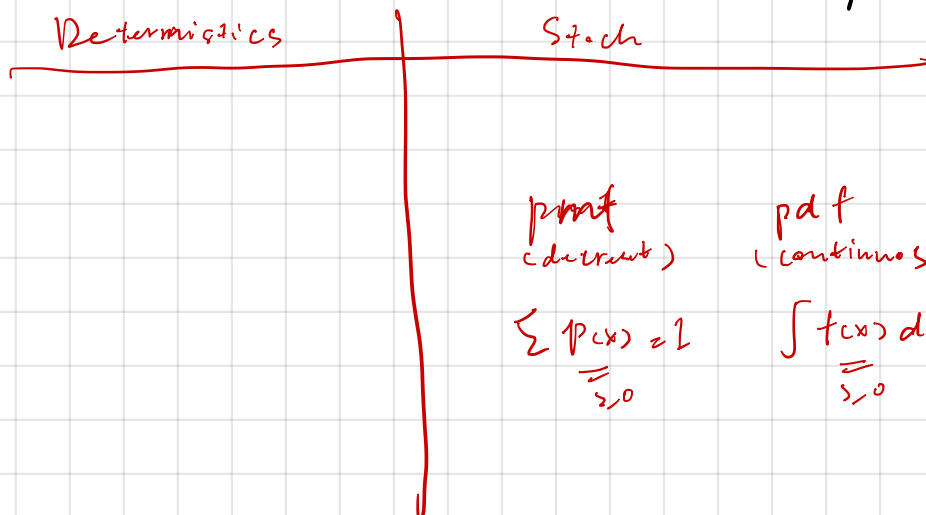
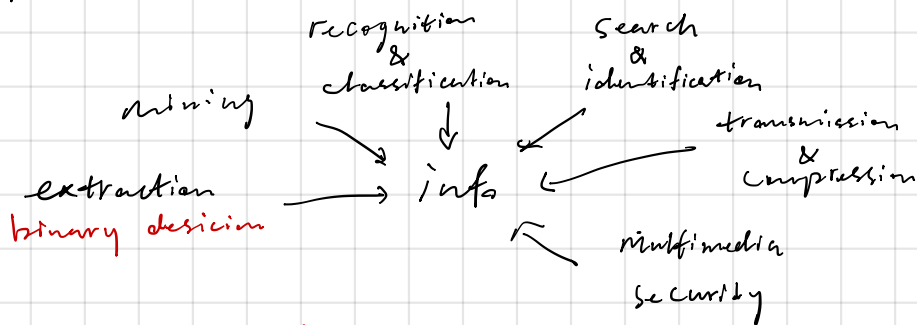


main probs of interest



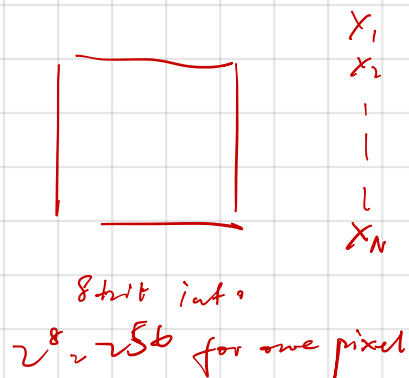
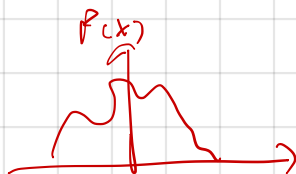
information theory & machine learning.

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

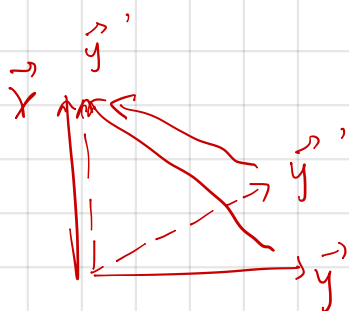
$$X \sim f(x)$$

$$\mu = E(X) = \sum x p(x)$$

$$\sigma^2 = \sum (x-\mu)^2 p(x)$$



(μ, σ^2) . most unpredictable distributions



$$\langle \vec{x}, \vec{y} \rangle \rightarrow \frac{\pi}{2}, d(x,y) \rightarrow \max \quad \langle \vec{x}, \vec{y} \rangle \rightarrow 0 \quad d(x,y) \rightarrow \min$$

$$\cos \langle \vec{x}, \vec{y} \rangle \rightarrow 0 \quad \cos \langle \vec{x}, \vec{y} \rangle \rightarrow 1$$

