

# Parametric Information Bottleneck <sup>1</sup>

## SAIL@UNIST Group Meeting

T.T. Nguyen and J. Choi

SAIL@UNIST

November 3, 2018



---

<sup>1</sup>v4: 18/11/02

# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

- Classification and Adversarial Robustness
- Learning dynamics



# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

- Classification and Adversarial Robustness
- Learning dynamics



Statistical Artificial Intelligence  
Laboratory @UNIST

# Introduction

- Deep neural networks (DNNs) = flexible modeling capability
  - Multi-layered neural networks
- Learning principle = a principled way of exploiting a model for certain tasks
- The Maximum Likelihood Estimate (MLE) principle:
  - A de-facto learning principle for DNNs
  - Maximizing the model's likelihood of seeing the training data

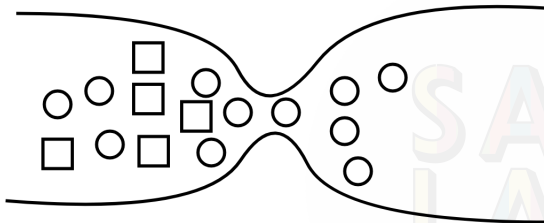
# MLE + DNNs = ?

- MLE: ignores the special topology of DNNs during the learning
  - generic to many models, not dedicatedly tailored for DNNs
  - MLE sees the entire neural architecture  $f(\mathbf{x})$  as a whole, without taking advantage of the hierarchical structure of the neural function  $f$
- We need a clever way to exploit the topology of DNNs for learning!  
→ Parametric Information Bottleneck (PIB)!

# Parametric Information Bottleneck: A quick look

Informally, PIB principle is assuring that each layer of DNNs maximally preserves the information relevant to the task while it is being compressed!

"I'm letting only the relevant information past and squeeze out the irrelevant one"



PIB = an efficient way to induce layer-wise relevance-compression trade-offs to DNNs!

# More about PIB

Though being conceptually simple, the layer-wise compression-relevance trade-offs pose challenges theoretically and empirically:

- 1** Is it possible to obtain the optimal compression-relevance trade-offs simultaneously at all layers?
- 2** Computing compression and relevance in DNNs is highly intractable

The main contributions in PIB is efficiently inducing layer-wise compression-relevance trade-offs to DNNs by addressing 1. and proposing an approximate mutual information to 2.

# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

- Classification and Adversarial Robustness
- Learning dynamics



Statistical Artificial Intelligence  
Laboratory @UNIST



# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

- Classification and Adversarial Robustness
- Learning dynamics

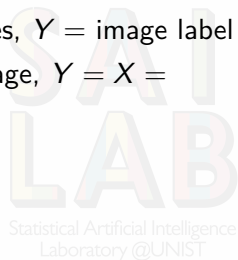


Statistical Artificial Intelligence  
Laboratory @UNIST

# Information Bottleneck Principle: Problem Setting

Learning problem setting:  $X$  = input random variable (RV),  $Y$  = output random variable

- E.g., supervised learning:  $X$  = input images,  $Y$  = image label
- E.g., unsupervised learning:  $X$  = input image,  $Y = X$  = reconstruction image



# Information Bottleneck Principle

- A principled way of extracting **relevant information** in one variable,  $X$  about another variable,  $Y$ .
- Encode  $X$  into an intermediate representation,

$$X \xrightarrow{p(z|x)} Z,$$

in such a way that  $Z$  preserves as much of **relevant information** about  $Y$  as possible.

# Information Bottleneck Principle

IB = optimizing the compression-relevance trade-off (the following problems are equivalent):

- Compression-Relevance Function

[Tishby et al., 1999, Slonim and Weiss, 2002]:

$$\arg \min_{P(Z|X): I(Z; Y) \geq D} I(Z; X)$$

where  $D$  is a positive number specifying the minimum relevance.

