## 1 Definitions Ahoy!

## 1.1 MLE/MAP

1. MLE: Finds the best parameters for a specific dataset,  $\mathcal{D}$ . Specifically, we want to find the parameters  $\hat{\theta}_{MLE}$  that maximize the likelihood for  $\mathcal{D}$ .

$$\hat{\theta}_{MLE} = \operatorname*{argmax}_{\theta} p(\mathcal{D} \mid \theta)$$

2. **MAP:** Finds the best parameters given  $\mathcal{D}$  and a prior belief about the parameters. Specifically, we want to find the parameters  $\hat{\theta}_{MAP}$  that maximize the posterior distribution  $p(\theta \mid \mathcal{D})$  of parameters  $\theta$ .

$$\begin{split} \hat{\theta}_{MAP} &= \operatorname*{argmax}_{\theta} p(\theta \mid \mathcal{D}) \\ &= \operatorname*{argmax}_{\theta} \frac{p(\mathcal{D} \mid \theta) p(\theta)}{\operatorname{Normalizing Constant}} \\ &= \operatorname*{argmax}_{\theta} p(\mathcal{D} \mid \theta) p(\theta) \\ &= \operatorname*{argmin}_{\theta} - \log \left( p(\mathcal{D} \mid \theta) p(\theta) \right) \\ &= \operatorname*{argmin}_{\theta} - \log p(\mathcal{D} \mid \theta) - \log p(\theta) \end{split}$$

3. MLE and MAP for conditional likelihood: When we want to predict the output y given the input x using our supervised dataset, we have to reformulate the MLE and MAP optimizations to use the conditional likelihood (and conditional posterior) instead:

$$\begin{split} \hat{\theta}_{MLE} &= \operatorname*{argmax}_{\theta} p(\mathcal{D} \mid \theta) \\ &= \operatorname*{argmax}_{\theta} \prod_{i=1}^{N} p\left(y^{(i)} \mid x^{(i)}, \theta\right) \\ &= \operatorname*{argmin}_{\theta} - \log \prod_{i=1}^{N} p\left(y^{(i)} \mid x^{(i)}, \theta\right) \\ &= \operatorname*{argmin}_{\theta} - \sum_{i=1}^{N} \log p\left(y^{(i)} \mid x^{(i)}, \theta\right) \end{split}$$

$$\begin{split} \hat{\theta}_{MAP} &= \operatorname*{argmax}_{\theta} p(\theta \mid \mathcal{D}) \\ &= \operatorname*{argmax}_{\theta} \left( \prod_{i=1}^{N} p\left(y^{(i)} \mid x^{(i)}, \theta\right) \right) p(\theta) \\ &= \operatorname*{argmin}_{\theta} - \log \prod_{i=1}^{N} p\left(y^{(i)} \mid x^{(i)}, \theta\right) - \log p(\theta) \\ &= \operatorname*{argmin}_{\theta} - \sum_{i=1}^{N} \log p\left(y^{(i)} \mid x^{(i)}, \theta\right) - \log p(\theta) \end{split}$$

## 2 Anybody have a MAP?

Imagine you are a data scientist working for an advertising company. The advertising company has recently run an ad and they want you to estimate its performance. The ad was shown to N people.  $y^{(i)}=1$  if person i clicked on the ad and 0 otherwise. Thus  $\sum_{i}^{N}y^{(i)}=N_{1}$  people decided to click on the ad. Assume that the probability that the i-th person clicks on the ad is  $\phi$  and the probability that the i-th person does not click on the ad is  $1-\phi$ .

$$p(\mathcal{D} \mid \theta) = p\left(y^{(1)}, y^{(2)}, ..., y^{(N)} \mid \phi\right) = \prod_{i=1}^{N} p\left(y^{(i)} \mid \phi\right) = \phi^{N_1} (1 - \phi)^{N - N_1}$$

1. Calculate  $\hat{\phi}_{MLE}$ .

Note

2. Your coworker tells you that  $\phi \sim \text{Beta}(\alpha, \beta)$ . That is:

$$p(\phi) = \frac{\phi^{\alpha - 1} (1 - \phi)^{\beta - 1}}{B(\alpha, \beta)}$$

Note that  $B(\alpha, \beta)$  is not a function of  $\phi$  and can be treated as a constant. Formulate the optimization of the log posterior,  $\operatorname{argmin}_{\phi} - \log p(\phi \mid \mathcal{D})$ , in terms of  $N, N_1, \phi, \alpha$ , and  $\beta$ .

Now, calcula	te $\phi_{MAP}$ .					
Calculate $\hat{\phi}_{M}$	e time, so you, som $IAP$ .		, decide to set	a - 0 + 1 - 1	, where $\rho = 100^{-4}$	0   1 -
How do $\hat{\phi}_{ML}$	$_{E}$ and $\hat{\phi}_{MAP}$ differ	:? Argue which	n estimate you	think is better	r.	

## 3 Conceptual MLE/MAP Questions

1.	When calculating the MAP estimates, we rely on the Bayes formula and then argue we can ignore $p(\mathcal{D})$ . Why do we usually ignore calculating $p(\mathcal{D})$ ?
2.	As the amount of data increases, how are MLE and MAP affected?
3.	Can MLE and MAP estimates be the same? If so, when?