# Loss Entropy:

## KL divergence:

Coin 1
$$\begin{cases} \% & \text{heads } (p_1) \\ \% & \text{tails } (p_r) \end{cases}$$

Coin 2
$$\begin{cases} p & \text{heads} \\ q & \text{tails} \end{cases}$$

$$\frac{P(\text{observations} \mid \text{real coin})}{P(m \mid \text{coin } 2)} \qquad i \left( \frac{p_1^{N_H} \, p_r^{N_T}}{q_1^{N_H} \, q_r^{N_T}} \right)^{\frac{1}{N}}$$

$$\frac{1}{N} \cdot \log \left( \frac{p_1^{N_H} \, p_r^{N_T}}{q_1^{N_H} \, q_r^{N_T}} \right)^{\frac{1}{N}} = \frac{1}{N} \log p_1^{N_H} + \frac{1}{N} \log p_r^{N_T}$$

$$- \frac{1}{N} \log q_1^{N_H} - \frac{1}{N} \log q_r^{N_T} \quad z \quad p_1 \log p_1 + p_r \log p_r$$

$$- p_1 \log q_1 - p_r \log q_r = p_1 \log \frac{p_1}{q_1} + p_r \log \frac{p_r}{q_r} \; ?$$

$$D_{KL}(P \| Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

Now **cross Entropy Loss:**

True class distribution ↑      predicted class distribution ↑

$$D_{KL}\left( \overbrace{P^*(y \mid x_i)}^{} \; \middle\| \; \overbrace{P(y \mid x_i ; \theta)}^{} \right) = \sum P^*(y \mid x_i) \log \frac{P^*(y \mid x_i)}{P(y \mid x_i ; \theta)}$$

$$= \underbrace{\sum_y p^*(y|x_i) \log p^*(y|x_i)}_{\text{Doesn't depend on } \theta} - \sum_y p^*(y|x_i) \log (y|x_i; \theta)$$

$$\Rightarrow \text{argmin} \; D_{KL}(p^* || P) \equiv \underset{\theta}{\text{arg min}} - \sum_y P^*(y|x_i) \log P(y|x_i; \theta)$$

$$\underset{\theta}{\text{argmin}} \; H(P^*, P)$$