- Lagrangian multiplier [Tishby et al., 1999]:

$$\arg \min_{P(Z|X)} I(Z; X) - \beta I(Z; Y)$$

where  $\beta$  is a (positive) Lagrangian multiplier.

# Information Bottleneck Principle

- **Compression** (measure) = representation complexity =  $I(Z; X)$
- **Relevance** = predictive power =  $I(Z; Y)$
- Mutual Information:

$$I(Z; Y) = \int p(\mathbf{z}, \mathbf{y}) \log \frac{p(\mathbf{z}, \mathbf{y})}{p(\mathbf{z})p(\mathbf{y})} d\mathbf{z} d\mathbf{y},$$

Intuitively,  $I(Z; Y)$  = the amount of information that  $Z$  contains about  $Y$ .

- Compression-relevance trade-off:
  - You might lose relevant information as well if you compress (the representation) too much!
  - But if you allow too much information (in the representation), you might let past irrelevant information (which is not useful for the task performance)

# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

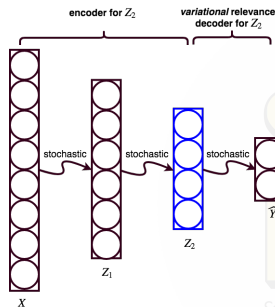
- Classification and Adversarial Robustness
- Learning dynamics



# Parametric Information Bottleneck

- PIB = Layer-wise application of IB to DNNs
- Topology of DNNs = Markov chain structure:

$$Y \rightarrow X \rightarrow Z_I \rightarrow Z_{I+1} \rightarrow \hat{Y}$$



- Encoder  $P(Z_I|X)$  is induced by the network architecture and Bayes' Rule, e.g.,  $P(Z_1|X) = \sigma(W_1X + b_1)$  in sigmoid-activated fully-connected neural network

# PIB

PIB = layer-wise multi-objective Information Bottleneck:

$$\forall 1 \leq l \leq L, \min_{P(Z_l|X)} \mathcal{L}_l[P(Z_l|X)] \quad (1)$$

where  $\mathcal{L}_l[P(Z_l|X)] := I(Z_l; X) - \beta_l I(Z_l; Y)$



# PIB

Unfortunately, we cannot achieve the layer-wise optimality for simultaneously all layers <sup>2</sup>

**Theorem 3.1** (Conflicting Information Optimality). *Given  $Y \rightarrow X \rightarrow Z_1 \rightarrow Z_2, Z_2 \not\perp Z_1, \beta_1 > 0$ , and  $\beta_2 > 0$ , then  $\mathcal{L}_1$  and  $\mathcal{L}_2$  defined by the layer-wise multi-objective Information Bottleneck are conflicting, i.e., there does not exist a single solution that minimizes  $\mathcal{L}_1$  and  $\mathcal{L}_2$  simultaneously.*

→ The layers to compromise their information optimality → we propose two simple compromised strategies: *JointPIB* and *GreedyPIB*

---

<sup>2</sup>Detailed proof at <https://arxiv.org/abs/1712.01272>

# PIB: compromised optimality

## ■ JointPIB:

$$\mathcal{L}^{joint} := \sum_{l=0}^L \gamma_l \tilde{\mathcal{L}}_l$$

- GreedyPIB: applies PIB progressively in a greedy manner. In other words, *GreedyPIB* tries to obtain the conditional optimality of a current layer which is conditioned on the achieved conditional optimality of the previous layers.

# Intractability of Mutual Information

- Compression and Relevance (which are mutual information by nature) of high-dimensional random variables in PIB are highly intractable
- We propose *Variational Relevance* and *Variational Compression* to address that issue for DNNs

SAI  
LAB

Statistical Artificial Intelligence  
Laboratory @UNIST

# Relevance in PIB

## ■ Relevance:

$$I(Z_I; Y) = H(Y) - H(Y|Z_I)$$

$$H(Y|Z_I) = - \int p(\mathbf{y}, \mathbf{z}_I) \log p(\mathbf{y}|\mathbf{z}_I) d\mathbf{y} d\mathbf{z}_I$$

