

Problem Set 7

Issued: Tuesday, April 7, 2015

Due: Tuesday, April 14, 2015

Problem 7.1 (practice)

Consider a Gaussian random variable $x \sim \mathcal{N}(0, \sigma_x^2)$, where σ_x^2 is fixed.

- (a) Calculate the differential entropy $h(x)$.

Now consider a model $p_{y|x}(y | x) \sim \mathcal{N}(x, \sigma_{y|x}^2)$, where $\sigma_{y|x}^2$ is fixed.

- (b) Calculate the mutual information between x and y .

Problem 7.2

Suppose we have observations of the form

$$y \triangleq \sqrt{\rho}x + z,$$

where x and z are independent Gaussian random variables each with zero-mean and unit-variance, and where $\rho > 0$ is some constant.

- (a) Determine the Bayes least-squares estimate $\hat{x}_{\text{BLS}}(y)$ of x based on the observed value y , and compute the corresponding mean-square estimation error $\lambda_{\text{BLS}}(\rho)$, expressing your answer as a function of ρ .

Hint: You may find it convenient to recall that the sum of two independent Gaussian random variables is also a Gaussian random variable.

- (b) Show that

$$h(y|x) = h(y - ax|x),$$

for any constant a .

- (c) Show that

$$I(x; y) = h(y) - h(z),$$

and then use this result to express

$$I(\rho) \triangleq I(x; y).$$

as a function of ρ .

(d) Establish that

$$\frac{d}{d\rho} I(\rho) = \frac{1}{2} \lambda_{\text{BLS}}(\rho).$$

Problem 7.3

Consider N flips of a biased coin: i.e., let \mathbf{y} be a vector of N independent identically distributed Bernoulli(x) observations, where $x \in [0, 1]$. That is, $y_i = 1$ with probability x , and $y_i = 0$ otherwise. Consider a symmetric Beta(θ, θ) prior on x :

$$p_x(x) = c(\theta) x^{\theta-1} (1-x)^{\theta-1}, \quad 0 \leq x \leq 1,$$

where $c(\theta) = (\Gamma(2\theta))/(\Gamma(\theta))^2$ and Γ is the Gamma function. Note that the uniform prior corresponds to $\theta = 1$.

A few facts about Gamma function and Beta function are summarized as follows:

- $\Gamma(z) \triangleq \int_0^\infty x^{z-1} e^{-x} dx$
- $\Gamma(z) = (z-1)\Gamma(z-1)$
- $\text{Beta}(x, y) \triangleq \int_0^1 t^{x-1} (1-t)^{y-1} dt = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$

- (a) For $N = 5$, plot, as a function of θ over the range $0 < \theta \leq 1$, the normalized mutual information $I(\mathbf{y}; x)/N$ obtained with the Beta prior. Plot additional curves for $N = 10, 15$, and 20 .

Hint: In MATLAB, the `quad()` function can perform the necessary integration

- (b) For the Beta prior, determine the form of the associated predictor (mixture) distribution $q(y_N | y_{N-1}, \dots, y_1)$ as a function of the Beta prior parameter θ . How does varying θ affect the predictor for small N ? For large N ?

Now consider the uniform prior on x .

- (c) Fix N and ϵ . Let q^* be the uniform mixing distribution (in Bayesian terminology, the marginal probability of \mathbf{y} under the uniform prior). Show that

$$D(p_{\mathbf{y}|x}(\mathbf{y} | x) \| q^*) \leq \log(N+1)$$

for all x . Then determine for what fraction ρ of values of $x \in [0, 1]$ we have

$$D(p_{\mathbf{y}|x}(\mathbf{y} | x) \| q^*) \geq (1 - \epsilon) \log(N+1).$$

Plot the associated ρ versus ϵ curve for $N = 5$.

To determine the fraction ρ , you may want to just look at a finite set of values of x (say M uniformly spaced samples between 0 and 1, where M is reasonably large) and evaluate $D(\cdot\|\cdot)$ for each of those, keeping track of the fraction that are above the threshold. This gives an approximation to what we want.

- (d) Now vary N , and examine the progression of ρ versus ϵ curves as N gets larger. What appears to happen as $N \rightarrow \infty$? (You can confirm this analytically by finding the limits as $N \rightarrow \infty$.)

