

Lecture 10 Mutual information

$$I(X;Y) = D_{KL}(p(x,y) \parallel p(x)p(y)) =$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}.$$

$$I(X;Y) \geq 0 ; = 0 \text{ iff } p(x,y) = p(x)p(y).$$

Note that

$$\begin{aligned} I(X;Y) &= \sum_{x,y} p(x|y) p(y) \log p(x|y) - \\ &\quad - \sum_{x,y} p(x,y) \log p(x) = \\ &= \underbrace{\sum_y p(y) \sum_x p(x|y) \log p(x|y)}_{-H(X|Y)} + H(X). \end{aligned}$$

$$\text{Likewise, } I(X;Y) = H(Y) - H(Y|X).$$

Thus, $I(X;Y)$ quantifies the reduction in uncertainty about X after observing Y , and vice versa.

Recall that $H(X|Y) \leq H(X)$, consistent with $0 \leq I(X;Y) = H(X) - H(X|Y)$.

One can show that

$$I(X;Y) = H(X) + H(Y) - \underbrace{H(X,Y)}_{H(X) + H(Y|X)} = \\ = H(Y) - H(Y|X).$$

Likewise, $I(X;Y) = \underbrace{H(X,Y)}_{H(X) + H(Y|X)} - H(X|Y) - H(Y|X) = \\ = H(X) - H(X|Y).$

Conditional mutual information

$$I(X;Y|Z) \equiv \sum_z p(z) \left[\sum_{x,y} p(x,y|z) \log \frac{p(x,y|z)}{p(x|z)p(y|z)} \right] \\ \stackrel{\textcircled{E}}{=} \sum_z p(z) \left[\sum_x p(x|z) \log p(x|z) \right] - \sum_z p(z) \left[\sum_y p(y|z) \times \right. \\ \left. \times \log p(y|z) \right] + \sum_z p(z) \left[\sum_{x,y} p(x,y|z) \log p(x,y|z) \right] = \\ = H(X|Z) + H(Y|Z) - \underbrace{H(X,Y|Z)}_{H(Y|Z) + H(X|Y,Z)} = H(X|Z) - \underbrace{H(X|Y,Z)}_{H(X,Y,Z) - H(Z,Y)}, \\ \text{see below}$$

$$H(X|Y,Z) = H(X,Y,Z) - H(Y) - H(Z|Y) = \\ = H(X,Y,Z) - \cancel{H(Y)} - H(Z,Y) + \cancel{H(Y)}$$

$$\textcircled{=} H(Y, Z) - \underbrace{H(Y) + H(Z)}_{-I(Y; Z)} + H(Y) + H(X, Z) - \underbrace{H(X, Y, Z)}_{I(Y; X, Z)} \textcircled{=}$$

$$\textcircled{=} I(Y; X, Z) - I(Y; Z) \quad \equiv$$

Thus, conditional MI is the extra info that $X|Z$ tells us about Y , excluding what Z alone has told us about Y .

Finally, we can write

$$\textcircled{=} I(Y; X, Z) = I(X; Y|Z) + I(Y; Z)$$

$X \leftrightarrow Y$ yields

$$I(Y, Z; X) = I(\underset{z_2}{\cancel{z}}; \underset{z_1}{\cancel{z}}) + I(Y; X|Z).$$

$$\begin{matrix} & z_1 \\ & \downarrow \\ z_2 & z_1 \end{matrix} \quad \begin{matrix} X \\ \downarrow \\ z_2 \\ \downarrow \\ z_1 \end{matrix}$$

~~Iff~~ $I(z_1, z_2; X) = I(z_1; X) + I(z_2; X|z_1).$

This generalizes to

$$I(z_1, \dots, z_n; X) = \sum_{n=1}^N I(z_n; X|z_1, \dots, z_{n-1}). \quad \equiv$$

