## Introduction to Machine Learning

Information Theory Mutual Information under Reparametrization (Deep-Dive) × 0 0 × × ×



## Learning goals

• Understand why MI is invariant under certain reparametrizations

## **MUTUAL INFORMATION PROPERTIES**

• MI is invariant w.r.t. injective reparametrizations that are in  $\mathcal{C}^1$  :

Let  $f, g : \mathbb{R}^d \to \mathbb{R}^d \in \mathcal{C}^1$  be injective transformations and X, Y be continuous random variables in  $\mathbb{R}^d$  then by the change of variables the joint and marginal densities of  $\tilde{X} = f(X), \tilde{Y} = g(Y)$ 

$$egin{aligned} & ilde{p}( ilde{x}, ilde{y}) = p(f^{-1}( ilde{x}),g^{-1}( ilde{y})) \cdot |J_{f^{-1}}( ilde{x})| \cdot |J_{g^{-1}}( ilde{y})|, \ & ilde{p}( ilde{x}) = p(f^{-1}( ilde{x})) \cdot |J_{f^{-1}}( ilde{x})|, \quad & ilde{p}( ilde{y}) = p(g^{-1}( ilde{y})) \cdot |J_{g^{-1}}( ilde{y})|. \end{aligned}$$

where p(x, y) is the joint density of X and Y and p(x), p(y) are the respective marginal densities. (*J* denotes the Jacobian)

With this, it follows that

$$I( ilde{X}; ilde{Y}) = \int ilde{
ho}( ilde{x}, ilde{y}) \log\left(rac{ ilde{
ho}( ilde{x}, ilde{y})}{ ilde{
ho}( ilde{x}) ilde{
ho}( ilde{y})}
ight) d ilde{x}d ilde{y} = *$$

× × ×

## **MUTUAL INFORMATION PROPERTIES**

$$\begin{split} * &= \int p(f^{-1}(\tilde{x}), g^{-1}(\tilde{y})) \cdot |J_{f^{-1}}(\tilde{x})| \cdot |J_{g^{-1}}(\tilde{y})| \\ &\quad \cdot \log \left( \frac{p(f^{-1}(\tilde{x}), g^{-1}(\tilde{y})) \cdot |J_{f^{-1}}(\tilde{x})| \cdot |J_{g^{-1}}(\tilde{y})|}{p(f^{-1}(\tilde{x}))|J_{f^{-1}}(\tilde{x})| \cdot p(g^{-1}(\tilde{y}))|J_{g^{-1}}(\tilde{y})|} \right) d\tilde{x} d\tilde{y} \\ &= \int p(f^{-1}(f(x)), g^{-1}(g(y))) \cdot |J_{f^{-1}}(f(x))| \cdot |J_{g^{-1}}(g(y))| \\ &\quad \cdot \log \left( \frac{p(f^{-1}(f(x)), g^{-1}(g(y)))}{p(f^{-1}(f(x)))p(g^{-1}(g(y)))} \right) |J_{f}(x)| \cdot |J_{g}(y)| dx dy \\ &= \int p(x, y) \cdot |J_{f^{-1}}(f(x))J_{f}(x)| \cdot |J_{g^{-1}}(g(y))J_{g}(y)| \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy \\ &= \int p(x, y) \cdot \log \left( \frac{p(x, y)}{p(x)p(y)} \right) dx dy = I(X; Y). \end{split}$$

(The fourth equality holds by the inverse function theorem)

X

XX