

25 Asymptotics of Inference and Universality

In these notes, we develop the asymptotics of inference. In particular, we examine the behavior of model capacity in the regime of a large number of samples, and identify the asymptotic least-informative prior.

Finally, we apply these results to the problem of universal inference, and sharply characterize when such inference is and is not possible.

25.1 Asymptotics of Model Capacity

To characterize the asymptotic behavior of the model capacity, we use the Bayesian terminology and analyze the mutual information between the hidden random variable \mathbf{x} and the observed variable \mathbf{y} , which consists of N i.i.d. random variables distributed according to $p_{\mathbf{y}|\mathbf{x}}(\cdot; x)$. We assume that \mathbf{x} is distributed according to the prior $p_{\mathbf{x}}(\cdot)$.

As the number of samples changes, the likelihood model of the full sequence \mathbf{y} also changes. Naturally, we would expect the capacity of the model and the least informative prior to depend on the number of samples. However, this creates an inconvenience of having to change the least informative prior as the amount of available data changes. In practice, for large N , we could use as a prior a probability distribution that is equal to the limit of the least informative prior when $N \rightarrow \infty$. Such priors are called Jeffreys' priors. Our goal in this section is to characterize the dependency of the model capacity on the number of samples N and to determine the Jeffreys' prior for a broad class of problems.

We first represent the mutual information as the mean divergence between the likelihood model and the marginal distribution of \mathbf{y} ,

$$I(\mathbf{x}; \mathbf{y}) = \mathbb{E}_{p_{\mathbf{x}}(\cdot)} \left[\mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y},x}(\mathbf{y}, x)}{p_{\mathbf{y}}(\mathbf{y})p_{\mathbf{x}}(x)} \middle| x = x \right] \right] \quad (1)$$

$$= \mathbb{E}_{p_{\mathbf{x}}(\cdot)} \left[\mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}}(\mathbf{y})} \middle| x = x \right] \right] \quad (2)$$

$$= \mathbb{E}_{p_{\mathbf{x}}(\cdot)} [D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x) \| p_{\mathbf{y}}(\cdot))] \quad (3)$$

and analyze the behavior of $D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x) \| p_{\mathbf{y}}(\cdot))$ as $N \rightarrow \infty$.

25.1.1 Model Capacity for Finite Model Classes

If \mathbf{x} is a random variable defined over a finite alphabet \mathcal{X} ,

$$D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x)||p_{\mathbf{y}}(\cdot)) = \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|x}(\mathbf{y}|x)}{p_{\mathbf{y}}(\mathbf{y})} \right] \quad (4)$$

$$= \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{\sum_{x'} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x') p_{\mathbf{x}}(x')} \right] \quad (5)$$

$$= \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{1}{p_{\mathbf{x}}(x) + \sum_{x' \neq x} p_{\mathbf{x}}(x') \frac{p_{\mathbf{y}|x}(\mathbf{y}|x')}{p_{\mathbf{y}|x}(\mathbf{y}|x)}} \right]. \quad (6)$$

Since \mathcal{X} is finite, the sum in the denominator contains a finite number of terms. For any sequence \mathbf{y} generated i.i.d. from the distribution $p_{\mathbf{y}|\mathbf{x}}(\cdot|x)$, each term in the sum decays exponentially with the number of samples N :

$$\frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x')}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)} \doteq e^{-ND(p_{\mathbf{y}|\mathbf{x}}(\cdot|x')||p_{\mathbf{y}|\mathbf{x}}(\cdot|x))}. \quad (7)$$

Therefore, the whole sum decays exponentially. For large N ,

$$D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x)||p_{\mathbf{y}}) \cong \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{1}{p_{\mathbf{x}}(x)} \right] = -\log p_{\mathbf{x}}(x) \quad (8)$$

and

$$I(\mathbf{x}; \mathbf{y}) = \mathbb{E}_{p_{\mathbf{x}}(\cdot)} [D(p_{\mathbf{y}|\mathbf{x}}(\cdot|\mathbf{x})||p_{\mathbf{y}}(\cdot))] \cong \mathbb{E}_{p_{\mathbf{x}}(\cdot)} [-\log p_{\mathbf{x}}(\mathbf{x})] = H(p_{\mathbf{x}}). \quad (9)$$