[chain rule for mutual information]

generalized correlation coefficient

Consider $\begin{pmatrix} x \\ y \end{pmatrix} \sim N(\vec{0}, \sigma^2 \begin{pmatrix} 1 & p \\ p & 1 \end{pmatrix})$ jointly gaussian ($d=2$)

Recall that if $X \sim N(\mu, \sigma^2)$,

$$h(x) = \frac{1}{2} + \frac{1}{2} \log(2\pi\sigma^2).$$

Likewise, if $\vec{X} \sim N(\vec{\mu}, \Sigma)$ d -dimensional gaussian

$$h(\vec{X}) = \frac{d}{2} + \frac{d}{2} \log(2\pi) + \frac{1}{2} \underbrace{\log |\Sigma|}_{\det(\Sigma)}.$$

In our case,

$$x \sim N(0, \sigma^2) \text{ and } y \sim N(0, \sigma^2), \text{ s.t.}$$

$$h(x) = h(y) = \frac{1}{2} + \frac{1}{2} \log(2\pi\sigma^2).$$

Moreover,

$$h(x, y) = 1 + \log(2\pi) + \frac{1}{2} \log[\sigma^4(1-p^2)].$$

$$\begin{aligned} \text{Then } I(x; y) &= h(x) + h(y) - h(x, y) = \\ &= \log(2\pi\sigma^2) - \log(2\pi) - \underbrace{\frac{1}{2} \log \sigma^4}_{\log \sigma^2} - \frac{1}{2} \log(1-p^2) \end{aligned} \quad \textcircled{11}$$

$$\therefore \frac{1}{2} \log \frac{1}{1-p^2}.$$

Special cases: $\stackrel{(a)}{=} p=1 \Rightarrow X=Y \Rightarrow I(x; y)=\infty$
Observing Y tells us X exactly

(b) $p=0 \Rightarrow X \perp Y \Rightarrow I(X;Y)=0$
 Observing Y tells us nothing
 about X

(c) $p=-1 \Rightarrow X = -Y \Rightarrow I(X;Y)=\infty$
 Observing Y tells us X exactly,
 as in (a)

Normalized MI

Recall that

$$I(X;Y) = \begin{cases} H(X) - H(X|Y) \leq H(X) \\ H(Y) - H(Y|X) \leq H(Y) \end{cases}$$

So, $0 \leq I(X;Y) \leq \min\{H(X), H(Y)\}$.

Thus, we can define

$$NMI(X,Y) = \frac{I(X;Y)}{\min(H(X), H(Y))}$$

clearly, $0 \leq NMI(X,Y) \leq 1$.

If $NMI(X,Y)=0 \Rightarrow I(X;Y)=0 \Rightarrow X \perp Y$.

If $NMI(X,Y)=1$ and $H(X) < H(Y)$,
 [$H(Y) < H(X)$ argued]
 in the same way

we have $I(X;Y) = H(X) - H(X|Y) = H(X) \Rightarrow$
 $\Rightarrow H(X|Y) = 0$, Y determines X
 completely

NMI may be tricky to compute for
 continuous variables \Rightarrow dependence on bin
 sizes.

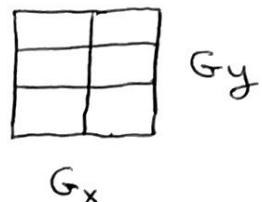
Maximal information coefficient (MIC)

