

Lecture 13.

HGR ^{Renyi}
Maximal Correlation.

Given 2 discrete random variables X, Y , want to measure the correlation between X and Y , how can we do?

Example: Pearson correlation coefficient:

$$\gamma(X, Y) \triangleq \frac{\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}$$

- If $\gamma(X, Y) = 0 \nRightarrow X, Y$ are indep.

A good correlation measurement $\rho(X, Y)$ should satisfy:

- (1) Commutable $\rho(X, Y) = \rho(Y, X)$
- (2) $0 \leq \rho(X, Y) \leq 1$
- (3) $\rho(X, Y) = 0$ if and only if X, Y are indep.
- (4) $\rho(X, Y) = 1$, if $Y = \xi(X)$ or $X = \eta(Y)$, for some deterministic functions ξ, η .
- (5) For one-to-one functions ξ, η , $\rho(\xi(X), \eta(Y)) = \rho(X, Y)$

Example: Pearson correlation coefficient $\gamma(X, Y)$

$$\gamma(X, Y) \triangleq \frac{\mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} \quad -1 \leq \gamma(X, Y) \leq 1$$

$|\gamma(X, Y)|$ satisfies (1), (2), not (3), (4), (5) (HW)

$|\gamma(X, Y)| = 1$, iff $X = aY + b$ or $Y = cX + d$

Example: mutual information $I(X;Y) \stackrel{\Delta}{=} \sum_{x,y} P_{XY}(x,y) \log \frac{P_{XY}(x,y)}{P_X(x)P_Y(y)}$
 joint dist. of X, Y

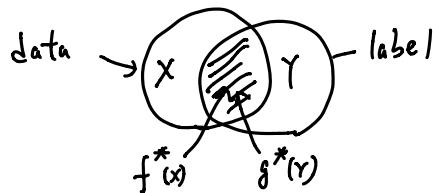
$I(X;Y)$ satisfies (1), (3), (5) but (2), (4)

- Can we find some $\rho(X,Y)$ satisfying (1)-(5) ?

Definition: The HGR maximal correlation $\rho(X,Y)$ is defined as

$$\rho(X,Y) \stackrel{\Delta}{=} \max_{\substack{f:X \rightarrow \mathbb{R}, g:Y \rightarrow \mathbb{R} \\ \mathbb{E}[f(x)] = \mathbb{E}[g(Y)] = 0 \\ \mathbb{E}[f^2(x)] = \mathbb{E}[g^2(Y)] = 1}} \mathbb{E}[f(x) \cdot g(Y)]$$

- $\rho(X,Y) = \max_{f,g} \gamma(f(x), g(Y))$
- Data X, Y , want find features of X and Y such that these features are the most related \Rightarrow HGR maximal correlation functions



$\rho(X,Y)$ satisfies (1), (2), (5), (4), (3)

Check (5): for one-to-one functions ξ, η , $\rho(\xi(X), \eta(Y)) = \max_{f,g} \gamma(\underbrace{f(\xi(X))}_{f'}, \underbrace{g(\eta(Y))}_{g'})$

$$= \max_{f',g'} \gamma(f'(x), g'(\eta(Y))) = \rho(X,Y)$$

Check (4): If $Y = \xi(X)$, let $K(Y) : Y \mapsto \mathbb{R}$ be a one-to-one function such that $\mathbb{E}[K(Y)] = 0$, $\mathbb{E}[K^2(Y)] = 1$, then take $f(x) = K(\xi(x))$, $g(Y) = K(Y)$

$$\Rightarrow \underbrace{\mathbb{E}[f(x) \cdot g(Y)]}_{\leq 1} = \mathbb{E}\left[\underbrace{k(\xi(x))}_{Y} \cdot k(Y)\right] = \mathbb{E}[k^2(Y)] = 1 \Rightarrow \rho(X,Y) = 1$$

Definition: The canonical dependence matrix $\tilde{B} \in \mathbb{R}^{|\mathcal{Y}| \times |\mathcal{X}|}$ is defined as

$$\tilde{B}(y, x) \triangleq \frac{P_{XY}(x, y) - P_X(x)P_Y(y)}{\sqrt{P_X(x) \cdot P_Y(y)}}, \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}$$

Example: X, Y are binary, $P_{XY}(x, y) = \begin{cases} \frac{1}{3} & \text{if } (x, y) = (0, 0) \\ \frac{1}{6} & \cdots \\ \frac{1}{6} & (0, 1) \\ \frac{1}{3} & (1, 0) \\ \frac{1}{3} & (1, 1) \end{cases}$

