

Inductive Mutual Information Estimation: A Convex Maximum-Entropy Copula Approach

Yves-Laurent Kom Samo

KXY Technologies Inc.

✉ ml@kxy.ai 📠 @Dr_YLKS 🐦 @Dr_YLKS



March 20, 2021



Overview

- 1 Overview
- 2 Our Approach
- 3 Empirical Evaluation



Motivation

Whenever there is **structure in data**, there is **mutual information**. Yet mutual information (MI) estimation is a notoriously hard problem to solve (McAllester and Stratos (2020)).

- Mutual information (MI) is only a weak property of a joint distribution $P_{x,y}$.
- But virtually all MI estimators assume that we have enough data to fully characterize $P_{x,y}$.
- As a result, existing MI estimators are data-inefficient.



Contributions

Our estimator (MIND) is the only MI estimator that works for small and large sample sizes and/or input dimensions.

Key Properties:

- Data-efficiency and low-variance
- Consistency
- Marginal (or scale) invariance



The MIND Estimator

- Define $\mathbf{u}_z = (F_1(z_1), \dots, F_d(z_d))$, where F_i is the CDF of z_i . u_z has uniform marginals and has CDF the copula of $\mathbf{z} = (z_1, \dots, z_d)$. \mathbf{u}_z is \mathbf{z} 's copula-uniform dual representation.
- $I(\mathbf{y}; \mathbf{x}) = I(\mathbf{u}_y; \mathbf{u}_x) = h(\mathbf{u}_y) + h(\mathbf{u}_x) - h(\mathbf{u}_y, \mathbf{u}_x)$.
- Estimate $h(\mathbf{u}_z)$ as the solution to the max-ent problem:

$$\begin{cases} \max_{P \in \mathcal{D}_d} h(P) \\ \text{s.t. } E_P[\phi_m(\mathbf{u})] = \alpha_m := E_{P_{u_z}}[\phi_m(\mathbf{u})] \end{cases}, \quad (\text{A-MIND})$$

over the space \mathcal{D}_d of continuous distribution on $[0, 1]^d$.

- ϕ_m is chosen so that $E_P(\phi_m(\mathbf{u}))$ reveals association between coordinates of \mathbf{u} . E.g. $E_P(\phi_s(u, v))$ is the Spearman rank correlation between u and v when $\phi_s(u, v) = 12uv - 3$.



Solving the Maximum-Entropy Problem

- **Solution:** $h_{AM}(\mathbf{u}_z) := -\alpha_m^T \boldsymbol{\theta}^*$, with $\boldsymbol{\theta}^*$ solution to the CVX problem

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} -\alpha_m^T \boldsymbol{\theta} + \int_{[0,1]^d} e^{\boldsymbol{\theta}^T \phi_m(\mathbf{u})} d\mathbf{u}.$$

- **Copula Entropy Error:**

$$h_{AM}(\mathbf{u}_z) - h(\mathbf{u}_z) = KL(P_{\mathbf{u}_z} \| P_{AM}(\mathbf{u}_z)).$$

- **Mutual Information Error:**

$$I_{AM}(\mathbf{y}; \mathbf{x}) - I(\mathbf{y}; \mathbf{x}) = KL(P_{\mathbf{u}_y} \| P_{AM}(\mathbf{u}_y)) + KL(P_{\mathbf{u}_x} \| P_{AM}(\mathbf{u}_x)) \\ - KL(P_{\mathbf{u}_{x,y}} \| P_{AM}(\mathbf{u}_{x,y})).$$

Can be null even when m is small!



Sample Estimator Properties

- **Constraint Estimator:** $\hat{\alpha}_{m,n} = \frac{1}{n} \sum_{i=1}^n \phi_m \left(\frac{\text{rg}(\mathbf{z}_i)}{n+1} \right)$, where $(\mathbf{z}_1, \dots, \mathbf{z}_n)$ are i.i.d samples and $\text{rg}(\mathbf{z}_i)$ are coordinatewise ranks. $\hat{\alpha}_{m,n}$ is consistent and has $O(1/n)$ MSE rate.
- **Consistency:** Overall, MIND is a consistent (in m and n) estimator of $I(\mathbf{y}; \mathbf{x})$ so long as $(\phi_m)_m$ are universal approximators of continuous functions on $[0, 1]^d$ for any d (e.g. polynomials of degree m).
- **Data-Efficiency:** We may achieve perfect estimation even for a small m (Corollary 3.1) and MSE rate is $O(1/n)$.



Synthetic Benchmark

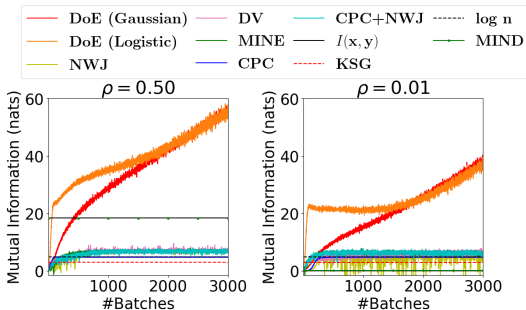


Figure: Estimation of the mutual information between two 128-dimensional vectors $\mathbf{x} = (x_1, \dots, x_d)$ and $\mathbf{y} = (y_1, \dots, y_d)$ from a draw of 1000 i.i.d. samples. (x_i, y_i) are i.i.d. Gaussians with mean zero, unit marginal variance, and correlation ρ .



Real-Life Application

$I(\mathbf{y}; \mathbf{x})$	\bar{R}^2	$RMSE$	d	n
1.50	0.95	18,531	80	1460

Table: Mutual information and highest performances achievable in the 'House Prices: Advanced Regression Techniques' Kaggle challenge.

$$\bar{R}^2(P_{\mathbf{y},\mathbf{x}}) = 1 - e^{-2I(\mathbf{y};\mathbf{x})}$$

$$RMSE(P_{\mathbf{y},\mathbf{x}}) = e^{-I(\mathbf{y};\mathbf{x})} \sqrt{\text{Var}(\mathbf{y})}$$



Reference

Reference

- Csiszar, I. (1975). “I-divergence geometry of probability distributions and minimization problems.” In The Annals of Probability (Ann. Probab.).
- Belghazi, M. I. et al. (2018). “Mutual Information Neural Estimation.” In International Conference on Machine Learning (ICML).
- McAllester, D. Stratos, K. (2020). “Formal limitations on the measurement of mutual information.” In International Conference on Artificial Intelligence and Statistics (AISTATS).

Code: <https://github.com/kxytechnologies/kxy-python>

