

Theory

1.1 Information Entropy

Definition 1.1. **Accuracy gain** from splitting R into R_1 and R_2 based on loss $L(R)$: $L(R) - \frac{|R_1|L(R_1) + |R_2|L(R_2)}{|R_1| + |R_2|}$

Definition 1.2. Given a random variable $X \sim p$, the **entropy** measures the amount of randomness/uncertainty in an arbitrary realization of X .

$$H(X) := \mathbb{E}_{X \sim p}[-\log_2 p(X)] \quad (1.1)$$

Definition 1.3. Given joint distribution $(X, Y) \sim p(X, Y)$, the **entropy of joint distribution** is defined as

$$H(X, Y) := \mathbb{E}_{(X, Y) \sim p(X, Y)}[-\log_2 p(X, Y)] = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log_2 p(x, y) \quad (1.2)$$

Definition 1.4. Given two random variables X and Y , the **conditional entropy of Y conditioned on specific realization of X** is defined to be

$$H(Y|X = x) := \mathbb{E}_{y \sim p(y|X=x)}[-\log_2 p(y|X = x)] = - \sum_{y \in \mathcal{Y}} p(y|X = x) \log_2 p(y|X = x) \quad (1.3)$$

The **expected conditional entropy**¹ is defined as

$$H(Y|X) = \mathbb{E}_{X \sim p(x)}[H(Y|X)] = \mathbb{E}_{X \sim p(x)}[\mathbb{E}_{y \sim p(y|X=x)}[-\log_2 p(y|X = x)]] = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x) \quad (1.4)$$

$$= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \text{Im}(Y)} p(y|X = x) \log_2 p(y|X = x) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log_2 p(y|X = x) = - \mathbb{E}_{(X, Y) \sim p(x, y)}[\log_2 p(Y|X)] \quad (1.5)$$

Proposition 1.1. For every $X \in \Delta(\mathcal{X})$, $H(X) \geq 0$.

Proposition 1.2 (Chain Rule). $H(X, Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$

Proposition 1.3. If $X \perp Y$, then knowing X does not provide extra information (i.e. reduce entropy) of Y . That is $H(Y|X) = H(Y)$.

Proposition 1.4. Y becomes deterministic by knowing Y , that is, $H(Y|Y) = 0$.

Proposition 1.5. By knowing X , the uncertainty about Y is reduced: $H(Y|X) \leq H(Y)$.

Definition 1.5. The **information gain** in Y due to X , or **mutual information** of X and Y is defined to be

$$IG(Y|X) := H(Y) - H(Y|X) \quad (1.6)$$

When X is completely uninformative about Y : $H(Y|X) = H(Y)$, then $IG(Y|X) = 0$.

When X is completely information about Y : $H(Y|X) = 0$ (deterministic), then $IG(Y|X) = H(Y)$.

Proposition 1.6 (Symmetry of Information Gain).

$$IG(Y|X) := H(Y) - H(Y|X) = H(X, Y) - H(X|Y) - H(Y|X) \quad (1.7)$$

$$= H(Y|X) + H(X) - H(X|Y) - H(Y|X) = H(X) - H(X|Y) = IG(X|Y) \quad (1.8)$$

BVD: Deterministic

$$\mathbb{E}_{x, \mathcal{D}} \left[(h_{\mathcal{D}}(x) - f(x))^2 \right] = \mathbb{E}_{x, \mathcal{D}} \left[(h_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x])^2 \right] + \mathbb{E}_x \left[(\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x] - f(x))^2 \right] \quad (1.9)$$

BVD: Stochastic Let $(\mathbf{x}^{(i)}, y^{(i)}) \in \mathcal{X} \times \mathbb{R}$ denote one training instance such that $(\mathbf{x}^{(i)}, y^{(i)}) \stackrel{i.i.d.}{\sim} p_{\text{sample}}$, where $p_{\text{sample}} \in \Delta(\mathcal{X} \times \mathbb{R})$.

Fixing $N \in \mathbb{N}$, one can construct a new distribution $p_{\text{dataset}} \in \Delta(\mathcal{X} \times \mathbb{R})^N$ such that $(\mathbf{x}^{(i)}, y^{(i)})_{i=1}^N =: \mathcal{D} \sim p_{\text{dataset}}$. Given a (random) training set \mathcal{D} , a (random) classifier function $h_{\mathcal{D}} \in \mathcal{H}$ is generated.

For every *query point* $\mathbf{x} \in \mathcal{X}$, the prediction $h_{\mathcal{D}}(\mathbf{x})$ is therefore random.