## ■ Relevance decoder:

$$p(\mathbf{y}|\mathbf{z}_I) = \int p_D(\mathbf{x}, \mathbf{y}) \frac{p(\mathbf{z}_I|\mathbf{x})}{p(\mathbf{z}_I)} d\mathbf{x}$$

# Variational Relevance

- It follows from Jensen's inequality that:

$$\begin{aligned}
 H(Y|Z_I) &= - \int p(\mathbf{y}|\mathbf{z}_I) p(\mathbf{z}_I) \log \boxed{p(\mathbf{y}|\mathbf{z}_I)} d\mathbf{y} d\mathbf{z}_I \\
 &\quad \text{relevance decoder} \downarrow \\
 &\leq - \int p(\mathbf{y}|\mathbf{z}_I) p(\mathbf{z}_I) \log \boxed{q(\mathbf{y}|\mathbf{z}_I)} d\mathbf{y} d\mathbf{z}_I \\
 &\quad \text{variational relevance decoder} \downarrow \\
 &= - \mathbb{E}_{(\mathbf{x}, \mathbf{y})_D} \mathbb{E}_{\mathbf{z}_I | \mathbf{x}} \log \boxed{q(\mathbf{y}|\mathbf{z}_I)} =: \tilde{H}(Y|Z_I) \quad (2) \\
 &\quad \text{variational relevance decoder} \downarrow
 \end{aligned}$$

where  $q(\mathbf{y}|\mathbf{z}_I)$  is any probability distribution.

- In PIB, we set  $q(Y|Z_I) = P(\hat{Y}|Z_I)$ : re-use the the higher-level network architecture to define variational relevance decoder
- $\tilde{H}(\hat{Y}|Z_I) =:$  variational conditional relevance (VCR)

# Variational Relevance

Variational Conditional Relevance (VCR) generalizes MLE to all layers in DNNs!

**Theorem 3.2** (Information on the extreme layers). *The VCR of the lowest-level (so-called **super**) layer (i.e.,  $l = 0$ ) is the negative log-likelihood (NLL) function of the neural network, i.e.,*

$$\tilde{H}(Y|Z_0) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y})_D} [\log p(\hat{\mathbf{y}}|\mathbf{x})]. \quad (18)$$

*Similarly, the VCR of the highest-level layer (i.e.,  $l = L$ ) equals that of the **compositional** layer  $Z = (Z_1, Z_2, \dots, Z_L)$ , a composite of all hidden layers; in addition, their VCR is an upper bound on the NLL:*

$$\tilde{H}(Y|Z_L) = \tilde{H}(Y|Z) \geq -\mathbb{E}_{(\mathbf{x}, \mathbf{y})_D} [\log p(\hat{\mathbf{y}}|\mathbf{x})]. \quad (19)$$



# Variational Compression

- Avoid directly estimating  $I(Z_I; X)$  by instead resorting to its upper bound  $I(Z_I; Z_{I-1})$
- Approximate  $I(Z_I; Z_{I-1})$  using a mean-field (factorized) variational distribution  $r(\mathbf{z}_I) = \prod_{i=1}^{n_I} r(z_{I,i})$ :

$$\begin{aligned} I(Z_I; X) &\leq I(Z_I; Z_{I-1}) \\ &= \int p(\mathbf{z}_I | \mathbf{z}_{I-1}) p(\mathbf{z}_{I-1}) \log \frac{p(\mathbf{z}_I | \mathbf{z}_{I-1})}{p(\mathbf{z}_I)} d\mathbf{z}_I d\mathbf{z}_{I-1} \\ &\leq \int p(\mathbf{z}_I | \mathbf{z}_{I-1}) p(\mathbf{z}_{I-1}) \log \frac{p(\mathbf{z}_I | \mathbf{z}_{I-1})}{r(\mathbf{z}_I)} d\mathbf{z}_I d\mathbf{z}_{I-1} \\ &= \mathbb{E}_{\mathbf{z}_{I-1}} \sum_{i=1}^{n_I} D_{KL} [p(z_{I,i} | \mathbf{z}_{I-1}) || r(z_{I,i})] \\ &=: \tilde{I}(Z_I; Z_{I-1}) \end{aligned} \tag{3}$$

