

4/19/2019 Discussion # 3

Entropy:

$\log$  (base 2)

Definition:

If a random variable (RV)  $X$  has  $k$  different values,  $a_1, a_2, \dots, a_k$ , its entropy is given by

$$H[X] = - \sum_{k=1}^k p(X=a_k) \log p(X=a_k)$$

Example:

$$k=2. \quad p(X=0) = p$$

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

aside:  $-\log a = \log \frac{1}{a}$

$$H(X) = \mathbb{E}[\log \frac{1}{p(X)}] = p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p} = H_2(p) \\ = H_b(p)$$



↑  
they are the  
same thing,  
just different  
notation.

Conditional entropy:

Definition:

Given two RVs  $X$  and  $Y$

$$H[Y|X] = \sum_k P(X=a_k) H[Y|X=a_k]$$

Information Gain: (Mutual information)

$$\text{Gain} = H[Y] - H[Y|X] = I(X;Y)$$

Properties:

$$0 \leq H(x) \leq \log |K| \quad K \text{ is a set}$$

$|K|$ : # of elements in set  $K$ .  
(Cardinality).

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

$$= I(Y;X)$$

$$0 \leq H(Y|X) \leq H(Y)$$

$$\text{when } H(Y|X) = H(Y) \quad X \perp\!\!\!\perp Y$$

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

**Note:**  $H(X) \geq H(X|Y)$  didn't mean  $H(X) \geq H(X|Y=y)$

Decision Tree algorithm:

Iterative Dichotomiser 3: (ID3):

Select the feature with largest information gain.

pb 1:

$$a) P(Y=1) = \frac{5}{8} \quad P(Y=0) = \frac{3}{8}$$

$$H_b\left(\frac{3}{8}\right) = -\frac{5}{8} \log \frac{5}{8} - \frac{3}{8} \log \frac{3}{8} = 0.9544.$$

$$b) H(Y|X_1)$$

$$P(X_1=0) = \frac{1}{2}, \quad P(X_1=1) = \frac{1}{2}$$

$$\begin{array}{c} X_1 \\ \diagup \quad \diagdown \\ 1 \quad 0 \\ \vdots \quad \vdots \\ 0 \quad 0 \end{array} \quad \begin{array}{l} P(Y=1|X=1) = 1 \\ P(Y=0|X=1) = 0 \\ P(Y=1|X=0) = \frac{1}{4} \\ P(Y=0|X=0) = \frac{3}{4} \end{array}$$

$$H(Y|X_1) = P(X_1=0)H(Y|X_1=0) + P(X_1=1)H(Y|X_1=1)$$

$$= \frac{1}{2} H_b\left(\frac{1}{4}\right) + \frac{1}{2} H_b(1)$$

$$\approx 0.4056.$$

$$c) H(Y|X_2), \quad H(Y|X_3), \quad H(Y|X_4)$$

$$H(Y|X_2) = \frac{2}{8} H_b(1) + \frac{6}{8} H_b\left(\frac{1}{2}\right) = 0.75$$

$$H(Y|X_3) = \frac{1}{2} H_b\left(\frac{3}{4}\right) + \frac{1}{2} H_b\left(\frac{1}{2}\right) = 0.9056$$

$$H(Y|X_4) = \frac{2}{8} H_b(1) + \frac{6}{8} H_b\left(\frac{1}{2}\right) = 0.75$$

$$(d) IC(Y; x_1) = 0.9544 - 0.4056 = 0.5488$$

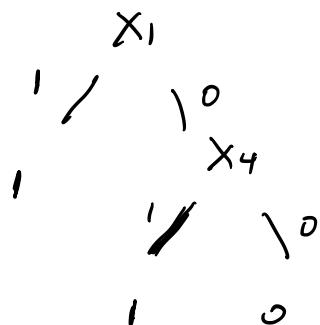
$$I(Y; X_2) = 0.2044$$

$$I(Y;X_3) = 0.1407$$

$$I(Y; X_4) = 0.2044$$

(e) after asking  $x_1$

sample:	$X_2$	$X_3$	$X_4$	$Y$
1	0	0	1	1
2	0	0	0	0
5	0	1	0	0
6	0	1	0	0



Pb (2)

