

Assignment-01

CS772

MLE for λ :

given $p(x|\lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$

for $x_i \Rightarrow \log(p(x_i|\lambda)) = x_i \log(\lambda) - \lambda - \log(x_i!)$

for $x = (x_1, \dots, x_N) \Rightarrow \log(p(x|\lambda))$

for MLE, $NLL(\lambda) = \log(\lambda) \sum_{i=1}^N x_i - N\lambda - \sum_{i=1}^N \log(x_i!)$

$\frac{\partial NLL(\lambda)}{\partial \lambda} = 0$
 $\Rightarrow \frac{\sum x_i}{\lambda} = N \Rightarrow \lambda = \frac{\sum_{i=1}^N x_i}{N}$ (MLE estimate)

for MAP estimate,

$L(\lambda) = \text{MLE} + \log(p(\lambda))$
 $= (\sum x_i) \log \lambda - N\lambda - \sum \log(x_i!)$
 $+ (\alpha-1) \log(\lambda) - \beta\lambda + \log\left(\frac{\beta^\alpha}{\Gamma(\alpha)}\right)$

$\frac{\partial L(\lambda)}{\partial \lambda} = 0 \Rightarrow \frac{\sum x_i}{\lambda} - N + \frac{\alpha-1}{\lambda} - \beta = 0$

$\Rightarrow \lambda_{\text{MAP}} = \frac{\sum x_i + (\alpha-1)}{N + \beta}$ (a kind of sum of obs. in prior)

Symbolically # obs. in prior

Since, γ -distribution is conjugate prior to Poisson, the posterior is δ -dist. as shown:

Posterior \propto Prior \cdot Likelihood

$$P(\lambda|x) \propto P(\lambda) P(x|\lambda)$$

$$\frac{\beta^\alpha \lambda^{\alpha-1} \exp(-\beta\lambda)}{T(\alpha)} \quad \frac{\lambda^{\sum x_i} \exp(-N\lambda)}{\prod \pi(x_i)!}$$

from x_i 's being i.i.d

$$\therefore P(\lambda|x) = \frac{\beta^\alpha}{T(\alpha)} \cdot \frac{\lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda)}{\prod \pi(x_i)!}$$

Normalizing constant

$$\text{where } K = \int_0^\infty \frac{\beta^\alpha}{T(\alpha) \prod \pi(x_i)!} \lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda) d\lambda$$

$$= \frac{\beta^\alpha}{T(\alpha) \prod \pi(x_i)!} \int_0^\infty \lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda) d\lambda$$

$$(*) \quad P(\lambda|x) = \frac{\lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda)}{\int_0^\infty \lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda) d\lambda}$$

Observe that $\text{Gamma}(\lambda', \sum x_i + \alpha, N + \beta)$ is the form of numerator i.e.

$$\int_0^\infty \text{Gamma}(\lambda', \sum x_i + \alpha, N + \beta) d\lambda = 1$$

$$\Rightarrow \int_0^\infty \lambda^{\sum x_i + \alpha - 1} \exp(-N\lambda - \beta\lambda) d\lambda = \frac{T(\sum x_i + \alpha)}{(N + \beta)}$$

Putting Back in (*) we see,

$$\text{Full Posterior distribution} \rightarrow P(\lambda|x) \sim \text{Gamma}(\lambda', \sum x_i + \alpha, N + \beta)$$

(C) Mode of posterior = $\frac{\alpha + \sum_{i=1}^N x_i - 1}{N + \beta}$ (from given properties)

Prior mode = $\frac{\alpha - 1}{\beta}$ and $\lambda_{MLE} = \frac{\sum x_i}{N}$

So, clearly $\frac{N}{N + \beta} \lambda_{MLE} = \frac{\sum x_i}{N + \beta}$
 $+ \frac{\beta}{N + \beta} \text{Mode}_{prior} = \frac{\alpha - 1}{N + \beta}$

$\frac{N}{N + \beta} \lambda_{MLE} + \frac{\beta}{N + \beta} \text{Mode}_{prior} = \text{Mode}_{posterior}$
 Linear combination

