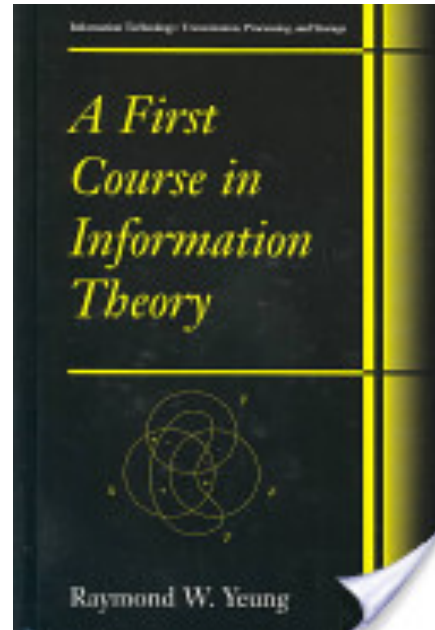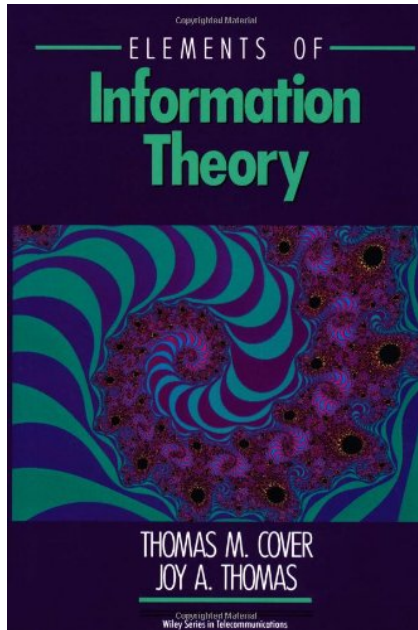# Analysis and processing of information

## S. Voloshynovskiy

# Recall of Information Theory

# Recommended books

# Course Outline

- Enropy, conditional entropy, joint entropy

- KLD

- Mutual information

- Particularites of IT measures for continuos random variables

**Chain rule for probability: joint probility** $p(x_1, x_2, ..., x_n)$

$$p(x_1, x_2, ..., x_n) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) ... p(x_n | x_{n-1}, ..., x_2, x_1) =$$
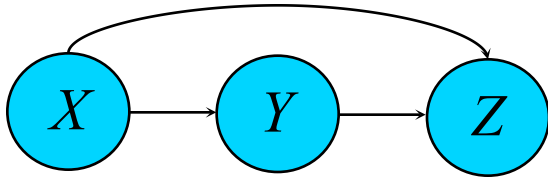$$= \prod_{i=1}^{n} p(x_i | x_{i-1}, ..., x_2, x_1)$$

$$p(x_1, x_2) = p(x_1) p(x_2 | x_1) \xrightarrow{\textbf{independence}} p(x_1, x_2) = p(x_1) p(x_2)$$

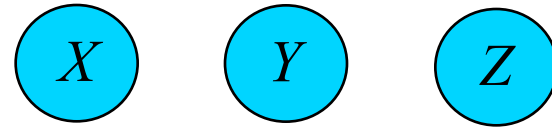$$p(x_1, x_2, ..., x_n) = p(x_1) p(x_2) p(x_3) ... p(x_n) = \prod_{i=1}^{n} p(x_i)$$

**For the independent R.V.s** $\{X_i\}$

$$p(x_1, x_2, ..., x_n) = p(x) p(x) p(x) ... p(x) = \prod_{i=1}^{n} p(x) = (p(x))^n$$

# Recall of probability



$$p(x,y,z) = p(x)p(y|x)p(z|x,y)$$

$$p(x,y,z) = p(x)p(y)p(z)$$

$$p(x,y,z) = p(x)p(y|x)p(z|y)$$

$$p(x,y,z) = p(z)p(y|z)p(x|y)$$

$$p(x,y,z) = p(z|x,y)p(x)p(y)$$

$$p(x,y,z) = p(z)p(x|z)p(y|z)$$

# Entropy

Recall from: Elements of Information Theory

**Definition (entropy):** *Entropy* of discrete r.v. $X \in \mathcal{X}$ $\quad \mathcal{X} = \{x_1, x_2, ..., x_N\}$

$$H(X) = -\sum_{x \in \mathcal{X}} p_X(x) \log_2 p_X(x) = -E_{p_X}\left[\log p_X(x)\right] \quad \text{(bits)}$$

$$E_{p_X}[x] = \sum_{i=1}^{N} x_i p(x_i)$$

$$H(X) = -\sum_{i=1}^{n} \Pr\{X = x_i\} \log_2 \Pr\{X = x_i\}$$

# Entropie: discrete random variable

**Discrete random variable:** $X \in \{0,1\}$

$\Pr\{X = 0\} = p$

$\Pr\{X = 1\} = 1 - p$

$$H(X) = -p \log_2 p - (1-p)\log_2(1-p) := H(p) \text{ or } H_2(p)$$

The function $H(p)$ is:

- Symmetric wrt $p = 0.5$ ;
- Maximum at $p = 0.5 \left(H(p) = 1\right)$.