$$(a) \quad \text{Var}(v) = \sum_{x_i \in V} (x_i - u)^2$$

$$\frac{\partial}{\partial u} \text{Var}(v) = -2 \sum_{x_i \in V} (x_i - u) = 0.$$

$$\sum_{x_i \in V} (x_i - u) = 0.$$

$$\sum_{x_i \in V} x_i - \sum_{x_i \in V} u = 0$$

$$|V| u = \sum_{x_i \in V} x_i$$
$$u = \frac{\sum_{x_i \in V} x_i}{|V|}$$

Pb (3)

$$(a) \underline{C} = \underline{\underline{A}} \underline{\underline{z}}$$

$$C_{11} = \sum_{j=1}^n a_{1j} z_{j1}$$

$$C_{22} = \sum_{j=1}^n a_{2j} z_{j2}$$

⋮

$$C_{nn} = \sum_{j=1}^n a_{nj} z_{j2} \quad \text{change of variable.}$$

$$\text{Tr}(\underline{C}) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} z_{ji} = \sum_{r=1}^n \sum_{s=1}^n a_{rs} z_{sr}$$

$$\frac{d \text{Tr}(\underline{\underline{A}} \underline{\underline{z}})}{d z_{ij}} = \frac{d (\sum_{r=1}^n \sum_{s=1}^n a_{rs} z_{sr})}{d z_{ij}} = a_{ji} \quad \text{only exist when } s=i, r=j$$

$$\therefore \frac{d \text{Tr}(\underline{\underline{A}} \underline{\underline{z}})}{d \underline{\underline{z}}} = A^T$$

$$(b) \quad \frac{d \operatorname{Tr}(\underline{\underline{Z}} \underline{\underline{A}} \underline{\underline{Z}}^T)}{d \underline{\underline{Z}}} = \underline{\underline{Z}} \underline{\underline{A}}^T + \underline{\underline{Z}} \underline{\underline{A}}$$

$$\text{Let } \underline{\underline{C}} = \underline{\underline{A}} \underline{\underline{Z}}^T \quad D = \underline{\underline{Z}} \underline{\underline{A}} \underline{\underline{Z}}^T = \underline{\underline{Z}} \underline{\underline{C}}$$

$$C_{11} = \sum_{t=1}^n a_{1t} z_{1t}$$

$$C_{12} = \sum_{t=1}^n a_{1t} z_{2t}$$

$$C_{21} = \sum_{t=1}^n a_{2t} z_{1t}$$

$$\vdots$$

$$C_{rs} = \sum_{t=1}^n a_{rt} z_{st}$$

$$d_{11} = \sum_{s=1}^n z_{1s} C_{s1} = \sum_{s=1}^n z_{1s} \sum_{t=1}^n a_{st} z_{1t}$$

$$= \sum_{s=1}^n \sum_{t=1}^n z_{1s} a_{st} z_{1t}$$

$$d_{22} = \sum_{s=1}^n \sum_{t=1}^n z_{2s} a_{st} z_{2t}$$

$$\operatorname{Tr}(D) = \sum_{r=1}^n \sum_{s=1}^n \sum_{t=1}^n z_{rs} a_{st} z_{rt}$$

$$= \sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss} + \sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}$$

$$\therefore \operatorname{Tr}(D) = \sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss} + \sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}$$

$$\text{Tr}(D) = \sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss} + \sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}$$

