

Setting:

$$y = \beta_1 x_1 + \dots + \beta_p x_p + \sigma \varepsilon$$

for $\varepsilon \sim N(0,1)$

parameters $\beta \in \mathbb{R}^p$, $\sigma > 0$

n observations

$(x_1, y_1) \dots (x_n, y_n)$

But $p \gg n$

(more parameters than observations)

If you use OLS, it is not well-defined

Assumptions: β is sparse.

In other words, a lot of β 's are zero, and we need to have this assumption to have identifiability.

From Bayesian perspective:

where is the uncertainty?

- ① number of β 's with entries 0?
- ② which β 's are zero?
- ③ the value of β 's with non-zero entry?

Encoding the belief for ③
By introducing

$$\gamma_1, \gamma_2, \dots, \gamma_p \in \{0, 1\}$$

$$\gamma_j \in 1 \Leftrightarrow \beta_j \neq 0 \Leftrightarrow "x_j \text{ is important}"$$

$$\gamma_j \in 0 \Leftrightarrow \beta_j = 0 \Leftrightarrow "x_j \text{ is not important}"$$

\Leftrightarrow means if and only if

$$p_j | r_j \sim \delta_0 (1 - r_j) + r_j \text{Unif}[-T, T]$$

for T large.

that is to say, if $r_j = 0$, then it is just a point mass 0.

if $r_j = 1$, you know p_j not zero, then you impose a super flat prior on it.



for the belief for ②

$$r_1, r_2, \dots, r_p \mid \theta \sim \text{Bernoulli}(\theta)$$

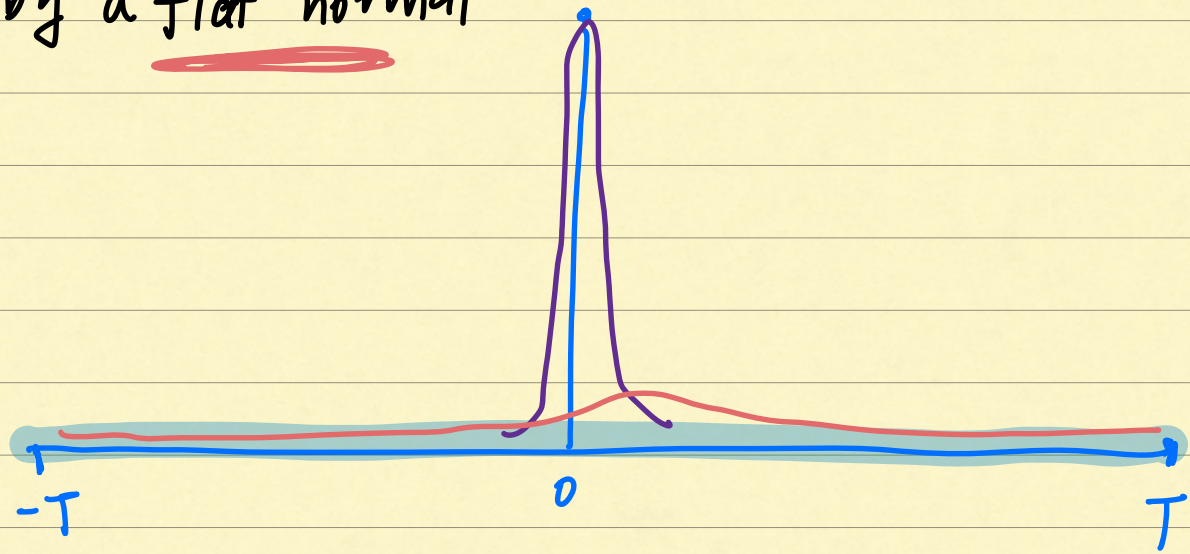
for the belief for ①

$\theta \sim \text{Beta}(a, b)$, and you would choose a and b to give small expectation of θ as you believe B 's are sparse.

