

Information Flow in Deep Neural Networks

Mohammadamin Kiani Mohammad Mohammadian Diba Hadi

Sharif University of Technology
Information Theory, Statistics, and Learning

Instructor: Prof. Yassaee
Project Mentor: Mr. Hadavi

Sep 4, 2025

- 1 Introduction & Background
- 2 Proposed Framework: Noisy DNNs
- 3 Mutual Information Estimation
 - Theoretical Guarantees for SP Estimator
- 4 Compression and Clustering Connection
- 5 Conclusions
- 6 Our Work & Empirical Results
 - Theoretical Stability Guarantees
 - Empirical Results
- 7 Empirical Results

- Treat the whole layer T , as a single random variable, characterized by its encoder, $P(T|X)$, and decoder, $P(Y|T)$ distributions.
- The plane of the mutual information values of any other variable with the input variable X and the desired output variable Y .
- **Rational:**
 - invariance to invertible re-parameterization $I(X; Y) = I(\psi(X); \psi(y))$
 - data processing inequality $I(X; Y) \geq I(X; Z)$
- Information path: layers are mapped to K monotonic connected points in the plane.
 - $I(X; Y) \geq I(T_1; y) \geq \dots \geq I(T_k; Y) \geq I(\hat{Y}; Y)$
 - $H(X) \geq I(X; T_1) \geq \dots \geq I(X; T_k) \geq I(X; \hat{Y})$

Motivation: Understanding DNN Internals

- **Goal:** Characterize information learned in DNN hidden layers.
- **Mutual Information (MI)** $I(\mathbf{X}; \mathbf{T}_\ell)$ as a measure:
 - \mathbf{X} : Input, \mathbf{T}_ℓ : Layer ℓ representation.
 - Principled: Invariant to invertible ops, meaningful units (bits/nats).
- **Information Bottleneck (IB) Theory Context:**
 - Learn **compressed** \mathbf{T}_ℓ of \mathbf{X} , **informative** about target \mathbf{Y} .
 - Suggests training phases: Fitting ($I(\mathbf{T}_\ell; \mathbf{Y}) \uparrow$) and Compression ($I(\mathbf{X}; \mathbf{T}_\ell) \downarrow$).
 - IB Tradeoff: optimal trade-off between compression of X and prediction of Y
 - $T(X) = \min_{p(t|x), p(y|t), p(t)} I(X; T) - \beta I(T; Y)$
 - β determines the level of relevant information captured by the representation T

The Problem with $I(\mathbf{X}; \mathbf{T}_\ell)$ in Deterministic DNNs

- Deterministic DNN: $\mathbf{T}_\ell = f_\ell(\dots f_1(\mathbf{X}))$.
- For common continuous nonlinearities:
 - Continuous $P_X \implies I(\mathbf{X}; \mathbf{T}_\ell) = \infty$.
 - Discrete P_X (e.g., dataset) $\implies I(\mathbf{X}; \mathbf{T}_\ell) = H(\mathbf{X})$ (constant).
- **Consequence:** True $I(\mathbf{X}; \mathbf{T}_\ell)$ is often vacuous for observing "compression dynamics."
- Observed changes in $I(\mathbf{X}; \text{Bin}(\mathbf{T}_\ell))$ (binned MI) in prior work likely stem from the [binning approximation](#), not true MI changes.

A Rigorous Framework: Noisy DNNs

- To make $I(\mathbf{X}; \mathbf{T}_\ell)$ well-defined & sensitive to parameters, map $\mathbf{X} \mapsto \mathbf{T}_\ell$ must be **stochastic**.
- **Proposed Model:** Add i.i.d. Gaussian noise at each hidden layer output.

$$\mathbf{T}_\ell = f_\ell(\mathbf{T}_{\ell-1}) + \mathbf{Z}_\ell, \quad \mathbf{Z}_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I}_{d_\ell})$$

- $\mathbf{S}_\ell = f_\ell(\mathbf{T}_{\ell-1})$ is the "signal" part before noise.
- $\mathbf{T}_\ell = \mathbf{S}_\ell + \mathbf{Z}_\ell$.
- This makes $I(\mathbf{X}; \mathbf{T}_\ell)$ finite and dependent on network weights.
- Data Processing Inequality holds: $\mathbf{X} - \mathbf{T}_1 - \dots - \mathbf{T}_L$.

