$$F(w) = \frac{1}{N} \sum_{i=1}^{N} f_i(w) + h(w) \to \min_{w}$$

$f_i \in C^2$ вып., $h \in C$ вып.

$$\begin{cases} (1 - \frac{\mu}{L})^k c = \varepsilon \\ k \log(1 - \frac{\mu}{L}) + \log c = \varepsilon \\ k \approx O\left(\frac{L}{\mu} \log \frac{1}{\varepsilon}\right) \end{cases}$$

| Метод | Число итераций |
|---|---|
| prox - GD | $O\left(N \frac{L}{\mu} \log \frac{1}{\varepsilon}\right)$ |
| acc. prox- GD | $O\left(N \sqrt{\frac{L}{\mu}} \log \frac{1}{\varepsilon}\right)$ |
| SAG / SVRG | $O\left(\left(N + \frac{L}{\mu}\right) \log \frac{1}{\varepsilon}\right)$ |
| ? | $O\left(\left(N + \sqrt{\frac{L}{\mu}}\right) \log \frac{1}{\varepsilon}\right)$ |

<u>SDCA</u> стохастический двойственный покоординатный подъём

$$F(w) = \frac{1}{N} \sum_{i=1}^{N} \psi_i(a_i^T w) + \frac{\lambda}{2} \|w\|^2 \to \min_{w} \quad, \quad \psi_i : \mathbb{R} \to \mathbb{R}$$

линейная регрессия: $\psi_i(z) = \frac{1}{2}(y_i - z)^2$, $a_i = x_i$

логистическая регрессия: $\psi_i(z) = \log(1 + \exp(-z))$, $a_i = y_i x_i$

SVM: $\psi_i(z) = \max(0, 1 - z)$, $a_i = y_i x_i$

$$P(w, z) = \begin{cases} \frac{1}{N} \sum_{i=1}^{N} \psi_i(z_i) + \frac{\lambda}{2} \|w\|^2 \to \min_{z, w} \\ Aw = z \end{cases}$$

$$L(w, z, \mu) = \frac{1}{N} \sum_{i=1}^{N} \psi_i(z_i) + \frac{\lambda}{2} \|w\|^2 + \frac{1}{N} \mu^T(z - Aw)$$

$$q(\mu) = \inf_{w,z} L(w,z,\mu)$$

$$\nabla_w L = \lambda w - \frac{1}{N} A^T \mu = 0, \quad w = \frac{1}{\lambda N} A^T \mu = \frac{1}{\lambda N} \sum_{i=1}^{N} \mu_i \cdot a_i$$

отн. $z_i$ : $\varphi_i(z_i) + \mu_i z_i \to \min_{z_i}$

$$\varphi_i^*(u) = \max_z (uz - \varphi_i(z))$$

$$q(\mu) = \frac{1}{N} \sum_{i=1}^{N} \min_{z_i}(\varphi_i(z_i) + \mu_i z_i) + \frac{\lambda}{2} \left\| \frac{1}{\lambda N} A^T \mu \right\|^2 -$$

$$-\lambda \left( \frac{1}{\lambda N} A^T \mu \right)^T \frac{1}{\lambda N} A^T \mu = \frac{1}{N} \sum_{i=1}^{N} \left( -\max_{z_i}(-\mu_i z_i - \varphi_i(z_i)) \right) -$$

$$- \frac{\lambda}{2} \left\| \frac{1}{\lambda N} A^T \mu \right\|^2 = \frac{1}{N} \sum_{i=1}^{N} (-\varphi_i^*(-\mu_i)) - \frac{\lambda}{2} \left\| \frac{1}{\lambda N} A^T \mu \right\|^2 \to \max_\mu$$


Схема SDCA

$$\mu^{(0)} ; \quad w^{(0)} = \frac{1}{\lambda N} A^T \mu^{(0)}$$

для $k = 0, 1, 2, \ldots$

    $i_k \sim \text{Unif}(1 \ldots N)$

    Найти $\delta \mu_{i_k}$ :

$$\frac{1}{N}\left(-\varphi_{i_k}^*(-\mu_{i_k}^{(k)} - \delta\mu_{i_k})\right) - \frac{\lambda}{2}\left\| w^{(k)} + \frac{1}{\lambda N} a_{i_k}^T \delta\mu_{i_k} \right\|^2 \to \max_{\delta\mu_{i_k}}$$

$$\mu^{(k+1)} = \mu^{(k)} + \delta\mu_{i_k} \cdot e_{i_k}$$

$$w^{(k+1)} = w^{(k)} + \frac{1}{\lambda N} a_{i_k}^T \delta\mu_{i_k}$$


Утв. $\varphi_i \in C_L^{1,1} \Rightarrow$ для $\mathbb{E} F(w_k) - F_{opt} \leq \varepsilon$

в SDCA достаточно сделать $O\left((N + \frac{L}{\mu}) \log \frac{1}{\varepsilon}\right)$ итер.

