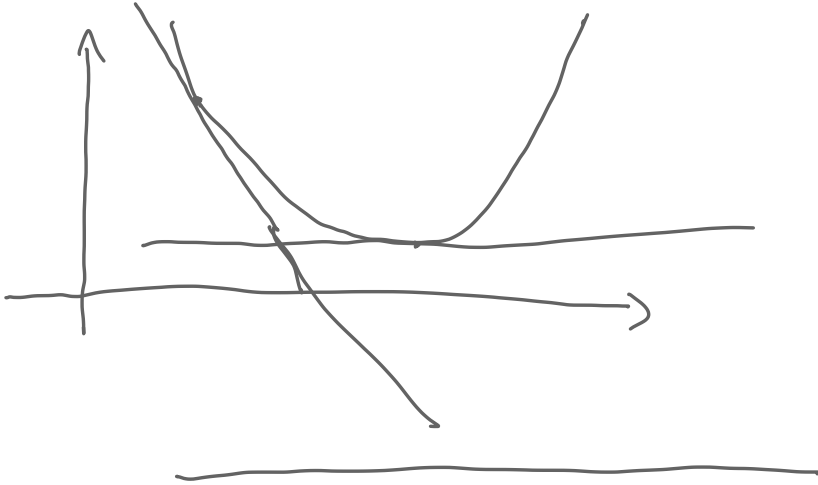


OPTIMAL LECTURE 5, 16/2/2021

- Project
- Exam



$$x_{k+1} = x_k - \eta H_F^{-1}(x_k) g_k$$

- CLASSIC GRADIENT METHODS
 - NOISE VARIANCE REDUCTION METHODS
- OPTIMIZATION FOR NEURAL NETWORKS

NOISE VARIANCE REDUCTION METHODS

$M \quad M_G \rightarrow$ NOISE OF GRADIENT

$$g_t = \frac{1}{n_m} \sum_{i=1}^{n_m} \underline{\nabla f(x_t, \xi_i)}$$

$$x_i \sim \text{i.i.d.} \quad \text{Var}(x_i) = \sigma^2$$

$$\text{Var}\left(\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{\sigma^2}{n}$$

$$\text{Var}\left(\nabla f(x_t, \xi)\right) = \sigma^2$$

$$\text{Var}(g_t) = \frac{\sigma^2}{n_m}$$

$$\mathbb{E}[F(w_{k+1}) - F^*] \leq \frac{\alpha LM}{2c\mu} + (1 - \alpha c\mu)^k \times$$

$$\left[F(w_1) - F^* - \frac{\alpha L \sigma^2}{2c\mu} \right]$$

$$\alpha \leq \frac{\mu}{L M_G}$$

$$M_G = \mu_G + M_V$$

$$\mathbb{E}[F(w_{k+1}) - F^*] \leq \frac{\alpha LM}{2c\mu} + (1 - \alpha c\mu)^k \times$$

$$\left[F(w_1) - F^* - \frac{\alpha LM}{2c\mu} \right]$$

$$\alpha \leq \frac{\mu}{L \mu_G}$$

$$\mu_G = \mu_G + \mu_V$$

Mini Batch

$$g_k = \frac{1}{n_m} \sum \nabla f(x_k, \xi_i)$$

$$\mu, \mu_0$$

$$\rightarrow \mathbb{E}[g_k]^T \nabla F(x_k) \geq \mu \|\nabla F(x_k)\|^2$$

$$\rightarrow \|\mathbb{E}[g_k]\|^2 \leq \mu_0 \|\nabla F(x_k)\|^2$$

$$\mathbb{E}[g_k] = \nabla F(w_k) \Rightarrow \nabla F(x_k)^T \nabla F(w_k) = \|\nabla F(w_k)\|^2 \geq$$

$$1) = \frac{1}{|S|} \sum_{x_i \in S} \nabla f(w_k, x_i) \quad \mu \|\nabla F(w_k)\|^2 \Rightarrow \mu \leq 1$$

$$2) = \frac{1}{n_m} \sum_{i=1}^{n_m} \mathbb{E}[\nabla f(w_k, \xi_i)] = \frac{1}{n_m} n_m \nabla F(w_k) = \nabla F(w_k)$$

$$\mathbb{E}[F(w_{k+1}) - F^*] \leq \frac{\alpha LM}{2c\mu} + (1 - \alpha c\mu)^k \times$$