X	0	1	
	$\frac{1}{3}$	$\frac{1}{6}$	
Y	0	$\frac{1}{6}$	$P_Y(0) = \frac{1}{2}$
	1	$\frac{1}{3}$	$P_Y(1) = \frac{1}{2}$
	$P_X(0)$	$P_X(1)$	
	$\frac{1}{2}$	$\frac{1}{2}$	

\tilde{B} is a 2×2 matrix

$$\tilde{B} = \begin{bmatrix} 0 & 1 \\ 0 & \tilde{B}(0, 0) = \frac{\frac{1}{3} - \frac{1}{2} \cdot \frac{1}{2}}{\sqrt{\frac{1}{2} \cdot \frac{1}{2}}} = \frac{\frac{1}{6}}{\frac{1}{2}} = \frac{1}{6} & \tilde{B}(0, 1) = \frac{\frac{1}{3} - \frac{1}{2} \cdot \frac{1}{2}}{\sqrt{\frac{1}{2} \cdot \frac{1}{2}}} = \frac{-\frac{1}{6}}{\frac{1}{2}} = -\frac{1}{6} \\ 1 & \tilde{B}(1, 0) = -\frac{1}{6} & \tilde{B}(1, 1) = \frac{1}{6} \end{bmatrix}$$

$$\tilde{B} = \begin{bmatrix} \frac{1}{6} & -\frac{1}{6} \\ -\frac{1}{6} & \frac{1}{6} \end{bmatrix}$$

Definition: The information vectors ϕ, ψ associated to functions f, g are defined:

$$\phi(x) = \underbrace{\sqrt{P_X(x)} \cdot f(x)}_{\text{1x1-dimensional vector}}, \quad \psi(y) = \underbrace{\sqrt{P_Y(y)} \cdot g(y)}_{\text{1y1-dim vector}}$$

$$\phi = \begin{bmatrix} \phi(1) \\ \phi(2) \\ \vdots \\ \phi(|\mathcal{X}|) \end{bmatrix}, \quad \psi = \begin{bmatrix} \psi(1) \\ \psi(2) \\ \vdots \\ \psi(|\mathcal{Y}|) \end{bmatrix}$$

Whenever functions f, g are given, we have the corresponding ϕ, ψ .

- property : norm-square: $\|\phi\|^2 = \sum_x \phi^2(x) = \sum_x P_X(x) f^2(x) = \mathbb{E}[f^2(x)]$, $\|\psi\|^2 = \mathbb{E}[g^2(y)]$

$$\text{inner product: } \underline{\sqrt{P_X}} = \begin{bmatrix} \sqrt{P_X(1)} \\ \sqrt{P_X(2)} \\ \vdots \\ \sqrt{P_X(|\mathcal{X}|)} \end{bmatrix}, \quad \langle \phi, \underline{\sqrt{P_X}} \rangle = \sum_x \phi(x) \cdot \sqrt{P_X(x)} = \sum_x P_X(x) \cdot f(x) = \mathbb{E}[f(x)]$$

$$\underline{\sqrt{P_Y}} = \begin{bmatrix} \sqrt{P_Y(1)} \\ \vdots \\ \sqrt{P_Y(|\mathcal{Y}|)} \end{bmatrix}, \quad \langle \psi, \underline{\sqrt{P_Y}} \rangle = \mathbb{E}[g(y)]$$

Thm: The HGR maximal correlation $\rho(x, y)$ is the largest singular value of \tilde{B}

$$\text{Proof: Note that } \mathbb{E}[f(x) \cdot g(y)] = \sum_{x, y} P_{XY}(x, y) \cdot f(x) \cdot g(y) = \sum_{x, y} \left(\frac{P_{XY}(x, y)}{\sqrt{P_X(x)P_Y(y)}} \cdot (\sqrt{P_X(x)} \cdot f(x)) \cdot (\sqrt{P_Y(y)} \cdot g(y)) \right)$$

$$= \frac{P_X(x)P_Y(y)}{\sqrt{P_X(x)P_Y(y)}} \cdot (\cancel{\sqrt{P_X(x)}f(x)}) (\cancel{\sqrt{P_Y(y)}g(y)})$$

$$\sum_{x, y} \frac{P_X(x)P_Y(y)}{\sqrt{P_X(x)P_Y(y)}} \cdot (\sqrt{P_X(x)} f(x)) \cdot (\sqrt{P_Y(y)} g(y)) = \sum_{x, y} P_X(x) f(x) \cdot P_Y(y) g(y)$$