Нет необходимости подбирать оптимальную длину шага $\alpha_k$. Все расчёты на итерации произ-водятся аналитически.

$$F(w) = \frac{1}{N} \sum_{i=1}^{N} \varphi_i(a_i^T w) + h(w) \to \min_w$$

$$\begin{cases} \frac{1}{N} \sum_{i=1}^{N} \varphi_i(z_i) + h(w) \to \min_w \\ Aw = z \end{cases}$$

$$L(z, w, \mu) = \frac{1}{N} \sum_{i=1}^{N} \varphi_i(z_i) + h(w) + \frac{1}{N} \mu^T (z - Aw)$$

$$q(\mu) = \frac{1}{N} \sum_{i=1}^{N} \left( -\varphi_i^*(-\mu_i) \right) + \min_w \left( h(w) + \frac{1}{N} \mu^T A w \right) =$$

$$= \frac{1}{N} \sum_{i=1}^{N} \left( -\varphi_i^*(-\mu_i^{\circledast}) \right) - \underbrace{\max_w \left( -w^T \frac{A^T \mu}{N} - h(w) \right)}_{h^*\left( -\frac{1}{N} A^T \mu \right)} \to \max_w$$

$$h(w) = \|w\|_2$$

$$h^*(u) = \max_w \left( w^T u - \|w\|_2 \right) = \left\{ \max_{w_i} \left( w_i u_i - |w_i| \right) \right\}$$

$$wu - |w| \to \max$$
$$w \geq 0 \quad , \quad (u - 1)w$$
$$w \leq 0 \quad , \quad (u + 1)w$$



$$\max_w \{ wu - |w| \} = 0$$
при $|u| \leq 1$, $\infty$ иначе ,

$$h^*(u) = [\, |u| \leq 1 \,] = \delta_{|u| \leq 1}(u)$$

$$F_k(w) = F(w) + \frac{\varkappa}{2} \|w - y_{k-1}\|^2 \to \min_w$$

$$L_{F_k} = L + \varkappa \qquad \text{cond } F = \frac{L}{\mu}$$
$$\mu_{F_k} = \mu + \varkappa \qquad \text{cond } F_k = \frac{L + \varkappa}{\mu + \varkappa}$$

Схема   Acc. SDCA

$$R^2 = \max_i \|x_i\|^2$$

$$\mathcal{x} = \frac{R^2 L}{N} - \mu \;;\; \eta = \sqrt{\frac{\mu}{\mu + \mathcal{x}}} \;;\; \beta = \frac{1 - \eta}{1 + \eta}$$

$$y_1 = w_1 = 0, \quad \alpha_1 = 0, \quad z_1 = (1 + \eta^{-2})(P(0) - q(0))$$

для $k = 2, 3, \ldots$

$$F_k(w) = F(w) + \frac{\mathcal{x}}{2}\|w - y_{k-1}\|^2$$

$$(w_k, \alpha_k) = prox\text{-}SDCA\left(F_k, \alpha_{k-1}, \underbrace{\frac{\eta}{2(1 + \eta^{-2})}}_{= \varepsilon_k}, z_{k-1}\right)$$

$$y_k = w_k + \beta(w_k - w_{k-1})$$

$$z_k = (1 - \eta/2) z_{k-1}$$

$$\widetilde{O}\left(\left(N + \sqrt{\frac{N R^2 L}{\mu}}\right) \log(1/\varepsilon)\right)$$
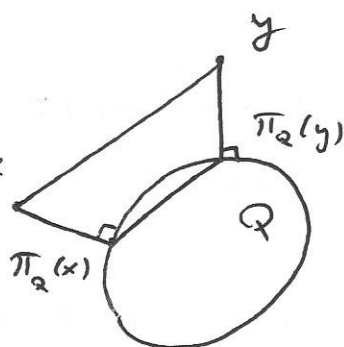
число итераций

