

Solution to Homework 3

Shoeb Mohammed and Zhuo Chen

February 21, 2016

1 MAP and MLE parameter estimation

$\mathcal{D} = \{x^{(i)} | 1 \leq i \leq m\}$ where $x^{(i)} \sim \text{i.i.d } \text{Ber}(\theta)$

1.1

If m_1 are number of heads, m_0 are number of tails and $m_0 + m_1 = m$ then the likelihood and MLE for θ are

$$p(\mathcal{D}|\theta) = \theta^{m_1}(1 - \theta)^{m_0} \quad (1)$$

$$\begin{aligned} \theta_{MLE} &= \text{argmax}_{\theta} \theta^{m_1}(1 - \theta)^{m_0} \\ &= \text{argmax}_{\theta} m_1 \log \theta + m_0 \log(1 - \theta) \end{aligned} \quad (2)$$

θ_{MLE} satisfies (first derivative of the likelihood equals zero)

$$\frac{m_1}{\theta_{MLE}} - \frac{m_0}{1 - \theta_{MLE}} = 0 \quad (3)$$

Thus,

$$\begin{aligned} \theta_{MLE} &= \frac{m_1}{m_0 + m_1} \\ &= \frac{m_1}{m} \end{aligned} \quad (4)$$

1.2

The prior is

$$\begin{aligned} p(\theta) &= \text{Beta}(\theta|a, b) \\ &\propto \theta^{(a-1)}(1 - \theta)^{(b-1)} \end{aligned} \quad (5)$$

Thus, the posterior is

$$\begin{aligned} p(\theta|\mathcal{D}) &= \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \\ &= \frac{p(\mathcal{D}|\theta)p(\theta)}{\sum_{\theta'} p(\mathcal{D}|\theta')p(\theta')} \\ &\propto \theta^{m_1+a-1} \theta^{m_0+b-1} \end{aligned} \quad (6)$$

Thus,

$$\begin{aligned}\theta_{MAP} &= \operatorname{argmax}_{\theta} \theta^{m_1+a-1} (1-\theta)^{m_0+b-1} \\ &= \operatorname{argmax}_{\theta} (m_1+a-1) \log \theta + (m_0+b-1) \log(1-\theta)\end{aligned}\tag{7}$$

Equation 7 is similar to MLE estimation. Thus,

$$\begin{aligned}\theta_{MAP} &= \frac{m_1+a-1}{m_0+m_1+a+b-2} \\ &= \frac{m_1+a-1}{m+a+b-2}\end{aligned}\tag{8}$$

It is clear from equations 8 and 4 that $\theta_{MAP} = \theta_{MLE}$ when $a = b = 1$.

2 Logistic regression and Gaussian Naive Bayes

2.1

$$\begin{aligned}p(y=1|x) &= g(\theta^T x) \\ p(y=0|x) &= 1 - g(\theta^T x)\end{aligned}\tag{9}$$

2.2

With naives Bayes assumption,

$$\begin{aligned}p(y=1|x) &= \frac{p(x|y=1)p(y=1)}{p(x)} \\ &= \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)} \\ &= \frac{p(x|y=1)\gamma}{p(x|y=1)\gamma + p(x|y=0)(1-\gamma)} \quad \text{since } y \sim \operatorname{Ber}(\gamma) \\ &= \frac{\prod_{j=1}^d \mathcal{N}(\mu_j^1, \sigma_j^2)\gamma}{\prod_{j=1}^d \mathcal{N}(\mu_j^1, \sigma_j^2)\gamma + \prod_{j=1}^d \mathcal{N}(\mu_j^0, \sigma_j^2)(1-\gamma)} \quad \text{since } p(x_j|y=1) \sim \mathcal{N}(\mu_j^1, \sigma_j^2) \text{ and } p(x_j|y=0) \sim \mathcal{N}(\mu_j^0, \sigma_j^2) \\ &= \frac{\mathcal{N}(\mu^1, \Sigma)\gamma}{\mathcal{N}(\mu^1, \Sigma)\gamma + \mathcal{N}(\mu^0, \Sigma)(1-\gamma)} \quad \text{where } \mu_0 = (\mu_1^0 \cdots \mu_d^0)^T, \mu_1 = (\mu_1^1 \cdots \mu_d^1)^T, \Sigma = \operatorname{diag}(\sigma_1^2 \cdots \sigma_d^2)\end{aligned}\tag{10}$$

