

## Lecture 10: State Estimation of Hidden Markov Process

*Lecturer: Jiantao Jiao, Tsachy Weissman Scribe: Alon Devorah, David Hallac, Kevin Shutzberg*

In this lecture, we complete and expand the recursive algorithm for state estimation of a hidden markov process through a memoryless noisy channel.

## 1 Recap

Let there be a Markov Process  $\{X_n\}_{n \geq 1}$  sent through a memoryless channel, and measurements  $\{Y_n\}_{n \geq 1}$ . Then the stochastic system is characterized by:

$$P_x, \{P_{X_t|X_{t-1}}\}_{t \geq 2}, \{P_{Y_t|X_t}\}_{t \geq 1}$$

Let

$$\begin{aligned}\alpha_t(X_t) &= p(x_t|y^t) \\ \beta_t(X_t) &= p(x_t|y^{t-1}) \\ \gamma_t(X_t) &= p(x_t|y^n)\end{aligned}$$

In lecture 9, we found that:

$$\alpha_t(X_t) = \frac{\beta_t(X_t)p(y_t|x_t)}{\sum_{\tilde{x}_t} \beta_t(X_t)p(y_t|\tilde{x}_t)} \quad (1)$$

$$\beta_{t+1}(X_t) = \sum_{x_t} \alpha_t(x_t)p(x_{t+1}|x_t) \quad (2)$$

$$\gamma_t(X_t) = \sum_{X_{t+1}} \frac{\gamma_{t+1}\alpha_{t+1}p(x_{t+1}|x_t)}{\beta_{t+1}(X_{t+1})} \quad (3)$$

### 1.1 Finding $P(x^n|y^n)$

**Claim:**

$$X^{t-1} - (X_t, Y^n) - X_{t+1}^n$$

ie given all observations  $Y^n$ ,  $X^n$  is still a Markov Process.

**Proof:** We prove using the "Markov Lemma", described in Lecture 9.

$$\begin{aligned}P(x^n, y^n) &= \prod_{i=1}^n p(x_i|x_{i-1})p(y_i|x_i) \\ &= \prod_{i=1}^t p(x_i|x_{i-1})p(y_i|x_i) \prod_{i=t+1}^n p(x_i|x_{i-1})p(y_i|x_i)\end{aligned}$$

These two terms can be seen as two general funtions:

$$\begin{aligned} & \phi_1(x^t, y^t) \phi_2(x_t^n, y_{t+1}^n) \\ & \phi_1(x^{t+1}, x_t, y^n) \phi_2(x_t, x_{t+1}^n, y^n) \end{aligned}$$

.. Where we've trivially expanded the vectors  $Y$  to go to  $n$  as this includes any subset of  $Y^n$ .  
Looking at the posterior:

$$\begin{aligned} p(x^n|y^n) &= p(x_n|y^n)p(x_{n-1}|x_n, y^n)p(x_{n-2}|x_{n-1}, y^n), \dots, p(x_1|x_2, y^n) \\ &= p(x_n|y^n) \prod_{t=1}^{n-1} p(x_t|x_{t+1}, y^t) \end{aligned}$$

We let  $y^n \rightarrow y^t$  because we already have clean state information  $x_{t+1}$

$$\begin{aligned} p(x_t|x_{t+1}, y^t) &= \frac{p(x_t, x_{t+1}|y^t)}{p(x_{t+1}|y^t)} = \frac{p(x_t|y^t)p(x_{t+1}|x_t, y^t)}{p(x_{t+1}|y^t)} \\ &= \frac{\alpha_t(x_t)p(x_{t+1}|x_t)}{\beta_{t+1}(x_{t+1})} \\ \Rightarrow p(x^n|y^n) &= \gamma_n(x_n) \prod_{t=1}^{n-1} \frac{\alpha_t(x_t)p(x_{t+1}|x_t)}{\beta_{t+1}(x_{t+1})} \\ \log p(x^n|y^n) &= \sum_{t=1}^n g_t(x_t, x_{t+1}) \end{aligned}$$

Where we define  $g$  as:

$$\begin{aligned} g_t(x_t, x_{t+1}) &:= \log \frac{\alpha_t(x_t)p(x_{t+1}|x_t)}{\beta_{t+1}(x_{t+1})}, \quad t \in \{1, \dots, n-1\} \\ g_t(x_t, x_{t+1}) &:= \log \gamma_n(x_n), \quad t = n \end{aligned}$$

- if  $g$  was a funtcion of  $x_t$  only, we could simply use a greedy algorithm and maximize each term.
- since  $g_t$  is a function of both  $x_t, x_{t+1}$

**Definition 1. Entropy:** Let  $U$  a discrete R.V. taking values in  $\mathcal{U}$ . The **entropy** of  $U$  is defined by:

(4)

**Note:** The entropy  $H(U)$  is not a random variable. In fact it is not a function of the object  $U$ , but rather a functional (or property) of the underlying distribution  $P_U^{(u)}, u \in \mathcal{U}$ . An analogy is  $E[U]$ , which is also a number (the mean) corresponding to the distribution.