The model capacity is equal to

$$C_N = \max_{p_{\mathbf{x}}} I(\mathbf{x}; \mathbf{y}) \cong \max_{p_{\mathbf{x}}} H(p_{\mathbf{x}}) = \log |\mathcal{X}| \quad (10)$$

and is achieved for $p_{\mathbf{x}}(x) = |\mathcal{X}|^{-1}$. Since neither the model capacity nor the prior that achieves capacity depends on N , this is also the Jeffreys' prior:

$$p_{\mathbf{x}}^{\infty}(x) = \frac{1}{|\mathcal{X}|}. \quad (11)$$

25.1.2 Model Capacity for Continuous Model Classes

The analysis for a continuous parameter \mathbf{x} is a bit more involved. We present it here for completeness, but note that the approximations taken in the derivations require particular care in specifying the conditions on the likelihood model to ensure that such approximations are valid.

We start by manipulating the log ratio of the likelihood and the marginal distribution of \mathbf{y} :

$$D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x)||p_{\mathbf{y}}(\cdot)) = \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}}(\mathbf{y})} \right] \quad (12)$$

$$= \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)} \right] + \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)}{p_{\mathbf{y}}(\mathbf{y})} \right] \quad (13)$$

where \hat{x}_N is the ML estimate of the parameter \mathbf{x} based on the observations \mathbf{y} .

Using the Taylor series expansion of the log-likelihood function around the ML estimate \hat{x}_N , we obtain

$$\frac{1}{N} \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x) = \frac{1}{N} \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N) + \frac{1}{N} (x - \hat{x}_N) \left. \frac{\partial \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{\partial x} \right|_{\hat{x}_N} \quad (14)$$

$$+ \frac{1}{2N} (x - \hat{x}_N)^2 \left. \frac{\partial^2 \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{\partial x^2} \right|_{\hat{x}_N} + o((x - \hat{x}_N)^2) \quad (15)$$

$$= \frac{1}{N} \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N) \quad (16)$$

$$+ \frac{1}{2N} (x - \hat{x}_N)^2 \sum_{n=1}^N \left. \frac{\partial^2 \log p_{y|x}(y_n|x)}{\partial x^2} \right|_{\hat{x}_N} + o((x - \hat{x}_N)^2), \quad (17)$$

which implies

$$\frac{1}{N} \log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)} = \frac{1}{2} (x - \hat{x}_N)^2 \frac{1}{N} \sum_{n=1}^N \left. \frac{\partial^2 \log p_{y|x}(y_n|x)}{\partial x^2} \right|_{\hat{x}_N} + o((x - \hat{x}_N)^2). \quad (18)$$

Since

$$\hat{x}_N \xrightarrow{\text{w.p.1}} x \quad (19)$$

$$\frac{1}{N} \sum_{n=1}^N \left. \frac{\partial^2 \log p_{y|x}(y_n|x)}{\partial x^2} \right|_{\hat{x}_N} \xrightarrow{\text{w.p.1}} \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\left. \frac{\partial^2 \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{\partial x^2} \right|_{\hat{x}_N} \right] = -J_{\mathbf{y}}(\hat{x}_N), \quad (20)$$

for sufficiently large N

$$\frac{1}{N} \log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)} \cong -\frac{1}{2} J_{\mathbf{y}}(\hat{x}_N) (x - \hat{x}_N)^2. \quad (21)$$

This implies that the first term in the right-hand side of (13) is equal to

$$\mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)} \right] \cong \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[-\frac{1}{2} (N J_{\mathbf{y}}(\hat{x}_N)) (x - \hat{x}_N)^2 \right] \quad (22)$$

$$\cong \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|x)} \left[-\frac{1}{2} (N J_{\mathbf{y}}(x)) (x - \hat{x}_N)^2 \right] \quad (23)$$

$$\cong -\frac{1}{2} (N J_{\mathbf{y}}(x)) \text{var}(\hat{x}_N) \cong -\frac{1}{2}. \quad (24)$$

