# BART from the ground up    p266-277

Assume:

$$\vec{y} = f(x) + \vec{\epsilon}$$

$$\epsilon_1, \ldots, \epsilon_n \overset{iid}{\sim} N(0, \sigma^2)$$

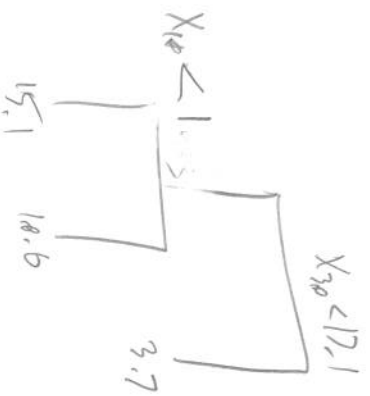$X$ n×p real value or categorical , $\vec{y} \in \mathbb{R}^n$ (for now)    cols

We approximate $f$ by sum of m tree models

$$\hat{\vec{y}} = \sum_{t=1}^{m} g_t(X) + \vec{\epsilon}$$

Each $g_t$ is a tree w/ binary splits h>2 like the CART setup:

$X_{10} < 1...$
$X_{30} < 17.1$
15.1    10.6    3.7

---

Note how the $g_t$'s are identified by structure and split rules, which we call $T_t$, and a collection of rules n with each comp_ztions of X which we denote.

$$M_t = \langle M_{t1}, M_{t2}, \ldots, M_{tb} \rangle$$

where $b_t$ is the number of nodes in the t-th tree

$$\hat{\vec{y}} = \sum g_t(X | T_t, M_t) + \vec{\epsilon}$$

Sum of trees: now majority voting.

Goal: to be able to compute $P(\hat{y} | x^*)$, thereby if course make predictions

$$\hat{y} = E[\hat{Y} | x^*]$$

Same as regression

We employ a Bayesian setup:

$$P(T_1, M_1, \ldots, T_m, M_m, \sigma^2 \mid y, x) \propto \underbrace{P(y \mid T_1, M_1, \ldots, T_m, M_m, \sigma^2, x)}_{\sigma} P(T_1, M_1, \ldots, T_m, M_m, \sigma^2 \mid x)$$

L1b

$$P(T_1, M_1, \ldots, T_m, M_m, \sigma^2 \mid x) \quad \text{prior}$$

Since we impose conditions on $x$ (fixed design), we're going to suppress this notation from now on.

Let's look at the prior form

$$P(T_1, M_1, \ldots, T_m, M_m, \sigma^2) = P(T_1, M_1, \ldots, T_m, M_m) P(\sigma^2) \qquad \text{Assume } \sigma^2 \perp \text{tree structure}$$

$$= P(\sigma^2) \prod_{t=1}^{m} P(T_t, M_t) \qquad \text{Assume trees iid of each other}$$

$$= P(\sigma^2) \prod_{t=1}^{m} P(M_t \mid T_t) P(T_t) \qquad \text{Bayes Rule}$$

$$= \underbrace{P(\sigma^2)}_{III} \underbrace{\prod_{t=1}^{m} P(T_t)}_{I} \underbrace{\prod_{\ell=1}^{b_t} P(M_{t\ell} \mid T_t)}_{II} \qquad \text{Assume each node iid of each other}$$
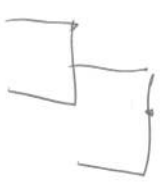
three things to specify

$P(T) \hookrightarrow$ physical structure
$\hookrightarrow$ split rule

Structure: begin with •

$$P(\text{split}) = \alpha(1 + \text{depth})^{-\beta}, \quad \alpha \in (0,1), \quad \beta \in [0,\infty)$$

For each split rule:

Pick $j \in \{1,...,p\}$ uniformly

If $x_j \in \{1,...,k\}^h$ pick a random class $\in 2^{\{1,...,k\}}$

$\Rightarrow x_j \in$ subset

If $x_j \in \mathbb{R}^n$, pick $i \in \{1,...,n\}$ rule is

$x_j < X_{ij}$

---

$$P(\mu_\ell | T) = N(m_n, \sigma_n^2)$$

How do we pick these two hyperparams?
We use a half empirical Bayes idea:

Pick $m_n, \sigma_n^2$ s.t. $[y^{(i)}, y^{(0)}]$ spans 95%
of the possible spread. Assume mean is 4

center:

$E[Y|x] = E[\sum g_t] = n E[g_1] = n E[m_1] = m m_n \qquad mm = \frac{y^{(i)} - y^{(0)}}{2}$

$Var[Y|x] = $

$= m \sigma_n$

Back to Soln 101:

$$\left[ M_n \pm 1.96 \, \sigma_n \sqrt{m} \right] = \left[ y_{(i)}, y_{(b)} \right]$$

Before we solve, why not:

$$y_i := \frac{y - y_{(i)}}{y_{(b)} - y_{(i)}} - \frac{1}{2} \in \left[ -\frac{1}{2}, \frac{1}{2} \right]$$
from $[0,1]$

$$\Rightarrow M_n = 0 \Rightarrow \sigma_n = \sqrt{\frac{\frac{1}{2}}{1.96 \sqrt{m}}} \Rightarrow \sigma_n^2 \approx \frac{1}{15.4 M} \underbrace{\qquad}_{}$$ (K)

$$\Rightarrow P(M+\varepsilon | T) = N(0, \sigma_n^2)$$

$$P(\sigma^2) = Inv G\left(\frac{2}{2}, \frac{2}{2}\right) \leftarrow$$

To pick $(2,7)$ we use
a tilt Emp. Bayes idea.

Explicitly, turn to
Gibbs sampling ...

---

Let
$$M_m = y_{(c)} \overset{\Delta}{=} \frac{y_{(a)} + y_{(c)}}{2}$$

$$y_{(b)} - y_{(c)} = 1.96 \sigma_m \sqrt{m}$$
$$\Rightarrow \sigma_m = \frac{y_{(b)} - y_{(c)}}{1.96 \sqrt{m}}$$
$$m_{mm} - 1.96 \sigma_m \sqrt{m}$$
$$= y_{(b)} - y_{(c)}$$

$$\rightarrow \text{ or MASR for} \atop \text{in genome}$$

Lets say $P(\theta < S_o) = 0.9 = \textcircled{1}$  4ff pan
(extremey conservative)

We fix $\textcircled{2} = 3$ for large spread
and compute $\lambda$ from PDF.

---

$$\frac{10 \qquad 20}{}$$

How big should $m$ be?

— One of my projects: make a fully Bayesian sense
— Use X-validation
— Over one 200 sites $\lambda$ doesn't make
  such a difference...

Now, return to primer
  samples...

If you can't compute
$$P(T, M_1, ... , T_m, M_m, \sigma^2 | y)$$

Explicitly, turn to
Gibbs sampling...

$\nu = 3$

$P(\sigma^2)$
$\nu \sigma^2$
$\sigma^2$

$T_1, M_1 \mid \bar{R}_2, M_2, \dots \to \bar{T}_m, M_m, \sigma^2;$

$\vdots$

$T_m, M_m \mid T_1, M_1, \dots \to T_{m-1}, M_{m-1}, \sigma^2;$

$\sigma^2 \mid T_1, M_1, \dots \to T_m, M_m \checkmark$

How does $T_1, M_1$ depend on other trees?

formula:

$\vec{y} = \vec{g}_1 + \vec{g}_2 t + \dots + \vec{g}_m + \vec{\epsilon}$

$\implies \vec{g}_1 = \vec{y} - \underbrace{(\vec{g}_2 t + \dots + \vec{g}_m)}_{\text{contributing}\atop\text{other trees}} - \vec{\epsilon}$

$\underbrace{\qquad}_{R_1}$

what's left over is $R_1$

what's "remaining", "residual"
response"

$\longleftarrow$ plays the role of the new "y"

How does $\sigma^2$ depend on trees?

$\vec{\epsilon} = \vec{y} - \sum \vec{g}_t$

$\underbrace{\qquad}$

i.v. E for residuals

---

Now we have a backfitting Gibbs sampler:

$T_1, M_1 \mid R_1, \sigma^2$

$\vdots$

$T_m, M_m \mid R_m, \sigma^2$

$\sigma^2 \mid E$

We wish we could...