The entropy is maximal, if the symbols are equilikely.

# Joint entropy

**Definition (joint entropy):** *Joint entropy* of two discrete random variables $X$ and $Y$ is defined by:

$$H(X,Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log_2 p_{X,Y}(x,y) = -E_{p_{X,Y}}\left[\log p_{X,Y}(x,y)\right]$$

$$\mathcal{X} = \left\{x_1, x_2, ..., x_n\right\} \quad \mathcal{Y} = \left\{y_1, y_2, ..., y_m\right\}$$

$$H(X,Y) = H(Y,X)$$

# Conditional entropy

**Definition (conditional entropy):** *Conditional entropy* of r.v. $X$ given $Y$ is defined by:

$$H\left(X\middle|Y\right) = \sum_{y\in\mathcal{Y}} p_Y\left(y\right) H\left(X\middle|Y=y\right) = -\sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} p_{X,Y}\left(x,y\right)\log p_{X|Y}\left(x\middle|y\right)$$

$$\mathcal{X} = \left\{x_1, x_2, ..., x_n\right\} \quad \mathcal{Y} = \left\{y_1, y_2, ..., y_m\right\}$$

$$H\left(X\middle|Y=y\right) = -\sum_{i=1}^{n} p_{X|Y}\left(x_i\middle|y\right)\log p_{X|Y}\left(x_i\middle|y\right)$$

$$H\left(X\middle|Y\right) \neq H\left(Y\middle|X\right)$$

# Some properties of joint entropy

**Chain rule for 2 R.V.s**

$$H(X,Y) = H(X) + H(Y|X)$$

**Proof:**

$$-\log p_{X,Y}(x,y) = -\log p_{Y|X}(y|x) p_X(x) = -\log p_{Y|X}(y|x) - \log p_X(x).$$

$$H(X,Y) = E_{p_{X,Y}}\left[-\log p_{X,Y}(x,y)\right] = E_{p_{X,Y}}\left[-\log p_X(x)\right] + E_{p_{X,Y}}\left[-\log p_{Y|X}(y|x)\right] =$$
$$= H(X) + H(Y|X).$$

# Some properties of joint entropy

**General chain rule for entropy**

$$H\left(X_1, X_2, ..., X_n\right) = \sum_{i=1}^{n} H(X_i \mid X_{i-1}, ..., X_1)$$

**Proof:**

**=, if independent**

$$H\left(X_1, X_2\right) = H\left(X_1\right) + H\left(X_2 \mid X_1\right) \implies H\left(X_1, X_2\right) \leq H\left(X_1\right) + H\left(X_2\right)$$

$$H\left(X_1, X_2, X_3\right) = H\left(X_1\right) + H\left(X_2, X_3 \mid X_1\right) = H\left(X_1\right) + H\left(X_2 \mid X_1\right) + H\left(X_3 \mid X_2, X_1\right)$$

$$H\left(X_1, X_2, ..., X_n\right) = H\left(X_1\right) + H\left(X_2 \mid X_1\right) + ... + H\left(X_n \mid X_{n-1}, ..., X_1\right) =$$

$$= \sum_{i=1}^{n} H(X_i \mid X_{i-1}, ..., X_1)$$

# Conditional entropy

$$H\left(X_1, X_2, ..., X_n\right) = \sum_{i=1}^{n} H(X_i \underbrace{\left| X_{i-1}, ..., X_1\right.}_{})$$

# Conditional entropy

$$H\left(X_1, X_2, ..., X_n\right) \le \sum_{i=1}^{n} H(X_i)$$

**equality for independent** $\left\{X_i\right\}$

**Proof:**

**For independent** $\left\{X_i\right\}$

$$H\left(X_2 \middle| X_1\right) = H\left(X_2\right)$$

$$H\left(X_3 \middle| X_2, X_1\right) = H\left(X_3\right)$$

$$H(X_i \middle| X_{i-1}, ..., X_1) = H(X_i)$$

$$H\left(X_1, X_2, ..., X_n\right) = H\left(X_1\right) + H\left(X_2 \middle| X_1\right) + ... + H\left(X_n \middle| X_{n-1}, ..., X_1\right) =$$

$$H\left(X_1, X_2, ..., X_n\right) = H\left(X_1\right) + H\left(X_2\right) + ... + H\left(X_n\right) = \sum_{i=1}^{n} H\left(X_i\right)$$