Estimating $I(\mathbf{X}; \mathbf{T}_\ell)$ in Noisy DNNs

- Definition: $I(\mathbf{X}; \mathbf{T}_\ell) = h(\mathbf{T}_\ell) - \mathbb{E}_{\mathbf{X}}[h(\mathbf{T}_\ell | \mathbf{X} = \mathbf{x})]$.
- Direct computation of differential entropies $h(\cdot)$ is intractable.
- **Sample Propagation (SP) Estimator Idea:**
 - PDF of \mathbf{T}_ℓ : $p_{\mathbf{T}_\ell}(\mathbf{t}) = (p_{\mathbf{S}_\ell} * \phi_\beta)(\mathbf{t})$ (convolution).
 - $p_{\mathbf{S}_\ell}$: PDF of signal \mathbf{S}_ℓ . ϕ_β : PDF of noise \mathbf{Z}_ℓ .
 - Estimate $h(p_{\mathbf{T}_\ell})$ via $h(\hat{p}_{\mathbf{S}_\ell} * \phi_\beta)$, using empirical $\hat{p}_{\mathbf{S}_\ell}$ from samples of \mathbf{S}_ℓ .
 - Similar approach for conditional term $h(p_{\mathbf{T}_\ell | \mathbf{X} = \mathbf{x}})$.

The Sample Propagation (SP) Estimator

Given dataset $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^n$:

- ① **Estimate** $h(\mathbf{T}_\ell)$: Use $h(\hat{p}_{\mathcal{S}_\ell} * \phi_\beta)$.
 - Samples $\mathcal{S}_\ell = \{\mathbf{s}_{\ell,i}\}_{i=1}^n$ from $f_\ell(\dots f_1(\mathbf{x}_i))$.
 - This is entropy of a Gaussian mixture centered at $\mathbf{s}_{\ell,i}$.
- ② **Estimate** $h(\mathbf{T}_\ell | \mathbf{X} = \mathbf{x})$: Use $h(\hat{p}_{\mathcal{S}_\ell^{(\mathbf{x})}} * \phi_\beta)$.
 - For each \mathbf{x} , n_x samples $\mathcal{S}_\ell^{(\mathbf{x})}$ by passing \mathbf{x} multiple times through $f_\ell(\dots f_1(\cdot))$.

SP Estimator \hat{I}_{SP} :

$$\hat{I}_{SP} = h(\hat{p}_{\mathcal{S}_\ell} * \phi_\beta) - \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{X}} h(\hat{p}_{\mathcal{S}_\ell^{(\mathbf{x})}} * \phi_\beta)$$

Entropies of Gaussian mixtures often computed via Monte Carlo Integration.

Theoretical Guarantees: Preliminaries

Focus: Estimating $h(P_S * \phi_\beta)$ from N i.i.d. samples $\mathcal{S}_N = \{\mathbf{S}_i\}$ from P_S .

Estimator: $h(\hat{P}_{\mathcal{S}_N} * \phi_\beta)$.

- **Minimax Absolute-Error Risk** over distribution class \mathcal{F} for S :

$$R^*(N, \beta, \mathcal{F}) = \inf_{\hat{h}} \sup_{P_S \in \mathcal{F}} \mathbb{E} \left| h(P_S * \phi_\beta) - \hat{h}(\mathcal{S}_N, \beta) \right|$$

Measures worst-case error for the *best possible* estimator.

Theoretical Guarantees: Preliminaries

Focus: Estimating $h(P_S * \phi_\beta)$ from N i.i.d. samples $\mathcal{S}_N = \{\mathbf{S}_i\}$ from P_S .
Estimator: $h(\hat{P}_{\mathcal{S}_N} * \phi_\beta)$.

- **Minimax Absolute-Error Risk** over distribution class \mathcal{F} for S :

$$R^*(N, \beta, \mathcal{F}) = \inf_{\hat{h}} \sup_{P_S \in \mathcal{F}} \mathbb{E} \left| h(P_S * \phi_\beta) - \hat{h}(\mathcal{S}_N, \beta) \right|$$

Measures worst-case error for the *best possible* estimator.