$\mu_{prior} = \frac{\alpha}{\beta}$, $\mu_{posterior} = \frac{\alpha + \sum x_i}{N + \beta}$, $\lambda_{MLE} = \frac{\sum x_i}{N}$

$\therefore \frac{N}{N + \beta} \lambda_{MLE} + \frac{\beta}{\beta + N} \mu_{prior} = \mu_{posterior}$
 Linear combination

P2. From given formula,
 $\Sigma_{n+1} = \left(\beta \sum_{n=1}^{N+1} x_n x_n^T + \lambda I \right)^{-1} = \left(\beta \sum_{n=1}^N x_n x_n^T + \lambda I + x_{N+1} x_{N+1}^T \right)^{-1}$
 $= \left(\Sigma_N^{-1} + x_{N+1} x_{N+1}^T \right)^{-1}$

Putting it in given identity,

$$\left(\Sigma^{-1} + x_{n+1} x_{n+1}^T \right)^{-1} = \Sigma^{-1} - \frac{(\Sigma^{-1} x_n)(x_n^T \Sigma^{-1})}{(1 + x_n^T \Sigma^{-1} x_n)}$$

$$\text{Now, } \sigma_{n+1}^2(x^*) = \beta^{-1} + x_*^T \Sigma_{n+1}^{-1} x_*$$

$$= \beta^{-1} + x_*^T \left(\Sigma_n^{-1} - \frac{(\Sigma_n^{-1} x_n)(x_n^T \Sigma_n^{-1})}{1 + x_n^T \Sigma_n^{-1} x_n} \right) x_*$$

$$\sigma_{n+1}^2(x^*) = \underbrace{\beta^{-1} + x_*^T \Sigma_n^{-1} x_*}_{\sigma_n^2(x^*)} - \frac{x_*^T (\Sigma_n^{-1} x_n)(x_n^T \Sigma_n^{-1}) x_*}{1 + x_n^T \Sigma_n^{-1} x_n}$$

Relation for impact of data:

$$\sigma_{n+1}^2(x^*) = \sigma_n^2(x^*) - \underbrace{\frac{x_*^T (\Sigma_n^{-1} x_n)(x_n^T \Sigma_n^{-1}) x_*}{1 + x_n^T \Sigma_n^{-1} x_n}}_{\Delta_n}$$

Now, Σ_n^{-1} is a p.s.d. matrix (covariance matrix)

$$\therefore \forall v \in W, v^T \Sigma_n^{-1} v \geq 0$$

$$\text{or, } x_n^T \Sigma_n^{-1} x_n \geq 0 \Rightarrow \text{denominator} \geq 1$$

Now if $(x_*^T \Sigma_n^{-1} x_n)$ is q then

$$x_n^T \Sigma_n^{-1} x_* \text{ is } q^T = x_n^T \Sigma_n^{-1} x_*$$

$$\therefore \Delta_n = \frac{(q q^T)}{1 + x_n^T \Sigma_n^{-1} x_n} \quad \text{Also, } q \text{ is scalar, (from dimension)}$$

Symmetric matrix

$$\therefore q q^T = q^2$$

$$\therefore \Delta_n = \frac{q^2}{1 + x_n^T \Sigma_n^{-1} x_n} \geq 0$$

So, since $\Delta_n \geq 0$

$$\therefore \sigma_{n+1}^2(x^*) = \sigma_n^2(x^*) - \Delta_n \leq \sigma_n^2(x^*)$$

Thus, variance decreases for a new observation with data (as expected).