# Properties of entropy

$$H(X) \leq H(X,Y) \leq H(X,Y,Z) \leq \ldots$$

$$H(X) = H(X,Y) \Leftrightarrow Y = f(X) \qquad H(X,Y) = H(X) + \underbrace{H(Y|X)}_{0}$$

$$H(X,Y) = H(X,Y,Z) \Leftrightarrow Z = f(X,Y)$$



$H(X)$      $H(X,Y)$      $H(X,Y,Z)$

# Properties of entropy

$$H(X) \geq H(X|Y) \geq H(X|Y,Z) \geq \ldots$$

$$H(X) = H(X|Y) \Leftrightarrow X \perp Y$$
$$H(X|Y) = H(X|Y,Z) \Leftrightarrow X \perp Z|Y \ldots$$

# Relative entropy

**Definition (relative entropy):** *Relative entropy or Kullback-Leibler distance* **between pmfs** $p(x)$ **and** $q(x)$**:**

$$D(p\|q) = \sum_{x \in \mathcal{X}} p(x) \log_2 \frac{p(x)}{q(x)} = E_p\left[\log_2 \frac{p(x)}{q(x)}\right]$$

with the conventions:

$$0 \log_2 \frac{0}{q(x)} = 0, \; \forall q(x)$$

$$p(x) \log_2 \frac{p(x)}{0} = +\infty, \; \forall p(x) > 0.$$

$$D(p\|q) \neq D(q\|p)$$

# Relative entropy

**Example:**

$$x \in \mathcal{X} = \{0,1\}; \; p(0) = a; \; p(1) = 1 - a; \; q(0) = b; \; q(1) = 1 - b.$$

$$D\left(p\middle\|q\right) = \sum_{x \in \mathcal{X}} p\left(x\right) \log_2 \frac{p\left(x\right)}{q\left(x\right)} = a \log_2 \frac{a}{b} + (1-a) \log_2 \frac{1-a}{1-b};$$

$$D\left(q\middle\|p\right) = \sum_{x \in \mathcal{X}} q\left(x\right) \log_2 \frac{q\left(x\right)}{p\left(x\right)} = b \log_2 \frac{b}{a} + (1-b) \log_2 \frac{1-b}{1-a}.$$

$$a = \frac{1}{4}; b = \frac{1}{8}, \; D\left(p\middle\|q\right) = \frac{1}{4} \log_2 \frac{8}{4} + \left(1 - \frac{1}{4}\right) \log_2 \left(\frac{1 - \frac{1}{4}}{1 - \frac{1}{8}}\right) = 0.0832 \text{ bit;}$$

$$\boxed{D\left(p\middle\|q\right) \neq D\left(q\middle\|p\right)}$$

$$D\left(q\middle\|p\right) = \frac{1}{8} \log_2 \frac{4}{8} + \left(1 - \frac{1}{8}\right) \log_2 \left(\frac{1 - \frac{1}{8}}{1 - \frac{1}{4}}\right) = 0.0696 \text{ bit.}$$

# Mutual Information

**Definition (mutual information):** *Mutual information* between two r.v. $X$ and $Y$ is defined by:

$$I(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x) p_Y(y)}$$

$$I(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)} \frac{\cancel{p_Y(y)}}{\cancel{p_Y(y)}}$$

$$I(X;Y) = I(Y;X)$$

$$I(X;Y) \geq 0$$

$$\Rightarrow I(X;Y) = 0 \Leftrightarrow X \perp Y \quad \text{(if independent)}$$

