③ Entropy

⊙  $$H(q) = -\sum_{i=1}^{K} P_{qk} \log_2 P_{qk}$$

binary split:

$$H(q) = -P \log_2 P - (1-P) \log_2 (1-P) \quad; \quad 0 < P < 1$$

Need to prove: $-P \log P - (1-P) \log(1-P) \leq 1$ ?
(assume log base 2)

Balanced probabilities = High entropy (lowest information gain)
   Max entropy at $P = \frac{1}{2}$

$$-\left(\frac{1}{2}\right) \log\left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log\left(\frac{1}{2}\right) \overset{?}{\leq} 1$$

Also, → $\lim\limits_{P \to \frac{1}{2}^{\pm}} \left(-P \log P - (1-P) \log(1-P)\right) = 1 \quad \leq 1 \quad \checkmark$

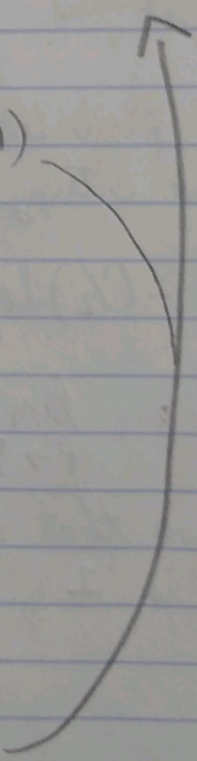∴ This shows that given a binary data split, the MAXIMUM entropy is 1

③ⓑ $n = \#$ of branches

| $n$ | $max(H)$ |
|---|---|
| 2 | 1 |
| 3 | 1.58 |
| 4 | 2 |
| 8 | 3 |
| 16 | 4 |

Same as $\log(n)$

The maximum entropy for multiway branching, given 'n' branches, is $\log(n)$.

Need to prove:

$$-\sum_{i=1}^{n} \left(\frac{1}{n}\right)\log_2\left(\frac{1}{n}\right) \overset{?}{=} \log_2(n)$$

$$= -\sum_{i=1}^{n} \frac{1}{n}\left(\overset{0}{\log_2(1)} - \log_2(n)\right)$$

$$= +\sum_{i=1}^{n} +\frac{1}{n}\log_2(n)$$

$$= n\left(\frac{1}{n}\log(n)\right)$$

$$= \log_2(n) \overset{\checkmark}{=} \log_2(n).$$

(4) $Gain(q, V) = b(q) - \sum_{i=1}^{|V|} \frac{N_i}{N_q} b(i)$

Show that maximizing Gain minimizes impurity (b) measure over the $|V|$ children.

$b(q)$ is always positive or $0$, so to maximize $Gain()$, we must minimize $\sum_{i=1}^{|V|} \frac{N_i}{N_q} b(i)$, which is subtracted from $b(q)$.

$\sum_{i=1}^{|V|} \frac{N_i}{N_q} b(i)$ : In order to minimize, $\left(\frac{N_i}{N_q}\right)$ and $(b(i))$ need to be minimized.

→ $\frac{N_i}{N_q}$ → $N_q$ is constant. A feature with $N_i$ values closer to $0$ or $N_q$ should be selected for, because this results more homogenous nodes leading to LOWER impurity.

→ $b(i)$ → High purity = Low impurity = low $\sum_{i=1}^{|V|} \frac{N_i}{N_q} b(i)$ so maximizing purity, and therefore Gain, can only occur through Minimizing the Impurity Measure.

⑤ $Gini(q) = \sum_{k \neq k'} P_k P_{k'}$

$$= \sum_{k=1}^{M} \left( P_k \times \sum_{k' \neq k} P_{k'} \right)$$

$$\longrightarrow (1-P_k)$$

$$= \sum_{k=1}^{M} \left( P_k \times (1-P_k) \right)$$

$$= \sum_{k=1}^{M} P_k (1-P_k) \quad \checkmark$$

$$\therefore \sum_{k=1}^{M} P_k (1-P_k) == \sum_{k \neq k'} P_k P_{k'} \quad \text{where } M > 2$$