Here's a handy approximation for the entropy of a binomial random variable n with parameters N and p , valid for large N :

$$H(n) \approx \frac{1}{2} (1 + \log(2\pi Np(1-p))) .$$

Problem 7.4

- (a) Determine a conjugate prior family for exponential models of the form

$$p_{y|x}(y|x) = xe^{-xy}.$$

where $y \geq 0$.

- (b) Determine a conjugate prior family for Poisson models of the form

$$p_{y|x}(y|x) = \frac{x^y e^{-x}}{y!} \quad y = 0, 1, 2, \dots$$

Problem 7.5

Let $\mathbf{z} = [z_1, \dots, z_N]^T$ be a vector of N i.i.d. variables distributed according to

$$p_{\mathbf{z}}(z; x) = \begin{cases} x/3, & z = 0, \\ 2x/3, & z = 1, \\ 1 - x, & z = 2, \end{cases}$$

where $x \in [0, 1]$ is an unknown parameter. We do not observe \mathbf{z} directly, but instead get access to the corresponding sequence of measurements $\mathbf{y} = [y_1, \dots, y_N]^T$ such that

$$y_n = g(z_n) = \begin{cases} 0, & z_n = 0, 2, \\ 1, & z_n = 1, \end{cases}$$

for $n = 1, \dots, N$.

In a particular run of the experiment, we observe that $1/4$ of all elements in the observed sequence \mathbf{y} are ones, and the rest are zeros (i.e., $\sum_{n=1}^N y_n = N/4$).

- (a) Find the ML estimate \hat{x}_{ML} of the parameter x for that particular run.

Let $\mathcal{P}_z \triangleq \{p_z(\cdot; x) : x \in [0, 1]\}$ be the set of distributions $p_z(\cdot; x)$ parameterized by x as defined above, and

$$\hat{\mathcal{P}}^z(\mathbf{y}) \triangleq \left\{ \hat{p}_z : \sum_{z: g(z)=y} \hat{p}_y(z) = \hat{p}_y(y; \mathbf{y}) \quad \text{for all } y \in \mathcal{Y} \right\}$$

be the set of distributions \hat{p}_z that are consistent with the empirical distribution $\hat{p}_y(\cdot; \mathbf{y})$ for the observed sequence \mathbf{y} .

- (b) Draw \mathcal{P}_z and $\hat{\mathcal{P}}^z(\mathbf{y})$ on the probability simplex. Mark clearly the points of intersection of the sets with each other and with the simplex boundaries. You can assume that N is large enough so that the set of empirical distributions is dense.

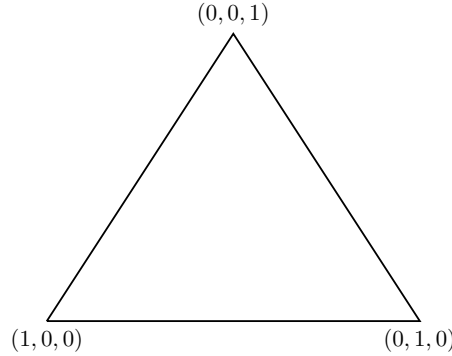


Figure 1: Probability simplex for \mathbf{z} , where a point (p_0, p_1, p_2) corresponds to a distribution $(p_z(0), p_z(1), p_z(2))$.

- (c) In this part, we estimate the parameter x via alternating projections. Recall that in iteration i of the algorithm ($i = 1, 2, \dots$), the E-step constructs a distribution

$$\hat{p}_z^{(i)} = \arg \min_{\hat{p}_z \in \hat{\mathcal{P}}^z(\mathbf{y})} D(\hat{p}_z \| p_z(\cdot; x^{(i-1)})) .$$

The M-step of the algorithm then determines the next estimate of the parameter

$$x^{(i)} = \arg \min_x D(\hat{p}_z^{(i)} \| p_z(\cdot; x)) ,$$

which in turn defines a distribution $p_z^{(i)} = p_z(\cdot; x^{(i)}) \in \mathcal{P}_z$.