Define $\text{MIC}(X, Y) = \max_G \frac{|I(X; Y)|_G}{\log G_{\min}}$,

where $G = \text{set of 2D grids}$, and

$$G_{\min} = \min \{G_x, G_y\}$$

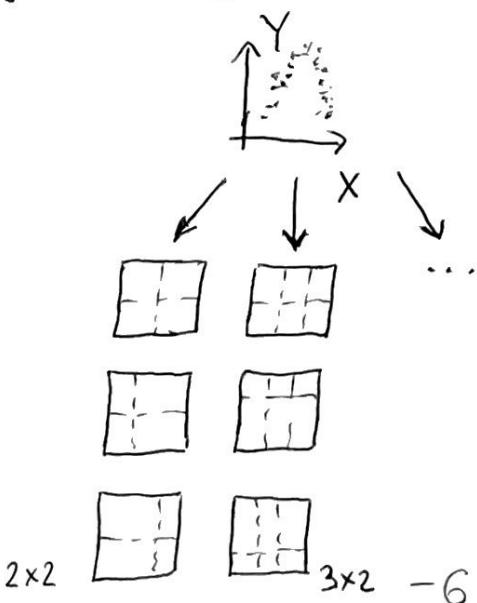
$G_x \times G_y$ grid



$\log G_{\min}$ is the max entropy value for $H(X)$ or $H(Y)$ [the min of the two, actually] \Rightarrow this ensures $0 \leq \text{MIC}(X, Y) \leq 1$.

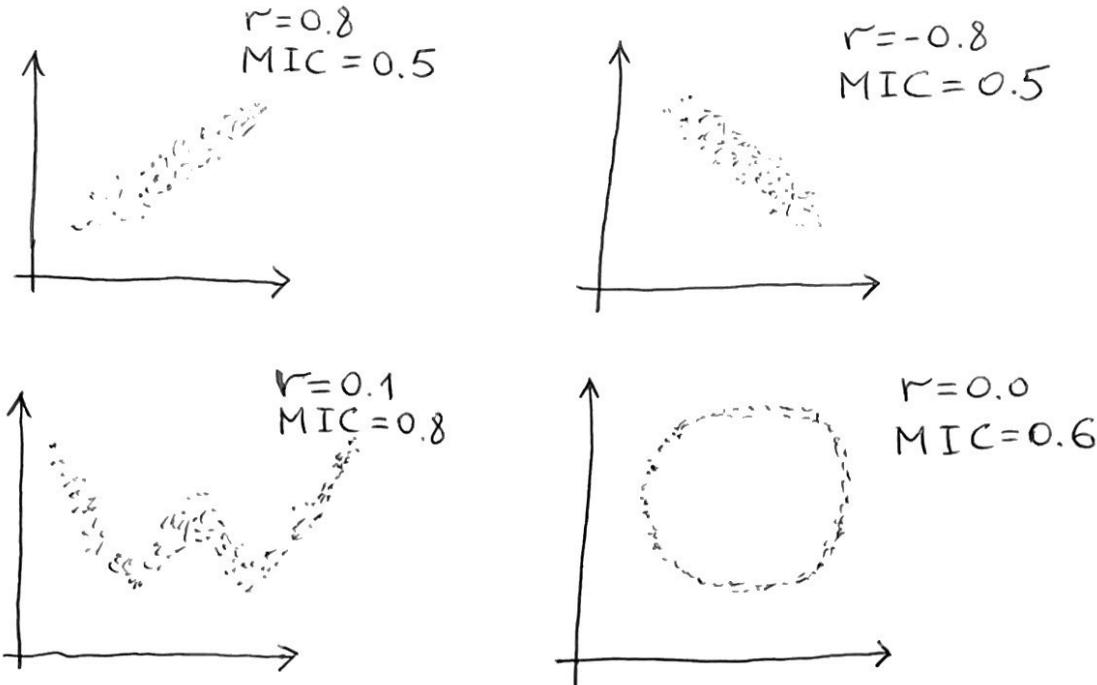
Finally, $|I(X; Y)|_G$ is $I(X; Y)$ computed on a given grid G .

Empirically, $G_x G_y \leq N^{0.6}$.
↑ sample size



MIC is the max value on all these grids

Ex.



Note that MIC is not restricted to linear relations $\Rightarrow \text{MIC} \approx 1.0$ for any non-linear relationship.

