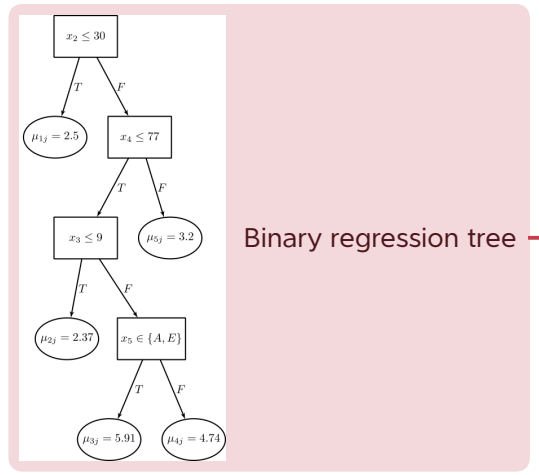


Nicholas, 2018

Bayesian additive regression tree (BART)



Building a sum of trees
 $y = \sum_{j=1}^m g(x; T_j, M_j) + \epsilon, \epsilon \sim N(0, \sigma^2)$

Prior of tree

- Prior of tree at depth d nonterminal
 $d_j = \frac{\alpha}{(1 + d)^\beta}, \alpha \in (0, 1), \beta \geq 0$
- The distribution on the splitting variable assignments at each interior node is the uniform prior on available variables
- The distribution on the splitting rule assignment in each interior node, conditional on the splitting variable is the uniform prior on the discrete set of available splitting values

Prior of mean

- $y \rightarrow [y_{\min} = -0.5, y_{\max} = 0.5] \Rightarrow \mu_{ij} \sim N(0, \sigma_\mu^2), \sigma_\mu = \frac{1}{2k\sqrt{m}}$

Prior of σ

- Chipman (2010)
 $(\nu, q) = (3, 0.9), P(\sigma < \hat{\sigma}) = q$
- $\sigma^2 \sim \frac{\nu\lambda}{\chi_\nu^2}$

Number of trees
 $m = 200$

regularization of prior

- The terminal node parameters of every tree should be independent.
- There should not be any correlation between trees
- Variance should differ
- Approximate normal distribution
- To decide the number of trees

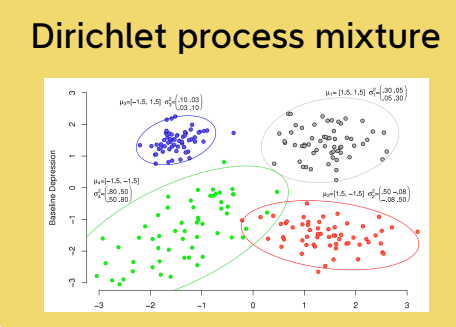
Posterior sampling

- backfitting MCMC algorithm
 $p((T_1, M_1), \dots, (T_m, M_m), \sigma | y)$
- Gibbs sampling
 $T_j | R_j, \sigma$
 $M_j | T_j, R_j, \sigma$
- Particle Gibbs algorithms

Partial residual
 $R_j = y - \sum_{k \neq j} g(x; T_k, M_k)$

Draw σ from full conditional
 $\sigma | T_1, \dots, T_m, M_1, \dots, M_m, y$

Prior independence

$$p((T_1, M_1), \dots, (T_m, M_m), \sigma) = \left[\prod_j p(T_j, M_j) \right] p(\sigma)$$
$$= \left[\prod_j p(M_j | T_j) p(T_j) \right] p(\sigma)$$
$$p(M_j | T_j) = \prod_i p(\mu_{ij} | T_j)$$


Dirichlet process mixture

Dirichlet distribution

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{k=1}^K x_k^{\alpha_k - 1}$$

Normalizing constant

$$B(\alpha) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}, \quad \alpha = (\alpha_1, \dots, \alpha_K)$$

Note

- G_0 is a base distribution
- α is a positive scaling parameter
- G is a random probability measure that has the same support as G_0

Dirichlet process: distribution over distributions

$$G \sim DP(\alpha, G_0), X_n | G \sim G$$

Example

Consider Gaussian G_0

Comparison

G_0 is continuous, so the probability that any two samples are equal is precisely zero.
However, G is a discrete distribution, made up of a countably infinite number of point masses (Dirac-delta).
Therefore, there is always a non-zero probability of two samples colliding

Stick breaking viewpoint

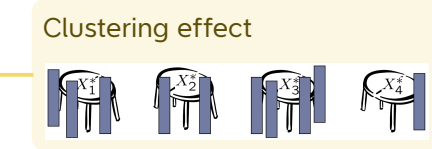
$$V_1, V_2, \dots, V_{\infty} \sim \text{Beta}(1, \alpha)$$
$$f(V_i = v_i) = \alpha(1 - v_i)^{\alpha-1}$$
$$X_1^*, X_2^*, \dots, X_{\infty}^* \sim G_0$$
$$z_i(x) = v_i \prod_{j=1}^i (1 - v_j)$$
$$G = \sum_{i=1}^{\infty} z_i(x) \delta_{X_i^*}$$