Suppose we initialize the algorithm with $x^{(0)} = 3/4$.

- (i) Determine $\hat{p}_z^{(1)}$.
- (ii) Determine $p_z^{(1)}$.
- (iii) Will this process converge? If yes, determine the limit distribution $p_z^{(\infty)}$.
- (d) Now suppose $x^{(0)} = 0$. Will the algorithm converge? If yes, determine the limit distribution $p_z^{(\infty)}$.

Problem 7.6

Let us develop and analyze the Arimoto-Blahut algorithm for computing model capacity. To start, let \mathcal{X} and \mathcal{Y} be arbitrary finite alphabets. Let $\{q_{y|x}(\cdot|x), x \in \mathcal{X}\}$ denote the class of models whose capacity we wish to compute. Recall that the model capacity is given by

$$C \triangleq \max_{p_x \in \mathcal{P}^{\mathcal{X}}} I(x; y) \quad (1)$$

where

$$I(x; y) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p_x(x) q_{y|x}(y|x) \log \frac{q_{y|x}(y|x)}{q_y(y)}, \quad \text{with} \quad q_y(y) = \sum_{x' \in \mathcal{X}} p_x(x') q_{y|x'}(y|x'),$$

and with $\mathcal{P}^{\mathcal{X}}$ denoting the set of all distributions on \mathcal{X} .

Next, let

$$\varphi(p_x, p_{x|y}) \triangleq \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p_x(x) q_{y|x}(y|x) \log \frac{p_{x|y}(x|y)}{p_x(x)}$$

for any $p_x \in \mathcal{P}^{\mathcal{X}}$ and $p_{x|y} \in \mathcal{P}^{\mathcal{X}|\mathcal{Y}}$, where $\mathcal{P}^{\mathcal{X}|\mathcal{Y}}$ is the set of all conditional distributions for x given y .

- (a) Show that for any $p_x \in \mathcal{P}^{\mathcal{X}}$,

$$I(x; y) = \max_{p_{x|y} \in \mathcal{P}^{\mathcal{X}|\mathcal{Y}}} \varphi(p_x, p_{x|y}) \quad (2)$$

and that the associated maximizing distribution is

$$p_{x|y}(x|y) = \frac{p_x(x) q_{y|x}(y|x)}{q_y(y)}.$$

Using (2) in (1), note that model capacity can be equivalently expressed as

$$C = \max_{p_x \in \mathcal{P}^{\mathcal{X}}} \max_{p_{x|y} \in \mathcal{P}^{\mathcal{X}|\mathcal{Y}}} \varphi(p_x, p_{x|y}). \quad (3)$$

The Arimoto-Blahut algorithm evaluates (3) via an alternating maximization procedure. Specifically, starting from some $p_x^{(0)}$, we compute, for $k = 1, 2, \dots$,

$$p_{x|y}^{(k)} = \arg \max_{p_{x|y} \in \mathcal{P}^{\mathcal{X}|\mathcal{Y}}} \varphi(p_x^{(k-1)}, p_{x|y}), \quad (4a)$$

$$p_x^{(k)} = \arg \max_{p_x \in \mathcal{P}^{\mathcal{X}}} \varphi(p_x, p_{x|y}^{(k)}). \quad (4b)$$

Hence, the estimate of capacity after the k th iteration is

$$C_k = \max_{p_x \in \mathcal{P}^{\mathcal{X}}} \varphi(p_x, p_{x|y}^{(k)}) = \varphi(p_x^{(k)}, p_{x|y}^{(k)}).$$

(b) Determine $\alpha_{k-1}(x)$ such that (4) corresponds to the following update step:

$$p_x^{(k)}(x) = \frac{p_x^{(k-1)}(x) e^{\alpha_{k-1}(x)}}{\sum_{x' \in \mathcal{X}} p_x^{(k-1)}(x') e^{\alpha_{k-1}(x')}}.$$

Hint: It may be convenient to optimize (4a) and (4b) separately.

