

# Entropy

The number of arrangements of  $n$  states  $\{1, \dots, s\}$  that yield a distribution  $\hat{p}$ :

$$C(\hat{p}) = \binom{n}{\hat{p}_1 n, \dots, \hat{p}_s n} = \frac{n!}{(\hat{p}_1 n)! \cdots (\hat{p}_s n)!}$$

**Stirling's Approximation:**

$$k! \approx k^k e^{-k} \sqrt{2\pi k}$$

**Shannon Entropy:**

- For  $p$  a pmf,  $H(p) = -\sum_{x=1}^s p(x) \log p(x)$
- For  $p$  a pdf,  $H(p) = -\int_{-\infty}^{\infty} p(x) \log p(x) dx$

**Entropy Approximation:**  $C(\hat{p}) \approx e^{nH(p)}$

**Maximum Entropy Principle:** Let  $p$  satisfy

1.  $\sum_{x=1}^s p(x) = 1$
2.  $\sum_{x=1}^s p(x) \mathcal{E}(x) \approx \theta$
3.  $p$  maximizes  $H(p)$

then  $p$  has the form

$$p(x) = \frac{1}{Z_\lambda} e^{\lambda \mathcal{E}(x)} = \frac{1}{\sum_{x=1}^s e^{\lambda \mathcal{E}(x)}} e^{\lambda \mathcal{E}(x)}$$

where  $\lambda$  is found via the constrain  $\sum p(x) \mathcal{E}(x) = \theta$ .

In the case of seeking  $\arg \max_p H(p)$  subject to constraints  $\sum p_x = 1$  and  $\sum p_x \mathcal{E}_i(x) = \theta_i$  for  $i = 1 : k$ ,  $p$  will have form

$$p(x) = \frac{1}{Z_\lambda} \exp\left[\sum_{i=1}^k \lambda_i \mathcal{E}_i(x)\right]$$

**Large Deviation Principle:** For  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p$  a pmf on  $\{1, \dots, s\}$ , if  $\mathcal{E} : \{1, \dots, s\} \rightarrow \mathbb{R}$  satisfies  $\frac{1}{n} \sum_{i=1}^n \mathcal{E}(X_k) = \theta$ , then  $\mathbb{E}_{\hat{p}}[\mathcal{E}(X)] = \theta$

**Observations:** For  $q$  a distribution on  $\{1, \dots, s\}$ ,

$$\mathbb{P}(\hat{p} = q) = \binom{n}{n_1, \dots, n_s} \prod_{x=1}^s p_x^{q_x \cdot n}$$

Further,

$$\frac{e^{-nD(q||p)}}{(n+1)^s} \leq \mathbb{P}(\hat{p} = q) \leq e^{-nD(q||p)}$$

**Kullback-Leibler Divergence:**

$$D(q || p) = -\sum_{x=1}^s q_x \log \frac{p_x}{q_x}$$

- $D(q || p) \geq 0$
- $D(q || p) = 0 \iff p = q$

**Convexity:**  $f$  convex if  $f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$  for  $\lambda \in [0, 1]$ .

- $f$  convex iff  $f''(x) \geq 0$
- $f$  concave iff  $f''(x) \leq 0$  iff  $-f$  is convex
- For  $x \in \mathbb{R}^s$ ,  $f$  convex iff  $h(\lambda) = f(\lambda x + (1-\lambda)y)$  convex

**Jensen's Inequality:** For  $g : \mathbb{R} \rightarrow \mathbb{R}$  convex,  $\mathbb{E}[g(X)] \geq g(\mathbb{E}[X])$

**Sanov's Theorem:** For  $B$  an open subset of the space of distributions on  $\{1, \dots, s\}$ ,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}(\hat{p} \in B) = -\inf_{q \in B} D(q || p)$$

Further, if  $p^* = \arg \min_{q \in B} D(q || p)$  and  $\hat{p} \in B$ ,  $\hat{p} \xrightarrow{\mathbb{P}} p^*$

**Exponential Families:**

$$p(x) = \frac{1}{Z(\lambda)} h(x) e^{\lambda \cdot T(x)}$$

where  $Z(\lambda)$  satisfies:

1.  $\frac{\partial}{\partial \lambda_k} \log Z_k \mathbb{E}_{T_k}(X)$
2.  $\frac{\partial^2}{\partial \lambda_k \partial \lambda_j} \log Z_k = \text{Cov}_{T_k(X), T_j(X)}$
3.  $\log Z_k$  is convex in  $\lambda$  and strictly convex unless  $\exists a \in \mathbb{R}^k$  such that  $\sum a_k T_k(x) = b$  for  $b$  constant.
4.  $\log Z(\lambda) - \sum \lambda_k \theta_k$  is convex in  $\lambda$  and minimized when  $\mathbb{E}[T(X)] = \theta_k$ .