# Mutual Information

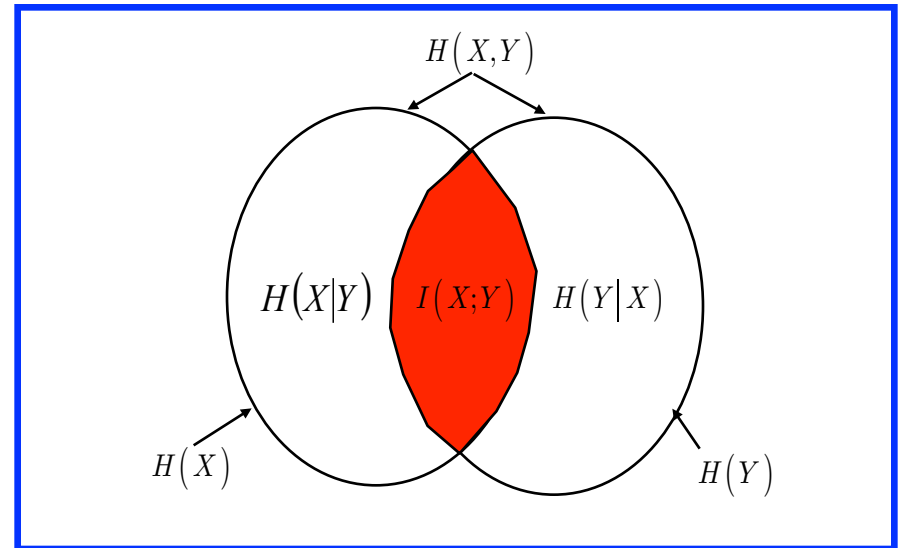## Relationship between mutual information and entropy

$$I(X;Y) = H(X) - H(X|Y)$$
$$I(X;Y) = H(Y) - H(Y|X)$$
$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

$$I(X;X) = H(X)$$

$$I(X;X) = H(X) - \underbrace{H(X|X)}_{0}$$



## Relationship between mutual information and KLD

$$I(X;Y) = D\big(p(x,y)\big\|p(x)p(y)\big) = E_{p(x,y)}\left[\log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)}\right]$$

# Mutual Information

**Conditional mutual information**

$$I\left(X;Y\big|Z\right) = H\left(X\big|Z\right) - H\left(X\big|Y,Z\right)$$

**Proof:**

$$I\left(X;Y\big|Z\right) = E_{p\left(x,y,z\right)}\left[\log\frac{p_{X,Y\big|Z}\left(x,y\big|z\right)}{p_{X\big|Z}\left(x\big|z\right)p_{Y\big|Z}\left(y\big|z\right)}\right] = E_{p\left(x,y,z\right)}\left[\log\frac{p_{X\big|Y,Z}\left(x\big|y,z\right)}{p_{X\big|Z}\left(x\big|z\right)}\right]$$

$$p_{X,Y\big|Z}\left(x,y\big|z\right) = p_{Y\big|Z}\left(y\big|z\right)p_{X\big|Y,Z}\left(x\big|y,z\right)$$

$$= \underbrace{E_{p\left(x,y,z\right)}\left[\log p_{X\big|Y,Z}\left(x\big|y,z\right)\right]}_{-H\left(X\big|Y,Z\right)} - \underbrace{E_{p\left(x,y,z\right)}\left[\log p_{X\big|Z}\left(x\big|z\right)\right]}_{H\left(X\big|Z\right)}$$

# Mutual Information

**Chain rule for mutual information**

$$I\left(X_1, X_2, ..., X_n; Y\right) = \sum_{i=1}^{n} I(X_i; Y \mid X_{i-1}, ..., X_1)$$

**Proof:**

$$I\left(X_1, X_2, ..., X_n; Y\right) = H\left(X_1, X_2, ..., X_n\right) - H\left(X_1, X_2, ..., X_n \mid Y\right)$$

$$\left\{ I\left(Z; Y\right) = H\left(Z\right) - H\left(Z \mid Y\right) \right\}$$

Chaine rule for entropy

$$\left\{ H(X_1, X_2, ..., X_n) = \sum_{i=1}^{n} H(X_i \mid X_{i-1}, ..., X_1) \right\}$$

$$= \sum_{i=1}^{n} H(X_i \mid X_{i-1}, ..., X_1) - \sum_{i=1}^{n} H(X_i \mid X_{i-1}, ..., X_1, Y) =$$

$$= \sum_{i=1}^{n} I(X_i; Y \mid X_{i-1}, ..., X_1) \quad \square$$

# Mutual Information

**Principe on non-creation of information by processing**

If $X \to Y \to Z$ form a Markov chain

$$I(X;Y) \geq I(X;Z) \text{ and } I(Y;Z) \geq I(X;Z)$$

$$p_{X,Y,Z}(x,y,z) = p_X(x)\, p_{Y|X}(y|x)\, p_{Z|Y}(z|y)$$

$X,Y,Z$ form a Markov chain

If $Z = g(Y)$, one has $I(X;Y) \geq I(X;g(Y))$

**Definition (« Data processing lemma»)**:

the amount of infromation can not be increased by any processing!

