

$$C(\beta_1, \beta_2, \dots, \beta_p) = \sum_i^N \left(y_i - \sum_j^p x_{ij} \beta_j \right)^2 + \lambda \sum_j^p |\beta_j|$$

$$\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p = \arg \min_{\beta_1, \beta_2, \dots, \beta_p} C(\beta_1, \beta_2, \dots, \beta_p)$$

this problem is trick because
of the $|\beta_j|$ term



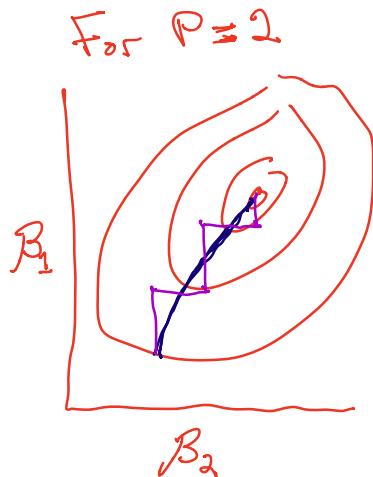
We will go over the derivation of the
coordinate descent algorithm for lasso and
then, at the end, you will write out
pseudocode in groups.

History: (according to Hastie)
Tibshirani's student

- PhD student Wenjiang Fu had main idea in 1997, called it "shooting alg."
- 2002 Tibshirani and Hastie try simulating similar but have bug + answer doesn't work
- 2006 they realised their mistake and published stuff with it. Advisor +, student +

$$C(\beta_1, \beta_2, \dots, \beta_p) = \sum_i^N \left(y_i - \sum_j^p x_{ij} \beta_j \right)^2 + \lambda \sum_j^p |\beta_j|$$

$$\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p = \underset{\beta_1, \beta_2, \dots, \beta_p}{\operatorname{argmin}} C(\beta_1, \beta_2, \dots, \beta_p)$$



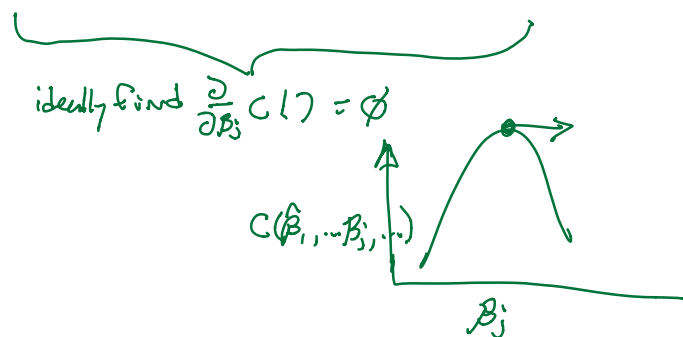
one approach
"coordinate descent"

initialize $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_p$

while not converged:

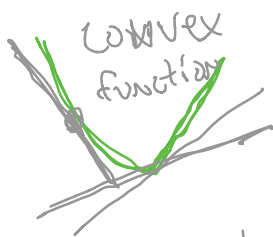
for j in P :

$$\hat{\beta}_j = \underset{\beta_j}{\operatorname{argmin}} C(\hat{\beta}_1, \hat{\beta}_2, \dots, \beta_j, \dots, \hat{\beta}_p)$$

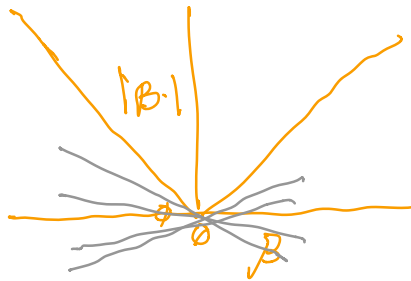


$$C(\hat{\beta}_1, \hat{\beta}_2, \dots, \beta_k, \dots, \hat{\beta}_p) = \sum_i^N \left(y_i - \sum_{j \neq k}^p x_{ij} \hat{\beta}_j - x_{ik} \beta_k \right)^2 + \lambda \sum_{j \neq k}^p |\hat{\beta}_j| + \lambda |\beta_k|$$

$$\frac{\partial C(\cdot)}{\partial \beta_k} = \sum_i^N 2 \left(y_i - \sum_{j \neq k}^p x_{ij} \hat{\beta}_j - x_{ik} \beta_k \right) x_{ik} + \lambda \frac{\partial}{\partial \beta_k} |\beta_k| \leftarrow \text{oops!!}$$



gradients/subgradients
provide a lower bound
to convex functions



$$\text{at } \beta_k < 0, \frac{\partial |\beta_k|}{\partial \beta_k} = -1$$

$$\text{at } \beta_k > 0, \frac{\partial |\beta_k|}{\partial \beta_k} = 1$$

$$\text{at } \beta_k = 0, \text{ Subgradient so } -1 \leq \frac{\partial |\beta_k|}{\partial \beta_k} \leq 1$$