**Jensen's Inequality:** Let  $Q$  denote a *convex* function, and  $X$  denote any random variable. Jensen's inequality states that

$$E[Q(X)] \geq Q(E[X]). \quad (5)$$

Further, if  $Q$  is strictly convex, equality holds iff  $X$  is deterministic.

*Example:*  $Q(x) = e^x$  is a convex function. Therefore, for a random variable  $X$ , we have by Jensen's inequality:

$$E[e^X] \geq e^{E[X]}$$

Conversely, if  $Q$  is a *concave* function, then

$$E[Q(X)] \leq Q(E[X]). \quad (6)$$

*Example:*  $Q(x) = \log x$  is a concave function. Therefore, for a random variable  $X \geq 0$ ,

$$E[\log X] \leq \log E[X] \quad (7)$$

## 1.2 Properties of Entropy

W.L.O.G suppose  $\mathcal{U} = \{1, 2, \dots, m\}$

1.  $H(U) \leq \log m$ , with equality iff  $P(u) = \frac{1}{m} \forall u$  (i.e. uniform).

**Proof:**

$$H(U) = E\left[\log \frac{1}{P(U)}\right] \quad (8)$$

$$\leq \log E\left[\frac{1}{P(U)}\right] \text{ (Jensen's inequality, since log is concave)} \quad (9)$$

$$= \log \sum_u P(U) \cdot \frac{1}{P(U)} \quad (10)$$

$$= \log m. \quad (11)$$

Equality in Jensen, iff  $\frac{1}{P(U)}$  is deterministic, iff  $p(u) = \frac{1}{m}$

2.  $H(U) \geq 0$ , with equality iff  $U$  is deterministic.

**Proof:**

$$H(U) = E\left[\log \frac{1}{P(U)}\right] \geq 0 \text{ since } \log \frac{1}{P(U)} \geq 0 \quad (12)$$

The equality occurs iff  $\log \frac{1}{P(U)} = 0$  with probability 1, iff  $P(U) = 1$  w.p. 1 iff  $U$  is deterministic.

3. For a PMF  $q$ , defined on the same alphabet as  $p$ , define

$$H_q(U) \triangleq \sum_{u \in \mathcal{U}} p(u) \log \frac{1}{q(u)}. \quad (13)$$

Note that this is the expected surprise function, but instead of the surprise associated with  $p$ , it is the surprise associated  $U$ , which is distributed according to PMF  $p$ , but incorrectly assumed to be having the PMF of  $q$ . The following result stipulates, that we will (on average) be more surprised if we had the wrong distribution in mind. This makes intuitive sense! Mathematically,

$$H(U) \leq H_q(U), \quad (14)$$

with equality iff  $q = p$ .

**Proof:**

$$H(U) - H_q(U) = E\left[\log \frac{1}{p(u)}\right] - E\left[\log \frac{1}{q(u)}\right] \quad (15)$$

$$H(U) - H_q(U) = E\left[\log \frac{q(u)}{p(u)}\right] \quad (16)$$

By Jensen's, we know that  $E \left[ \log \frac{q(u)}{p(u)} \right] \leq \log E \left[ \frac{q(u)}{p(u)} \right]$ , so

$$H(U) - H_q(U) \leq \log E \left[ \frac{q(u)}{p(u)} \right] \quad (17)$$

$$= \log \sum_{u \in \mathcal{U}} p(u) \frac{q(u)}{p(u)} \quad (18)$$

$$= \log \sum_{u \in \mathcal{U}} q(u) \quad (19)$$

$$= \log 1 \quad (20)$$

$$= 0 \quad (21)$$

Therefore, we see that

$$H(U) - H_q(U) \leq 0.$$

Equality only holds when Jensen's yields equality. That only happens when  $\frac{q(u)}{p(u)}$  is deterministic, which only occurs when  $q = p$ , i.e. the distributions are identical.

**Definition 2. Relative Entropy.** An important measure of distance between probability measures is relative entropy, or the Kullback-Leibler divergence:

$$D(p||q) \triangleq \sum_{u \in \mathcal{U}} p(u) \log \frac{p(u)}{q(u)} = E \left[ \log \frac{p(u)}{q(u)} \right] \quad (22)$$

Note that property 3 is equivalent to saying that the relative entropy is always greater than or equal to 0, with equality iff  $q = p$  (convince yourself).