$$= \mathbb{E}[f(x)] \cdot \mathbb{E}[g(y)]$$

$$\Rightarrow \mathbb{E}[f(x) \cdot g(y)] = \sum_{x,y} \underbrace{\frac{P_{xy}(x,y) - P_x(x) \cdot P_y(y)}{\sqrt{P_x(x) P_y(y)}}}_{\tilde{B}(y,x)} \cdot (\underbrace{\sqrt{P_x(x)} f(x)}_{\phi(x)}) \cdot (\underbrace{\sqrt{P_y(y)} g(y)}_{\psi(y)}) = \sum_{x,y} \tilde{B}(y,x) \cdot \phi(x) \cdot \psi(y) = \phi^T \tilde{B} \psi$$

$$\Rightarrow \rho(x, y) = \max_{\phi, \psi} \phi^T \tilde{B} \psi \quad \text{subject to: } \|\phi\|^2 = \|\psi\|^2 = 1 \quad (\Leftrightarrow \mathbb{E}[f^2(x)] = \mathbb{E}[g^2(y)] = 1)$$

$$\langle \phi, \sqrt{P_x} \rangle = \langle \psi, \sqrt{P_y} \rangle = 0$$

$$(\Leftrightarrow \mathbb{E}[fx] = \mathbb{E}[gy] = 0)$$

Fact 1: $\arg \max_{\|\phi\|^2 = \|\psi\|^2 = 1} \phi^T \tilde{B} \psi = \text{largest right / left singular vectors of } \tilde{B}$

Fact 2: singular vector are orthogonal with each other.

check: $\left[\begin{matrix} \tilde{B} & \sqrt{P_x} \end{matrix} \right] (y) = \sum_x \tilde{B}(y, x) \cdot \sqrt{P_x(x)} = \sum_x \frac{P_{xy}(x, y) - P_x(x) \cdot P_y(y)}{\sqrt{P_y(y)}} = \frac{P_y(y) - P_y(y)}{\sqrt{P_y(y)}} = 0$

$\uparrow \quad \uparrow$
 $|y| \times 1 \quad 1 \times 1$ right

$\Rightarrow \sqrt{P_x}$ is a \checkmark singular vector of \tilde{B} with singular value 0.

$$\text{since } \tilde{B} \cdot \sqrt{P_x} = 0$$

$\sqrt{P_y}$ is also a left singular vector of \tilde{B}

$\Rightarrow \rho(x, y) = \sigma_1$, with " $=$ " iff ϕ and ψ are the largest right/left singular vectors of \tilde{B}

\Rightarrow Let ϕ_i, ψ_i be the largest right/left singular vectors of \tilde{B} , then

$$f^*(x) = \frac{1}{\sqrt{P_x(x)}} \cdot \phi_i(x), \quad g^*(y) = \frac{1}{\sqrt{P_y(y)}} \cdot \psi_i(y)$$

- $\text{svd}(\tilde{B}) = [\sigma_i \cdot \phi_i \cdot \psi_i]$

$\uparrow \quad \uparrow \quad \uparrow$
 $\rho(x, y) \quad f(x) \quad g(y)$

Check (3): $\rho(x, y) = 0$ iff $\sigma_i = 0$ iff $\tilde{B} = 0$ iff $P_{xy}(x, y) = P_x(x) \cdot P_y(y)$
iff x, y are indep.

Power iteration / method

For a positive definite matrix M , eigenvalues $\lambda_1 > \lambda_2 > \dots > \lambda_m$
eigenvectors v_1, v_2, \dots, v_m

$$\Rightarrow M = \sum_{i=1}^m \lambda_i \cdot v_i \cdot v_i^T : \text{eigen-decomposition}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \ddots & 0 \\ 0 & & \lambda_m \end{bmatrix}, V = [v_1 \dots v_m] \Rightarrow M = V \Lambda \cdot V^T$$

The algorithm to compute v_1 ,

- (i) Choose $u^{(0)} \neq v_1$ power step
- (ii) Compute $u^{(1)} = M \cdot u^{(0)}$, $u^{(1)} \leftarrow \frac{u^{(1)}}{\|u^{(1)}\|}$ normalization
- $u^{(2)} = M \cdot u^{(1)}$
- \vdots
- $u^{(t)} = M \cdot u^{(t-1)}$, normalization
- $\Rightarrow u^{(t)} = \frac{M^t \cdot u^{(0)}}{\|M^t \cdot u^{(0)}\|}$

