

Brownbag on Neural Mutual Information estimation

Demystifying the "Correlation of the 21st Century." Terry Speed.

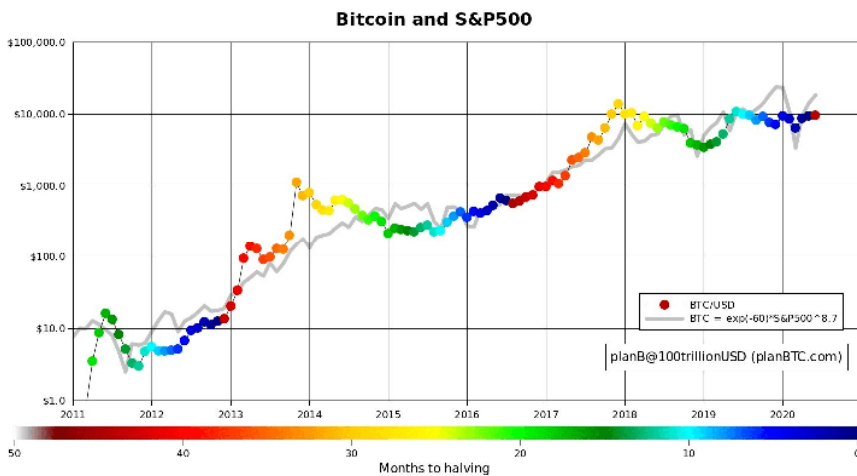
Matthias Hermann, HTWG-Konstanz, Institute for Optical Systems

9.11.2021

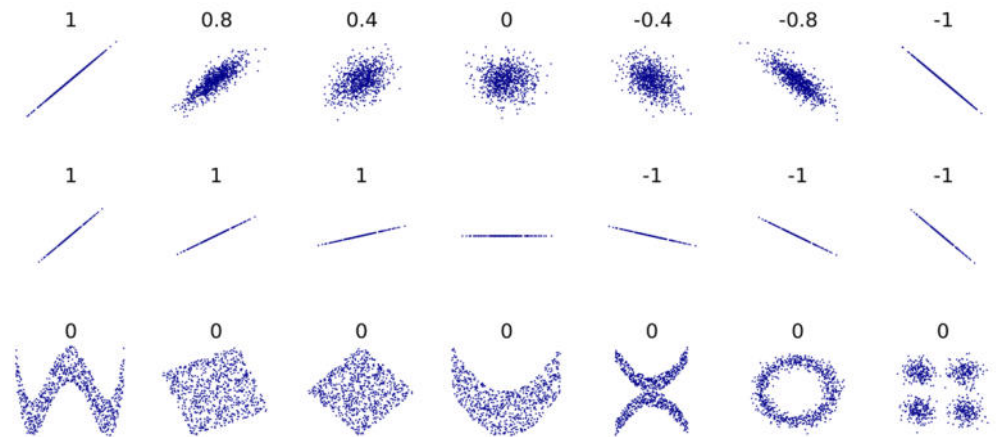
Agenda

- What is mutual information?
- Estimating mutual information
 - Histogram
 - Kruskov
- Mutual Information Neural estimation (MINE)
- Application to sensor registration

Linear dependence



Bitcoin X and S&P500 Y is highly correlated.



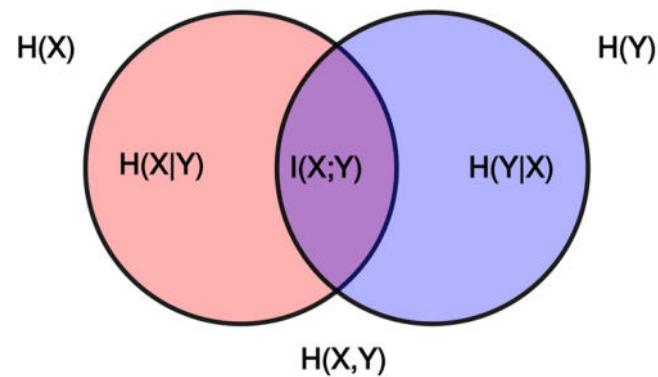
Pearson correlation coefficient measures linear relationships between X and Y. But fails with slope and non-linear relationships.

Non-linear dependence

- Mutual Information is given by

$$I(X, Y) = H[X] - H[X|Y].$$

- Specifies “how much (in bits) do we know about **X** given **Y**?”.