$$\frac{\partial C(\cdot)}{\partial \beta_k} = \sum_i^N 2(y_i - \sum_{j \neq k}^p \hat{\beta}_j x_{ij} - x_{ik} \beta_k) x_{ik} + \begin{cases} -\lambda, & \beta_k < 0 \\ [-\lambda, \lambda], & \beta_k = 0 \\ \lambda, & \beta_k > 0 \end{cases}$$

$$\underbrace{-\sum_i^N 2(y_i - \sum_{j \neq k}^p \hat{\beta}_j x_{ij}) x_{ik}}_{g_k} + \hat{\beta}_k \sum_i^N x_{ik}^2 + \begin{cases} -\lambda, & \beta_k < 0 \\ [-\lambda, \lambda], & \beta_k = 0 \\ \lambda, & \beta_k > 0 \end{cases}$$

$$\underline{\underline{\beta_k < 0}}$$

$$0 = -g_k + \hat{\beta}_k \sum_i^N x_{ik}^2 - \lambda$$

$$2 \hat{\beta}_k \sum_i^N x_{ik}^2 = g_k + \lambda$$

$$\hat{\beta}_k = (g_k + \lambda) / 2 \sum_i^N x_{ik}^2$$

and this is for when $\beta_k < 0$, so...

$$(g_k + \lambda) / \sum_i^N x_{ik}^2 < 0$$

or when $g_k + \lambda < 0$

$$\text{or } g_k < -\lambda$$

$$\hat{\beta}_k = \begin{cases} (g_k + \lambda) / \sum_i^N x_{ik}^2, & g_k < -\lambda \\ \emptyset, & -\lambda \leq g_k \leq \lambda \\ (g_k - \lambda) / \sum_i^N x_{ik}^2, & g_k > \lambda \end{cases}$$

$\beta_k > 0$ { look above and see it's
the same idea so

$$\hat{\beta}_k = (g_k - \lambda) / \sum_i^N x_{ik}^2$$

and this occurs at $\beta_k > 0 \Rightarrow g_k > \lambda$

$$\underline{\underline{\beta_k = \emptyset \dots}}$$

$$\text{so } \hat{\beta}_k = \emptyset$$

and this is for when

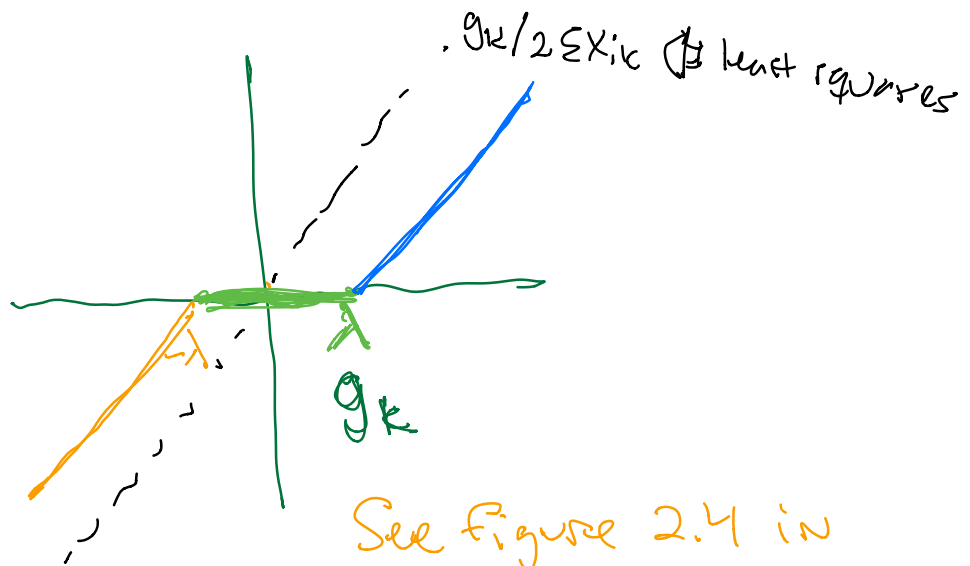
$$\emptyset \in \left[(g_k + \lambda) / 2 \sum_i^N x_{ik}, (g_k - \lambda) / 2 \sum_i^N x_{ik} \right]$$

which occurs when

$$-\lambda \leq g_k \text{ and } g_k \leq \lambda$$

$$\hat{\beta}_k = \begin{cases} (g_k + \lambda) / 2 \sum_i^N x_{ik}, & g_k < -\lambda \\ \emptyset, & -\lambda \leq g_k \leq \lambda \\ (g_k - \lambda) / 2 \sum_i^N x_{ik}, & g_k > \lambda \end{cases}$$

$$\text{where } g_k = \sum_i^N 2 \left(y_i - \sum_{j \neq k}^P \hat{\beta}_j \right) x_{ik}$$



See figure 2.4 in book.

"Soft thresholding function"