

Coordinate Optimization

- i) Choose coordinate j
- ii) Update value at coordinate j

$$x^{t+1} = x^t - \alpha_t \nabla_{j_t} f(x^t) \cdot e_{j_t}$$

Faster if updates are d times faster

$$\nabla_{j_t} f(x^t) \cdot e_{j_t} \in O(nd)$$

2.1.1

compute $\nabla_{j_t} S(x^t) \in O(n)$

- need $X \cdot w$ to compute $\nabla_{j_t} S(x^t)$

$$X \cdot w \in O(n \cdot d) \quad \text{to } \text{~~compute } \nabla_{j_t} S(x^t)~~$$

- can update Xw based on how w changes at each iteration $\wedge \in O(n)$

2.1.2

$$L_c = \max_{j_t} \{L_{j_t}\} \quad \text{since } L \text{ is an upper bound}$$

Sample Discrete(p)

$$\sum_i p_i = 1 \quad \text{e.g. } p = [0.2 \ 0.3 \ 0.5] \quad \text{returns } \begin{cases} 1 & \text{with prob. } 20\% \\ 2 & \text{" " } 30\% \\ 3 & \text{" " } 50\% \end{cases}$$

Proximal Gradient

For solving problems of the form:

$$\arg \min_{x \in \mathbb{R}^d} \underbrace{f(x)}_{\text{smooth}} + \underbrace{r(x)}_{\text{non-smooth}}$$

We use updates

$$x^{t+1} = \arg \min_{y \in \mathbb{R}^d} \left\{ \frac{1}{2} \|y - (x^t - \alpha_t \nabla f(x^t))\|^2 + \alpha_t r(y) \right\}$$

$$= \text{prox}_{\alpha_t r} [x^t - \alpha_t \nabla f(x^t)]$$

2.2.1

$$\begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix}_{n \times d} \begin{bmatrix} w_1 & \dots & w_k \\ | & & | \\ | & & | \end{bmatrix}_{d \times k} = \begin{bmatrix} \\ \\ \end{bmatrix}_{n \times k}$$

$X \qquad \qquad \qquad W \qquad \qquad \qquad Y$

$$f(w) = \sum_i -\log \frac{e^{w_1^T x_i}}{\sum_k e^{w_k^T x_i}}$$

$$g(w) = \dots$$

↙ Add L2-regularization to all and gradient (Use Frobenius norm)

2.2.2

Instead of findMin, use proxGradL1(S, w, λ, ...)

↳ returns $\arg \min_{w \in \mathbb{R}^d} f(w) + \lambda \|w\|_1$

- Just pass in the differentiable part of the function (and gradient for that part) - ~~the proxGradL1~~ ^{proxGradL1} does the regularization by applying the proximal operator for L1-regularization - which is an element-wise soft threshold

$$x_{\downarrow} = \frac{x_{\uparrow}}{|x_{\uparrow}|} \max \{0, |x_{\uparrow}| - \alpha_t\}$$

2.2.3

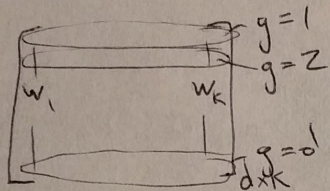
Mod. by the proximal operator to be the proximal operator for group L1 regularization

- which is a group-wise soft threshold

$$x_g = \frac{x_g}{\|x_g\|_2} \max \{ 0, \|x_g\|_2 - \lambda \}$$

- Add a parameter to identify which group each x_g is in.

- Declare the value in `softmaxClassifier.GL1` and pass it to `proxGrad.GL1`



0 rows in W correspond to unused features (columns of x)

▷ x_g = row of W

Stochastic Gradient

$$x^{t+1} = x^t - \frac{\alpha_t}{n} \sum_i \nabla^i f(x^t)$$

gradient wrt training example i

$$\text{Stochastic: } x^{t+1} = x^t - \alpha_t \nabla^{i_t} f(x^t)$$

↳ gradient wrt random training example i_t

⇒ $O(n)$ faster per iteration

2.3.3

AdaGrad - want different step sizes for different dimensions

$$x^{t+1} = x^t - \alpha_t D_t \nabla f(x^t)$$

Diagonal Matrix

$$D(i, i) = \frac{1}{\sqrt{\delta + \sum_{k=0}^t (\nabla_i f(x^k))^2}}$$

2.3.4

S.A.G. - keep memory of gradients wrt each training example
- update 1 gradient per iteration

$$G = \begin{bmatrix} g^1 \\ g^2 \\ \vdots \\ g^n \end{bmatrix}$$

At each iteration

$i \leftarrow$ random training index

$$g^i \leftarrow \nabla^i f(x^t)$$

$$x^{t+1} \leftarrow x^t - \frac{\alpha_t}{n} \sum_i g^i$$

\hookrightarrow Can initialize G with all zeros and still use n at every iteration