семинар

$$\frac{1}{N} \sum_{i=1}^{N} f_i(x)$$

$$f(x) = \mathbb{E}\, F(x, z) = \int_{\Omega} F(x, w)\, dP(w)$$

$$\min_{x \in Q} f(x), \quad Q - \text{вып.}, \quad f: Q \to \mathbb{R}$$

SVM : $\dfrac{1}{N} \sum\limits_{i=1}^{N} \max\{0; 1 - \langle a_i, x\rangle\}$

robust regression : $\dfrac{1}{N} \sum\limits_{i=1}^{N} |\langle a_i, x\rangle - b_i|$



$$\|\pi_Q(x) - \pi_Q(y)\| \leq \|x - y\|$$

$\min\limits_{x} f(x) \qquad \mathbb{E}(g_k | x_k) \in \partial f(x_k)$

$\{x_k\} \subseteq Q \qquad f(x) - f(x_k) \geq \langle G_k, x_* - x_k\rangle \quad \forall x \in Q$

$$G_k \in \partial f(x_k)$$

$$x_{k+1} = \pi_Q(x_k - \alpha_k g_k)$$

[4]

$$\|x_{n+1} - x^*\|^2 \le \|x_n - x^* - \alpha_n g_n\|^2 = \|x_n - x^*\|^2 - 2\alpha_n\langle g_n, x_n - x^*\rangle +$$
$$+ \alpha_n^2\|g_n\|^2 \ ; \quad \alpha_n\langle g_n, x_n - x^*\rangle \le \frac{1}{2}\|x_n - x^*\|^2 - \frac{1}{2}\|x_{n+1} - x^*\|^2 + \frac{\alpha_n^2}{2}\|g_n\|^2$$

телескопирующаяся сумма: $\sum_{n=1}^{T}(a_n - a_{n+1}) = a_1 - a_{T+1}$

$$\sum_{n=1}^{T}\alpha_n\underbrace{\langle g_n, x_n - x^*\rangle}_{\ge f(x_n) - f^* \ ?} \le \frac{1}{2}\underbrace{\|x_1 - x^*\|^2}_{R^2} + \sum_{n=1}^{T}\frac{\alpha_n^2}{2}\|g_n\|^2$$

$$\sum_{n=1}^{T}\mathbb{E}(\alpha_n\langle g_n, x_n - x^*\rangle) \le \frac{R^2}{2} + \ {}^{\#}\sum_{n=1}^{T}\mathbb{E}\frac{\alpha_n^2}{2}\|g_n\|^2$$

необходимо: $\alpha_n$ — детерминированные, тогда

$$\sum_{n=1}^{T}\alpha_n\,\mathbb{E}\langle g_n, x_n - x^*\rangle \le \frac{R^2}{2} + \sum_{n=1}^{T}\frac{\alpha_n^2}{2}\mathbb{E}\|g_n\|^2$$

$$\mathbb{E}\langle g_n, x_n - x^*\rangle = \mathbb{E}\,\mathbb{E}(\langle g_n, x_n - x^*\rangle\,|\,x_n) \stackrel{\#}{=} \mathbb{E}(\langle\mathbb{E}g_n, x_n - x^*\rangle\,|\,x_n) \ge$$
$$\stackrel{\#}{\ge} \mathbb{E}(\mathbb{E}(f(x_n) - f^*\,|\,x_n)) = \mathbb{E}(f(x_n) - f^*)$$

$$\frac{\sum_{n=1}^{T}\alpha_n\,\mathbb{E}(f(x_k) - f^*)}{\sum_{n=1}^{T}\alpha_n} \ \ge\ \mathbb{E}f\left(\underbrace{\sum_{n=1}^{T}\alpha_n x_n}_{\bar x}\right)^{-f^*} = \mathbb{E}f(\bar x) - f^*$$

йенсен

$$\{\mathbb{E}(f(x_n) - f^*\,|\,x_n) \ge \langle\mathbb{E}g_n, x_n - x^*\rangle\,|\,x_n)\}$$

$$\frac{\frac{R^2}{2}}{\sum_{n=1}^{T}\alpha_n} + \frac{\sum_{n=1}^{T}\alpha_n^2\mathbb{E}\|g_n\|^2}{\sum_{n=1}^{T}\alpha_n}$$

$$\alpha_n = \alpha: \quad \frac{R^2}{\alpha T} + \alpha\frac{\sum_{n=1}^{T}\mathbb{E}\|g_n\|^2}{T} \ \to\ \min_\alpha \quad (*)$$

$$\frac{R^2}{2} = \alpha^2\sum_{n=1}^{T}\mathbb{E}\|g_n\|^2, \quad \alpha = \frac{R}{\left(\sum_{n=1}^{T}\mathbb{E}\|g_n\|^2\right)^{\frac{1}{2}}}$$

$$\frac{R}{\sqrt{T}}\left(\frac{1}{T}\sum_{n=1}^{T}\mathbb{E}\|g_n\|^2\right)^{\frac{1}{2}}, \quad \alpha = \frac{R}{\sqrt{T}\,M}$$