# Mutual Information

**Proof**

Using the chain rule for the mutual information:

$$I\left(X;Y,Z\right) = I\left(X;Z\right) + I\left(X;Y\middle|Z\right) \quad (a)$$

$$= I\left(X;Y\right) + I\left(X;Z\middle|Y\right) \quad (b)$$

$$I\left(X;Z\middle|Y\right) = 0 \quad \text{(based on Markovianity)}$$

Thus:

$$I\left(X;Y\right) = I\left(X;Z\right) + I\left(X;Y\middle|Z\right) \quad \Rightarrow \quad \boxed{I\left(X;Y\right) \geq I\left(X;Z\right)}$$

# Continuous R.Vs.: Differential entropy

**Definition (entropy):** *Differential entropy* of r.v. $X$ with pdf $f_X(x)$

$$h(X) = -\int_{\mathcal{X}} f_X(x)\log_2 f_X(x)\,dx = E_{f_X}\left[-\log_2 f_X(x)\right]$$

$$P\left[x_k; \Delta x\right] = \int_{\Delta x} f_X(x)\,dx \cong f_X(x_k)\Delta x$$

# Continuous R.Vs.: Differential entropy

$$H\left(X;\Delta x\right) = -\sum_{i=1}^{N} P\left[x_k, \Delta x\right] \log_2 P\left[x_k, \Delta x\right] = -\sum_{i=1}^{N} f_X\left(x_k\right) \Delta x \log_2 f_X\left(x_k\right) \Delta x$$

$$= -\sum_{i=1}^{N} f_X\left(x_k\right) \Delta x \log_2 f_X\left(x_k\right) - \underbrace{\sum_{i=1}^{N} f_X\left(x_k\right) \Delta x}_{=1} \underbrace{\log_2 \Delta x}_{const}$$

$$\int_{-\infty}^{+\infty} f_X\left(x\right) dx = 1$$

In the limit $\Delta x$ tends to zero for large $N$. As a result, $\log_2 \Delta x$ tends to infinity.

$$H\left(X;\Delta x\right) = \underbrace{-\int f_X\left(x\right) \log_2 f_X\left(x\right) dx}_{h\left(X\right)} - \log_2 \Delta x.$$

# Continuous R.Vs.: Differential entropy

**Definition (Differential entropy of Gaussian r.v.):** *differential entropy of Gaussian r.v.* $X \sim \mathcal{N}\left(0, \sigma_X^2\right)$ is:

$$h\left(X\right) = -E_{f_X}\left[\log_2 f_X\left(x\right)\right] = \frac{1}{2}\log_2\left(2\pi e \sigma_X^2\right).$$

Proof:

$$h\left(X\right) = -\int_{\mathcal{X}} f_X\left(x\right)\ln f_X\left(x\right)dx \quad [nants] = -\int_{-\infty}^{\infty} f_X\left(x\right)\left[-\frac{x^2}{2\sigma_X^2} - \ln\sqrt{2\pi\sigma_X^2}\right]dx =$$

$$= \int_{-\infty}^{\infty} f_X\left(x\right)\frac{x^2}{2\sigma_X^2}\,dx + \ln\sqrt{2\pi\sigma_X^2}\int_{-\infty}^{\infty} f_X\left(x\right)dx$$

# Continuous R.Vs.: Differential entropy

$$h(X) = \int_{-\infty}^{\infty} f_X(x) \frac{x^2}{2\sigma_X^2} \, dx + \ln \sqrt{2\pi\sigma_X^2} \underbrace{\int_{-\infty}^{\infty} f_X(x) \, dx}_{= 1}$$

$$Var[X] = \int_{-\infty}^{\infty} x^2 f_X(x) \, dx = \sigma_X^2$$

$$\int_{-\infty}^{\infty} f_X(x) \frac{x^2}{2\sigma_X^2} \, dx = \frac{\sigma_X^2}{2\sigma_X^2} = \frac{1}{2}$$

$$h(X) = \frac{1}{2} + \ln \sqrt{2\pi\sigma_X^2} = \frac{1}{2} \ln e + \frac{1}{2} \ln 2\pi\sigma_X^2 = \frac{1}{2} \left[ \ln e + \ln 2\pi\sigma_X^2 \right] = \frac{1}{2} \ln 2\pi e \sigma_X^2 \left[ nants \right]$$

$$h(X) = \frac{1}{2} \log_2 2\pi e \sigma_X^2 \left[ bits \right]$$

# Differential entropy: properties

1. Translation does not change the entropy:

$$h\big(X + a\big) = h\big(X\big).$$

2. Impact of scaling on the differential entropy, if $X$ is scalar r.v.:

$$h\big(Xa\big) = h\big(X\big) + \log\big|a\big|,$$

determinant of $\mathbf{A}$

and if $X$ is a random:

$$h\big(\mathbf{A}\mathbf{X}\big) = h\big(\mathbf{X}\big) + \log\big|\det(\mathbf{A})\big|,$$