Suppose y is not deterministic in x , then the expected mean squared error when the model is applied on new instances sampled from p_{sample} is

$$\mathbb{E}_{\mathbf{x}, y, \mathcal{D}}[(h_{\mathcal{D}}(\mathbf{x}) - y)^2] = \mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\mathbf{x}, y}[(h_{\mathcal{D}}(\mathbf{x}) - y)^2 | \mathcal{D}]] = \mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\mathbf{x}, y}[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x] + \mathbb{E}_y[y|x] - y)^2 | \mathcal{D}]] \quad (1.10)$$

$$= \mathbb{E}_{\mathcal{D}}\{\mathbb{E}_x[\mathbb{E}_y[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])^2]] + 2\mathbb{E}_{x, y}[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])(\mathbb{E}_y[y|x] - y)] + \mathbb{E}_{x, y}[(\mathbb{E}_y[y|x] - y)^2]\} \quad (1.11)$$

$$= \mathbb{E}_{\mathcal{D}}\{\mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])^2] + 2\mathbb{E}_{x, y}[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])(\mathbb{E}_y[y|x] - y)] + \mathbb{E}_{x, y}[(\mathbb{E}_y[y|x] - y)^2]\} \quad (1.12)$$

$$(1.13)$$

By law of iterative expectation,

$$\mathbb{E}_{x, y}[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])(\mathbb{E}_y[y|x] - y)] = \mathbb{E}_x[\mathbb{E}_y[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])(\mathbb{E}_y[y|x] - y)]] = \mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])(\mathbb{E}_y[y|x] - \mathbb{E}_y[y])] = 0 \quad (1.14)$$

¹This is independent of specific realization of X

By dropping irrelevant expectation operators,

$$\Delta = \mathbb{E}_{\mathcal{D}}[\mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])^2]] + \underbrace{\mathbb{E}_{x,y}[(\mathbb{E}_y[y|x] - y)^2]}_{\text{Bayes Error } \varepsilon^2} = \mathbb{E}_{\mathcal{D}}[\mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_y[y|x])^2]] + \varepsilon^2 \quad (1.15)$$

$$= \mathbb{E}_{\mathcal{D}}[\mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x] + \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x] - \mathbb{E}_y[y|x])^2]] + \varepsilon^2 \quad (1.16)$$

Note that $\mathbb{E}_{\mathcal{D}}[\mathbb{E}_x[(h_{\mathcal{D}}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x])(\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x] - \mathbb{E}_y[y|x])]] = 0$ The first component reduced to zero after applying law of iterative expectation. Non-deterministic case

$$\mathbb{E}_{x,y,\mathcal{D}}[(h_{\mathcal{D}}(x) - y)^2] = \mathbb{E}_{x,\mathcal{D}}[(h_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x])^2] + \mathbb{E}_x[(\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)|x] - \mathbb{E}_y[y|x])^2] + \mathbb{E}_{x,y}[(\mathbb{E}_y[y|x] - y)^2] \quad (1.17)$$

2 Mathematics & Probability

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \det(\Sigma)^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$

$$\text{Var}(X) = \mathbb{E}[(X - \mu)(X - \mu)^T] \in \mathbb{R}^{d \times d}$$

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)^T] \in \mathbb{R}^{d \times d} \quad p(\theta | \text{data}) = \frac{p(\text{data} | \theta)p(\theta)}{p(\text{data})} \quad \theta^{\text{MAP}} = \underset{\theta}{\text{argmax}} p(\theta | \text{data}) = \underset{\theta}{\text{argmax}} p(\text{data} | \theta)p(\theta)$$

$$\theta^{\text{MAP}} = \underset{\theta}{\text{argmax}} p(X_1, \dots, X_N | \theta) p(\theta) = \underset{\theta}{\text{argmax}} p(\theta) \prod_{i=1}^N p(X_i | \theta) = \underset{\theta}{\text{argmax}} \log p(\theta) + \sum_{i=1}^N \log p(X_i | \theta)$$

Proposition 2.1 (Law of Total Expectation). $\mathbb{E}_Y[\mathbb{E}_{X|Y}[X|Y]] = \mathbb{E}[X]$.

Proof. $\mathbb{E}[\mathbb{E}[X|Y]] = \int \left[\int x p(x|y) dx \right] p(y) dy = \iint x p(x, y) dx dy = \mathbb{E}[X]$ ■

$$\underset{\mathbf{w}}{\text{minimize}} \mathcal{J}(\mathbf{w}) =: \frac{1}{2} \|\mathbf{t} - \mathbf{X}\mathbf{w}\|_2^2 \quad \mathcal{J}(\mathbf{w}) = \frac{1}{2} \|\mathbf{t}\|_2^2 + \frac{1}{2} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{t}^\top \mathbf{X} \mathbf{w}.$$

Theorem 2.1 (Bayes Optimal). $\underset{y}{\text{argmin}} \mathbb{E}[(y - t)^2 | \mathbf{x}] = \mathbb{E}[t | \mathbf{x}]$ where $t \sim p(t | \mathbf{x})$.

