

MLE estimate

Consistent, asymptotically normal ($\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \sim N(0, J^{-1} I J^{-1})$) as $n \rightarrow \infty$)

Efficient: ($\hat{\theta}_{ML}$ minimize $E((\hat{\theta}_{ML} - \theta_0)^2)$ namely variance as $n \rightarrow \infty$)

Equivariance: if $\hat{\theta}_n$ is MLE of θ , $g(\hat{\theta}_n)$ is MLE of $g(\theta)$

Cramer Rao lower bound
For unbiased estimator

$$E[(\hat{\theta} - \theta_0)^2] \geq \frac{1}{I_n(\theta_0)}$$

Where $I_n(\theta_0)$ is the Fisher information

$$Var[(\frac{\partial}{\partial \theta} \log P(X_{1,...,n} | \theta))] = I_n(\theta)$$

LS estimate

MLE when noise follows iid Gaussian with constant σ

$$RSS(\beta) = (y - X\beta)^T (y - X\beta)$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$$\hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2)$$

Optimality: smallest variance among all linear unbiased estimates

Expected prediction error:

$$E_D E_{Y|X} (\hat{f}(x) - Y)^2 = E_D \left(\hat{f}(x) - E(Y|X) \right)^2$$

$$+ E_{Y|X} (Y - E(Y|X))^2$$

$$= E_D \left(\hat{f}(x) - E_D \hat{f}(x) \right)^2$$

$$+ E_D \left(E_D \hat{f}(x) - E(Y|X) \right)^2$$

$$+ E_{Y|X} (Y - E(Y|X))^2$$

Expected risk

$$\text{Total: } R(f) = E_{X,Y} [Q(Y, f(X))] = \int_Y \int_X Q(Y, f(X)) P(X, Y) dX dY$$

$$\text{Conditional: } R(f, X) = \int_Y Q(Y, f(X)) P(Y|X) dY$$

Connect to test error: The test error can be used to approximate the expected risk. The expected value of the test error is the expected

Ridge and Lasso

$$\text{Reg: } \lambda \beta^T \beta; \lambda ||\beta||_1$$

$$\hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

both assume

$$Y|(X, \beta) \sim N(X^T \beta, \sigma^2 I)$$

Difference:

1. **Prior** on β

$$\beta^{\text{ridge}} \sim N(0, \frac{\sigma^2}{\lambda} I)$$

$$P(\beta^{\text{Lasso}}) = \frac{\lambda}{4\sigma^2} \exp\left(-\frac{|\beta|\lambda}{2\sigma^2}\right)$$

2. **Cost function**

Ridge is strictly convex while

Lasso is just convex. Not

differentiable when $|\beta|=0$

Spline regression

$$RSS(f, \lambda)$$

$$= \sum_{i=1}^n (y_i - f(x_i))^2$$

$$+ \lambda \int (f''(x))^2 dx$$

Gaussian process

Assumption:

$$1. P(Y|X, \beta, \sigma) \propto$$

$$\exp\left(-\frac{1}{2\sigma^2} (Y - X^T \beta)^2\right)$$

$$2. P(\beta|\Lambda) \propto \exp\left(-\frac{1}{2} \beta^T \Lambda \beta\right)$$

Posterior distribution of β

$N(\beta|\mu_\beta, \Sigma_\beta)$ where

$$\mu_\beta = (X^T X + \sigma^2 \Lambda)^{-1} X^T y$$

$$\Sigma_\beta = \sigma^2 (X^T X + \sigma^2 \Lambda)^{-1}$$

$$\text{Cov}[y] = X \Lambda^{-1} X^T + \sigma^2 I$$

distribution of y

$$p\left(\begin{bmatrix} y \\ y_{n+1} \end{bmatrix} \middle| \begin{bmatrix} x_{n+1}, X, \sigma \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} y \\ y_{n+1} \end{bmatrix} \middle| \begin{bmatrix} C_n & k \\ k^T & c \end{bmatrix}, \begin{bmatrix} C_n & k \\ k^T & c \end{bmatrix}\right),$$

where $C_n = K + \sigma^2 I$,
 $c = k(x_{n+1}, x_{n+1}) + \sigma^2$,
 $k = k(x_{n+1}, X)$,
and $K = k(X, X)$.