Thm: $u^{(t)} \xrightarrow{t \rightarrow \infty} v_1$

Proof: Let $u^{(0)} = \sum_{i=1}^m \alpha_i v_i$, then $M^t \cdot u^{(0)} = \sum_{i=1}^m \alpha_i \underbrace{M^t \cdot v_i}_{\lambda_i^t \cdot v_i} = \sum_{i=1}^m \alpha_i \lambda_i^t v_i$

since λ_1 is the largest among all λ_i 's

$$\Rightarrow \frac{\lambda_1^t}{\lambda_1^t} \rightarrow 0, t \rightarrow \infty \quad \Rightarrow \frac{M^t \cdot u^{(0)}}{\|M^t \cdot u^{(0)}\|} \rightarrow v_1, t \rightarrow \infty$$

largest right singular vector of \tilde{B}

Fact 3: ϕ_1 is also the largest eigenvector of $\tilde{B}^T \tilde{B}$
 ϕ_1 is .. of $\tilde{B} \cdot \tilde{B}^T$

The algorithm to compute ϕ_1, ψ_1 ,

- (i) choose $\phi^{(0)} \neq \phi_1$, and $\phi^{(0)} \perp \sqrt{P_X}$
- (ii) Iterating: $\begin{cases} \psi^{(i)} = \tilde{B} \cdot \phi^{(i)} \\ \phi^{(i+1)} = \tilde{B}^T \psi^{(i)} \end{cases}$, for $i = 0, 1, \dots, t$

(iii) Normalization

$$\phi^{(0)} \rightarrow \psi^{(0)} = \tilde{B} \cdot \phi^{(0)} \rightarrow \phi^{(1)} = \tilde{B}^T \cdot \psi^{(0)} \rightarrow \phi^{(1)} = \tilde{B} \cdot \phi^{(0)} \rightarrow \dots$$

For ϕ : $\phi^{(i+1)} = \tilde{B}^T \cdot \tilde{B} \cdot \phi^{(i)} \Rightarrow \phi^{(i+1)}$ converge to largest eigenvector of $\tilde{B}^T \cdot \tilde{B}$

We can write $f^{(i)}(x) = \frac{1}{\sqrt{P_x(x)}} \phi^{(i)}(x)$, $g^{(i)}(y) = \frac{1}{\sqrt{P_y(y)}} \cdot \phi^{(i)}(y)$

Then, the updating step $\psi^{(i)} = \tilde{B} \cdot \phi^{(i)}$ can be expressed in terms of functions as:

$$\begin{aligned} g^{(i)}(y) &= \frac{1}{\sqrt{P_y(y)}} \cdot \psi^{(i)}(y) = \frac{1}{\sqrt{P_y(y)}} \cdot \sum_x \tilde{B}(y, x) \cdot \phi^{(i)}(x) = \frac{1}{\sqrt{P_y(y)}} \cdot \sum_x \frac{P_{xy}(x|y) - P_x(x)P_y(y)}{\sqrt{P_x(x)P_y(y)}} \cdot (\cancel{\sqrt{P_x(x)}} f^{(i)}(x)) \\ &= \frac{1}{\sqrt{P_y(y)}} \cdot \left(\sum_x P_x(x|y) \cdot f^{(i)}(x) - \sum_x \cancel{P_x(x)P_y(y)} f^{(i)}(x) \right) \\ &= \sum_x P_{x|y}(x|y) \cdot f^{(i)}(x) = \mathbb{E}[f^{(i)}(x) | Y=y] \end{aligned}$$

Alternative conditional Expectation (ACE)

(i) Choose zero-mean $f^{(0)}(x)$, $\mathbb{E}[f^{(0)}(x)] = 0$

(ii) Iterating : $\begin{cases} g^{(i)}(y) = \mathbb{E}[f^{(i)}(x) | Y=y] \\ f^{(i+1)}(x) = \mathbb{E}[g^{(i)}(Y) | X=x] \end{cases}$, for $i=0 \dots t$

(iii) Normalization.

Compute from samples : $(x_1, y_1), \dots, (x_n, y_n)$

$$\begin{aligned} \text{for } y_1, \dots, y_n, \quad y_1 = y_2 = \dots = y_k = y \\ \Rightarrow \mathbb{E}[f(x) | Y=y] = \frac{1}{k} \sum_{i=1}^k f(x_i) \end{aligned}$$