

## 1 Jensen's inequality

$f(x)$  を convex downward とする .  $\forall \lambda_n \geq 0, x_n (n = 1, \dots, N) | \sum_{n=1}^N \lambda_n = 1$  のとき ,

$$\sum_{n=1}^N \lambda_n f(x_n) \geq f\left(\sum_{n=1}^N \lambda_n x_n\right) \quad (1.1)$$

が成り立つ . これを Jensen's inequality という .

*Proof.* convex downward の定義より ,

$$\begin{aligned} \forall \lambda, x_1, x_2 | 0 \leq \lambda \leq 1 \Rightarrow \\ \lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2) \end{aligned} \quad (1.2)$$

なので ,  $N=2$  のときには Jensen's inequality は明らかに成り立つ .

また ,  $N=k$  のときに Jensen's inequality が成り立つとすると ,

$$\begin{aligned} \sum_{n=1}^{k+1} \lambda_n f(x_n) &= \sum_{n=1}^k \lambda_n f(x_n) + \lambda_{k+1} f(x_{k+1}) \\ &\quad \left( A = \sum_{n=1}^k \lambda_n \text{ として } \right) \\ &= A \sum_{n=1}^k \frac{\lambda_n}{A} f(x_n) + \lambda_{k+1} f(x_{k+1}) \\ &\geq A \sum_{n=1}^k \frac{\lambda_n}{A} f(x_n) + \lambda_{k+1} f(x_{k+1}) \\ &\geq A f\left(\sum_{n=1}^k \frac{\lambda_n}{A} x_n\right) + \lambda_{k+1} f(x_{k+1}) \quad \left( \because \sum_{n=1}^k \frac{\lambda_n}{A} = 1, \text{ 上の仮定} \right) \\ &\geq f\left(\sum_{n=1}^{k+1} \lambda_n x_n\right) \quad \left( \because A + \lambda_{k+1} = 1, \text{ 凸関数の定義} \right) \end{aligned} \quad (1.3)$$

よって帰納的に , 任意の整数  $N > 2$  について Jensen's inequality が成り立つ . Q.E.D.

□

直感的には , 任意の凸関数上の点をとったときに , その点を結んでできる多角形の内部の点は凸関数上の点よりも上にあるということで理解できる .

## 2 Laplace's method

Laplace's method (smoothing) とは , ある分布  $q(\theta)$  を Gaussian で近似する手法である . ここでは , Laplace's method による分布の積分の近似を示す . [2]

パラメータ  $\theta$  の分布に関する積分を考える ( $n$  は標本数)

$$\int \exp(nq(\theta)) d\theta \quad (2.1)$$

を考える． $\theta$  の mode を  $\hat{\theta}$  とし，その近傍でテイラー展開する．( $n$  が十分に大きければ  $\theta$  は mode に集中すると仮定して近似している)． $\frac{\partial q(\theta)}{\partial \theta} \Big|_{\hat{\theta}} = 0$  であるので，

$$q(\theta) = q(\hat{\theta}) - \frac{1}{2}(\theta - \hat{\theta})^T J(\hat{\theta})(\theta - \hat{\theta}) + o((\theta - \hat{\theta})^2) \quad (2.2)$$

ここで，

$$J(\hat{\theta}) = - \frac{\partial^2 q(\theta)}{\partial \theta \partial \theta^T} \Big|_{\hat{\theta}} \quad (2.3)$$

とした．

$$\int \exp \left\{ -\frac{1}{2}(\theta - \hat{\theta})^T J(\hat{\theta})(\theta - \hat{\theta}) \right\} d\theta = (2\pi)^{p/2} n^{-p/2} |J(\hat{\theta})|^{-1/2} \quad (2.4)$$

より，結局，

$$\int \exp(nq(\theta)) d\theta \approx (2\pi)^{p/2} n^{-p/2} |J(\hat{\theta})|^{-1/2} \exp(nq(\hat{\theta})) \quad (2.5)$$