## Source Coding

**Source Code:**  $C : \{1, \dots, t\} \rightarrow \{0, 1\}^*$

**Prefix Code:** a code  $C$  for which  $C(x)$  is not a prefix of  $C(y)$  for  $x \neq y$

**Kraft-McMillan Inequality:** For all prefix codes  $C$ ,

$$\sum_{x=1}^t 2^{-|C(x)|} \leq 1$$

and for any code lengths  $\ell_1, \dots, \ell_t$  satisfying

$$\sum_{x=1}^t 2^{-\ell_x} \leq 1$$

there exists a prefix code  $C$  with  $|C(x)| = \ell_x$

**Optimal Coding:** Let  $\vec{X} \sim p$ . For the optimal code  $C^* = \arg \min_C \text{prefix} \mathbb{E}_p |C(X)|$ ,

$$H(p) \leq \mathbb{E}_p |C(X)| \leq H(p) + 1$$

# Statistical Inference

**Unbiased Estimator:**  $\mathbb{E}[\hat{\theta}] = \theta$

**Consistent Estimator:**

- Almost sure consistency:  $\mathbb{P}(\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta) = 1$
- Consistent in probability:  $\forall \varepsilon > 0, \mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0$
- Consistent in mean square:  $\mathbb{E}[(\hat{\theta}_n - \theta)^2] \rightarrow 0$ .

**Mean Square Error:**

$$\text{MSE}(\hat{\theta}_n) = \mathbb{E}[(\hat{\theta}_n - \theta)^2] = \text{Var}[\hat{\theta}] + \text{Bias}(\hat{\theta})^2$$

**Theorem:**  $\text{MSE}[\hat{\theta}_n] \rightarrow 0 \implies \mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0$

**Bias:**  $\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta$

**Variance:**  $\text{Var}[\hat{\theta}] = \mathbb{E}[\hat{\theta}^2] - \mathbb{E}[\hat{\theta}]^2$

**Kernel Density Estimation:** for a kernel  $k$  satisfying

1.  $k(x) \geq 0$
2.  $\int xk(x) dx = 0$
3.  $\int x^2k(x) dx = 1$

we define the kernel density estimator

$$\hat{f}_{n,w}(x; X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n k_w(x - X_i) = \frac{1}{n} \sum_{i=1}^n \frac{1}{w} k\left(\frac{x - X_i}{w}\right)$$

**Convolution:** Let  $Z \sim f$  and  $Y \sim g$  be independent. Then

$$Z + Y \sim (f \star g)(x) = \int_{\mathbb{R}} f(t)g(x - t) dt$$

**Integrated Square Error:**

$$\text{ISE}(\hat{f}) = \int_{\mathbb{R}} \left| \hat{f}_n(x; X_{1:n}) - f(x) \right|^2 dx$$

**Mean Integrated Square Error:**

$$\text{MISE}(\hat{f}) = \mathbb{E}[\text{ISE}(\hat{f})] = \int_{\mathbb{R}} \mathbb{E} \left| \hat{f}_n(x; X_{1:n}) - f(x) \right|^2 dx$$

**Asymptotics:** For  $f, k$  smooth, as  $w \rightarrow 0$ ,

$$\text{MISE}(\hat{f}_{n,w}) = \alpha w^4 + \frac{\beta}{nw} + \text{error}$$

**Silverman's Rule of Thumb:** For parameters  $\alpha, \beta$  unknown, choose the kernel bandwidth  $w \propto n^{-1/5}$

**Cross-Validation Estimator:**

$$\hat{f}_{n,w}^{(i)}(x; X_{1:n}) = \frac{1}{n} \sum_{j \neq i} \hat{f}_{n-1,w}(X_j)$$

