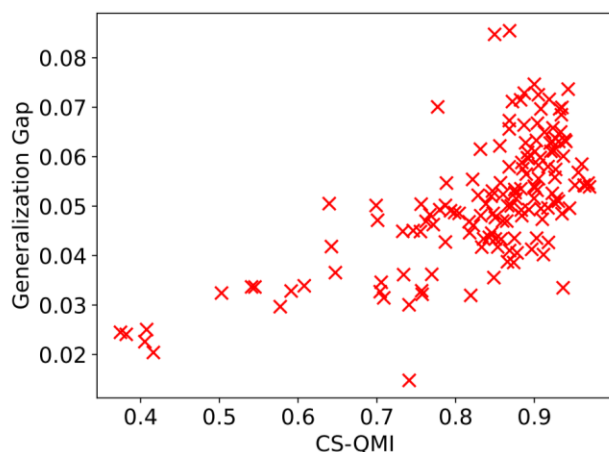
Figure 1: Information plane diagrams on California Housing and Beijing PM2.5 datasets.

no compression ( $r = 0$ )

maximum compression  
( $r = 1$ )

# Information Bottleneck in DNNs

- Generalization



$I_{CS}(\mathbf{x}; \mathbf{t})$  (left) and  $I_{CS}(\mathbf{x}; \mathbf{t}|\mathbf{y})$  (right) with respect to the generalization gap in California housing w.r.t. 100 individually trained NNs.

Both  $I(X; T|Y)$  and  $I(X; T)$  correlate well with empirical generalization error gap.

Does Compression imply generalization?