We can also use (21) to show that

$$\frac{1}{N} \log \frac{p_{\mathbf{x}|\mathbf{y}}(x|\mathbf{y})}{p_{\mathbf{x}|\mathbf{y}}(\hat{x}_N|\mathbf{y})} = \frac{1}{N} \log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|x)}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)} + \frac{1}{N} \log \frac{p_{\mathbf{x}}(x)}{p_{\mathbf{x}}(\hat{x}_N)} \quad (25)$$

$$\cong -\frac{1}{2} J_{\mathbf{y}}(\hat{x}_N)(x - \hat{x}_N)^2 + \frac{1}{N} \log \frac{p_{\mathbf{x}}(x)}{p_{\mathbf{x}}(\hat{x}_N)}. \quad (26)$$

As $N \rightarrow \infty$, the effect of the second term is clearly removed and

$$\frac{1}{N} \log \frac{p_{\mathbf{x}|\mathbf{y}}(x|\mathbf{y})}{p_{\mathbf{x}|\mathbf{y}}(\hat{x}_N|\mathbf{y})} \cong -\frac{1}{2} J_{\mathbf{y}}(\hat{x}_N)(x - \hat{x}_N)^2. \quad (27)$$

By exponentiating both sides of the equation, we realize that

$$p_{\mathbf{x}|\mathbf{y}}(x|\mathbf{y}) \propto e^{-\frac{1}{2} N J_{\mathbf{y}}(\hat{x}_N)(x - \hat{x}_N)^2}, \quad (28)$$

i.e., as $N \rightarrow \infty$, the posterior probability distribution of the hidden parameter \mathbf{x} can be approximated by a Gaussian distribution with mean \hat{x}_N and variance $(N J_{\mathbf{y}}(\hat{x}_N))^{-1}$:

$$p_{\mathbf{x}|\mathbf{y}}(x|\mathbf{y}) = \frac{(N J_{\mathbf{y}}(\hat{x}_N))^{1/2}}{(2\pi)^{1/2}} e^{-\frac{1}{2} N J_{\mathbf{y}}(\hat{x}_N)(x - \hat{x}_N)^2}. \quad (29)$$

Now we analyze the asymptotic behavior of the second term in (13):

$$\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{x}_N)}{p_{\mathbf{y}}(\mathbf{y})} = \log \frac{p_{\mathbf{x}|\mathbf{y}}(\hat{x}_N|\mathbf{y})}{p_{\mathbf{x}}(\hat{x}_N)} \quad (30)$$

$$\cong \frac{1}{2} \log \frac{N}{2\pi} + \log \sqrt{J_{\mathbf{y}}(\hat{x}_N)} - \log p_{\mathbf{x}}(\hat{x}_N) \quad (31)$$

$$\cong \frac{1}{2} \log \frac{N}{2\pi} + \log \frac{\sqrt{J_{\mathbf{y}}(x)}}{p_{\mathbf{x}}(x)}. \quad (32)$$

We made use of $\hat{x}_N \xrightarrow{\text{w.p.1}} x$ to arrive at the last approximation. Substituting into (13), we obtain

$$D(p_{\mathbf{y}|\mathbf{x}}(\cdot|x) \| p_{\mathbf{y}}(\cdot)) \cong -\frac{1}{2} + \frac{1}{2} \log \frac{N}{2\pi} + \log \frac{\sqrt{J_{\mathbf{y}}(x)}}{p_{\mathbf{x}}(x)} \quad (33)$$

$$= \frac{1}{2} \log \frac{N}{2\pi e} + \log \frac{\sqrt{J_{\mathbf{y}}(x)}}{p_{\mathbf{x}}(x)}. \quad (34)$$

Here we assume natural logarithms.