# Contents

## 1 Introduction

## 2 Information Bottleneck Principle

## 3 Parametric Information Bottleneck

- Layer-wise Multi-Objective IB
- Approximate Mutual Information

## 4 Experiments

- Classification and Adversarial Robustness
- Learning dynamics





# Classification and Adversarial Robustness

- PIB offers a DNN model with competitive classification performance and robustness against adversarial attacks

Model	Classification		Adv. Robustness (%)	
	MNIST (Error %)	CIFAR10 (Accuracy %)	Targeted	Untargeted
DET	1.73	53.91	00.00	00.00
VIB Alemi et al. (2017)	1.45	54.41	83.70	93.10
SFNN Raiko et al. (2015)	1.44	55.94	83.00	95.20
GreedyPIB	1.54	<b>57.61</b>	83.21	94.30
JointPIB	<b>1.36</b>	55.62	<b>84.16</b>	<b>96.00</b>

Table 1: The performance of PIB for classification and adversarial robustness on MNIST and CIFAR10 in comparison with MLE and a partially information-theoretic treatment VIB.

# Learning dynamics

- Closer look at the inside of learning: PIB preserves relevant information throughout all layers during the learning

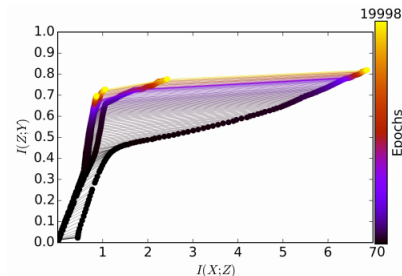


Figure 1. Learning dynamics of SFNN

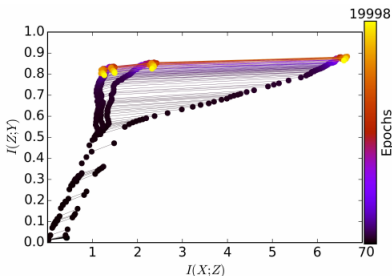


Figure 2. Learning dynamics of JointPIB

# Learning dynamics

- GreedyPIB progressively preserves relevant information throughout all layers in a greedy manner.

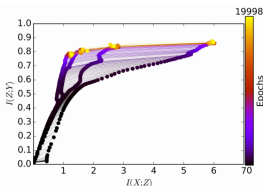


Figure 3. Learning dynamics of GreedyPIB at  $l = 1$ .

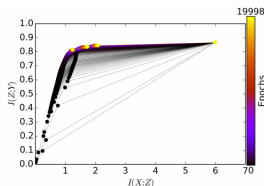


Figure 4. Learning dynamics of GreedyPIB at  $l = 2$ .

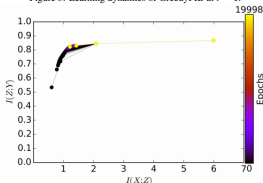


Figure 5. Learning dynamics of GreedyPIB at  $l = 3$ .

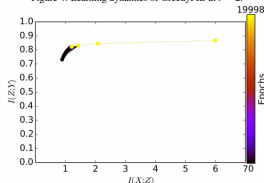


Figure 6. Learning dynamics of GreedyPIB at  $l = 4$ .

Thank you for listening!



Statistical Artificial Intelligence  
Laboratory @UNIST

# References

 Slonim, N. and Weiss, Y. (2002).

Maximum likelihood and the information bottleneck.  
*In Advances in Neural Information Processing Systems 15 [Neural Information Processing Systems, NIPS 2002, December 9-14, 2002, Vancouver, British Columbia, Canada], pages 335–342.*

 Tishby, N., Pereira, F. C., and Bialek, W. (1999).

The information bottleneck method.  
*In Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing.*