

1 Non-smooth hyper-parameter learning

2 Clarice Poon

3 March 30, 2022

Consider

$$\begin{aligned} \min_y h(y), \quad \text{where } h(y) = f(y, x(y)), \\ \text{s.t. } x(y) \in \operatorname{argmin}_x g(y, x). \end{aligned}$$

In general, we have

$$\partial_x g(y, x(y)) = 0,$$

and by the implicit function theorem, provided that $\partial_x^2 g(y, x(y))$ is invertible, $y \mapsto x(y)$ is differentiable with

$$x'(y) = -\partial_{xx} g(y, x(y))^{-1} \partial_{yx} g(y, x(y)).$$

One can then evaluate the gradient of h as

$$\nabla h(y) = \partial_y f(y, x(y)) - \partial_x f(y, x(y))^\top \partial_{xx} g(y, x(y))^{-1} \partial_{yx} g(y, x(y)). \quad (1)$$

Quasi-Newton? One approach is to find p_k such that

$$\partial_{yx} g(y, x(y)) \approx \partial_{yy} g(y, x(y)) p_k$$

4 and evaluate $\nabla h(y) \approx \partial_y f(y, x(y)) - \partial_x f(y, x(y))^\top p_k$.

Can we do a quasi-Newton approach? We know that

$$\partial_x g(y, x(y)) = 0$$

Let $x_k \triangleq x(y_k)$. Then

$$\begin{aligned} \partial_x g(y_k, x(y_{k-1})) &= \partial_x g(y_k, x_{k-1}) - \partial_x g(y_k, x_k) \\ &= \partial_{xx} g(y_k, x_k)(x_{k-1} - x_k) + o(\|x_{k-1} - x_k\|) \end{aligned}$$

Suppose we find B_k such that it minimises

$$\min_B \|\partial_x g(y_k, x(y_{k-1})) - B(x_{k-1} - x_k)\|$$

and treat it as an approximation to $\partial_{xx} g(y_k, x(y_k))^{-1}$. The idea is to compute

$$\nabla h(y_k) \approx \partial_y f(y_k, x(y_k)) - \partial_x f(y_k, x(y_k))^\top B_k \partial_{yx} g(y_k, x(y_k)).$$

One possible update of B_k is as $\tau_k \text{Id} + u_k u_k^\top$. Define $s_k = x_{k-1} - x_k$ and $z_k = \partial_y g(y_k, x(y_{k-1}))$. We want to find diagonal + rank-1 matrix B to minimise

$$\min_B \|Bz_k - s_k\|$$

5 i) Define $\tau_k = \langle s_k, z_k \rangle / \|z_k\|^2$ and project onto $[\tau_{\min}, \tau_{\max}]$. Note that before
6 projection, $\tau_k = \text{argmin}_\tau \|\tau z_k - s_k\|$.

7 ii) Let $B_0 \triangleq \gamma \tau_k \text{Id}$ where $\gamma \in (0, 1)$.

iii) If $\langle s_k - B_0 z_k, z_k \rangle \leq 10^{-8} \|z_k\| \|s - B_0 z_k\|$ the $U_k = 0$. Else:

$$U_k = \frac{(s_k - B_0 z_k)(s_k - B_0 z_k)^\top}{\langle s_k - B_0 z_k, z_k \rangle}.$$

8 iv) Let $B_k = B_0 + U_k$.