Sampling from DP

$$X_1, X_2, \dots, X_{N-1} \sim \begin{cases} X_1 \text{ (new draw from } G_0 \text{ with probability } \frac{\alpha}{\alpha+1}) \\ \text{with probability } \frac{1}{N-1} \end{cases}$$
$$P(X_1, \dots, X_N) = P(X_1)P(X_2|X_1) \dots P(X_N|X_1, \dots, X_{N-1})$$
$$= \alpha^K \frac{\Gamma(\sum_{k=1}^K \alpha_m(X_k) + 1)}{\alpha(1 + \alpha) \dots (\alpha + 1 + \alpha)} \prod_{k=1}^K G_0(X_k)$$

Chinese restaurant process: posterior of DP

Consider a restaurant with infinitely many tables, where the X_k 's represent the patrons of the restaurant. From the above conditional probability distribution, we can see that a customer is more likely to sit at a table if there are already more people sitting there. However, with probability proportional to α , the customer will sit at a new table.



Background

Given a data set, and are told that it was generated from a mixture of Gaussian distributions. But no one has any idea how many Gaussians produced the data.

Dirichlet process mixture model

Assumption
A finite mixture model assumes that the data come from a mixture of a finite number of distributions

Finite mixture models

$$\pi \sim \text{Dirichlet}(\alpha K, \dots, \alpha K)$$
$$c_n \sim \text{Multinomial}(\pi)$$
$$\eta_n \sim G_0$$
$$y_n | c_n, \eta_1, \dots, \eta_K \sim F(\cdot | \eta_{c_n})$$

Illustration

Infinite mixture models

Take limit as K goes to ∞
Note: the N data points still come from at most N different components (Blei et al 2003)

Inference for DPMM

Procedure (infer)

- Take a random guess
- Calculate mean for each cluster — Let #clusters grow or shrink
- Compute cluster assignment for each sample
- Update cluster by sampling from the distribution

Procedure (generate)

- Generate a new label assignment
- For unique assignment

Comparison between EM and Gibbs

- Gibbs often faster to implement
- EM easier to diagnose convergence
- EM can be parallelized
- Gibbs is more widely applicable

AFT model

Main model

Expression

- Hazard function
 $\lambda(t/x) = \exp(\beta'x)\lambda_0(\exp(\beta'x)t)$
- Survival function
 $S_i(t) = S_0(e^{x_i^T \beta} t)$

Common distributions of AFT models

- Weibull distribution and exponential AFT model
 $S(t) = \exp(-\lambda t^\alpha)$
- Log-normal AFT model
- Log-logistic distribution

MLE estimation

The log likelihood function is given by

$$\log L(\beta) = \sum_{i=1}^n \log \lambda(t_i | x_i) + \sum_{i=1}^n \log S(t_i | x_i)$$

Estimation of AFT models

AFT models are fitted by using maximum likelihood estimation (MLE) method. The likelihood of observed data is given by

$$L(\beta) = \prod_{i=1}^n \lambda(t_i | x_i) S(t_i | x_i)^{1 - \delta_i}$$

Advantage

However, since the PH model specifies the effect of covariates x is not multiplicative on the hazard function, it is not easy to interpret. For example, the estimate of relative risk is $\exp(\beta'x)$, which is not a probability. So, if the ordinary linear regression model can handle external observations, x

Discussion

Statistical issues

Illustration

Nonparametric AFT model
 $\log(T) = m(A, x) + W, A \in \{0, 1\}$

Model components

- DPMM
 $M \sim \text{Gamma}(\Psi_1, \Psi_2)$
 $G | M \sim \text{DPMM}(M, G_0), G_0 = N(0, \sigma^2)$
 $\pi_i | G \sim G$
- Prior specification
 $\Psi_1 = 2, \Psi_2 = 0.1$
- BART
 $m \sim \text{BART}(\alpha, \beta, k, J)$
- Prior specification
 $\alpha = 0.95, \beta = 2, k = 2, J = 200$
- Assumptions
 $y = \hat{\mu}_{AFT} + \epsilon, \epsilon \sim N(0, \sigma_{AFT}^2)$
- Centered transformation
 $y_i^r = y_i \exp\{-\hat{\mu}_{AFT}\}$
- Prior of terminal nodes
 $\mu_{j,l} \sim N(0, \frac{4\sigma_{AFT}^2}{Jk^2})$
- Prior of BART
 $m(A, x) \sim N(0, \frac{4\sigma_{AFT}^2}{k^2})$
- Nonparametric AFT model
 $\log T_i = m(A_i, x_i) + W_i$
- Prior specification
 $P(\text{Var}(W | G, \sigma) \leq \delta_W^2) = 0.5 \rightarrow \kappa = \delta_W^2 / 4, \sigma_\tau^2 = \kappa$
- Posterior inference — Blocked Gibbs sampling

Weibull distribution

$$f(t; \lambda, \alpha) = \alpha \lambda t^{\alpha-1} \exp(-\lambda t^\alpha)$$

Hazard function

$$h(t) = \lambda \alpha t^{\alpha-1}$$

Cumulated hazard function

自由主题