

today: • CRF
• variance reduced SGD

$$\psi_i(\tilde{y}) \triangleq \phi(x^{(i)}, y^{(i)}) - \phi(x^{(i)}, \tilde{y})$$

CRF: log-loss $s(x, y; w) = -\log p(y|x; w)$ $p(y|x; w) \propto \exp(s(x, y; w))$

suppose a MAF $s(x, y; w) = \sum_{c \in \mathcal{C}} s_c(x, y; w) = \sum_c \langle w, \phi_c(x, y) \rangle$

optimization objective:

primal $\max_{\tilde{y}} \psi_i(\tilde{y}) - w^T \phi_i(\tilde{y})$ dual

$$\Delta_i = [\psi_i(\tilde{y})]_{\tilde{y} \in \mathcal{Y}_i}$$

SVM struct: $\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_i H_i(w)$

$\max_{\alpha_i \in \Delta_i, s_i} -\frac{\lambda \|w(\alpha)\|^2}{2} + \frac{1}{n} \sum_i \alpha_i^T \alpha_i$

CRF: $\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_i -\log p(y^{(i)}|x^{(i)}; w)$
 $\log(\sum_{\tilde{y}} \exp(-w^T \psi_i(\tilde{y})))$

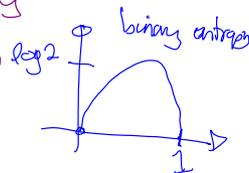
$\max_{\alpha_i \in \Delta_i, s_i} -\frac{\lambda \|w(\alpha)\|^2}{2} + \frac{1}{n} \sum_i H_i(\alpha_i)$

$\triangleq -\sum_{\tilde{y}} \alpha_i(\tilde{y}) \log \alpha_i(\tilde{y})$

KKT: $w(\alpha) = \frac{1}{\lambda n} \sum_i \sum_{\tilde{y} \in \mathcal{Y}_i} \alpha_i(\tilde{y}) \psi_i(\tilde{y})$
 $= \frac{1}{\lambda n} \sum_i \sum_{c \in \mathcal{C}} \sum_{\tilde{y} \in \mathcal{Y}_i} \alpha_i(\tilde{y}) \phi_{i,c}(\tilde{y})$

of optimality

$\alpha_i^*(\tilde{y}) = p(\tilde{y}|x^{(i)}; w^*)$



$\nabla_w \text{CRF-primal}(w) = \lambda w + \frac{1}{n} \sum_i \nabla_w (-\log p(y^{(i)}|x^{(i)}; w))$
 $\log Z_i(w) - w^T \phi_i(y^{(i)})$

$\alpha_i^* \in \text{interior of } \Delta_i$

unlike sparse solution in structured SVM

$\mathbb{E}_{\tilde{y}|x^{(i)}, w} [\psi_i(\tilde{y})] - \psi_i(y^{(i)})$

$\mathbb{E}_{\tilde{y}|x^{(i)}, w} [-\psi_i(\tilde{y})]$ to compute this;

$$\nabla_w = 0 \Rightarrow w^* = \frac{1}{n} \sum_i \frac{\mathbb{E}_{\tilde{y}|x^{(i)}, w^*} [\psi_i(\tilde{y})]}{\sum_{\tilde{y}} p(\tilde{y}|x^{(i)}, w^*) \psi_i(\tilde{y})}$$

$$\begin{aligned} \psi_i(\tilde{y}) &= \sum_c \psi_{i,c}(\tilde{y}_c) \\ \mathbb{E}[\psi_i(\tilde{y})] &= \sum_c \mathbb{E}[\psi_{i,c}(\tilde{y}_c)] \\ &\stackrel{C}{\approx} \sum_{\tilde{y}_c} p(\tilde{y}_c|x^{(i)}, w) \psi_{i,c}(\tilde{y}_c) \end{aligned}$$

need to do marginal inference
 \rightarrow need marginalization oracle

[vs a maximization oracle for SVM struct]

Optimization for CRF:

primal objective is smooth and strongly convex [vs non-smooth for SVM struct]

- for a while, batch L-BFGS was method of choice [batch \Rightarrow slow for large n]
- [Collins & al, JMLR 2008] is online exponentiated gradient

block-coordinate method on dual; exponentiated gradient step on block dual function
 \downarrow

$$\alpha_i(\tilde{y})^{(t+1)} \propto \alpha_i(\tilde{y})^{(t)} \exp(-\eta_t \nabla_{\alpha_i} D(\alpha^{(t)}))$$

\rightarrow get linear convergence rate with cheap $O(1)$ updates (like SGD) [vs. $O(n)$]

low hanging fruit: SDCA \rightarrow have some convergence properties but with cheap line-search and AF AIK, have not been tried

Variance reduced SGD

minimize $f(w) = \frac{1}{n} \sum_i f_i(w)$ where f is μ -strongly convex
 L -smooth (i.e. ∇f is L -Lipschitz)

batch gradient method $w_{t+1} = w_t - \gamma \nabla f(w_t)$ $\approx \exp(-\rho t)$

L Cauchy 1847

$$\frac{1}{n} \sum_i \nabla f_i(w_t)$$

$$\gamma = \frac{1}{L}; \quad f(w_t) - f^* \leq (1-\rho)^t (f(w_0) - f^*)$$

$$\rho \approx \frac{\mu}{L} = \frac{1}{K} \quad K \leq \frac{L}{\mu} \quad \text{condition \#}$$