5

① стох. субгр. метод

$$\alpha_k = \alpha = \frac{R}{M\sqrt{T}} \quad , \quad (*) = \frac{MR}{\sqrt{T}}$$

пример: $f(x) = \frac{1}{N} \sum_{i=1}^{N} |\langle a_i, x \rangle - b_i|$

$i_k$ равномерно из $\{1 \ldots N\}$

$g_k = \text{sgn}(\langle a_{i_k}, x_k \rangle - b_{i_k}) a_{i_k}$

$\mathbb{E}(g_k | x_k) \in \partial f(x_k)$

$\|g_k\|^2 \leq \|a_{i_k}\|^2, \quad \mathbb{E}\|g_k\|^2 \leq \frac{1}{N} \sum_{i=1}^{N} \|a_i\|^2 = M$

$$(*) = \frac{\left( \frac{1}{N} \sum_{i=1}^{N} \|a_i\|^2 \right)^{\frac{1}{2}} R}{\sqrt{T}} \leq \max_{1 \leq i \leq N} \|a_i\|^2 R / \sqrt{T}$$

② стох. субгр. метод с адаптивными длинами шагов

$\alpha_k$ — случайные, зависят от $(g_1 \ldots g_k)$

$$\langle g_k, x_k - x^* \rangle \leq \frac{1}{2\alpha_k} \|x_k - x^*\|^2 - \frac{1}{2\alpha_k} \|x_{k+1} - x^*\|^2 + \frac{\alpha_k}{2} \|g_k\|^2$$

$$\sum_{k=1}^{T} \frac{1}{2\alpha_k} \|x_k - x^*\|^2 - \sum_{k=1}^{T} \frac{1}{2\alpha_k} \|x_{k+1} - x^*\|^2 =$$

$$= \sum_{k=1}^{T} \frac{1}{2\alpha_k} \|x_k - x^*\|^2 - \sum_{k=2}^{T+1} \frac{1}{2\alpha_{k+1}} \|x_k - x^*\|^2 =$$

$$= \frac{1}{2\alpha_1} \|x_1 - x^*\|^2 - \frac{1}{2\alpha_T} \|x_{T+1} - x^*\|^2 + \sum_{k=2}^{T} \left( \frac{1}{2\alpha_k} - \frac{1}{2\alpha_{k-1}} \right) \underbrace{\|x_k - x^*\|^2}_{\leq R^2} \quad (≣)$$

пусто мн-во $Q$-ограничено: $\|x_k - x^*\| \leq R$

и $\{\alpha_k\}$ монотонно убывают

$$(≣) \frac{R^2}{2\alpha_1} + \frac{R^2}{2}\left( \frac{1}{\alpha_T} - \frac{1}{\alpha_1} \right) = \frac{R^2}{2\alpha_T}$$

$$\sum_{k=1}^{T} \langle g_k, x_k - x^* \rangle \leq \frac{R^2}{2\alpha_T} + \sum_{k=1}^{T} \frac{\alpha_k}{2} \|g_k\|^2$$

$$\overline{x}_T = \frac{1}{T} \sum_{u=1}^{T} x_u$$

$$\mathbb{E}\, f(\overline{x}_T) - f^* \leq \mathbb{E}\left( \frac{R^2}{2T\alpha_T} + \frac{1}{2T} \sum_{u=1}^{T} \alpha_u \|g_u\|^2 \right) \quad (\star)$$

$$\alpha_u \equiv \alpha: \quad \frac{R^2}{\alpha} + \alpha \sum_{u=1}^{t} \|g_u\|^2 \longrightarrow \min_{\alpha}$$

$$\alpha_u = \frac{R}{\left( \sum_{u=1}^{T} \|g_u\|^2 \right)^{\frac{1}{2}}}$$

Выбираем $\quad \alpha_u = \dfrac{R}{\left( \sum_{s=1}^{u-1} \|g_u\|^2 \right)^{\frac{1}{2}}}$

$$(\star) = \frac{R}{T} \left( \sum_{s=1}^{T-1} \|g_s\|^2 \right)^{\frac{1}{2}} + \frac{R}{T} \sum_{u=1}^{T} \frac{\|g_u\|^2}{\left( \sum_{s=1}^{u-1} \|g_s\|^2 \right)^{\frac{1}{2}}}$$

