

Today:

Revisiting an older method for minimization that regained popularity.

Randomized Coordinate Descent:

$$w_{k+1} = w_k - \gamma [\nabla f(w)]_{i_k}$$

where  $i_k$  is picked:

- uniformly at random
- random with importance weights
- greedily.

Q: How does RCD compare with GD?

We will assume the following coordinate-wise smoothness:

Property: A function  $f(w)$  is  $\beta_i$ -coordinate-wise smooth if:

$$\forall a, x, i \quad f(x + ae_i) \leq f(x) + a[\nabla f(x)]_i + \frac{\beta_i a^2}{2}$$

$$\text{Let } \beta = \max_i \beta_i.$$

Further let us assume that  $f(w)$  is  $\mu$ -PL:

$$\forall w, \quad \|\nabla f(w)\|^2 \geq \mu/2 (f(w) - f^*)$$

Then, we can show the following guarantees for RCD with uniform sampling

Thm: If  $f$  is  $\beta$ -coordinate-wise smooth and  $\mu$ -PL then RCD with stepsize  $\frac{1}{L}$ :

$$w_{k+1} = w_k - \frac{1}{\beta} [\nabla f(w_k)]_{i_k}$$

$$i_k \sim \text{unif}(1, \dots, n)$$

obtains the following rate:

$$\mathbb{E} f(w_k) - f^* \leq \left(1 - \frac{\mu}{d\beta}\right)^k (f(x_0) - f^*)$$

Proof:

By plugging

$$w_{k+1} = w_k - \frac{1}{\beta} [\nabla f(w_k)]_{i_k}$$

In the coordinate-wise smoothness property, we get

$$f(w_{k+1}) \leq f(w_k) - \frac{1}{2\beta} |\nabla f(w_k) \cdot i_k|^2$$

By taking expectation <sup>(w.r.t  $i_k$ )</sup>, we have:

$$\begin{aligned} E f(w_{k+1}) &\leq f(w_k) - \frac{1}{2\beta} E |\nabla f(w_k) \cdot i_k|^2 \\ &= f(w_k) - \frac{1}{2\beta} \sum_{i=1}^d \frac{1}{d} |\nabla f(w_k) \cdot i|^2 \\ &= f(w_k) - \frac{1}{2\beta d} \|\nabla f(w_k)\|^2 \end{aligned}$$

We will now apply the PL condition to get:

$$E f(w_{k+1}) \leq f(w_k) - \frac{\mu}{2\beta d} (f(w_k) - f^*)$$

$\Rightarrow$

$$\begin{aligned} E f(\omega_{k+1}) - f^* &\leq f(\omega_k) - f^* - \frac{\mu}{2\beta_0 L} (f(\omega_k) - f^*) \\ &= \left(1 - \frac{\mu}{\beta_0 L}\right) (f(\omega_k) - f^*) \end{aligned}$$

Applying expectations and recursively expanding yields the result.



Unfortunately, as it stands RCD with uniform sampling is not clearly better than UCD which achieves a rate

$$f(\omega_k) - f^* \leq \left(1 - \frac{\mu}{L}\right)^k (f(x_0) - f^*)$$

which requires  $O(d)$  less iterations than RCD for the same accuracy  $\epsilon$ .

To improve RCD we can employ importance sampling:

$$\Pr(i_k = i) = \frac{\beta_i}{\sum_{j=1}^d \beta_j}$$

That is, each coordinate is sampled proportionally to its "effect" on  $f(w)$ .

For this weighted sampling we get:

$$w_{k+1} = w_k - \frac{1}{\beta_i} [\nabla f(w_k)]_{i_k}$$

and

$$\begin{aligned} E f(w_{k+1}) &\leq f(w_k) - \frac{1}{\beta_i} E | [\nabla f(w_k)]_{i_k} | \\ &= f(w_k) - \frac{1}{\beta_i} \frac{\beta_i}{\sum \beta_i} \frac{1}{d} \|\nabla f(w_k)\|^2 \end{aligned}$$

$$\leq f(\omega_k) - \frac{1}{\bar{\beta} \cdot d} \|\nabla f(\omega_k)\|^2$$

The above imply the following Theorem:

Thm. RCD with importance sampling according to

$$\Pr(i_k = i) = \frac{\beta_i}{\sum_j \beta_j}$$

yields the following rate:

$$\mathbb{E} f(\omega_k) - f^* \leq \left(1 - \frac{h}{d\bar{\beta}}\right)^k (f(\omega_0) - f^*)$$

Remark:  $\bar{\beta}$  can be significantly smaller than  $\beta$ .

Example:

If we are aiming for  $\epsilon$ -accuracy then:

$$T_{\epsilon}^{\text{importance}} = O\left(\frac{d \sum_i p_i}{n \cdot \mu} \log\left(\frac{f_0 - f^*}{\epsilon}\right)\right)$$

$$T_{\epsilon}^{\text{uniform}} = O\left(\frac{d p_{\max}}{\mu} \log\left(\frac{f_0 - f^*}{\epsilon}\right)\right)$$

Hence importance sampling can be extremely helpful.

Main Question:

- How easy are  $p_i$  to compute?
- Similar sampling for SGD?



