

Example of MSE/MAP ESTIMATION

Note Title

1/18/2013

Consider a Bernoulli distribution

$$p(x_i|u) = u^{x_i} (1-u)^{1-x_i} \quad 0 \leq u \leq 1$$

This describes a distribution over binary values 0, 1 that could represent Heads/Tails, Yes/No, or any other two-state event.

We will rewrite this likelihood function in a more general way, to make it easier to generalize later to distributions over arbitrary # of discrete states.

Let $u = u_1$ and $1-u = u_2$, so $u_1 + u_2 = 1$

$$\text{Let } z_{i1} = \begin{cases} 1 & \text{if } x_i = 1 \\ 0 & \text{if } x_i = 0 \end{cases} \quad z_{i2} = \begin{cases} 1 & \text{if } x_i = 0 \\ 0 & \text{if } x_i = 1 \end{cases}$$

These z_{ij} are binary indicator variables used to form a "1 of K" representation. [Chris Bishop's PRML book discusses this representation in more detail, if you have access to the book]

example $x = \begin{matrix} x_1 & x_2 & x_3 & x_4 & x_5 \\ \{H, T, T, H, H\} \end{matrix}$

z_{ik}	$k=1$	$k=2$
$i=1$	1	0
$i=2$	0	1
$i=3$	0	1
$i=4$	1	0
$i=5$	1	0
	<hr/>	<hr/>
	3	2

Note for each row, only one variable is 1 and the other is 0. In a general 1 of K situation, one variable is 1 and the other $K-1$ variables are 0.

Furthermore, note sum down each column is # samples taking value K

$$\sum z_{i1} = N_1 \quad \sum z_{i2} = N_2$$

Using notation from previous page, we can now write our bernoulli distribution in more general form as:

$$P(x_i | u) = u_1^{z_{i1}} u_2^{z_{i2}}$$

Now, given N samples $X = \{x_1, x_2, \dots, x_N\}$, form the joint likelihood function for parameters $u = \{u_1, u_2\}$

$$L(u) = P(X|u) = \prod_{i=1}^N \prod_{k=1}^2 u_k^{z_{ik}} = \prod_{i=1}^N u_1^{z_{i1}} u_2^{z_{i2}}$$

Now solve for MLE estimate

Lagrange multiplier
to enforce constraint
 $u_1 + u_2 = 1$

$$\log L = \sum_{i=1}^N z_{i1} \log u_1 + z_{i2} \log u_2 + \lambda (1 - u_1 - u_2)$$

$$\frac{\partial \log L}{\partial u_1} = \sum_{i=1}^N \frac{z_{i1}}{u_1} - \lambda = 0 \Rightarrow \sum z_{i1} = N_1 = \lambda u_1$$

$$\frac{\partial \log L}{\partial u_2} = \sum_{i=1}^N \frac{z_{i2}}{u_2} - \lambda = 0 \Rightarrow \sum z_{i2} = N_2 = \lambda u_2$$

$$\frac{\partial \log L}{\partial \lambda} = 1 - u_1 - u_2 = 0 \Rightarrow u_1 + u_2 = 1$$

sum both sides $N_1 + N_2 = \lambda \underbrace{(u_1 + u_2)}_1$

$$\begin{aligned} \text{So } u_1 &= N_1 / \lambda = N_1 / (N_1 + N_2) \\ u_2 &= N_2 / \lambda = N_2 / (N_1 + N_2) \end{aligned}$$

NOTE: The MLE estimates are just the relative frequency of counts of values occurring in the sample data.

To compute MAP estimates, recall Bayes rule

$$\overset{\text{posterior}}{P(u|x)} = \frac{\overset{\text{likelihood}}{P(x|u)} \overset{\text{prior}}{P(u)}}{\underset{\text{evidence}}{P(x)}}$$

Since the denominator $P(x)$ does not depend on u , we can ignore it if we only care about $\arg\max P(u|x)$.

$$\text{so } \hat{u}_{\text{MAP}} = \arg\max_u P(u|x) = \arg\max_u P(x|u)P(u)$$

From MLE derivation on previous page, we know that the joint likelihood $P(x|u)$ is

$$P(x|u) = \prod_{i=1}^n u_1^{z_{i1}} u_2^{z_{i2}} = u_1^{\sum z_{i1}} u_2^{\sum z_{i2}} = u_1^{n_1} u_2^{n_2}$$

with $u_1 + u_2 = 1$, note: our previous MLE derivation could have been streamlined by writing the joint likelihood function in this way.

We would like to multiply this by a "conjugate prior" that does not greatly complicate the form of the product. The beta distribution is a conjugate prior for the bernoulli distribution.

$$\text{Beta}(u|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \underset{\text{our } u_1}{u^{a-1}} \underset{\text{our } u_2}{(1-u)^{b-1}}$$

So let our prior be

$$P(u) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} u_1^{a-1} u_2^{b-1}$$

normalizing constant. Does not depend on u parameters, so we will drop it from now on.

-4-

Now form the posterior

$$p(u|x) \propto p(x|u)p(u) = u_1^{N_1} u_2^{N_2} u_1^{a-1} u_2^{b-1} \\ = u_1^{N_1+a-1} u_2^{N_2+b-1}$$

Now want $\operatorname{argmax}_{u_1, u_2} p(u|x)$, subject to $u_1 + u_2 = 1$

$$\log p(u|x) = (N_1 + a - 1) \log u_1 + (N_2 + b - 1) \log u_2 + \lambda(1 - u_1 - u_2)$$

$$\frac{\partial \log p}{\partial u_1} = \frac{(N_1 + a - 1)}{u_1} - \lambda = 0$$

$$\frac{\partial \log p}{\partial u_2} = \frac{N_2 + b - 1}{u_2} - \lambda = 0$$

$$\frac{\partial \log p}{\partial \lambda} = 1 - u_1 - u_2 = 0 \\ \Rightarrow u_1 + u_2 = 1$$

NOTE! in the homework you will fill in the details, for a more general distribution

Solving for u_1 and u_2 in a similar manner to our MLE derivation, we find

$$u_1 = \frac{N_1 + a - 1}{N_1 + N_2 + a + b - 2} \quad u_2 = \frac{N_2 + b - 1}{N_1 + N_2 + a + b - 2}$$

Intuitive interpretation: $a-1$ and $b-1$ are "virtual" sample counts given to us prior to observing the actual data, so we pool together both the real and virtual samples.