To compute the model capacity, we evaluate

$$I(\mathbf{x}; \mathbf{y}) = \mathbb{E}_{p_{\mathbf{x}}(\cdot)} [D(p_{\mathbf{y}|\mathbf{x}}(\cdot|\mathbf{x}) \| p_{\mathbf{y}}(\cdot))] \cong \frac{1}{2} \log \frac{N}{2\pi e} + \mathbb{E}_{p_{\mathbf{x}}(\cdot)} \left[\log \frac{\sqrt{J_{\mathbf{y}}(x)}}{p_{\mathbf{x}}(x)} \right] \quad (35)$$

$$= \frac{1}{2} \log \frac{N}{2\pi e} - D(p_{\mathbf{x}}(\cdot) \| q(\cdot)), \quad (36)$$

where

$$q(x) = \frac{\sqrt{J_y(x)}}{\int_x \sqrt{J_y(x)} dx}. \quad (37)$$

The model capacity is equal to

$$C_N = \max_{p_x} I(\mathbf{x}; \mathbf{y}) \cong \frac{1}{2} \log \frac{N}{2\pi e}. \quad (38)$$

The Jeffreys' prior is equal to $q(\cdot)$:

$$p_x^\infty(x) = \frac{\sqrt{J_y(x)}}{\int_x \sqrt{J_y(x)} dx}. \quad (39)$$

We note that $p_x^\infty(x)$ might be an improper prior if the integral of $\sqrt{J_y(x)}$ does not exist.

25.1.3 Vector Parameters

We can also extend the analysis to the case of a vector parameter \mathbf{x} whose dimensionality is k . The derivation uses the Fisher information matrix in place of the Fisher information function. It is easy to show that (21) becomes

$$\frac{1}{N} \log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{\mathbf{x}}_N)} \cong -\frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}}_N)^T J_y(\hat{\mathbf{x}}_N)(\mathbf{x} - \hat{\mathbf{x}}_N). \quad (40)$$

Since the covariance of the ML estimator is equal to $(NJ_y(\mathbf{x}))^{-1}$, the first term in the right-hand side of (13) becomes

$$\mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|\mathbf{x})} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})}{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{\mathbf{x}}_N)} \right] \cong \mathbb{E}_{p_{\mathbf{y}|\mathbf{x}}(\cdot|\mathbf{x})} \left[-\frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}}_N)^T (NJ_y(\mathbf{x}))(\mathbf{x} - \hat{\mathbf{x}}_N) \right] = -\frac{k}{2}. \quad (41)$$

Similarly,

$$p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) \cong \frac{(N)^{k/2}}{(2\pi)^{k/2}} (\det(J_y(\hat{\mathbf{x}}_N)))^{1/2} e^{-\frac{1}{2}(\mathbf{x} - \hat{\mathbf{x}}_N)^T (NJ_y(\hat{\mathbf{x}}_N))(\mathbf{x} - \hat{\mathbf{x}}_N)} \quad (42)$$

is a multivariate Gaussian distribution. The second term in the right-hand side of (13) therefore becomes

$$\mathbb{E}_{p_{\mathbf{x}}(\cdot)} \left[\log \frac{p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\hat{\mathbf{x}}_N)}{p_{\mathbf{y}}(\mathbf{y})} \right] \cong \frac{k}{2} \log \frac{N}{2\pi} + \log \frac{\sqrt{\det(J_y(\mathbf{x}))}}{p_{\mathbf{x}}(\mathbf{x})}. \quad (43)$$

Combining the two terms and maximizing, we arrive at the model capacity:

$$C_N = \max_{p_x} I(\mathbf{x}; \mathbf{y}) \cong \frac{k}{2} \log \frac{N}{2\pi e}, \quad (44)$$

and the Jeffreys' prior:

$$p_{\mathbf{x}}^\infty(x) = \frac{\sqrt{\det(J_y(\mathbf{x}))}}{\int_{\mathbf{x}} \sqrt{\det(J_y(\mathbf{x}))} d\mathbf{x}}. \quad (45)$$

25.2 Asymptotics of Inference: Universal Inference

The concept of *universal inference* reflects a notion that after we observe a sufficient amount of data, we could perform inference and make predictions about future observations as well as if we knew the true probability distribution that generates the data.

A formal definition is as follows.