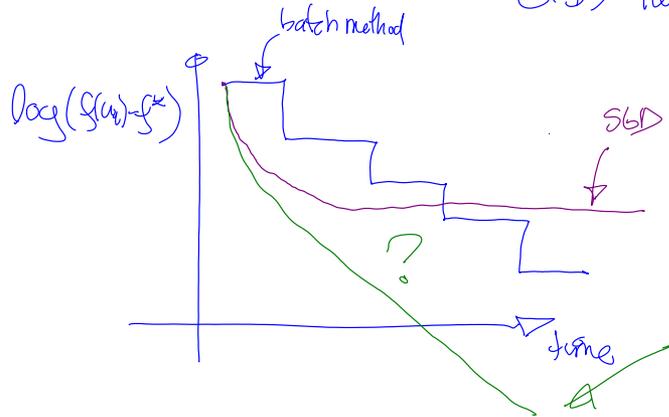
stochastic gradient method
[Robbins & Monro 1951]
incremental gradient method

$$w_{t+1} = w_t - \gamma_t \nabla f_{i_t}(w_t)$$

i_t ^{with} $\{1, \dots, n\}$

$\gamma_t = \text{const}$ \rightarrow linear rate up to level of radius γ
 $\gamma_t \sim \frac{1}{\sqrt{t}}$ $\rightarrow \tilde{O}(\frac{1}{\sqrt{t}})$ rate sublinear

$O(1)$ iteration



variance reduced SGD methods

SAG (stochastic average gradient)

[Le Roux, Schmidt & Bach 2012]

standard batch gradient: $w_{t+1} = w_t - \gamma \frac{1}{n} \sum_i \nabla f_i(w_t)$

SAG is $\left[\begin{array}{l} \text{pick } i_t \text{ ; update } g_{i_t}^{(t+1)} = \nabla f_{i_t}(w_t) \text{ ; } g_j^{(t+1)} = g_j^{(t)} \text{ for } j \neq i_t \\ \text{and then } w_{t+1} = w_t - \gamma \sum_i g_i^{(t+1)} \end{array} \right.$

store these "stale" gradients

$O(1)$ iteration (but $O(n)$ storage)

big surprise \rightarrow this converges linearly and fast?

incremental randomized gradient method (FRG) (Bach et al. 2017) (Liu et al. 2017) (Liu et al. 2018) (Liu et al. 2019)

increment aggregated gradient method (IAG) [Blatt et al 2007] where you cycle deterministically through i_t

Convergence rate: fhm: with $\delta_i = \frac{1}{16L}$ where $L = \max_i (\text{Lipschitz constant of } \nabla f_i)$

$$\mathbb{E}[f(w_n)] - f^* \leq \left(1 - \min\left\{\frac{\mu}{16L}, \frac{1}{8n}\right\}\right)^n C_0$$

ie. $\rho_{SAG} = \min\left\{\frac{1}{16k_{SAG}}, \frac{1}{8n}\right\}$ compare with $\rho_{grad} \approx \frac{1}{k_{grad}}$

example: $n = 700k$, $L = 0.25$, $\mu = \frac{1}{n}$ ($\Rightarrow k = \frac{n}{4}$)

Rate comparison

- Assume that $N = 700000$, $L = 0.25$, $\mu = 1/N$:
 - Gradient method has rate $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.99998$.
 - Accelerated gradient method has rate $(1 - \sqrt{\frac{\mu}{L}}) = 0.99761$.
 - SAG (N iterations) has rate $(1 - \min\left\{\frac{\mu}{16L}, \frac{1}{8N}\right\})^N = 0.88250$.
 - Fastest possible first-order method: $\left(\frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}\right)^2 = 0.99048$.

Pointers

- Online exponentiated gradient for CRF paper:
 - Collins, M., Globerson, A., Koo, T., Carreras, X., and Bartlett, P. L. Exponentiated gradient algorithms for conditional random fields and max-margin Markov networks. [JMLR, 9:1775-1822, 2008](#) (include a comparison with L-BFGS and also gives the dual of the CRF objective)
- stochastic average gradient (SAG):
 - original NIPS 2012 paper:
 - N. Le Roux, M. Schmidt, F. Bach
A Stochastic Gradient Method with an Exponential Convergence Rate for Finite Training Sets. [NIPS 2012](#) | [code](#)
 - massive journal version:
 - M. Schmidt, N. Le Roux, F. Bach.

Minimizing Finite Sums with the Stochastic Average Gradient

[Mathematical Programming](#), 162:83-162, 2017 | [arxiv](#)

- SAGA paper -- unbiased version of SAG (with simpler proof) as well as describe the related methods of SDCA and SVRG (will be covered next class) [see references therein for SDCA and SVRG...]
 - A. Defazio, F. Bach and S. Lacoste-Julien
SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives
[NIPS 2014](#)
 - for an even more bare bone proof, see:
T. Hofmann, A. Lucchi, S. Lacoste-Julien, and Brian McWilliams
Variance Reduced Stochastic Gradient Descent with Neighbors,
[NIPS 2015](#)