Note that for step iii) the choice of U_k is precisely finding $U_k = uu^\top$ such that

$$B_0 z_k + u \langle u, z_k \rangle - s_k = 0.$$

• **[ToDo:**

10 **1. If we repeatedly updated B_k with rank-1 matrices, show that B_k**
11 **converges to $\partial_{xx} g(y_*, x(y_*))^{-1}$.**

12 **2. Suppose that $f = g$ and at iteration k , we have an approximate**
13 **solution $\hat{x}_k \approx x(y_k)$. Possible ways of computing $\nabla h(y_k)$ include**

14 **i) $p_1 = \partial_y f(y_k, \hat{x}_k)$**

15 **ii) $p_2 = \partial_y f(y_k, \hat{x}_k) - \partial_x f(y_k, \hat{x}_k)^\top B_k \partial_{yx} f(y_k, \hat{x}_k)$**

16 **iii) $p_3 = \partial_y f(y_k, \hat{x}_k) + \partial_x f(y_k, \hat{x}_k)^\top \partial_y \hat{x}_k$ where we obtain $\partial_y \hat{x}_k$ via**
17 **autodiff.**

18 **The first and 3rd option have been analysed recently (Ablin**
19 **et al). For the second approach, can we bound the difference**
20 **between taking approximation B_k and the true Hessian?**

21 **]**

A differentiable approach to nonsmooth bilevel programming One example is where g is nonsmooth is when y correspond to a hyperparameter λ and x is the lasso regression coefficients:

$$f(\lambda, \beta) \triangleq \|A_{\text{test}} \beta - y\|^2 \quad \text{and} \quad g(\lambda, \beta) \triangleq \frac{1}{2} \|A_{\text{train}} \beta - b\|_2^2 + \lambda \|\beta\|_1.$$

The difficulty is in this case is that since g is non-smooth, the formula (1) cannot be used. One alternative is to consider instead

$$f(\lambda, (u, v)) \triangleq \|A_{\text{test}} uv - y\|^2 \quad \text{and} \quad g(\lambda, (u, v)) \triangleq \|A_{\text{train}} uv - b\|_2^2 + \lambda \|u\|^2 / 2 + \lambda \|v\|^2 / 2.$$

The advantage with this approach is that g is a smooth function and one can show that the Hessian of g is invertible when $\beta \triangleq u(\lambda) \odot v(\lambda)$ is a nondegenerate solution, that is,

$$\max_{i \notin \text{Supp}(\beta)} |A_{\text{train}}^\top (A_{\text{train}}\beta - b)|_i < 1.$$

22 Things to do

23 • Check properties of the Hessian of g .

24 • Acceleration using support pruning.

The square root lasso The square root lasso is

$$\min_{\beta \in \mathbb{R}^n} \|X\beta - y\|_2 + \lambda \|\beta\|_1.$$

One interesting aspect of this is that when $y = X\beta_0 + w$, the minimiser β satisfies

$$\|\beta - \beta_0\| \lesssim \|w\|$$

25 for some constant λ . This is remarkable since, for the Lasso, to achieve this kind
26 of error bound, one would require that $\lambda \sim \|w\|$ and some knowledge of the noise
27 level is required.

One remark is that the square root lasso is equivalent to

$$\min_{\sigma > 0} \min_{\beta} \frac{1}{2\sigma} \|X\beta - y\|_2^2 + \frac{\sigma}{2} + \lambda \|\beta\|_1,$$

and we can therefore write this in the bilevel formulation with

$$f(\sigma, \beta) = \frac{1}{2\sigma} \|X\beta - y\|_2^2 + \frac{\sigma}{2} + \lambda \|\beta\|_1$$

and

$$g(\sigma, \beta) = \frac{1}{2\sigma} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1.$$

One question I have is what happens when we consider

$$\begin{aligned} f(\sigma, \beta) &= \frac{1}{2\sigma} \|A_{\text{test}}\beta - y\|_2^2 + \frac{\varepsilon\sigma}{2} + \iota_{\sigma > 0} \\ g(\sigma, \beta) &= \frac{1}{2\sigma} \|A_{\text{train}}\beta - y\|_2^2 + \|\beta\|_1 \end{aligned}$$