$$\mu_{y_{n+1}} = K^{-1} C_n^{-1} y$$

$$\sigma_{y_{n+1}}^2 = c - K^{-1} C_n^{-1} k$$

Model averaging

Bias $[\hat{f}(x)] =$

$$\frac{1}{B} \sum_{i=1}^B E_D(\hat{f}_i(x)) - E(Y|X) =$$

$$\frac{1}{B} \sum_{i=1}^B \text{Bias}[\hat{f}_i(x)]$$

$$V[\hat{f}(x)] = \frac{1}{B^2} \sum_{i=1}^B V[\hat{f}_i(x)] +$$

$$\frac{1}{B^2} \sum_{i \neq j} \text{Cov}(\hat{f}_i(x), \hat{f}_j(x))$$

K-fold CV: tend to underfit

LOOCV: unbiased but the variance can be very large

due to highly correlated training sets

Bootstrap: too optimistic due to overlap of training and test set

$$\hat{R}^* = \frac{1}{B} \sum_{b=1}^B \sum_{i=1}^n l(y_i, \hat{f}^{*b}(x_i))$$

The chance of a sample to have appeared in the bootstrap is $1 -$

$$\left(1 - \frac{1}{n}\right)^n \approx 0.632$$

The “Jackknife” method

With LOOCV

$$\text{bias}^{\text{JK}} = (n-1)(\tilde{S}_n - \hat{S}_n)$$

$$\tilde{S}^{\text{JK}} = \hat{S}_n - \text{bias}^{\text{JK}}$$

$$\text{with } \tilde{S}_n = \frac{1}{n} \sum_{i=1}^n \hat{S}_{n-1}^{-i}$$

$$E[\text{bias}[\hat{S}_n]] = \frac{a_1}{n} + \frac{a_2}{n^2} + \dots$$

$$E[\text{bias}^{\text{JK}}] = \frac{a_1}{n} + O(n^{-2})$$

With bootstrap (debiasing)

$$\text{bias}_{\text{boot}} = \frac{1}{B} \sum_{b \leq B} \hat{S}^*(b) - \hat{S}_n$$

$$\bar{S} = \hat{S}_n - \text{bias}$$

$$= 2 \hat{S}_n - \frac{1}{B} \sum_{b \leq B} \hat{S}^*(b)$$

Gradient Descent procedure

Algorithm: Initialize α , threshold θ , $\eta(0)$, $k \leftarrow 0$
Repeat

$$\alpha \leftarrow \alpha - \eta(k) \nabla J(\alpha)$$

$$k \leftarrow k + 1$$

$$\text{Until } |\eta(k) \nabla J(\alpha)| < \theta$$

Optimal learning rate

$$H = \frac{\partial^2 J}{\partial \alpha_i \partial \alpha_j}$$

$$J(a(k+1))$$

$$\approx J(a(k)) + \nabla J^T(a(k)) (\alpha(k+1) - a(k))$$

$$- a(k) + 1/2 ((\alpha(k+1) - a(k))^T H (\alpha(k+1) - a(k)))$$

$$= J(a(k)) - \eta(k) \nabla J^T(a(k)) \nabla J$$

$$+ 1/2 \eta^2(k) \nabla J^T(a(k)) H \nabla J$$

$$\frac{\partial}{\partial \eta} J(a(k+1)) = 0 \Rightarrow$$

$$\eta^{\text{opt}} = \frac{\|\nabla J\|^2}{\nabla J^T H \nabla J}$$

Newton’s algorithm

$$\frac{\partial}{\partial a(k+1)} J(a(k+1)) = 0 \Rightarrow$$

$$a(k+1) = a(k) - H^{-1} \nabla J$$

Perceptron

$$J_p(a) = \sum_{x \in X_{mc}} (-y a^T x)$$

Update rule: $a(k+1) =$

$$a(k) + \eta(k) \sum_{x \in X_{mc}} y x$$

WINNOW algorithm

Performs better when many dimensions are irrelevant.