$$x \sim p(x) \quad y \sim p(y)$$

$$E_{p(x)} E_{p(y)} \|x - y\| = ? \quad \text{we don't know.}$$

SVD

Single value decomposition

$$A \approx M \times N \xRightarrow{\text{SVD}} A = U \Sigma V^T$$

\uparrow \uparrow \uparrow
 $m \times m$ $m \times n$ $n \times n$

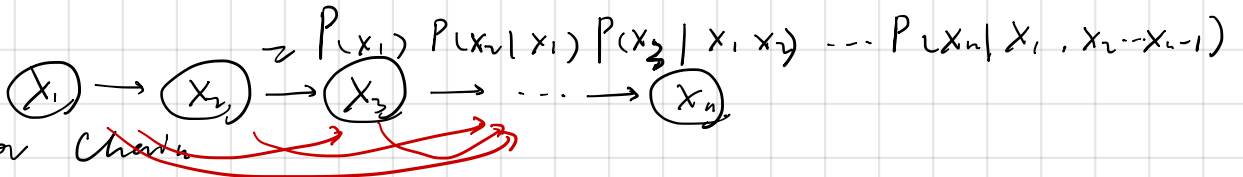
$$\Sigma = \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_r \\ & & & & 0 \end{bmatrix}$$

Information Theory

Chain Rule $P(x, y) = P(y|x) P(x) = P(x|y) P(y)$

Bayes Rule $P(y|x) = \frac{P(x|y) P(y)}{P(x)}$

$$\begin{aligned} P(x_1, x_2, x_3, \dots, x_n) &= P(x_1 | x_2, x_3, \dots, x_n) P(x_2, x_3, \dots, x_n) \\ &= P(x_1 | x_2, x_3, \dots, x_n) P(x_2 | x_3, \dots, x_n) P(x_3, \dots, x_n) \\ &\quad \vdots \\ &= P(x_1 | x_2, x_3, \dots, x_n) P(x_2 | x_3, \dots, x_n) \dots P(x_{n-1} | x_n) P(x_n) \\ &= P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots P(x_n | x_1, x_2, \dots, x_{n-1}) \end{aligned}$$



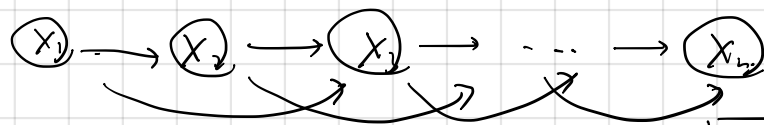
first order:

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1) P(x_2 | x_1) P(x_3 | x_2) \dots P(x_n | x_{n-1})$$



second order:

$$P(x_1, x_2, x_3, \dots, x_n) = P(x_1) P(x_2 | x_1) P(x_3 | x_1, x_2) \dots P(x_n | x_{n-1}, x_{n-2})$$

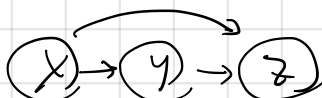


$$= \prod_{i=2}^n P(x_i | x_{i-1}, x_{i-2})$$

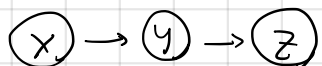
forward and backward

$$P(x_2 | x_0, x_1) = P(x_2 | x_1) \quad \text{base case.}$$

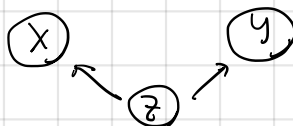
$$P(x_1 | x_1, x_0) = P(x_1)$$



$$P(x, y, z) = P(x) P(y|x) P(z|x, y) \quad P(x, y, z) = P(x) P(y) P(z)$$



$$P(x, y, z) = P(x) P(y|x) P(z|y) \quad P(x, y, z) = P(z) P(y|z) P(x|y)$$



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

$$X \sim P(X), \quad X = [x_1, x_2, \dots, x_n], \quad n \text{ samples}$$

$$E(X) = \sum_{i=1}^n P(X=x_i) x_i \quad X \in \mathcal{X}$$

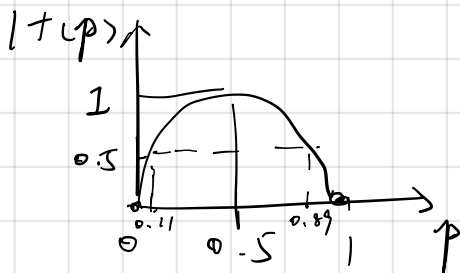
$$H(X) = - \sum_{i=1}^n P(X=x_i) \log_2 P(X=x_i) = - \sum_{x \in \mathcal{X}} P(x) \log_2 P(x) \\ = E[-\log_2 P(x)]$$