Step 1:

$$\begin{aligned}
 \frac{d\text{Tr}(D)}{z_{ii}} &= \frac{d}{dz_{ii}} \left( \underbrace{\sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss}}_{r=i, s=i} + \underbrace{\sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}}_{r=i, s=i, t \neq i} \right) \\
 &= 2z_{ii} a_{ii} + \sum_{t=1, t \neq i}^n a_{it} z_{it} + \sum_{s=1, s \neq i}^n z_{is} a_{si} \\
 &= 2z_{ii} a_{ii} + \sum_{t=1, t \neq i}^n a_{it} z_{it} + \sum_{t=1, t \neq i}^n z_{it} a_{ti} \quad \checkmark \text{change of variable.} \\
 &\stackrel{(1)}{=} z_{ii} a_{ii} + z_{ii} a_{ii} + \sum_{t=1, t \neq i}^n a_{it} z_{it} + \sum_{t=1, t \neq i}^n z_{it} a_{ti} \\
 &\stackrel{(2)}{=} \sum_{t=1}^n a_{it} z_{it} + \sum_{t=1}^n a_{ti} z_{it} \\
 &= \sum_{t=1}^n z_{it} (a_{it} + a_{ti})
 \end{aligned}$$

$$\therefore \frac{d\text{Tr}(D)}{z_{ii}} = \sum_{t=1}^n z_{it} (a_{it} + a_{ti})$$

diagonal term:

$$\text{Tr}(D) = \sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss} + \sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}$$

**Step 2:**

$$\begin{aligned}
 \frac{d\text{Tr}(D)}{z_{ij}} &= \frac{d}{dz_{ij}} \left( \underbrace{\sum_{r=1}^n \sum_{s=1}^n z_{rs}^2 a_{ss}}_{r=i, s=j} + \underbrace{\sum_{r,s,t,t+s}^n z_{rs} a_{st} z_{rt}}_{\substack{r=i, s=j, t \neq j \\ r=i, t=j, s \neq j}} \right) \\
 &= 2z_{ij} a_{jj} + \sum_{t, t \neq j}^n a_{jt} z_{it} + \sum_{s, s \neq j}^n z_{is} a_{sj} \\
 &= 2z_{ij} a_{jj} + \sum_{t, t \neq i}^n a_{jt} z_{it} + \sum_{t, t \neq j}^n z_{it} a_{tj} \\
 &= \sum_{t=1}^n a_{jt} z_{it} + \sum_{t=1}^n z_{it} a_{tj} \\
 &= \sum_{t=1}^n z_{it} (a_{jt} + a_{tj}) \\
 \therefore \frac{d\text{Tr}(D)}{z_{ij}} &= \sum_{t=1}^n z_{it} (a_{jt} + a_{tj})
 \end{aligned}$$

Combine step 1 and 2:

$$\frac{d\text{Tr}(D)}{z_{ii}} = \sum_{t=1}^n z_{it} (a_{it} + a_{ti})$$

Combine step 1 and 2:

$$\frac{d \text{Tr}(D)}{z_{ii}} = \sum_{t=1}^n z_{it} (a_{it} + a_{ti})$$

$$\frac{d \text{Tr}(D)}{z_{ij}} = \sum_{t=1}^n z_{it} (a_{jt} + a_{tj})$$

$$\frac{d \text{Tr}(D)}{\Xi} = \Xi (A + A^\top) = \Xi A + \Xi A^\top$$

(c) Aside:  $\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$   
 $\text{Tr}(AB) = \text{Tr}[(AB)^\top]$

$$\begin{aligned} E_r(w) &= \text{Tr}[(XW - T)^\top (XW - T)] \\ &= \text{Tr}[W^\top X^\top X W - \underbrace{T^\top X W}_{\text{part A}} - \underbrace{W^\top X^\top T}_{\text{part B}} + \underbrace{T^\top T}_{\text{part C}}] \\ &= \text{Tr}[W^\top X^\top X W] - 2\text{Tr}[T^\top X W] + \text{Tr}[T^\top T] \end{aligned}$$

$$\nabla_w E_r(w) = (X^\top X + X^\top X) W - 2X^\top T = 0.$$

$$\Rightarrow (X^\top X) W = X^\top T$$

$$W = (X^\top X)^{-1} X^\top T$$