Search for 2 weight vectors a^+, a^- (for each class). If a

point is misclassified: $a_i^+ \leftarrow$

$$\alpha^{+x^{ki}} a_i^+; a_i^- \leftarrow \alpha^{-x^{ki}} a_i^-$$

(if $y_k = +1$); $a_i^+ \leftarrow$

$$\alpha^{-x^{ki}} a_i^+; a_i^- \leftarrow \alpha^{+x^{ki}} a_i^-$$

(if $y_k = -1$)

Fisher’s LDA

$$\tilde{m}_\alpha = w^T m_\alpha = w^T \frac{1}{|X_\alpha|} \sum_{x \in X_\alpha} x$$

$$\Sigma_\alpha = \Sigma_{x \in X_\alpha} (x - m_\alpha)(x - m_\alpha)^T$$

$$\Sigma_W = \Sigma_1 + \Sigma_2;$$

$$\tilde{\Sigma}_1 + \tilde{\Sigma}_2 = w^T \Sigma_W w$$

Fisher’s Separation Criterion

$$J(w) = \frac{\|\tilde{m}_1 - \tilde{m}_2\|^2}{\tilde{\Sigma}_1 + \tilde{\Sigma}_2}$$
$$= \frac{w^T (m_1 - m_2)(m_1 - m_2)^T w}{w^T \Sigma_W w}$$

$$= \frac{w^T \Sigma_B w}{w^T \Sigma_W w}$$

$$\frac{d}{dw} J(w) = 0 \Rightarrow \Sigma_B w$$

$$= \Sigma_W w \frac{w^T \Sigma_B w}{w^T \Sigma_W w} \Rightarrow \Sigma_W^{-1} \Sigma_B w$$

$$= \lambda w \text{ with } \lambda = \frac{w^T \Sigma_B w}{w^T \Sigma_W w}$$

$$\hat{w} = \Sigma_W^{-1} (m_1 - m_2)$$

complexity: $O(nd^2)$

Logistic regression

$$R(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$P(Y = y|x) = \frac{1}{1 + \exp(-y w^T x)}$$

Kernel

$$\text{RBF: } k(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

$$\text{Sigmoid: } k(x, z) = \tanh(k_{xz} - b)$$

$$k(x, z) = \exp(-\alpha \|x - z\|)$$

$$k((x, h), (y, h')) = p(h|x)p(h'|y)$$

Properties:

$$k(x, z) = k_1(\phi(x), \phi(z))$$

$$k(x, z) = \exp(k_1(x, z))$$

$$k(x, z) = f(x)f(y) \text{ (f is real-valued function)}$$

$$k(x, z) = p(k_1(x, z)) \text{ (p(x) is polynomial with positive coefficients)}$$

SVM

Hard margin

$$\text{Minimize } \frac{1}{2} w^T w$$

$$\text{s.t } \forall i: z_i (w^T y_i + w_0) \geq 1$$

Lagrange Function:

$$L(w, w_0, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^n \alpha_i [z_i (w^T y_i + w_0) - 1]$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_{i \leq n} \alpha_i z_i y_i$$

$$\frac{\partial L}{\partial w_0} = 0 \Rightarrow - \sum_{i \leq n} \alpha_i z_i = 0$$

Dual problem:

$$\max W(\alpha) = \sum_{i=1}^n \alpha_i -$$

$$\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n z_i z_j \alpha_i \alpha_j y_i^T y_j \text{ s.t}$$

$$\forall i \alpha_i \geq 0 \wedge \sum_{i \leq n} \alpha_i z_i = 0$$

$$w^* = \sum_{i=1}^n \alpha_i^* z_i y_i$$

$$w_0^* = -\frac{1}{2} (\min_{i: z_i=1} w^{*T} y_i + \max_{i: z_i=-1} w^{*T} y_i)$$

KKT conditions require

$$\alpha_i [z_i (w^T y_i + w_0) - 1] = 0$$

$$\alpha_i \geq 0$$

$$z_i (w^T y_i + w_0) - 1 \geq 0$$

$$\text{Soft margin}$$

$$\text{Min } \frac{1}{2} w^T w + C \sum_{i=1}^n \mathcal{E}_i$$

$$\text{s.t } \forall i: z_i (w^T y_i + w_0) \geq 1 - \mathcal{E}_i$$

$$\mathcal{E}_i \geq 0$$

$$L(w, w_0, \mathcal{E}, \alpha, \beta)$$

$$= \frac{1}{2} w^T w + C \sum_{i=1}^n \mathcal{E}_i$$

$$- \sum_{i=1}^n \alpha_i [z_i (w^T y_i + w_0) - 1 + \mathcal{E}_i]$$

$$- \sum_{i=1}^n \beta_i \mathcal{E}_i$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow C - \alpha_i - \beta_i = 0$$