**Stone's Theorem:** For

$$\hat{w}_n = \arg \min_w \int \hat{f}_{n,w}^2(x) - \frac{2}{n} \sum_{i=1}^n \hat{f}_{n-1,w}^{(i)}(X_i) dx$$

we have

$$\text{ISE}(\hat{f}_{\hat{w}_n}) \xrightarrow{a.s.} \inf_w \text{ISE}(\hat{f}_{w,n}, f)$$

**Maximum Likelihood Estimation:** Sample  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} p_\theta$  for  $\theta$  unknown. Then

$$\hat{\theta} = \arg \max_{\theta} p_\theta(X_1 = x_1, \dots, X_n = x_n) = \arg \min_{\theta} D(\hat{p} \| p_\theta)$$

## Classification

Let  $X \in \mathbb{R}^d$  be a RV with classes  $Y \in \{1, \dots, c\}$ . Define  $\pi_i = \mathbb{P}(Y = i)$  and  $f_i(x)$  the class conditioned density.

Then for any set  $A$ ,

$$\mathbb{P}_X(A) = \sum_{i=1}^c \pi_i \mathbb{P}(A)$$

$$f_x(A) = \sum_{i=1}^c \pi_i f_i(A)$$

**Bayes' Classification Rule:**

$$h^*(x) = \arg \max_{i=1:c} \mathbb{P}(Y = i | X = x)$$

**Neyman-Pearson Classification:** Fix  $t \in (0, \infty)$ . Then

$$h_t(x) = \begin{cases} 1 & \text{if } \frac{f_2(x)}{f_1(x)} < t \\ 2 & \text{if } \frac{f_2(x)}{f_1(x)} > t \end{cases}$$

**Note:** In the case  $t = \frac{\pi_1}{\pi_2}$ , then NP is equivalent to Bayes'.

**Detection Rate:**  $\mathbb{P}(h(X) = 2 | Y = 2)$

**False Alarm Rate:**  $\mathbb{P}(h(X) = 2 | Y = 1)$

**Theorem:** Fix  $t \in (0, \infty)$  and choose any classifier  $h$ . If  $\text{FAR}(h) \leq \text{FAR}(h_t)$ , then  $\text{DR}(h) \leq \text{DR}(h_t)$

## Generative Classifiers

*Motivation:* From Bayes,

$$h^*(x) = \arg \max_{c \in \{1, \dots, s\}} \mathbb{P}(Y = c | X = x) = \arg \max_c \frac{\pi_c f_c(x)}{\mathbb{P}(X = x)}$$

so we can estimate  $\hat{\pi}_c \approx \pi_c$  and then it suffices to estimate  $f_c(x)$ .

**Naive Bayes:** Assume  $f_c(x_1, \dots, x_d) = \prod_{j=1}^d f_{c,j}(x_j)$ . Then instead of needing to find  $(f_c)^s$  with  $f_c: \mathbb{R}^d \rightarrow \mathbb{R}$ , it suffices to find  $f_{c,j}: \mathbb{R} \rightarrow \mathbb{R}$ .

**Quadratic Discriminant Analysis:** Assume  $f_c(x) \sim \mathcal{N}(\mu_c, \Sigma_c)$ . Then

$$f_c(x, \mu_c, \Sigma_c) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{(x - \mu_c) \Sigma_c^{-1} (x - \mu_c)}{2} \right)$$

for  $x, \mu_c in \mathbb{R}^d$  and  $\Sigma_c \in \mathbb{R}^{d \times d}$ .

**Linear Discriminant Analysis:** Assume  $f_c(x) \sim \mathcal{N}(\mu_c, \Sigma)$ .

**Discriminative Classifiers:**

*Motivation:*

$$h^*(x) = \arg \max_{c \in \{1, \dots, s\}} \mathbb{P}(Y = c \mid X = x) = \arg \max_c r_c(x)$$

**Linear Regression:** Assume  $s = 2$  so  $r_1(x) + r_2(x) = 1$ . Then

$$\log \frac{r_2(x)}{r_1(x)} = \alpha + \beta x \implies r_2 = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

and it suffices to use MLE to estimate  $\alpha, \beta$ .

**Softmax:** Assume  $s > 2$ . Then model

$$\log \frac{r_k(x)}{r_1(x)} = \alpha_k + \beta_k x \implies r_k(x) = \frac{e^{\alpha_k + \beta_k x}}{1 + \sum_{k=n}^s e^{\alpha_k + \beta_k x}}$$