$$X \in \{0, 1\} \quad X \in \mathcal{X}$$

$$P(X=0) = p$$

$$P(X=1) = 1-p$$

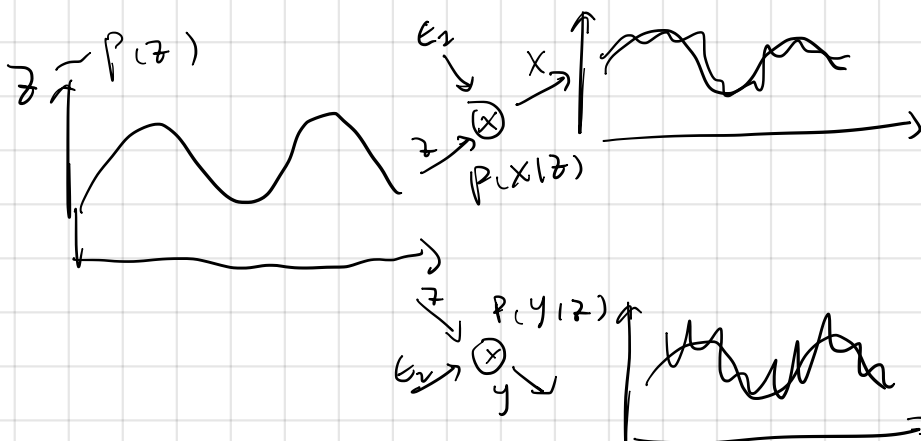
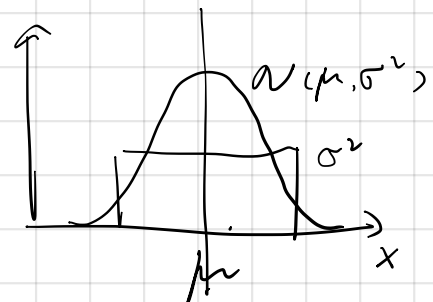
$$H(X) = -p \log_2 p - (1-p) \log_2 (1-p) = H(p)$$



$$-\log_2 p$$

* Uniform distribution is the most unexpected given the fixed domain $[a, b]$ so the entropy of it is the highest

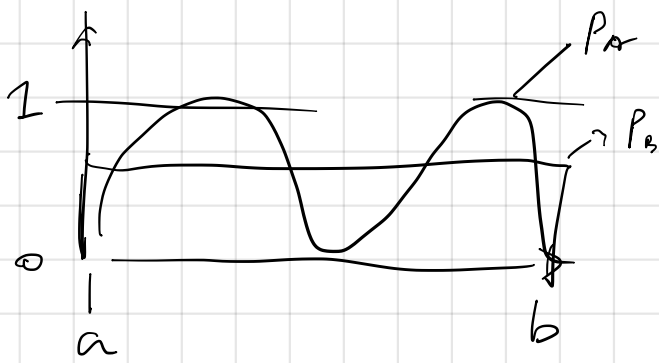
* Gaussian distribution has the highest entropy given the fixed σ^2



$$x - z = \epsilon_1$$

$$\epsilon_1 \perp \epsilon_2 \text{ independent}$$

$$y - z = \epsilon_2$$



$$H(P_A) < H(P_B)$$

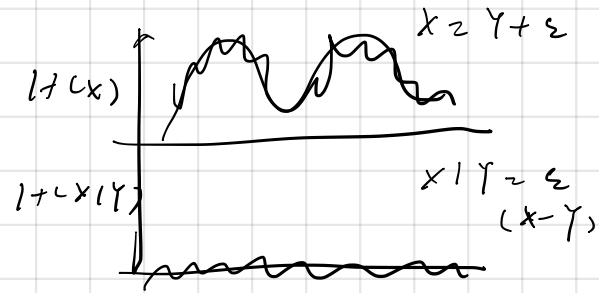
Q. ?

⚠ n independent $X_i \perp X_j$

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i)$$

Conditioning kills entropy

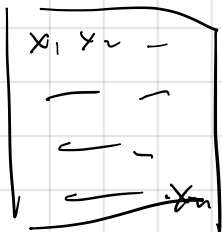
$$H(X|Y) \leq H(X)$$



$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$$

using 1st Markov chain to prove this

$$H(X_1, X_2, \dots, X_n) \text{ bits} \leq 8 \cdot n \text{ bits}$$



$$H(X_n) = \log_2 256 = 8 \text{ bits}$$

n

given $P(x, y, z)$. to know $P(x)$

marginalisation)

$$P(x) = \sum_y \sum_z P(x, y, z)$$

$$P(x) P(y|x) P(z|y, x)$$

independence & correlation
for all the moments, just second moment of the distribution