Dual problem is same, only different in $\forall i \ C \geq \alpha_i \geq 0 \wedge \sum_{i \leq n} \alpha_i z_i = 0$

KKT conditions require

$$\alpha_i [z_i (w^T y_i + w_0) - 1 + \varepsilon_i] = 0$$

$$\varepsilon_i (C - \alpha_i) = 0$$

Multiclass

$$(w_{z_i}^T y_i + w_{z_i,0}) - \max_{z \neq z_i} w_z^T y_i + w_{z,0} \geq m$$

Structured SVMs

Joint feature map $\Psi(z, y)$

Scoring function $f_w(z, y) = w^T \Psi(z, y)$

Classification: $h(y) = \operatorname{argmax}_{z \in \mathbb{K}} f_w(z, y)$

Objective:

$$\text{Min } \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

s.t $\xi_i \geq 0$

$$\begin{aligned} & w^T \Psi(z_i, y_i) - w^T \Psi(z, y_i) \\ & \geq \Delta(z, z_i) - \varepsilon_i \quad \forall y_i \in \mathcal{Y}, z \neq z_i \\ & \Leftrightarrow w^T \Psi(z_i, y_i) - \max_{z \neq z_i} [\Delta(z, z_i) \end{aligned}$$

$$+ w^T \Psi(z, y_i)] \geq -\varepsilon_i \quad \forall y_i \in \mathcal{Y}$$

Lagrangian: $\frac{1}{2} w^T w +$

$$C \sum_{i=1}^n \varepsilon_i -$$

$$\sum_{i=1}^n \sum_{z_j \in \mathbb{K}_i} a_{ij} [w^T \Psi_i(z_j) -$$

$$\Delta_i(z_j) + \varepsilon_i] - \sum_{i=1}^n \beta_i \varepsilon_i \text{ where}$$

$$\mathbb{K}_i = \mathbb{K} \setminus \{z_i\}$$

Ensemble methods

Bagging

For $b=1$ to B do

Z^{*b} = b -th bootstrap sample from Z

Construct classifier c_b based on Z^{*b}

end for

return $\hat{c}_B(x) = \operatorname{sgn}(\sum_{i=1}^B c_i(x))$

Random Forest

For $b=1$ to B do

Z^{*b} = b -th bootstrap sample from Z

Repeat

1. Select m variables at random

from p variables ($m \approx \sqrt{p}$)

2. pick the best variable/split-point among selected variables

3. split the node into two

daughter node

Until node size n_{\min} is reached

End for

Return the ensemble of

trees $\{\hat{c}_b(x)\}_{b=1}^B$

Adaboost

Initialize $w_1 = 1/n$

For $b = 1$ to B do

1. Fit a classifier $c_b(x)$ using weight w_i

2. Compute $\epsilon_b \leftarrow$

$$\sum_{i=1}^n w_i^{(b)} \mathbb{I}_{\{c_b(x_i) \neq y_i\}} / \sum_{i=1}^n w_i^{(b)}$$

3. Compute $\alpha_b \leftarrow \log \frac{1 - \epsilon_b}{\epsilon_b}$

4. Set $\forall i: w_i \leftarrow$

$$w_i \exp(\alpha_b \mathbb{I}_{\{c_b(x_i) \neq y_i\}})$$

End for

Return $\hat{c}_B(x) =$

$$\operatorname{sgn}(\sum_{i=1}^n \alpha_b c_b(x))$$

Comparison

(a) same training set on every

iteration; varies the training

sets using resampling.