4. If  $X_1, X_2, \dots, X_n$  are independent random variables, then

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i) \quad (23)$$

**Proof:**

$$H(X_1, X_2, \dots, X_n) = E \left[ \log \frac{1}{p(x_1, x_2, \dots, x_n)} \right] \quad (24)$$

$$= E [-\log p(x_1, x_2, \dots, x_n)] \quad (25)$$

$$= E [-\log p(x_1)p(x_2) \dots p(x_n)] \quad (26)$$

$$= E \left[ -\sum_{i=1}^n \log p(x_i) \right] \quad (27)$$

$$= \sum_{i=1}^n E [-\log p(x_i)] \quad (28)$$

$$= \sum_{i=1}^n H(X_i). \quad (29)$$

Therefore, the entropy of independent random variables is the sum of the individual entropies. This is also intuitive, since the uncertainty (or surprise) associated with each random variable is independent.

**Definition 3. Conditional Entropy of  $X$  given  $Y$**

$$H(X|Y) \triangleq \mathbb{E} \left[ \log \frac{1}{P(X|Y)} \right] \quad (30)$$

$$= \sum_{x,y} Pr[x,y] \frac{1}{\log P(x|y)} \quad (31)$$

$$= \sum_y P(y) \left[ \sum_x P(x|y) \frac{1}{\log P(x|y)} \right] \quad (32)$$

$$= \sum_y P(y) H(X|y). \quad (33)$$

*Note:* The conditional entropy is a functional of the joint distribution of  $(X, Y)$ . Note that this is also a number, and denotes the “average” surprise in  $X$  when we observe  $Y$ . Here, by definition, we also average over the realizations of  $Y$ . Note that the conditional entropy is NOT a function of the random variable  $Y$ . In this sense, it is very different from a familiar object in probability, the conditional expectation  $E[X|Y]$  which is a random variable (and a function of  $Y$ ).

5.  $H(X|Y) \leq H(X)$ , equal iff  $X \perp Y$

**Proof:**

$$H(X) - H(X|Y) = \mathbb{E} \left[ \log \frac{1}{P(X)} \right] - \mathbb{E} \left[ \log \frac{1}{P(X|Y)} \right] \quad (34)$$

$$= \mathbb{E} \left[ \log \frac{P(X|Y)}{P(X)} \right] = \mathbb{E} \left[ \log \frac{P(X,Y)}{P(X)P(Y)} \right] \quad (35)$$

$$= \sum_{x,y} P(x,y) \log \frac{P(x,y)}{P(x)P(y)} \quad (36)$$

$$= D(P_{x,y} || P_x \times P_y) \quad (37)$$

$$\geq 0 \quad \text{equal iff } X \perp Y. \quad (38)$$

The last step follows from the non-negativity of relative entropy. Equality holds iff  $P_{x,y} \equiv P_x \times P_y$ , i.e.  $X$  and  $Y$  are independent.

**Definition 4. Joint Entropy of  $X$  and  $Y$**

$$H(X, Y) \triangleq \mathbb{E} \left[ \log \frac{1}{P(X, Y)} \right] \quad (39)$$

$$= \mathbb{E} \left[ \log \frac{1}{P(X)P(Y|X)} \right] \quad (40)$$

6. Chain rule for entropy:

$$H(X, Y) = H(X) + H(Y|X) \quad (41)$$

$$= H(Y) + H(X|Y) \quad (42)$$

7. Sub-additivity of entropy

$$H(X, Y) \leq H(X) + H(Y), \quad (43)$$

with equality iff  $X \perp Y$  (follows from the property that conditioning does not increase entropy)

**Definition 5. *Mutual information between  $X$  and  $Y$***

We now define the mutual information between random variables  $X$  and  $Y$  distributed according to the joint PMF  $P(x, y)$ :

$$I(X, Y) \triangleq H(X) + H(Y) - H(X, Y) \quad (44)$$

$$= H(Y) - H(Y|X) \quad (45)$$

$$= H(X) - H(X|Y) \quad (46)$$

$$= D(P_{x,y} || P_x \times P_y) \quad (47)$$

The mutual information is a canonical measure of the information conveyed by one random variable about another. The definition tells us that it is the reduction in average surprise, upon observing a correlated random variable. The mutual information is again a functional of the joint distribution of the pair  $(X, Y)$ . It can also be viewed as the relative entropy between the joint distribution, and the product of the marginals.