と近似できる．

### 3 BIC

モデルの事後分布に基づくモデルの評価基準である BIC(Baysian Information criterion) を示す．一般に尤度を用いてモデルの比較を行うことはできない．なぜなら，尤度にはパラメータの次元などによるバイアスが含まれているからである [2]．モデルを複雑にすればするほど尤度を大きくすることができる問題からこのことは明らかである．BIC では事後分布を考えることによりこの問題を解決してる．ただ，BIC では簡単にモデルの評価を行うことができるが，パラメータ事前分布が広がっており，そのヘッセ行列が非退化であるという仮定が多くの場合に妥当でないという問題もある [1]．

$P$  次元のパラメータ  $\theta$  を持つモデル  $M$  のもとで観測データ  $X = \{x_1, \dots, x_N\}$  が観測されたときを考える．このとき，モデルの事後分布はベイズの定理より，

$$p(M|X) = \frac{p(X|M)p(M)}{\int p(X|M)p(M)dM} \quad (3.1)$$

となる．この事後確率を最大化するようなモデルを選ぶことが BIC の目標である．ここで， $p(M)$  は一定とすると，分母は定数であるので， $p(X|M)$  を最大化すれば良いことがわかる．これをパラメータ  $\theta$  を用いてあらわすと，

$$p(X|M) = \int p(X|\theta)p(\theta|M)d\theta \quad (3.2)$$

となる．この周辺分布  $p(X|M)$  を Laplace's method(section2) 用いて近似する．

$L(\theta) = \ln p(X|\theta)$ ,  $\pi(\theta = p(\theta))$  として，ラプラス近似を行う．

$$\begin{aligned}
p(X|M) &= \int \exp(L(\boldsymbol{\theta}))\pi(\boldsymbol{\theta})d\boldsymbol{\theta} \\
&\approx \int \exp\left\{L(\hat{\boldsymbol{\theta}}) - \frac{N}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T J(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\right\} \times \left\{\pi(\hat{\boldsymbol{\theta}}) + (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \frac{\partial \pi(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right\} d\boldsymbol{\theta} \\
&\approx \exp(L(\hat{\boldsymbol{\theta}}))\pi(\hat{\boldsymbol{\theta}}) \int \exp\left\{-\frac{N}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T J(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\right\} d\boldsymbol{\theta} \\
&= \exp(L(\hat{\boldsymbol{\theta}}))\pi(\hat{\boldsymbol{\theta}})(2\pi)^{P/2}N^{-P/2}|J(\hat{\boldsymbol{\theta}})|^{-1/2}
\end{aligned} \tag{3.3}$$

ただし ,

$$J(\hat{\boldsymbol{\theta}}) = -\frac{1}{N} \frac{\partial^2 L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \tag{3.4}$$

として ,

$$\int (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \exp\left\{-\frac{N}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T J(\hat{\boldsymbol{\theta}})(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})\right\} d\boldsymbol{\theta} = 0 \tag{3.5}$$

を用いた . 対数を取って-2 を乗じて ,

$$-2 \ln(p(X|M)) = -2L(\hat{\boldsymbol{\theta}}) - 2 \ln(\pi(\hat{\boldsymbol{\theta}})) + P \ln N + \ln |J(\hat{\boldsymbol{\theta}})| \tag{3.6}$$

$N$  に関して  $O(1)$  の項を無視して ,

$$BIC = -2L(\hat{\boldsymbol{\theta}}) + P \ln N \tag{3.7}$$

を得る .

## References

- [1] Christopher Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 1st ed. 2006. Corr. 2nd printing 2011. Springer, Feb. 2010. ISBN: 9780387310732. URL: <http://amazon.co.jp/o/ASIN/0387310738/>.
- [2] 小西 貞則 and 北川 源四郎. 情報量規準 (シリーズ・予測と発見の科学). 朝倉書店, Sept. 2004. ISBN: 9784254127829. URL: <http://amazon.co.jp/o/ASIN/4254127820/>.