Data processing inequality (DPI)
 Suppose $X \rightarrow Y \rightarrow Z$ form a Markov chain of Y
 chain : $X \perp Z | Y$.
[potentially noisy] observe'n of a f'n
 of Y
noisy observe'n
 of a function
 of X

Chain rule: $I(X; Y, Z) = I(X; Z) + \underbrace{I(X; Y|Z)}_{\geq 0} =$

$$= I(X; Y) + \underbrace{I(X; Z|Y)}_{\leq 0 \text{ since } \cancel{X \perp Z | Y}}, \text{ or}$$

≤ 0 since ~~$X \perp Z | Y$~~

$$\underbrace{I(X; Z)}_{\leq I(X; Y)} \leq I(X; Y).$$

Similarly, $I(X; Z) \leq I(Y; Z)$

We cannot increase the amount of information we have about unknown X by further processing.

Sufficient statistics

Consider $\vec{\theta} \rightarrow D \rightarrow S(D)$ some function of the data

\uparrow \nwarrow
 unknown observed
 model data
 prms
 (hidden vars)

DPI gives: $I(\vec{\theta}; S(D)) \leq I(\vec{\theta}; D)$.

If $I(\vec{\theta}; S(D)) = I(\vec{\theta}; D)$,

$S(D)$ is a sufficient statistic of D ,
for the purpose of inferring $\vec{\theta}$.

Define minimal sufficient statistic

as the one that maximally compresses D
without hurting inference of $\vec{\theta}$.

Ex. (a) N Bernoulli trials:

$S(D) = \{N, N_1\} \rightarrow$ to find Bernoulli $\vec{\theta}$

\uparrow # pos. outcomes
 no need to keep the entire sequence
 of events

(b) N samples from a 1D gaussian:

$$N(\mu, \sigma^2)$$

To find $\vec{\theta} = \{\mu, \sigma^2\}$, just need
 $S(D) = \{N, \text{mean of data}, \text{var. of data}\}$.

Fano's inequality

Consider $Y \rightarrow X \rightarrow \hat{Y}$

\uparrow \uparrow \nwarrow
 true features $f(X)$, estimator
 class in the of Y (classifier)
 labels classifier

Define $E: \hat{Y} \neq Y$ misclassification event

$$P_e = P(E=1) \quad \text{prob. of error}$$

chain rule for entropy:

$$\begin{aligned} H(E, Y | \hat{Y}) &= H(Y | \hat{Y}) + \underbrace{H(E | Y, \hat{Y})}_{=0}, \quad E \text{ completely} \\ &\quad \text{defined by} \\ &= H(E | \hat{Y}) + H(Y | E, \hat{Y}) \quad Y \text{ and } \hat{Y} \end{aligned}$$

Now, $H(E | \hat{Y}) \leq H(E)$ conditioning reduces entropy

$$\begin{aligned} H(Y | E, \hat{Y}) &= \underbrace{P(E=0)}_{1-P_e} \underbrace{H(Y | \hat{Y}, E=0)}_{=0, \text{ no error}} + \\ &+ \underbrace{P(E=1)}_{P_e} \underbrace{H(Y | \hat{Y}, E=1)}_{\leq \log K} \leq P_e \log K. \end{aligned}$$

classes

Finally, DPI gives $I(Y; \hat{Y}) \leq I(Y; X)$, or

$$H(Y) - H(Y | \hat{Y}) \quad " \quad H(Y) - H(Y | X)$$

$$H(Y|X) \leq H(Y|\hat{Y}) \leq \underbrace{H(E)}_{\leq 1, \text{ entropy of } \uparrow \text{ a Bernoulli rv if } \log \rightarrow \log_2} + P_e \log K.$$

Thus, $H(Y|X) \leq 1 + P_e \log K$, or

$$P_e \geq \frac{H(Y|X) - 1}{\log K} \quad \text{Fano's inequality}$$

We want to minimize $H(Y|X)$ in order to ~~maximize~~ minimize the RHS \Rightarrow
 \Rightarrow maximize $I(Y; X)$.

In other words, pick features that have high MI with class labels Y.