- **Sample Complexity** $N^*(\eta, \beta, \mathcal{F})$: Smallest N for $R^* \leq \eta$. Minimum samples needed to achieve a target accuracy η .
- Classes considered: \mathcal{F}_d (distributions with support in $[-1, 1]^d$, e.g., for tanh layers), $\mathcal{F}_{d, \mu, K}^{(SG)}$ (subgaussian distributions, e.g., for ReLU layers with subgaussian inputs).

Guarantee 1: Sample Complexity

Statement (Simplified): For fixed noise $\beta > 0$, large dimension d , target error $\eta < \eta_0(\beta)$:

$$N^*(\eta, \beta, \mathcal{F}_d) \geq \Omega\left(\frac{2^{\gamma(\beta)d}}{d\eta}\right)$$

Guarantee 1: Sample Complexity

Statement (Simplified): For fixed noise $\beta > 0$, large dimension d , target error $\eta < \eta_0(\beta)$:

$$N^*(\eta, \beta, \mathcal{F}_d) \geq \Omega\left(\frac{2^{\gamma(\beta)d}}{d\eta}\right)$$

- **Core Implication:** Estimating $h(P_S * \phi_\beta)$ is **fundamentally hard in high dimensions**.
- The required number of samples N^* grows at least **exponentially with dimension d** .
- The exponent $\gamma(\beta)$ is monotonically decreasing in β .
 - Larger noise variance β^2 (larger β) \implies smaller $\gamma(\beta)$ \implies (relatively) less severe exponential dependence.
 - More noise "smooths" the distribution, making estimation easier.
- This is a lower bound, applying to *any* estimator, not just the SP one.

Guarantee 2: Risk of $h(\hat{P}_{S_N} * \phi_\beta)$

Statement (Simplified): For $P_S \in \mathcal{F}_{d,\mu,K}^{(SG)}$ (or \mathcal{F}_d):

$$\mathbb{E} \left| h(P_S * \phi_\beta) - h(\hat{P}_{S_N} * \phi_\beta) \right| \leq C(d, \mu, K, \beta) \cdot \frac{1}{\sqrt{N}}$$

Guarantee 2: Risk of $h(\hat{P}_{S_N} * \phi_\beta)$

Statement (Simplified): For $P_S \in \mathcal{F}_{d,\mu,K}^{(SG)}$ (or \mathcal{F}_d):

$$\mathbb{E} \left| h(P_S * \phi_\beta) - h(\hat{P}_{S_N} * \phi_\beta) \right| \leq C(d, \mu, K, \beta) \cdot \frac{1}{\sqrt{N}}$$

- **Core Implication:** The specific SP-based estimator $h(\hat{P}_{S_N} * \phi_\beta)$ achieves the **parametric rate of convergence** $O(1/\sqrt{N})$ with respect to sample size N .
- This is generally the best possible convergence rate for parametric estimation problems.
- However, the constant $C(d, \mu, K, \beta)$ can be (and often is) **exponential in dimension d** .
 - $C(d, \mu, K, \beta) \approx \left(\frac{1}{\sqrt{2}} + \frac{K}{\beta} \right)^d \times \text{polynomial factors in } d, \mu, K, 1/\beta$.
 - This reflects the curse of dimensionality.

Guarantee 3: MI Estimation Risk

The absolute-error risk of the full SP MI estimator \hat{I}_{SP} (using n samples for unconditional term, and $n_x = n$ for each of n conditional terms):

$$\sup_{P_X} \mathbb{E} \left| I(\mathbf{X}; \mathbf{T}_\ell) - \hat{I}_{SP} \right| \leq 2\Delta_{\beta, d_\ell}(n) + \frac{d_\ell \log(1 + 1/\beta^2)}{4\sqrt{n}}$$

- $\Delta_{\beta, d_\ell}(n)$ is the risk bound for estimating a single entropy term (i.e., $O(C(d_\ell)/\sqrt{n})$).

Guarantee 3: MI Estimation Risk

The absolute-error risk of the full SP MI estimator \hat{I}_{SP} (using n samples for unconditional term, and $n_x = n$ for each of n conditional terms):

$$\sup_{P_X} \mathbb{E} \left| I(\mathbf{X}; \mathbf{T}_\ell) - \hat{I}_{SP} \right| \leq 2\Delta_{\beta, d_\ell}(n) + \frac{d_\ell \log(1 + 1/\beta^2)}{4\sqrt{n}}$$

- $\Delta_{\beta, d_\ell}(n)$ is the risk bound for estimating a single entropy term (i.e., $O(C(d_\ell)/\sqrt{n})$).
- **Core Implication:** The overall MI estimation error also converges at the parametric rate $O(1/\sqrt{n})$.
- The bound depends on:
 - Layer dimension d_ℓ (through $\Delta_{\beta, d_\ell}(n)$ and the second term).
 - Noise variance β^2 (larger β can reduce the bound via both terms).
 - Number of samples n .
- The term $1/\beta^2$ relates to the Signal-to-Noise Ratio (SNR) between the signal \mathbf{S}_ℓ and noise \mathbf{Z}_ℓ .