(b) weights according to its

accuracy; same weight

(c) Both are good at reducing

variance but only Boosting tries

to reduce bias

PAC learning

Efficiently PAC learnable:

if an algorithm A receives as

input a sample Z of size $n \geq$

$\operatorname{poly}(\frac{1}{\epsilon}, \frac{1}{\delta})$, then A outputs

\hat{c} such that

$$P(R(\hat{c}) \leq \epsilon) \geq 1 - \delta$$

Axis-aligned rectangles

$$P(R(\hat{c}) \leq \epsilon) \geq P(\hat{R}IG)$$

$$\geq 1 - 4e^{-\frac{n\epsilon}{4}}$$

$$P(-\hat{R}IG) =$$

$P(\hat{R} \text{ doesn't intersect } T_i^\epsilon \text{ for some } i) \leq$

$$\sum_i P(\hat{R} \cap T_i^\epsilon = \emptyset) =$$

$$\sum_i \prod_{j \leq n} P(x_j \notin T_i^\epsilon) = 4 \left(1 - \right.$$

$$\left. \frac{\epsilon}{4} \right)^n \leq 4e^{-\frac{n\epsilon}{4}}$$

$$1 - 4e^{-\frac{n\epsilon}{4}} \geq 1 - \delta \Leftrightarrow n \geq \frac{4}{\epsilon} \ln \frac{4}{\delta}$$

$$\left[\frac{4}{\epsilon} \ln \frac{4}{\delta} \right] \text{ is polynomial in } \frac{1}{\epsilon} \text{ and } \frac{1}{\delta}$$

When $\hat{R}_n(\hat{c}) = 0$ and $|\mathcal{H}|$

is finite

$$n \geq \frac{1}{\epsilon} (\log |\mathcal{H}| + \log 1/\delta)$$

Then $P(R(\hat{c}) \leq \epsilon) \geq 1 - \delta$

Proof: $P(R(\hat{c}) > \epsilon) =$

$$P(\hat{R}_n(\hat{c}) = 0 \wedge R(\hat{c}) > \epsilon) \leq$$

$$\Pr(\exists h \in H: \hat{R}_n(h) = 0 \wedge$$

$$R(h) > \epsilon) \leq \sum_{h \in H} P(\hat{R}_n(h) =$$

$$0 \wedge R(h) >$$

$$\epsilon) = \sum_{h \in H} P(\hat{R}_n(h) =$$

$$0 | R(h) > \epsilon) P(R(h) > \epsilon) \leq$$

$$\sum_{h \in H} (R_n(h) = 0 | R(h) > \epsilon) \leq$$

$$\sum_{h \in H} \prod_{j \leq n} P(x_j \notin h\Delta c | \Pr(h\Delta c)) <$$

$$\epsilon) \leq |H| (1 - \epsilon)^n$$

Under the stochastic setting

$$P(R(\hat{c}) - \inf_{c \in C} R(c) \leq \epsilon) \geq$$

$$1 - \delta$$

If \mathcal{C} is finite,

$$P(R(\hat{c}) - \inf_{c \in C} R(c) > \epsilon) \leq$$

$$2|\mathcal{C}| \exp(-2n\epsilon^2)$$

VC_c is VC dimension of a

concept class \mathcal{C} (complexity

measure)