$$\begin{aligned}p(y=0|x) &= \frac{p(x|y=0)p(y=0)}{p(x)} \\ &= \frac{p(x|y=0)p(y=0)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)} \\ &= \frac{\mathcal{N}(\mu^0, \Sigma)(1-\gamma)}{\mathcal{N}(\mu^1, \Sigma)\gamma + \mathcal{N}(\mu^0, \Sigma)(1-\gamma)} \quad \text{where } \mu_0 = (\mu_1^0 \cdots \mu_d^0)^T, \mu_1 = (\mu_1^1 \cdots \mu_d^1)^T, \Sigma = \operatorname{diag}(\sigma_1^2 \cdots \sigma_d^2)\end{aligned}\tag{11}$$

2.3

Proof. With uniform class priors, equation 10 gives

$$\begin{aligned}
p(y = 1|x) &= \frac{\mathcal{N}(\mu^1, \Sigma)}{\mathcal{N}(\mu^1, \Sigma) + \mathcal{N}(\mu^0, \Sigma)} \\
&= \frac{1}{1 + \frac{\mathcal{N}(\mu^0, \Sigma)}{\mathcal{N}(\mu^1, \Sigma)}} \\
&= \frac{1}{1 + \frac{\exp(\frac{1}{2}(x-\mu^1)^T \Sigma^{-1}(x-\mu^1))}{\exp(\frac{1}{2}(x-\mu^0)^T \Sigma^{-1}(x-\mu^0))}} \\
&= \frac{1}{1 + \frac{\exp((x-\mu^1)^T \Lambda^2 (x-\mu^1))}{\exp((x-\mu^0)^T \Lambda^2 (x-\mu^0))}} \quad \text{where } \Lambda = \text{diag} \left(\frac{1}{\sqrt{2}\sigma_1} \cdots \frac{1}{\sqrt{2}\sigma_d} \right) \\
&= \frac{1}{1 + \frac{\exp((\Lambda(x-\mu^1))^T (\Lambda(x-\mu^1)))}{\exp((\Lambda(x-\mu^0))^T (\Lambda(x-\mu^0)))}} \\
&= \frac{1}{1 + \frac{\exp((\Lambda(z+a))^T (\Lambda(z+a)))}{\exp((\Lambda(z-a))^T (\Lambda(z-a)))}} \quad \text{where } a = \frac{\mu^0 - \mu^1}{2} \text{ and } z = x - \frac{\mu^0 + \mu^1}{2} \\
&= \frac{1}{1 + \exp((\Lambda(z+a))^T (\Lambda(z+a)) - (\Lambda(z-a))^T (\Lambda(z-a)))} \\
&= \frac{1}{1 + \exp(4(\Lambda a)^T (\Lambda z))} \\
&= \frac{1}{1 + \exp\left(4a^T \Sigma^{-1} \left(x - \frac{\mu^0 + \mu^1}{2}\right)\right)} \\
&= g(\theta^T x') \quad \text{where } \theta^T = [(\mu^0 - \mu^1)^T \Sigma^{-1}(\mu^0 + \mu^1), 2(\mu^1 - \mu^0)^T \Sigma^{-1}] \text{ and } x' = [1, x]
\end{aligned} \tag{12}$$

□

3 Softmax regression and OVA logistic regression

3.1 3.1

Implementing the loss function for softmax regression (naive version) —