In the remainder of the problem, we interpret the Arimoto-Blahut algorithm as an alternating divergence minimization procedure. To this end, we define the sets

$$\begin{aligned} \mathcal{P} &= \{p_x q_{y|x} : p_x \in \mathcal{P}^{\mathcal{X}}\} \\ \mathcal{Q} &= \{p_{x|y} q_{y|x} : p_{x|y} \in \mathcal{P}^{\mathcal{X}|\mathcal{Y}}\} \end{aligned}$$

Note that \mathcal{Q} is a set of measures (nonnegative functions), but not distributions, i.e., $p_{x|y}(x|y)q_{y|x}(y|x)$ is nonnegative but does not sum to one. We can still use divergence to compare measures, and the associated information geometry holds. In particular, the divergence of measure $\nu(x, y)$ from measure $\mu(x, y)$ is

$$D(\mu \parallel \nu) \triangleq \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} \mu(x, y) \log \frac{\mu(x, y)}{\nu(x, y)}.$$

(c) Show that the Arimoto-Blahut algorithm is an alternating divergence minimization procedure. Specifically, show that

(i) Eq. (4a) is a reverse I-projection from $p_x^{(k-1)} q_{y|x}$ into \mathcal{Q} .

(ii) Eq. (4b) is an I-projection from $p_{x|y}^{(k)} q_{y|x}$ into \mathcal{P} .

(d) Determine which, if either, of the sets \mathcal{P} and \mathcal{Q} are convex.

Problem 7.7 (practice)

Consider N light bulbs whose lifetime is uniformly distributed between 0 and x , where x is an unknown scalar parameter. All N light bulbs are installed in one room and turned on at the same time. You come into the room after time a , (a is a known constant), observe the state of each bulb and then immediately leave the room. The goal is to estimate the lifetime parameter x .

Let z_n be the lifetime of bulb n and y_n be the binary random variable that is equal to 1 if bulb n was still on when you entered the room ($n \in \{1, \dots, N\}$). We assume that all bulbs are independent of each other.

Suppose when you enter the room, the first $K > 0$ bulbs are on, and the rest are burnt out. In other words, you observe $y_n = 1$ for $n \in \{1, \dots, K\}$, and $y_n = 0$ for $n \in \{K + 1, \dots, N\}$.

- (a) Determine the probability of observing $\mathbf{y} = [y_1, \dots, y_N]$ defined above.
- (b) Determine the maximum likelihood estimate of the parameter x given the observed data \mathbf{y} .

In the remainder of this problem, we investigate the performance of the alternating projections algorithm that uses $\mathbf{z} = [z_1, \dots, z_N]$ as full data.

- (c) Determine function $f(\cdot)$ such that any empirical probability distribution $\hat{p}_{\mathbf{z}}(\cdot)$ consistent with the observed data must satisfy

$$\mathbb{E}_{\hat{p}_{\mathbf{z}}(\cdot)} [f(\mathbf{z})] = \frac{K}{N}.$$

Hint: It might be useful to think of the expectation as an integral, i.e., for any distribution $q(\cdot)$,

$$\mathbb{E}_q [f(\mathbf{z})] = \int_{\mathbf{z}} f(\mathbf{z}) q(\mathbf{z}) d\mathbf{z}.$$

- (d) Suppose we initialize the algorithm with $x^{(0)} > a$. Determine the empirical probability distribution $\hat{p}_{\mathbf{z}}^{(0)}(\cdot)$ that we obtain after performing the E-step of the algorithm, i.e.,

$$\hat{p}_{\mathbf{z}}^{(0)}(\cdot) = \arg \min_{\hat{p}_{\mathbf{z}}(\cdot) \in \hat{\mathcal{P}}^{\mathbf{z}}} D(\hat{p}_{\mathbf{z}}(\cdot) \| p_{\mathbf{z}}(\cdot; x^{(0)})),$$

where $\hat{\mathcal{P}}^{\mathbf{z}}$ is the set of all empirical probability distributions $\hat{p}_{\mathbf{z}}(\cdot)$ consistent with the observed data.

- (e) Determine the next value of the parameter $x^{(1)}$ that results from performing the M-step of the algorithm, i.e.,

$$x^{(1)} = \arg \min_x D(\hat{p}_z^{(0)}(\cdot) \| p_z(\cdot; x)).$$

- (f) Will the algorithm converge? If yes, will it converge to the correct estimate?