The Link: Compression & Geometric Clustering

- Consider $I(\mathbf{X}; \mathbf{T}_\ell) = I(\mathbf{S}_\ell; \mathbf{S}_\ell + \mathbf{Z}_\ell)$.
- This is MI over an AWGN-like channel:
 - "Input symbols": deterministic pre-noise activations $\mathcal{S}_\ell = \{\mathbf{s}_{\ell, \mathbf{x}}\}$.
 - $I(\cdot)$ measures **distinguishability** of $\mathbf{s}_{\ell, \mathbf{x}}$ after adding noise \mathbf{Z}_ℓ .
- **Hypothesis: Clustering drives compression.**
 - Training \implies representations $\mathbf{s}_{\ell, \mathbf{x}}$ of inputs \mathbf{x} from the *same class* may **cluster**.
 - Closer points in $\mathcal{S}_\ell \implies$ Gaussian components in $p_{\mathbf{T}_\ell}$ overlap more.
 - More overlap \implies harder to resolve inputs \implies **reduction in $I(\mathbf{X}; \mathbf{T}_\ell)$** .
- Paper argues: "Compression" in deterministic nets via binned MI was tracking this clustering.

Key Conclusions

- $I(\mathbf{X}; \mathbf{T}_\ell)$ in *deterministic* DNNs is often ill-defined for studying representation dynamics.
- *Noisy DNN* framework allows rigorous study of $I(\mathbf{X}; \mathbf{T}_\ell)$.
- Sample Propagation (SP) estimator developed for $I(\mathbf{X}; \mathbf{T}_\ell)$ in noisy DNNs.
- Detailed theoretical guarantees for SP estimator (risk, sample complexity, bias) show $O(1/\sqrt{N})$ rate but highlight the curse of dimensionality (exponential dependence on d).
- "Compression" ($I(\mathbf{X}; \mathbf{T}_\ell) \downarrow$) in noisy nets is linked to geometric clustering of representations.

Proposed Work: Input Perturbations

Problem: How does $I(\mathbf{X}; \mathbf{T}_\ell)$ react to distributional shift $P_X \rightarrow P_{X'}$?

Perturbation model: Control the shift the input distribution to another one by some metric

Shift metrics:

- Wasserstein $W_p(P_X, P_{X'})$, Total Variation $TV(P_X, P_{X'})$, $KL(P_X \| P_{X'})$

Question:

- **Q:** Robustness of Estimation — how close is the estimated $I(\mathbf{X}'; \mathbf{T}_\ell)$ with the real information of original dataset?

Shift Metrics: KL, TV, Wasserstein

KL divergence

$$D_{\text{KL}}(P\|Q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

- Asymmetric; $D_{\text{KL}} \geq 0$ with equality iff $P=Q$.
- **Likelihood-sensitive**: large when $q(x)=0$ where $p(x)>0$.
- Useful for modeling/calibration where densities are available.

Total Variation (TV)

$$\text{TV}(P, Q) = \frac{1}{2} \int |p(x) - q(x)| dx = \sup_A |P(A) - Q(A)|$$

- Symmetric, bounded: $0 \leq \text{TV} \leq 1$.
- **Support-overlap** measure; ignores geometry of x .

Wasserstein W_p

$$W_p(P, Q) = \left(\inf_{\gamma \in \Pi(P, Q)} \mathbb{E}_{(X, Y) \sim \gamma} \|X - Y\|^p \right)^{1/p}$$

- **Geometry-aware**: depends on ground metric $\|\cdot\|$.
- Stable under small input shifts; meaningful even with disjoint supports.