However, using MCMC is not feasible
challenge!

as $r_1, r_2, \dots, r_p \mid \theta$ needs updating at the same time and you have too many configurations.

challenge 2: the graph is not differentiable
 continuous relaxation, a modification
 of the graph, replace the point mass
 by a peaky normal, and uniform replaced
 by a flat normal



Your model becomes

$$\theta \sim \text{Beta}(a, b)$$

$$r_i \dots r_p | \theta \sim \text{Bernoulli}(\theta)$$

$$p_j | r_j, \sigma^2 \sim r_j \cdot N(0, v_1 \sigma^2) + (1 - r_j) N(0, v_0 \sigma^2)$$

$$v_1 \gg \gg \gg v_0$$

challenge 1:

Not to know for sure the posterior density for each β , but just find the MAP.

$$\operatorname{argmax}_{\beta} p(\beta, \theta, \sigma^2 | \mathbf{y})$$

$$p(\beta | \sigma^2, \theta) = \frac{1}{\pi} \left[\theta \cdot \text{Glab density} + (1-\theta) \cdot \text{spike density} \right]$$

But the addition is hard to work with because if you take log. doesn't factorize.

$$\log \frac{1}{\pi} \left[\theta \cdot \text{Glab density} + (1-\theta) \cdot \text{spike density} \right]$$
$$= \sum \log \left(\theta \frac{1}{\sqrt{\pi} \sigma} e^{-\beta_j^2 / 2\sigma^2} + (1-\theta) \frac{1}{\sqrt{\pi} \sigma} e^{-\beta_j^2 / 2\sigma^2} \right)$$

↓

Hard to optimize,

so you replace that with its lower bound

$$\mathbb{E}_r, \log p(\beta, \theta, \sigma^2, r | y)$$

you know

$$\log p(\beta, \theta, \sigma^2, r | y) \geq \mathbb{E}_r, \log p(\beta, \theta, \sigma^2, r | y)$$

By Jensen's inequality

thus, for each iteration

↓ optimize the $\mathbb{E}_r, \log p(\beta, \theta, \sigma^2, r | y)$
to get new $\beta, \theta, \sigma^2, r$

↓ take Expectation

$$p(r | \beta, \theta, \sigma^2, y) \propto \prod_{j=1}^p \left(\frac{1}{2\pi\sigma^2} e^{-\frac{r_j^2}{2\sigma^2}} \right) \left(\frac{1}{2\pi\sigma^2} e^{-\frac{\beta_j^2}{2\sigma^2}} \right)^{1-r_j}$$

$$\prod_{j=1}^n \theta^{r_j} (1-\theta)^{1-r_j}$$

thus

$$P(r_j=1 | \cdot) \propto \theta \cdot P_1(B_j)$$

$$P(r_j=0 | \cdot) \propto (1-\theta) \cdot P_0(B_j)$$

define $P_j^* \equiv \frac{\theta P_1(B_j)}{\theta P_1(B_j) + (1-\theta) P_0(B_j)}$

write out $P(\beta, \theta, \sigma^2, r | y)$

$$-\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum (y_i - x_i^T \beta)^2$$

$$+ \sum [r_j \log P_1(B_j) + (1-r_j) \log P_0(B_j)]$$

$$+ \sum [r_j \log \theta + (1-r_j) \log (1-\theta)]$$

$$+ \log p(\theta) + \log p(\sigma^2)$$

and $r \mid \beta, \sigma^2, \theta, y$ would be

$$\begin{aligned} & -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum (y_i - x_i^T \beta)^2 \\ & + \sum [p_j^* \log p_1(\beta_j) + (1-p_j^*) \log p_0(\beta_j)] \\ & + \sum p_j^* \log \theta + (1-p_j^*) \log (1-\theta) \\ & + \log p(\theta) + \log p(\sigma^2) \end{aligned}$$

So update that via EM.

$$r \mid \beta^{(t)}, \sigma^{2(t)}, \theta^{(t)}, y$$

and recall,

$$p_1(\beta_j) = \left(\frac{1}{2\pi\sigma^2 \cdot v_1} \right)^{-1/2} e^{-\beta_j^2 / 2v_1\sigma^2}$$

thus

$$\log p_1(\beta_j) = \frac{-\beta_j^2}{2v_1\sigma^2} - \frac{1}{2} \log \sigma^2$$

So for the optimization step,

LEMVS

the objective function is

$$\sum (y_i - x_i^T \beta)^2 - \frac{1}{2} \sum \beta_j^2 \left[\frac{p_j^*}{v_1} + \frac{1 - p_j^*}{v_0} \right]$$

it looks like a Ridge

But the penalty depends if
you come from spike, small
penalty; if you come from slab,
Big penalty.