Capturing non-linear dependencies

Given random variables \mathbf{X} and \mathbf{Y} , and function f

$$Y = f(X) + \sigma\epsilon$$

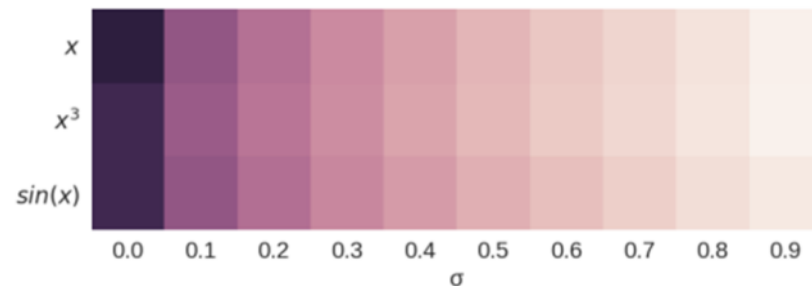
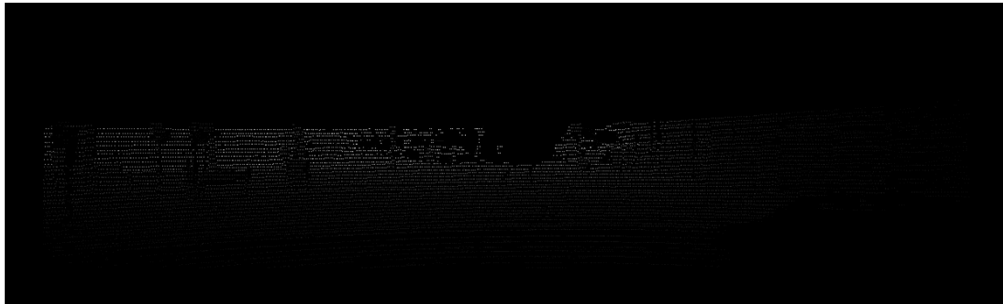


Fig. 1. Mutual information $I(\mathbf{X}, \mathbf{Y})$ measures the dependency of \mathbf{X} and \mathbf{Y} and is invariant to the deterministic nonlinear transformation f (“equitability”).

MI is defined for arbitrary variables

- Given random variables **ProjectedLidarPoints** and **CameraImage**

$$I(\text{ProjectedLidarPoints}, \text{CameraImage}) \approx 20.89 \text{ bits.}$$



Estimating mutual information

- Definition

$$I(X, Y) = KL(P_{XY} || P_X \otimes P_Y),$$

where P_{XY} is the joint distribution and $P_X \otimes P_Y$ is the product of their marginals.

$$I(X, Y) = \underbrace{H_{P_{XY}}[P_X \otimes P_Y]}_{\text{cross-entropy}} - \underbrace{H[P_{XY}]}_{\text{joint-entropy}} = \underbrace{H[P_X]}_{\text{entropy}} + \underbrace{H[P_Y]}_{\text{entropy}} - \underbrace{H[P_{XY}]}_{\text{joint-entropy}}$$

- Naive “Histogram” approach:

$$I(X, Y) \approx I_{\text{binned}}(X, Y) = \sum_{y \in Y} \sum_{x \in X} p_{(X, Y)}(x, y) \log \frac{p_{(X, Y)}(x, y)}{p_X(x)p_Y(y)}$$

- Can be problematic with empty bins!

Estimating mutual information

- Improving estimation by k-nearest neighbor statistics [2]:

$$I(X, Y) = \underbrace{H[P_X]}_{\text{entropy}} + \underbrace{H[P_Y]}_{\text{entropy}} - \underbrace{H[P_{XY}]}_{\text{joint-entropy}}$$

$$I(X, Y) \approx I_{knn}(X, Y) = \hat{H}[P_X] + \hat{H}[P_Y] - \hat{H}[P_{XY}]$$

$$\hat{H}[X] = -\underbrace{\psi(k)}_{\text{digamma}} + \underbrace{\psi(N)}_{\text{digamma}} + \log \underbrace{c_d}_{\text{volume } d\text{-ball}} + \frac{d}{N} \sum_i^N \log \underbrace{\epsilon(i)}_{\text{distance to } k\text{-th neighbor}}$$

Estimating mutual information

- K-nearest neighbor method is problematic for estimating joint entropy

$$\hat{H}[P_{XY}]$$

as choosing the same \mathbf{k} and computing

$$I_{knn}(X, Y) = \hat{H}[P_X] + \hat{H}[P_Y] - \hat{H}[P_{XY}]$$

would effectively use different scales for the joint and marginal space.