нер-во: $\quad \displaystyle\sum_{u=1}^{n} \frac{a_i}{\left( \sum_{j=1}^{k} a_j \right)^{\frac{1}{2}}} \leq 2 \left( \sum_{i=1}^{n} a_i \right)^{\frac{1}{2}}$

$$\mathbb{E}(\star) = \frac{R}{T} \mathbb{E}\left( \frac{1}{T} \sum_{u=1}^{T} \|g_u\|^2 \right)^{\frac{1}{2}} \leq \frac{R}{\sqrt{T}} \left( \frac{1}{T} \sum_{u=1}^{T} \underbrace{\mathbb{E}\|g_u\|^2}_{\leq M^2} \right)^{\frac{1}{2}} \leq \frac{MR}{\sqrt{T}}$$

③ Ada Grad

$$x_{u+1} = \pi_Q^{B_u}\left( x_u - B_u^{-1} g_u \right) = \operatorname*{argmin}_{x \in Q} \left\{ f(x_u) + \langle g_u, x - x_u \rangle + \right.$$

$$+ \frac{1}{2} \langle B_u(x - x_u), x - x_u \rangle \right\} \qquad // \text{раньше: } B_u = \frac{1}{\alpha_u} I$$

теперь: $B_u = \operatorname{diag}\{b_u\}$

$$\|x_{u+1} - x^*\|_{B_u}^2 = \|x_u - x^*\|_{B_u} - \dots$$

$$\sum_{u=1}^{T} \|x_u - x^*\|_{B_u}^2 - \sum_{u=1}^{T} \|x_{u+1} - x^*\|_{B_u}$$

$$\langle D s, s \rangle = \sum_j d_j s_j^2 \leq \|s\|_\infty^2\, t_2(D)$$

$$\langle \underbrace{(B_k - B_{u-1})}_{D}\underbrace{(x_u - x^*)}_{S}\,,\,\underbrace{x_u - x^*}_{S}\rangle$$

$$\bar{x}_T = \frac{1}{T}\sum_{u=1}^{T} x_u$$

$$\mathbb{E}\, f(\bar{x}_T) - f^* \le \frac{R_\infty^2\,\ell_2(B_T)}{2T} + \frac{1}{2T}\sum_{u=1}^{T}\langle B_u^{-1} g_u, g_u\rangle \;\underline{(*)}$$

$$B_u = \operatorname{diag}\Big\{\sum_{s=1}^{u-1} g_s g_s^T\Big\}^{\frac{1}{2}}$$

$$(*) = \frac{R_\infty}{\sqrt{T}}\sum_{i=1}^{n}\mathbb{E}\Big(\frac{1}{T}\sum_{u=1}^{T} g_{u i}^2\Big)$$

$$\|S\|_\infty \le \|S\|_2 \quad,\quad R_{\infty 2} = \sqrt{n}\,R_{\infty\infty} \quad,\quad R_2 \le \sqrt{n}\,R_\infty$$

## Ada Grad vs Ada Step Size

$$\frac{1}{n}\sum_{j=1}^{n}\sqrt{a_j} \;\le\; \frac{1}{\sqrt{n}}\Big(\sum_{j=1}^{n} a_j\Big)^{\frac{1}{2}} \quad,\quad \frac{1}{n}\sum_{j=1}^{n}\sqrt{a_j} \le \Big(\frac{1}{n}\sum_{j=1}^{n} a_j\Big)^{\frac{1}{2}}$$

Вогнутость корня

$$\sum_{j=1}^{n}\sqrt{a_j} \;\le\; \sqrt{n}\Big(\sum_{j=1}^{n} a_j\Big)^{\frac{1}{2}}$$

ada grad $\le$ ada step size $\le$ sgd

$$R_\infty = \|x - x_*\|_\infty = \max_j\{|x - x_*|_j\}$$

$$R_2 = \|x - x_*\|_2 = \Big(\sum_j (x - x_*)_j^2\Big)^{\frac{1}{2}}$$