Definition 1 (Universal Predictor). *Let $p_{\mathbf{y}}^N(\cdot; x)$ be an observation model, with $x \in \mathcal{X}$ denoting an unknown parameter. Let q_1, q_2, \dots, q_N be a sequence of predictors, where $q_n(\cdot; y^{n-1})$ is an estimate of the distribution for y_n based on $y^{n-1} = y^{n-1} = (y_1, \dots, y_{n-1})$. Then this sequence of predictors is universal if*

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{p_{\mathbf{y}}^N(\cdot; x)} \left[\log \frac{p_{y_n|y^{n-1}}(y_n|y^{n-1}; x)}{q_n(y_n; y^{n-1})} \right] = 0, \quad \text{for all } x \in \mathcal{X}, \quad (46)$$

where

$$p_{y_n|y^{n-1}}(y_n|y^{n-1}; x) = \frac{p_{y^n}(y^n; x)}{p_{y^{n-1}}(y^{n-1}; x)}. \quad (47)$$

Notice that the expectation in (46) corresponds to the approximation loss for the n th predictor, and that this loss is averaged over the N predictors, so any initial transient behavior has negligible effect with such a measure. Note too that when the model is i.i.d., i.e.,

$$p_{\mathbf{y}}^N(\mathbf{y}; x) = p_{y^N}(y^N; x) = \prod_{n=1}^N p_y(y_n; x), \quad (48)$$

then the numerator in the logarithm in (46) is simply $p_y(y_n; x)$.

Given this definition of universality, the following theorem characterizes when universal inference is possible. We state the theorem for the case of i.i.d. models, but it can be shown to hold more generally.

Theorem 1 (Universal Prediction Theorem). *For an i.i.d. source $\mathbf{y} = y^N = (y_1, \dots, y_N)$ whose distribution depends on an unknown parameter x , universal inference is possible if and only if*

$$\lim_{N \rightarrow \infty} \frac{C_N}{N} = 0, \quad (49)$$

where C_N is the capacity of the parameterized model.

Proof. We begin by showing the “if” part of the theorem. In particular, we show by construction that when $C_N/N \rightarrow 0$ for all $x \in \mathcal{X}$, there exists a universal predictor.

First, let $p_{y^n}(\cdot)$ denote the mixture distribution for the (partial) observations y^n that is induced by the least informative prior $p_*(\cdot)$ over \mathcal{X} , i.e.,

$$p_{y^n}(y^n) = \sum_{x \in \mathcal{X}} p_*(x) \prod_{i=1}^n p_y(y_i; x), \quad (50)$$

and, in turn, let

$$\delta(x) = \frac{1}{N} D(p_{\mathbf{y}}^N(\cdot; x) \| p_{\mathbf{y}}^N(\cdot)), \quad (51)$$

where $p_{\mathbf{y}}^N = p_{y^N}$.

Now

$$\frac{C_N}{N} = \frac{1}{N} \min_{q^N(\mathbf{y})} \max_{x \in \mathcal{X}} D(p_{\mathbf{y}}^N(\cdot; x) \| q^N(\cdot)) \quad (52)$$

$$= \frac{1}{N} \max_{x \in \mathcal{X}} D(p_{\mathbf{y}}^N(\cdot; x) \| p_{\mathbf{y}}^N(\cdot)) \quad (53)$$

$$= \frac{1}{N} D(p_{\mathbf{y}}^N(\cdot; x) \| p_{\mathbf{y}}^N(\cdot)) \quad (54)$$

$$= \max_{x \in \mathcal{X}} \delta(x), \quad (55)$$

where to obtain (52) we have used the definition of C_N as the minimax redundancy, where to obtain (53) we have used the redundancy-capacity theorem to allow interchanging the order of minimization and maximization, and the fact that the minimization is achieved by the mixture model derived from the least informative prior, and where in (55) we have used (51).