This is precisely the hyperparameter learning framework but with added regularisation on the parameter σ . Note that the outer problem can be written as an unconstrained smooth problem as follows: Let $z = \sqrt{\sigma}$ and $v = A_{\text{test}}\beta - y$, then

$$\begin{aligned} \min_{z \in \mathbb{R}} f(z, \beta(z^2)), \quad \text{where} \quad f(z, \beta) &\triangleq \min_{zv = A_{\text{test}}\beta - y} \frac{1}{2} \|v\|^2 + \frac{\varepsilon}{2} z^2, \\ \beta(z^2) &\triangleq \operatorname{argmin}_{\beta} g(z^2, \beta) \end{aligned}$$

Notice that the minimisation problem in f is convex wrt v , so by taking the convex dual,

$$f(z, \beta) = \max_{\alpha \in \mathbb{R}^m} -\frac{\varepsilon}{2} z^2 \|\alpha\|^2 + \frac{\varepsilon}{2} z^2 + \langle \alpha, -A_{\text{test}}\beta + y \rangle$$

The maximiser α is unique (as the problem is strongly concave) and

$$\partial_{\beta} f = -A_{\text{test}}^{\top} \alpha \quad \text{and} \quad \partial_z f = z \|\alpha\|^2.$$

28 Numerically, we can certainly handle this, the question is whether this kind of
29 regularisation is interesting in practice.

Let $F(\sigma) = f(\sigma, \beta(\sigma))$. Let's look at the optimality conditions

$$\begin{aligned} \partial_{\sigma} f &= \frac{-1}{\sigma^2} \|A_{\text{test}}\beta - y_{\text{test}}\|^2 + \frac{\varepsilon}{2} \\ \partial_{\beta} f &= \frac{1}{\sigma} A_{\text{test}}^{\top} (A_{\text{test}}\beta - y_{\text{test}}) \end{aligned}$$

Also, $\beta = \beta(\sigma)$ satisfies

$$A_{\text{train}}^{\top} A_{\text{train}} \beta = A_{\text{train}}^{\top} y - \sigma \text{sign}(\beta)$$

In general, $\sigma \mapsto \beta(\sigma)$ is differentiable almost everywhere with gradient

$$\beta'(\sigma) = -(A_{\text{train}}^{\top} A_{\text{train}})_{J,J}^{-1} \text{sign}(\beta).$$

where $J = \text{Supp}(\beta)$. So, when $F'(\sigma) = 0$, we have

$$\frac{-1}{\sigma^2} \|A_{\text{test}}\beta - y_{\text{test}}\|^2 - \frac{1}{\sigma} \langle (A_{\text{test}}\beta - y_{\text{test}}), A_{\text{test}}(A_{\text{train}}^{\top} A_{\text{train}})_{J,J}^{-1} \text{sign}(\beta) \rangle + \frac{\varepsilon}{2} = 0$$

which implies $\lambda = 1/\sigma$ satisfies, for $C \triangleq \langle (A_{\text{test}}\beta - y_{\text{test}}), A_{\text{test}}(A_{\text{train}}^{\top} A_{\text{train}})_{J,J}^{-1} \text{sign}(\beta) \rangle$,

$$\lambda = \frac{-C + \sqrt{C^2 + 2\varepsilon \|A_{\text{test}}\beta - y_{\text{test}}\|^2}}{2 \|A_{\text{test}}\beta - y_{\text{test}}\|^2}$$

NB: For the standard problem where $f(\sigma, \beta) = \frac{1}{2} \|A_{\text{test}}\beta - y\|^2$, then

$$F'(\sigma) = -\langle (A_{\text{test}}\beta - y_{\text{test}}), A_{\text{test}}(A_{\text{train}}^{\top} A_{\text{train}})_{J,J}^{-1} \text{sign}(\beta) \rangle.$$

Suppose $A_{\text{test}} = A_{\text{train}}$, then this says that

$$F'(\sigma) = \sigma \langle (A_{\text{train}}^{\top} A_{\text{train}})_{J,J}^{-1} \text{sign}(\beta), \text{sign}(\beta) \rangle > 0$$

30 which means that we optimise to $\sigma = 0$ as expected.