

## Exercise: BIC for a 2d discrete distribution

(Source: Jaakkola.)

Let  $x \in \{0, 1\}$  denote the result of a coin toss ( $x = 0$  for tails,  $x = 1$  for heads). The coin is potentially biased, so that heads occurs with probability  $\theta_1$ . Suppose that someone else observes the coin flip and reports to you the outcome,  $y$ . But this person is unreliable and only reports the result correctly with probability  $\theta_2$ ; i.e.,  $p(y|x, \theta_2)$  is given by

	$y = 0$	$y = 1$
$x = 0$	$\theta_2$	$1 - \theta_2$
$x = 1$	$1 - \theta_2$	$\theta_2$

Assume that  $\theta_2$  is independent of  $x$  and  $\theta_1$ .

1. Write down the joint probability distribution  $p(x, y|\theta)$  as a  $2 \times 2$  table, in terms of  $\theta = (\theta_1, \theta_2)$ .
2. Suppose have the following dataset:  $\mathbf{x} = (1, 1, 0, 1, 1, 0, 0)$ ,  $\mathbf{y} = (1, 0, 0, 0, 1, 0, 1)$ . What are the MLEs for  $\theta_1$  and  $\theta_2$ ? Justify your answer. Hint: note that the likelihood function factorizes,

$$p(x, y|\theta) = p(y|x, \theta_2)p(x|\theta_1) \quad (1)$$

What is  $p(\mathcal{D}|\hat{\theta}, M_2)$  where  $M_2$  denotes this 2-parameter model? (You may leave your answer in fractional form if you wish.)

3. Now consider a model with 4 parameters,  $\theta = (\theta_{0,0}, \theta_{0,1}, \theta_{1,0}, \theta_{1,1})$ , representing  $p(x, y|\theta) = \theta_{x,y}$ . (Only 3 of these parameters are free to vary, since they must sum to one.) What is the MLE of  $\theta$ ? What is  $p(\mathcal{D}|\hat{\theta}, M_4)$  where  $M_4$  denotes this 4-parameter model?
4. Suppose we are not sure which model is correct. We compute the leave-one-out cross validated log likelihood of the 2-parameter model and the 4-parameter model as follows:

$$L(m) = \sum_{i=1}^n \log p(x_i, y_i|m, \hat{\theta}(\mathcal{D}_{-i})) \quad (2)$$

and  $\hat{\theta}(\mathcal{D}_{-i})$  denotes the MLE computed on  $\mathcal{D}$  excluding row  $i$ . Which model will CV pick and why? Hint: notice how the table of counts changes when you omit each training case one at a time.

5. Recall that an alternative to CV is to use the BIC score, defined as

$$\text{BIC}(M, \mathcal{D}) \triangleq \log p(\mathcal{D}|\hat{\theta}_{MLE}) - \frac{\text{dof}(M)}{2} \log N \quad (3)$$

where  $\text{dof}(M)$  is the number of free parameters in the model, Compute the BIC scores for both models (use log base  $e$ ). Which model does BIC prefer?