$T_1 \mid R_1, \sigma^2, M_1 \quad \Longleftarrow$
$M_1 \mid R_1, \sigma^2, T_1 \quad \Longleftarrow$ Gibbs logic

We can marginalize out $M_1$ to find

$P(T_1 \mid R_1, \sigma^2) = \int P(T_1, M_1 \mid R_1, \sigma^2) \, dM_1$

$\propto \int P(R_1 \mid T_1, M_1, \sigma^2) \, P(T_1 \mid M_1, \sigma^2) \, dM_1$

$= P(T_1 \mid \sigma^2) \int P(R_1 \mid T_1, M_1, \sigma^2) \, P(M_1 \mid T_1, \sigma^2) \, dM_1$

$\phantom{=}\underbrace{\quad}_{\text{Assume Normal}} \quad \underbrace{\quad}_{\text{Normal}}$

So it can be done...

So this rng so to is:

$= P(T_1 \mid \sigma^2) P(R_1 \mid T_1, \sigma^2)$

$\underbrace{\qquad}_{\text{Lik.}}$

Lik.

Arrington

① $T_i$  $(b_i = 3)$

$M_{11}$
$\overrightarrow{Y_{11}}$  ← the remaining left over
$(n_1)$

$M_{12}$   $M_{13}$
$\overrightarrow{Y_{12}}$  $\overrightarrow{Y_{13}}$
$(n_{12})$  $(n_{13})$

$s.t$  $n = n_{11} + n_{12} + n_{13}$

$Y_{11k_1}, \dots, Y_{11k_{11}} \mid M_{11} \overset{iid}{\sim} N(M_{11}, \sigma^2)$  Comes from error

Recall  $Y_{11k_1} = M_{11} + \varepsilon_{11k_1}$

$M_{11}, \text{ the prior on } M_{11} \sim N(0, \sigma_m^2)$

So... easy enough!

$Y_{11k_1} \sim N(0, \sigma_m^2) * N(0, \sigma^2) = N(0, \sigma^2 + \sigma_m^2)$

$\text{Cov}[Y_{11k_1}, Y_{11k_2}] = \text{Cov}[M_{11} + \varepsilon_{11k_1}, M_{11} + \varepsilon_{11k_2}] = \sigma_m^2$

$\Rightarrow \overrightarrow{Y_{11}} \sim N_{n_{11}}(0, \underbrace{\sigma_m^2 J_{n_{11}} + \sigma^2 I_{n_{11}}}_{\Sigma_{11}})$  (equicorrelation)

not independent

---

$\frac{1}{(2\pi)^4} \prod_{\ell=1}^{b_i} |\Sigma_{i\ell}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\overrightarrow{Y_{i\ell}}^\top \Sigma_{i\ell}^{-1} \overrightarrow{Y_{i\ell}}\right)$

$\Rightarrow P(T_i \mid R_i, \sigma^2) \sim P(T_i \mid \sigma^2)$.

$\prod_{\ell=1}^{b_i} \frac{1}{(2\pi)^{n_{i\ell}/2} |\Sigma_{i\ell}|^{1/2}} \exp\left(-\frac{1}{2}\overrightarrow{Y_{i\ell}}^\top \Sigma_{i\ell}^{-1}\overrightarrow{Y_{i\ell}}\right)$   $P(R_i \mid T_i, \sigma^2)$

Since we have a $P$ marginal density, use
ex. $M-H$;

$\not\leq \ln(U(0,1)) < \ln \dfrac{P(T_i^* \mid T_i)}{P(T_i \mid T_i^*)} \dfrac{P(T_i^* \mid R_i, \sigma^2)}{P(T_i \mid R_i, \sigma^2)}$

$\dfrac{P(R_i \mid T_i^*, \sigma^2)}{P(R_i \mid T_i, \sigma^2)} \dfrac{P(T_i^*)}{P(T_i)}$

$\Rightarrow$ accept

So how to get $T_i^*$?

