

## Practice Set 2, Problem 1

$$y = \{y_1, y_2, \dots, y_n\}$$

$$P(y|\lambda) = \frac{\lambda^{y_n} e^{-\lambda}}{y_n!} \quad (\text{Poisson})$$

$$P(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \quad (\text{Gamma})$$

① MLE and MAP:

$$\hat{\lambda}_{MLE} = \underset{\lambda}{\operatorname{argmax}} \sum_{n=1}^N \log P(y_n|\lambda) \quad (\text{i.i.d. data})$$

$$\sum_{n=1}^N \log P(y_n|\lambda) = \left[ \sum_{n=1}^N y_n \log \lambda \right] - N \cancel{+} \text{Constant}$$

(terms that don't depend on  $\lambda$ )

Taking derivative and setting it to zero

$$\sum y_n \times \frac{1}{\lambda} - N = 0$$

$$\Rightarrow \boxed{\hat{\lambda}_{MLE} = \frac{\sum_{n=1}^N y_n}{N}}$$

MAP estimation will be almost identical with the extra  $\log P(\lambda)$  term.

$$\hat{\lambda}_{MAP} = \underset{\lambda}{\operatorname{argmax}} \left[ \sum_{n=1}^N \log P(y_n|\lambda) + \log P(\lambda) \right]$$



$$\log p(\lambda) = (\alpha - 1) \log \lambda - \beta \lambda + \text{Constant}$$

terms that  
don't depend on  
 $\lambda$

The MAP objective will be

$$\sum_{n=1}^N y_n \log \lambda - N \lambda + (\alpha - 1) \log \lambda - \beta \lambda$$

Maximizing w.r.t.  $\lambda$  will give the MAP solution

$$\hat{\lambda}_{\text{map}} = \frac{\sum_{n=1}^N y_n + \alpha - 1}{N + \beta}$$

(2) Posterior distribution of  $\lambda$

$$p(\lambda|y) = \frac{p(\lambda)p(y|\lambda)}{p(y)} = \frac{p(\lambda) \prod_{n=1}^N p(y_n|\lambda)}{p(y)}$$

Since the prior (gamma) and the likelihood (Poisson) ~~are~~ are conjugate, the posterior must be gamma!

Let's multiply the terms  $p(\lambda)$  and  $p(y_n|\lambda)$  are try to "identify" their gamma distributions parameters

$$p(\lambda|y) \propto p(y) \prod_{n=1}^N p(y_n|\lambda)$$



$$P(\lambda|y) \propto \lambda^{\alpha-1} e^{-\beta\lambda} \times \prod_{n=1}^N \frac{\lambda^{y_n} e^{-\lambda}}{y_n!}$$

$$\propto \lambda^{\sum_{n=1}^N y_n + \alpha - 1} e^{-(\beta + N)\lambda}$$

(ignoring the terms that don't depend on  $\lambda$ )

The above expression is clearly in form of a gamma distribution with

$$\text{Shape} = \sum_{n=1}^N y_n + \alpha$$

$$\text{rate} = \beta + N$$

[No need to worry about constant of proportionality. It must be a gamma]

Thus  $P(y|\lambda) = \text{Gamma}\left(\sum_{n=1}^N y_n + \alpha, \beta + N\right)$

③  $\hat{\lambda}_{\text{map}} = \frac{\sum_{n=1}^N y_n + \alpha - 1}{N + \beta} = \frac{\sum_{n=1}^N y_n}{N + \beta} + \frac{(\alpha - 1)}{(N + \beta)}$

Sum to 1  
(convex comb.)

$$= \left[ \frac{\sum_{n=1}^N y_n}{N} \right] \times \left[ \frac{N}{N + \beta} \right] + \left[ \frac{\alpha - 1}{\beta} \right] \times \left[ \frac{\beta}{N + \beta} \right]$$

MLE                      prior's mode

Posterior's mean:

$$\frac{\sum_{n=1}^N y_n + \alpha}{N + \beta} = \frac{\sum_{n=1}^N y_n}{N} \times \left[ \frac{N}{N + \beta} \right] + \left[ \frac{\alpha}{\beta} \right] \times \frac{\beta}{N + \beta}$$

MLE                      prior's mean



$$\begin{aligned}
 \textcircled{4} \quad P(y_*|y) &= \int P(y_*, \lambda | y) d\lambda \\
 &= \int P(y_* | \lambda) P(\lambda | y) d\lambda \\
 &\approx P(y_* | \hat{\lambda}_{MLE}) \quad (\text{if using MLE})
 \end{aligned}$$

$$\text{or} \quad P(y_* | \hat{\lambda}_{MAP}) \quad (\text{if using MAP})$$

In both these cases,  $P(y_*|y)$  is simply a Poisson with parameters  $\hat{\lambda}_{MLE}$  or  $\hat{\lambda}_{MAP}$ .

If using the full posterior, which is  $P(\lambda|y) = \text{Gamma}(\sum_{n=1}^N y_n + \alpha, \beta + N)$ ,  $P(y_*|y)$  will be

$$\begin{aligned}
 P(y_*|y) &= \int P(y_* | \lambda) P(\lambda | y) d\lambda \\
 &= \int \frac{\lambda^{y_*} e^{-\lambda}}{y_*!} \times \frac{(\beta + N)^{\sum_{n=1}^N y_n + \alpha}}{\int \lambda^{\sum_{n=1}^N y_n + \alpha} e^{-(\beta + N)\lambda} d\lambda} \lambda^{\sum_{n=1}^N y_n + \alpha - 1} e^{-(\beta + N)\lambda} d\lambda
 \end{aligned}$$

The above is actually a mixture of infinite many Poisson distributions. The result is actually not a Poisson but the "Negative Binomial" distribution.

(if you are interested in knowing more about this result, you may refer to the wikipedia article of NB distribution).

(This part was just for your "General Knowledge" :-)



## Practice Set 2 (Problem 2)

$$\nabla_w \text{NLL}(w) = - \left[ \sum_{n=1}^N y_n x_n - \frac{\exp(w^T x_n)}{1 + \exp(w^T x_n)} x_n \right]$$

$$\cancel{P(y_n=1/w, x_n)} = \mu_n$$

$$g = - \left[ \sum_{n=1}^N y_n x_n - \mu_n x_n \right] = - \sum_{n=1}^N (y_n - \mu_n) x_n$$

The expression above can't be written ~~as~~ by separating  $w$  on one side (like we did for linear regression). Therefore we can't find a closed form solution, and need iterative methods (e.g. gradient descent).

Intuitive meaning of the gradient's expression is

$$g = - \sum_{n=1}^N (y_n - \mu_n) x_n$$

↓  
Contribution  
of  $x_n$

$$\begin{array}{l} y_n = 0 \text{ but } \mu_n \rightarrow 1 \\ \text{or} \\ y_n = 1 \text{ but } \mu_n \rightarrow 0 \end{array}$$

if  $y_n - \mu_n$  difference is large, which will happen if there is a LARGE MISPREDICTION, then  $x_n$  will contribute more to the gradient (we actually discussed this also while discussing about gradient descent).