Next, note that the mixture-predictor

$$q_n(y_n; y^{n-1}) = p_{y_n|y^{n-1}}(y_n|y^{n-1}) = \frac{p_{y^n}(y^n)}{p_{y^{n-1}}(y^{n-1})}, \quad (56)$$

derived from (50) has average prediction loss given by

$$\delta_-(x) = \frac{1}{N} \mathbb{E}_{p_{\mathbf{y}}^N(\cdot; x)} \left[\log \frac{p_{\mathbf{y}}^N(\mathbf{y}; x)}{p_{\mathbf{y}}^N(\mathbf{y})} \right] \quad (57)$$

$$= \frac{1}{N} D(p_{\mathbf{y}}^N(\cdot; x) \| p_{\mathbf{y}}^N(\cdot)). \quad (58)$$

Thus, since for all $x \in \mathcal{X}$,

$$\frac{C_N}{N} = \max_{x \in \mathcal{X}} \delta(x) \geq \delta_-(x),$$

it follows that

$$\lim_{N \rightarrow \infty} \frac{C_N}{N} = 0 \quad \Rightarrow \quad \lim_{N \rightarrow \infty} \delta_-(x) = 0 \text{ for all } x \in \mathcal{X},$$

and thus the mixture-predictor derived from the least-informative prior is universal.

We now show the “only if” part of the theorem. In particular, we show that if a universal predictor exists, then $C_N/N \rightarrow 0$ as $N \rightarrow \infty$.

First, note that

$$\frac{C_N}{N} = \frac{1}{N} \min_{q^N} \max_{x \in \mathcal{X}} D(p_{\mathbf{y}}^N(\cdot; x) \| q^N(\cdot)) \quad (59)$$

$$\leq \max_{x \in \mathcal{X}} \frac{1}{N} D(p_{\mathbf{y}}^N(\cdot; x) \| q^N(\cdot)) \quad (60)$$

$$= \frac{1}{N} D(p_{\mathbf{y}}^N(\cdot; x_*) \| q^N(\cdot)) \quad (61)$$

where (60) holds for any particular choice of q^N , and where in (61) x_* is the maximizing x in (60) for the particular choice of q^N .

Now let $q_n(y_n; y^{n-1})$ for $n = 1, \dots, N$ denote our universal predictor, and note that its average loss $\delta_+(x)$ satisfies

$$\delta_+(x) \triangleq \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{p_{\mathbf{y}}^N(\cdot; x)} \left[\log \frac{p_{\mathbf{y}}(y_n; x)}{q_n(y_n; y^{n-1})} \right] \quad (62)$$

$$= \frac{1}{N} \sum_{n=1}^N D(p_{\mathbf{y}}^N(\cdot; x) \| q_{\mathbf{y}}^N(\mathbf{y})), \quad (63)$$

where in (63) we have

$$q_{\mathbf{y}}^N(\mathbf{y}) \triangleq \prod_{n=1}^N q_n(y_n; y^{n-1}). \quad (64)$$

Note that if each of the predictors are valid distributions (nonnegative and integrate to unity), (64) is also a valid distribution (nonnegative and integrates to unity), which can be verified by successive marginalization.

Finally, using (61) with (63) we see that

$$\max_{x \in \mathcal{X}} \delta_+(x) = \delta_+(x_*) \geq \frac{C_N}{N},$$

and thus

$$\lim_{N \rightarrow \infty} \delta_+(x) = 0 \text{ for all } x \in \mathcal{X} \quad \Rightarrow \quad \lim_{N \rightarrow \infty} \delta_+(x_*) = 0 \quad \Rightarrow \quad \lim_{N \rightarrow \infty} \frac{C_N}{N} = 0. \quad (65)$$

□

Some final comments. From our analysis of the model capacity, the condition of the theorem implies a condition of the distances from the likelihood models in the class and the marginal (mixture) probability distribution. We associate universal inference with those distances growing slower than the number of observations.

In both cases we analyzed in the previous section, the condition is satisfied, therefore we expect universal inference in the limit. But what if the data is generated by a distribution $q_{\mathbf{y}}(\cdot)$ that is not in the likelihood family? In this case, we can show that $D(q_{\mathbf{y}}(\cdot) \| p_{\mathbf{y}}(\cdot))$ grows linearly with N , implying that universal inference cannot be achieved. This again highlights the importance of choosing the model class that is large enough to include the true distribution.

25.3 Further Reading

Clarke & Barron. *Information-Theoretic Asymptotics of Bayes Methods*. In IEEE Transactions on Information Theory. This paper reviews and proves the relevant results on asymptotic behavior of model capacity.

Merhav & Feder. Universal Prediction. In IEEE Transactions on Information Theory. This paper discusses universal prediction.