- Kraskov's method [2] corrects for that by choosing \mathbf{k} dynamically!

Estimating mutual information in high dimensions

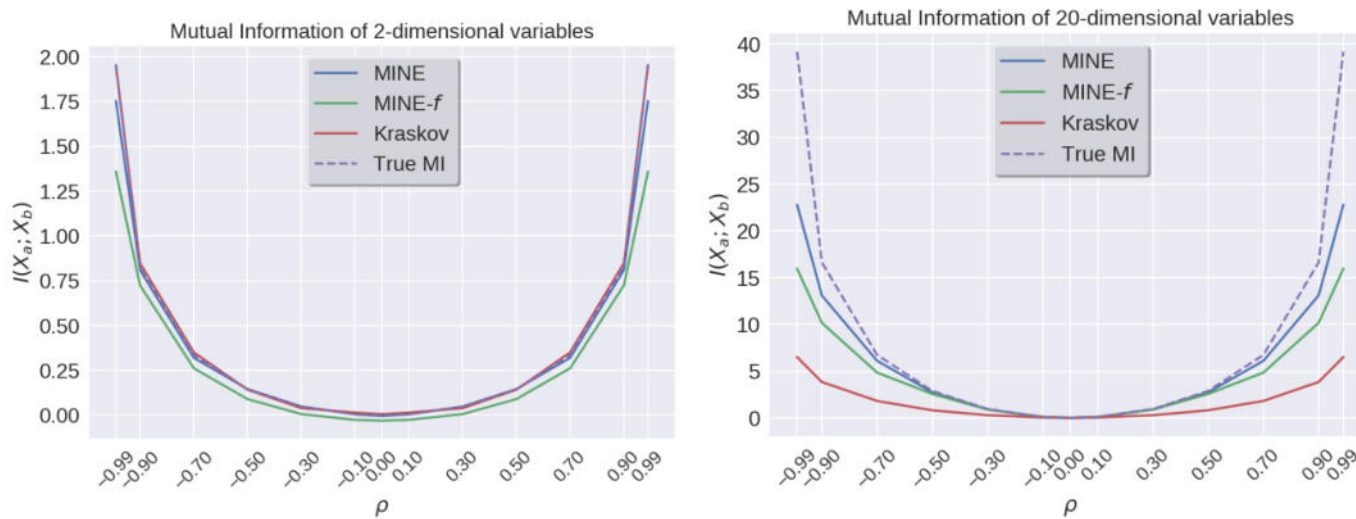


Fig. 2. Mutual information between two multivariate Gaussians with component-wise correlation $\rho \in (-1, 1)$ [1]. Kraskov's estimator [2] underestimates the true MI in high dimensions.

MINE - Neural mutual information estimation

- Maximizing the Donsker-Varadhan (DV) lower bound of the KL divergence

$$I(X, Y) = KL(P_{XY} || P_X \otimes P_Y)$$

$$\geq \sup_{f_\theta} DV(J, M; f_\theta) = E_{x \sim J}[f_\theta(x)] - \log(E_{y \sim M}[e^{f_\theta(y)}]),$$

with $J = P_{XY}$ and $M = P_X \otimes P_Y$ and function space $f \in F$.

- In the case of MINE [1], f_θ is a neural network and $\sup_{f_\theta} DV(\cdot; f_\theta)$ is computed by standard gradient ascent!
- The Auxiliary dataset $M = (X, Y^*)$ is constructed by sampling y^* without replacement (shuffling).

MINE - Neural mutual information estimation

- Maximizing the Donsker-Varadhan (DV) lower bound of the KL divergence

$$I(X, Y) = KL(P_{XY} || P_X \otimes P_Y)$$

$$\geq \sup_{f_\theta} DV(J, M; f_\theta) = E_{(x,y) \sim J} [f_\theta((x, y))] - \log(E_{(x,y^*) \sim M} [e^{f_\theta((x,y^*))}]),$$

with $J = P_{XY}$ and $M = P_X \otimes P_Y$ and function space $f \in F$.

- In the case of MINE [1], f_θ is a neural network and $\sup_{f_\theta} DV(\cdot; f_\theta)$ is computed by back-prob and standard gradient ascent!