Take $T_0$, and also as of it if they:
change
grow   25%
prune  25%
change  50%
swap  10%

| change | prob |
|---|---|
| grow | 25% |
| prune | 25% |
| change | 50% |
| swap | 10% |

Next in our non-parametric – Gibbs replacement:

$$\mu_1, R_1, T_1, \sigma^2 \begin{cases} \mu_{11} | \vec{r}_{11}, \sigma^2 \\ \mu_{12} | \vec{r}_{11}, \sigma^2 \\ \mu_{101} | \vec{r}_{11}, \sigma^2 \end{cases}$$

Sample due to conjugacy:

$$P(\mu_{11} | \vec{r}_{11}, \sigma^2) \propto P(\vec{r}_{11} | \mu_{11}, \sigma^2) P(\mu_{11} | \sigma^2)$$

$$= \left( \prod_{i=1}^{n_{11}} N(r_{11}, \sigma^2) \right) N(0, \sigma^2_n)$$

$$= N\left( \frac{\frac{n}{\sigma^2} \bar{r}_{11}}{\frac{1}{\sigma^2_n} + \frac{n}{\sigma^2}} , \frac{1}{\frac{1}{\sigma^2_n} + \frac{n}{\sigma^2}} \right)$$

$$\Rightarrow P\theta !$$

Gelman et al

Now we sample $\overline{T_2} | R_2 \sigma^2$, $\mu_2 | R_2, \sigma^2, \vec{r}_2, \cdots$

Now we sample due to conjugacy

$$\sigma^2 | E$$

Also sample due to conjugacy

$$P(\sigma^2 | E) \propto P(E | \sigma^2) P(\sigma^2)$$

$$= \left( \prod_{i=1}^{n} N(0, \sigma^2) \right) InvG\left( \frac{u}{2}, \right)$$

$$= InvG\left( \frac{u+n}{2}, \frac{2v + \sum \varrho_i^2}{2} \right)$$

Now we're good

$$P(T_1, \mu_{11}, \ldots, T_{n_1}, \mu_{n_1}, \sigma^2 | y)$$

We draw $N_G$ samples and burn the first B

$\hat{\theta}_{B+1}, \ldots, \hat{\theta}_{N_G - B}$   equal PDF

To get the best guess: avg of PDF

$$\hat{y} = \frac{1}{N_G - B} \sum_{g=B+1}^{N_G} \hat{y}_{(x^u)} \quad \underset{\text{OR}}{=} \quad \hat{y} = \hat{y}_{(50\%)}(x^u)$$

median

$$PPI_{y,95\%} = \left[ \hat{L}_{25\%}(Y^{u}), \hat{L}_{75\%}(Y^{u}) \right]$$

See 5.1 bullet $PB$ density

normal dens $(Y)$, $B$-fold

— pred dep function —

cross validation

p27-270 BART Probit

no rescaling y.

$$P(V=1 \mid x^{*}) = \mathbb{Z} \left( \sum_{b=1}^{\hat{m}} g_{b}(x^{**}) \right)$$

two versions of BART

$\hat{\tau}_{ij} \overset{iid}{\sim} N(\hat{b}_{b\tau})$  explicitly

defaults / cv

$\sigma = 1$

$v=3$, $g = 90\%$, $k=2$, $n=200$

augmentation idea

$$Z_1, \dots, Z_n \overset{iid}{\sim} N(\sum g_{b}, 1)$$

$N_0 = 1000$, $B=200$

s.t.

$Z_i \mid y_i=1 \sim \max \{ N(\sum g_b, 1), 0 \}$

$Z_i \mid y_i=0 \sim \min \{ N(\sum g_b, 1), 0 \}$

Competition

Lasso,

Grad. Boosang,

Random Forest,

Neural Nets

$(v, g) \in \{ (3, 90\%), (3, 99\%), (10, 75\%) \}$

$k \in \{2, 3, 5\}$

$m \in \{50, 200\}$

24 chains

Use $Z_i$'s as $y_i$'s and do

marginalia — within — plus backfitting for

$T_1, M_1, \dots, T_m, M_m$