Multi-class Classification $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b}$ (aka *logits*) Input dim = D , output dim = K , $\mathbf{W} \in \mathbb{R}^{K \times D}$

Pred.prob: $y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}$ $\mathcal{L}_{\text{CE}}(\mathbf{y}, \mathbf{t}) = -\sum_{k=1}^K t_k \log y_k = -\mathbf{t}^T (\log(\mathbf{y}))$ (*Softmax-cross-entropy*).

$$\frac{\partial \mathcal{L}_{\text{CE}}}{\partial \mathbf{w}_k} = \frac{\partial \mathcal{L}_{\text{CE}}}{\partial z_k} \cdot \frac{\partial z_k}{\partial \mathbf{w}_k} = (y_k - t_k) \cdot \mathbf{x}, \quad \mathbf{w}_k \leftarrow \mathbf{w}_k - \alpha \frac{1}{N} \sum_{i=1}^N \left(y_k^{(i)} - t_k^{(i)} \right) \mathbf{x}^{(i)}, \quad \mathbf{W} \leftarrow \mathbf{W} - \frac{\alpha}{N} (\mathbf{y} - \mathbf{t}) \mathbf{X}$$

(Forward)	(Backward)	(Forward: Reg)	(Backward I)	(Backward II)	$\mathbf{g}_i = \nabla \mathcal{L}(\mathbf{w}, \mathbf{x}_i, t_i)$
$\mathbf{z} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$	$\bar{\mathcal{L}} = 1$	$z = \mathbf{w}x + b$	$\overline{\mathcal{L}_{\text{reg}}} = 1$	$\bar{z} = \bar{y} \frac{dy}{dz}$	(GD)
$\mathbf{h} = \sigma(\mathbf{z})$	$\bar{\mathbf{y}} = \bar{\mathcal{L}}(\mathbf{y} - \mathbf{t})$	$y = \sigma(z)$	$\bar{\mathcal{R}} = \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{R}}$	$= \bar{y} \sigma'(z)$	$\mathbf{w} \leftarrow \mathbf{w} - \eta \sum_{i=1}^N g_i$
$\mathbf{y} = \mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)}$	$\overline{\mathbf{W}^{(2)}} = \bar{\mathbf{y}} \mathbf{h}^\top$	$\mathcal{L} = \frac{1}{2} (y - t)^2$	$= \overline{\mathcal{L}_{\text{reg}}} \lambda$	$\bar{w} = \bar{z} \frac{\partial z}{\partial w} + \bar{\mathcal{R}} \frac{d\mathcal{R}}{dw}$	(SGD)
$\mathcal{L} = \frac{1}{2} \ \mathbf{t} - \mathbf{y}\ ^2$	$\overline{\mathbf{h}^{(2)}} = \bar{\mathbf{y}}$	$\mathcal{R} = \frac{1}{2} w^2$	$\bar{\mathcal{L}} = \overline{\mathcal{L}_{\text{reg}}} \frac{d\mathcal{L}_{\text{reg}}}{d\mathcal{L}}$	$= \bar{z} x + \bar{\mathcal{R}} w$	$i \sim \mathcal{U}[1, N],$
	$\bar{\mathbf{h}} = \mathbf{W}^{(2)\top} \bar{\mathbf{y}}$	$\mathcal{L}_{\text{reg}} = \mathcal{L} + \lambda \mathcal{R}$	$= \overline{\mathcal{L}_{\text{reg}}}$	$\bar{b} = \bar{z} \frac{\partial z}{\partial b}$	$\mathbf{w} \leftarrow \mathbf{w} - \eta g_i$
	$\bar{\mathbf{z}} = \bar{\mathbf{h}} \circ \sigma'(\mathbf{z})$		$\bar{y} = \bar{\mathcal{L}} \frac{d\mathcal{L}}{dy}$	$= \bar{z}$	(mSGD)
	$\overline{\mathbf{W}^{(1)}} = \bar{\mathbf{z}} \mathbf{x}^\top$		$= \bar{\mathcal{L}}(y - t)$		$M \subset \{1, \dots, N\},$
	$\overline{\mathbf{b}^{(1)}} = \bar{\mathbf{z}}$				$\mathbf{w} \leftarrow \mathbf{w} - \eta \sum_{i \in M} g_i$

3 Misc

1. Activation functions $\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$ $\sigma(z) = \frac{1}{1 + \exp(-z)}$ $\text{ReLU}(z) = \max(0, z)$.
2. **Parametric Benefits** (i) Simpler (interpretability) (ii) Speed (iii) Less Data;
Drawbacks (i) Constrained (ii) Limited Complexity (iii) Poor fit.
3. **Non-parametric Benefits** (i) Flexibility (ii) Power (No prior assumptions) (iii) Performance;
Drawbacks (i) More data (ii) Slower (iii) Overfitting.
4. Decision of linear models: $\mathbf{W} \cdot \mathbf{x} + \mathbf{b} = \mathbf{0}$ (hyperplane).