P3

Class - distributions:

For real x_n : $P(x_n | y_n)$: Gaussian
params: Σ_n, μ_n

For binary x_n : $P(x_n | y_n = \frac{1}{0})$: Bernoulli
params: bias (p)

for V -valued x_n : $P(x_n | y_n = v')$: Multinoulli
Params: $P_{d,k}$ (prob. of x_n being as d given, y_n is k)

Formally,

for real x_n , Total params: $2DK$ (D values for each of K classes)
 $\cdot 2DK \rightarrow DK$ means and DK variances
 $+ (K-1) \rightarrow$ prior prob. $\text{sum}(y_n = K)$

$$\therefore P(x_n | y_n = k) = \mathcal{N}(\mu_{d,k}, \sigma_{d,k}^2)$$

Define an indicator variable,

$$z_k(y_n) = \begin{cases} 1 & \text{if } y_n = k \\ 0 & \text{o/w} \end{cases}$$

$$\therefore N_k = \sum_{n=1}^N z_k(y_n)$$

examples of class K .

\therefore MLE estimators are,

Params
for
 x_{nd}

$$\left\{ \begin{aligned} \mu_{d,k} &= \frac{\sum_{n=1}^N x_{nk} z_k(y_n)}{N_k} \\ \sigma_{d,k}^2 &= \frac{\sum_{n=1}^N (x_{nd} - \mu_{d,k})^2 z_k(y_n)}{N_k} \end{aligned} \right\}$$

and, for priors,

$$P(y_n = k | \pi) = \pi_k = \frac{N_k}{\sum_{k=1}^K N_k}$$

= fraction of examples in class K .

(*) If x_{nd} is binary,

$$P(x_{nd} = 1 | y_n = k) = p_{d,k}$$

\therefore there are a total of Dk parameters.

and also $(K-1)$ params for priors of $y(y_n = k)$

$$\text{Again, } z_k(y_n) = \begin{cases} 1 & \text{if } y_n = k \\ 0 & \text{o/w} \end{cases}$$

$$\text{So, } \left\{ \begin{aligned} N_k &= \sum_{n=1}^N z_k(y_n) \\ \pi_k &= \frac{N_k}{\sum_{k=1}^K N_k} \end{aligned} \right\} \text{ and MLE for } p_{d,k} \left\{ \begin{aligned} p_{d,k} &= \frac{\sum_{n=1}^N x_{nd} z_k(y_n)}{N_k} \end{aligned} \right.$$

So, as in 11 for case 1, the estimates are,

$$p_{d,k} = \frac{\sum_{n=1}^N x_{nd} z_k(y_n)}{N_k}, \quad \pi_k = \frac{N_k}{\sum_{k=1}^K N_k}$$

(*) x_{nd} is V -valued,

$$\therefore P(x_{nd} = i | y_{nd} = k) = p_{i,d,k} \quad (\text{multinomial dist.})$$

$$\text{with } \sum_{i=1}^V p_{i,d,k} = 1.$$

$$\therefore \text{total params} = \underbrace{DK(V-1)}_{\text{for distribution on } n} + \underbrace{(K-1)}_{\text{for prior}}$$

$$\text{Again, } \pi_k = \frac{N_k}{\sum_{k=1}^K N_k}, \quad P(y_n = k | \pi) = \pi_k$$

$$\text{Define, } z_{v,d}(x_n) = \begin{cases} 1 & \text{if } x_{nd} = v \\ 0 & \text{o/w} \end{cases}$$

$$p_{i,d,k} = \frac{\sum_{n=1}^N z_{v,d}(x_n) z_k(y_n)}{N_k}$$

MLE

estimate
(based on other
answers)

P4

$K=3$ is the best model because it has the max. marginal likelihood

$$P(y|K=3) = 2.5e-10 > P(y|K = \text{any other value})$$

$\therefore K=3$ is good since,

$$\frac{P(y|M_1)}{P(y|M_2)} = \frac{P(M_1|y)}{P(M_2|y)} \frac{P(M_1)}{P(M_2)}$$

assuming $P(M_1) = P(M_2)$ higher $P(y|M)$ means better model.

If we can include a new point, it's better to include at $x=-4$, since the var. of uncertainty is max. at $x=\pm 4$ for all $K=1,2,3$

at $K=3$, var at $x=-4$ is 7.180

(max. at x at all $x \in [-4, 4]$)

————— x —————