$$P(R(\hat{c}) - \inf_{c \in C} R(c) > \epsilon) \leq$$

$$9 n^{VC_c} \exp\left(-\frac{n\epsilon^2}{32}\right) \rightarrow$$

0 when $n \rightarrow$

∞ if VC_c is finite

Require uniform convergence

$$R(\hat{c}) - \inf_{c \in C} R(c) \leq 2 \sup_{c \in C} |\hat{R}_n(c)$$

$$- R(c)|$$

$$a. P(X \geq \epsilon) \leq \frac{E[X]}{\epsilon} \text{ when } X \geq 0$$

$$b. E[\exp(sX)] \leq \exp\left(\frac{s^2(b-a)^2}{8}\right)$$

$$\text{when } E[X] = 0, a \leq X \leq b, s > 0$$

$$c. P(S_n - ES_n \geq t) \leq$$

$$\exp\left(-\frac{t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

$$\text{when } S_n = \sum_{i=1}^n X_i, X_i \text{ is}$$

independent, $t > 0$

The Hoeffding bound

$$P(\tilde{S}_n - ES_n \geq \epsilon) \leq$$

$$\exp\left(-\frac{2n\epsilon^2}{\sum_{i=1}^n \frac{(b_i - a_i)^2}{n}}\right)$$

$$\text{Where } \tilde{S}_n = \frac{S_n}{n} \quad t = n\epsilon$$

Union Bound

$$P(\max_i X_i > \epsilon) \leq \sum_i P(X_i > \epsilon)$$

ERM for Hyperplanes

$$P(R(\hat{c}) - \inf_{c \in C} R(c) > \epsilon) \leq$$

$$(1 + 2^{\binom{n}{d}}) e^{2d\epsilon} e^{-\frac{n\epsilon^2}{2}}$$

$$\text{If } R(c^*) = 0 \text{ then } P(R(\hat{c}) \leq$$

$$\epsilon) \leq 2^{\binom{n}{d}} \exp(-\epsilon(n-d))$$

Nonparametric Bayesian

Methods

$$\text{Beta}(x|a, b) = \frac{1}{B(a, b)} x^{a-1} (1-x)^{b-1}$$

Stick-Breaking Process(GEM)

Repeatedly draw β_i from

$$\text{Beta}(1, \alpha) \quad \rho_k =$$

$$\left[\prod_{i=1}^{k-1} (1 - \beta_i) \right] \beta_k$$

Connect to DP

$$G(\theta) = \sum_{k=1}^{\infty} \rho_k \delta_{\theta_k}(\theta)$$

Chinese Restaurant Process

$P(\text{customer } n+1 \text{ joins table } t)$

$$= \frac{|t|}{\alpha + n}$$

$$\frac{\alpha}{\alpha + n}$$

$$\frac{\alpha}{\alpha + n}$$

$$\frac{\alpha}{\alpha + n}$$

DP mixture model

Probability of clusters

$$\rho \sim GEM(\alpha)$$

Center of clusters

$$\mu_k \sim N(\mu_0, \sigma_0)$$

Assignment of data points

$$z_i \sim \text{Categorical}(\rho)$$

Coordinates of data points

$$x_i \sim N(\mu_{z_i}, \sigma)$$

Asymptotics of

$$S(n) = \sum_{i=1}^n \frac{\alpha}{a+i-1} \rightarrow$$

$$a \ln(n) \text{ as } n \rightarrow \infty$$

$$S(n) < I(n) = \int_1^{n+1} \frac{\alpha}{a+x-1} dx$$

$$< \sum_{i=2}^{n+1} \frac{\alpha}{a+i-1}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$= S(n) - 1 + \frac{\alpha}{a+n}$$

$$P(z_i = k | z_{-i}, x, \alpha, \mu)$$

$$= \begin{cases} \frac{N_{k,-i}}{\alpha + N - 1} p(x_i | x_{-i,k}, \mu) \\ \frac{\alpha}{\alpha + N - 1} P(x_i | \mu) \end{cases}$$

```

1: for i = 1 to N in random order do
2:   Remove x_i's sufficient statistics from old cluster z_i;
3:   for k = 1 to K do
4:     Compute p_k(x_i) = p_k(x_i | x_{-i,k});
5:     Set N_{k,-i} = |x_{-i,k}|;
6:   Compute p(z_i = k | z_{-i}, x) = \frac{N_{k,-i}}{\alpha + N - 1};
7: end for
8: Compute p_*(x_i) = p(x_i | \mu);
9: Compute p(z_i = * | z_{-i}, x);
10: Normalize p(z_i | *);
11: Sample z_i ~ p(z_i | *);
12: Add x_i's sufficient statistics to new cluster z_i;
13: If any cluster is empty, remove it and decrease K;
14: end for

```

Extra

Taylor expansion

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}; \quad \frac{1}{1-x} = \sum_{n=0}^{\infty} x^n$$

$$\ln(1-x) = - \sum_{n=1}^{\infty} \frac{x^n}{n!}; \quad \ln(1+x)$$

$$= \sum_{n=1}^{\infty} (-1)^{n+1} \frac{x^n}{n!};$$

Gaussian

$$(2\pi)^{-\frac{k}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\text{AIC} : -2 \log(\hat{p}(X | \hat{\theta}_k)) + 2k$$

$$\text{BIC} : -2 \log(\hat{p}(X | \hat{\theta}_k, M_k)) + k' \log n$$