MINE – a proxy “classification” problem

- We need the auxiliary datasets $J = P_{XY}$ and $M = P_X \otimes P_Y$
- The dataset $J = P_{XY}$ is generated by concatenating the given training examples (\mathbf{X}, \mathbf{Y}) .
- The dataset $M = (X, Y^*)$ is constructed by sampling y^* without replacement (shuffling).

MINE - Application

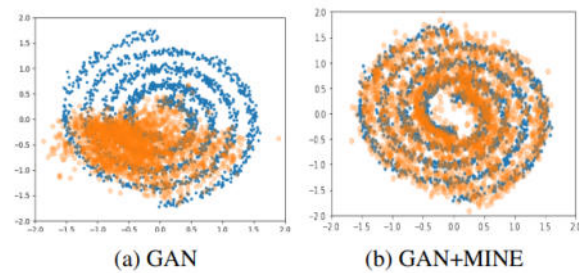


Figure 3. The generator of the GAN model without mutual information maximization after 5000 iterations suffers from mode collapse (has poor coverage of the target dataset) compared to GAN+MINE on the spiral experiment.

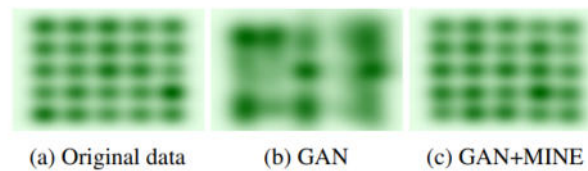


Figure 4. Kernel density estimate (KDE) plots for GAN+MINE samples and GAN samples on 25 Gaussians dataset.

Lidar-to-Camera registration I(Lidar, Camera)

- Problem definition:
 - Variable \mathbf{X} becomes the Lidar data
 - Variable \mathbf{Y} becomes the camera data
- Find unknown parameters rotation \mathbf{R} and translation \mathbf{t} , such that $\mathbf{C}_{coord} = [\mathbf{R}\mathbf{t}]\mathbf{L}_{coord}$, by Maximizing $I(\mathbf{X}, \mathbf{Y})$.

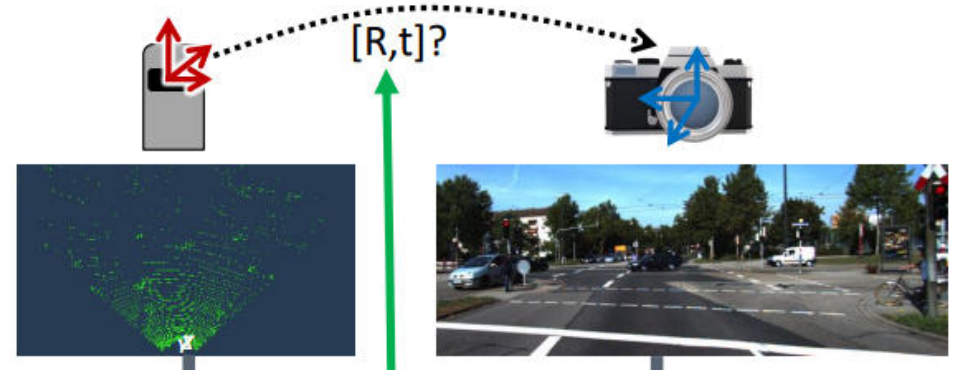


Fig. 3. Registering Lidar \mathbf{X} to camera \mathbf{Y} means finding the extrinsic calibration parameters \mathbf{R} and \mathbf{t} , where \mathbf{R} is a 3d rotation matrix and \mathbf{t} is a 3d translation vector.

Projecting Lidar data into image



Fig. 4. Visualizing projected Lidar data points in the image plane of the camera: Unregistered (left) and registered (right).

Thanks.

Matthias Hermann

Hochschule Konstanz
Institute for Optical Systems

matthias.hermann@htwg-konstanz.de

www.ios.htwg-konstanz.de



..novelty or surprise?

Berlyne D. E. (1960). Conflict, Arousal, and Curiosity.

Literature

- [1] Belghazi, M. I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Courville, A., & Hjelm, D. (2018, July). Mutual information neural estimation. In *International Conference on Machine Learning* (pp. 531-540). PMLR.
- [2] Kraskov, A., Stögbauer, H., & Grassberger, P. (2004). Estimating mutual information. *Physical review E*, 69(6), 066138.