Estimation Error under TV Shift

Theorem: Bounding the Estimator's Deviation under TV Shift

Let \hat{I}_{SP} be the estimator built using n samples from P_X . If the data distribution shifts to $P_{X'}$ such that $\text{TV}(P_X, P_{X'}) \leq \varepsilon$, the expected deviation is bounded by:

$$\mathbb{E} \left| I(\mathbf{X}'; \mathbf{T}_\ell) - \hat{I}_{SP} \right| \leq \underbrace{(\varepsilon \log(N_\ell - 1) + H_b(\varepsilon))}_{\text{Shift Error}} + \underbrace{\left(\frac{8cd_\ell + d_\ell \log(1 + 1/\beta^2)}{4\sqrt{n}} \right)}_{\text{Estimation Error}}$$

where the second term is the explicit risk bound from the paper.

Estimation Error under KL Shift

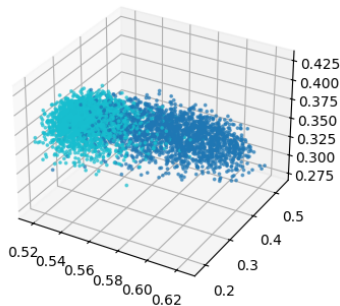
Theorem: Bounding the Estimator's Deviation under KL Shift

Let \hat{I}_{SP} be the estimator built using n samples from P_X . If the data distribution shifts to $P_{X'}$ such that $D_{KL}(P_{X'} \| P_X) \leq \varepsilon$, the deviation is bounded by:

$$\mathbb{E} \left| I(\mathbf{X}'; \mathbf{T}_\ell) - \hat{I}_{SP} \right| \leq \underbrace{\left(\sqrt{\frac{\varepsilon}{2}} \log(N_\ell - 1) + H_b \left(\sqrt{\frac{\varepsilon}{2}} \right) \right)}_{\text{Shift Error}} + \underbrace{\left(\frac{8cd_\ell + d_\ell \log(1 + 1/\beta^2)}{4\sqrt{n}} \right)}_{\text{Estimation Error}}$$

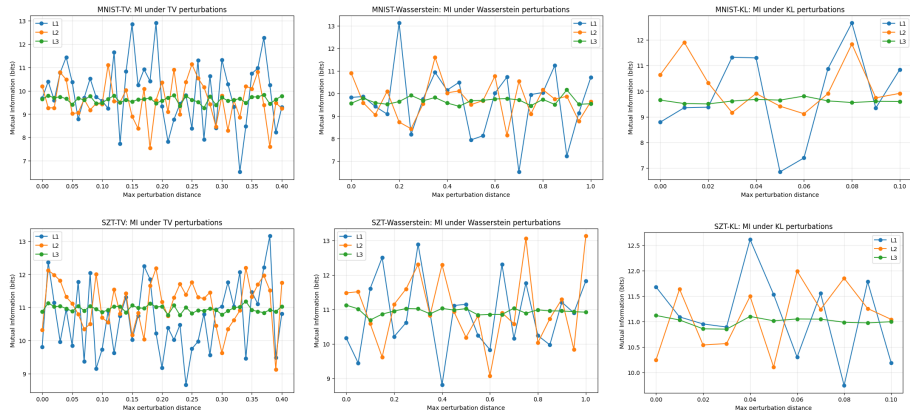
L5 Pre-noise Representation S_{L5}

L5 pre-noise S (final epoch)



- Pre-noise activations S_{L5} reveal emerging class-wise clusters as training proceeds.

MI under shift — Latest SZT & MNIST



• Each panel: $I(\mathbf{X}; \mathcal{T}_\ell)$ vs. perturbation; legend shows layers.

Key References

- [1] Goldfeld, Z., van den Berg, E., Greenewald, K., Melnyk, I., Nguyen, N., Kingsbury, B., & Polyanskiy, Y. (2019).
Estimating Information Flow in Deep Neural Networks.
- [2] Shwartz-Ziv, R., & Tishby, N. (2017).
Opening the black box of Deep Neural Networks via Information.
- [3] Tishby, N., & Zaslavsky, N. (2015).
Deep Learning and the Information Bottleneck Principle.
- [4] Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Lampinen, A., Teh, B. D., & Ganguli, S. (2018).
On the Information Bottleneck Theory of Deep Learning.
- [5] Goldfeld, Z., Greenewald, K., Polyanskiy, Y., & Weed, J. (2019).
Estimating Differential Entropy under Gaussian Convolutions.

Thank You!

Questions?