$$\left[F(w_1) - F^* - \frac{\alpha L \sigma^2}{2c\mu} \right]$$

$$\alpha \leq \frac{\mu}{L \sigma_G}$$

$$\sigma_G = \mu_G + \sigma_V$$

Mini Batch

$$g_k = \frac{1}{n_m} \sum \nabla f(x_k, \xi_k)$$

μ, μ_0

$$\mu = \mu_0 = 1$$

$$V(g_k) \leq \sigma + M_V \| \nabla F(w_k) \|^2$$

σ, σ_V functions of n_m

$$\sigma(n_m=1) = \sigma'$$

$$\sigma_V(n_m=1) = \sigma'_V$$

$$\sigma(n_m) = \frac{\sigma'}{n_m}$$

$$\sigma_V(n_m) = \frac{\sigma'_V}{n_m}$$

$$\mathbb{E}[F(w_{k+1}) - F^*] \leq \frac{\alpha L M^2}{2c n_m} + (1 - \alpha c)^k x$$

$$\left[F(w_1) - F^* - \frac{\alpha L M^2}{2c n_m} \right]$$

$$\alpha \leq \frac{\mu}{L \gamma_G}$$

$$\gamma_G = 1 + \gamma_V$$

RB vs SG ($n_m = 1$)

In terms of Iterations (K)
 the larger n_m , the better
 (the smaller the error)

What about time?

the time of 1 iteration is proportional
 to n_m

In the time RB does K iterations
 the SG does $n_m \times K$ iterations

MB

$$\mathbb{E}[F(w_{k+1}) - F^*] \leq \frac{\alpha_{SG} L M^1}{2c \cancel{\eta_m}} + (1 - \alpha_{SG}^{\cancel{\eta_m}})^K \times$$

$$\left[F(w_1) - F^* - \frac{\alpha_{SG} L M^1}{2c \cancel{\eta_m}} \right]$$

SG

$$\mathbb{E}[F(w_{n_m k+1}) - F^*] \leq \frac{\alpha_{SG} L M^1}{2c} + (1 - \alpha_{SG}^{\cancel{\eta_m}})^{n_m K} \times$$

$$\left[F(w_1) - F^* - \frac{\alpha_{SG} L M^1}{2c} \right]$$

$$\alpha \leq \frac{1}{L M_G} = \frac{1}{L(1 + \sigma_V)} \approx \frac{1}{L \sigma_V} =$$

$$= \frac{n_m}{L \sigma_V^c}$$

$$\Rightarrow \alpha_{SG}$$

$$I \text{ can take } \alpha_{MB} = n_m \alpha_{SG}$$

HB

$$E[F(w_{k+1}) - F^*] \leq \frac{\alpha_{SG} LM^1}{2c} + \boxed{\left(1 - n_m \alpha_{SG} C\right)^K x}$$

$$\left[F(w_1) - F_x - \frac{\alpha_{SG} LM^1}{2c}\right]$$

SG

$$E[F(w_{n_m k+1}) - F^*] \leq \frac{\alpha_{SG} LM^1}{2c} + \boxed{\left(1 - \alpha_{SG} C\right)^{n_m K} x}$$

$$\left[F(w_1) - F_x - \frac{\alpha_{SG} LM^1}{2c}\right]$$

HB

$$\left(1 - \underline{n_m} \alpha_{SG} C\right)^K$$

\leq

SG

$$\left(1 - \alpha_{SG} C\right)^{\underline{n_m K}}$$

$$n_m \alpha_{SG} C \ll 1$$

$$(1-x)^n \simeq 1 - nx$$

$$\simeq 1 - n_m \alpha_{SG} C K = 1 - \alpha_{SG} C n_m K$$

$$(1-x)^n = \sum_{i=0}^n \binom{n}{i} (-x)^{n-i}$$

$$\alpha \leq \frac{1}{L(1+\eta_V)} \approx \frac{1}{L\eta_V}$$

$$\eta_V = \frac{\sigma'_V}{n_m}$$

$$\alpha_{\text{MB}} \leq \frac{1}{L(1 + \frac{\sigma'_V}{n_m})} = \frac{n_m}{L(n_m + \sigma'_V)} \leq \frac{n_m}{L(1 + \sigma'_V)} = n_m \alpha_{\text{SG}}$$

$$\sigma'_V \gg 1$$

$$\frac{\sigma'_V}{n_m} \gg 1$$

$$\alpha_{\text{MB}} \leq n_m \alpha_{\text{SG}}$$

1st element : we have been
optimistic for MB

2) NOT TRUE Variance decreases as $\frac{1}{n_m}$

$$n_m = |S| \Rightarrow \text{no variance} \\ \sigma, \sigma_v = 0$$

Variance of $\frac{1}{n} \sum_i X_i$

without resample

$$\frac{\sigma^2}{n_m} \left(1 - \frac{n_m}{|S|} \right)$$

TOO PESSIMISTIC FOR THE OLS
(OLS HAS SOME ADVANTAGES
WE IGNORED)

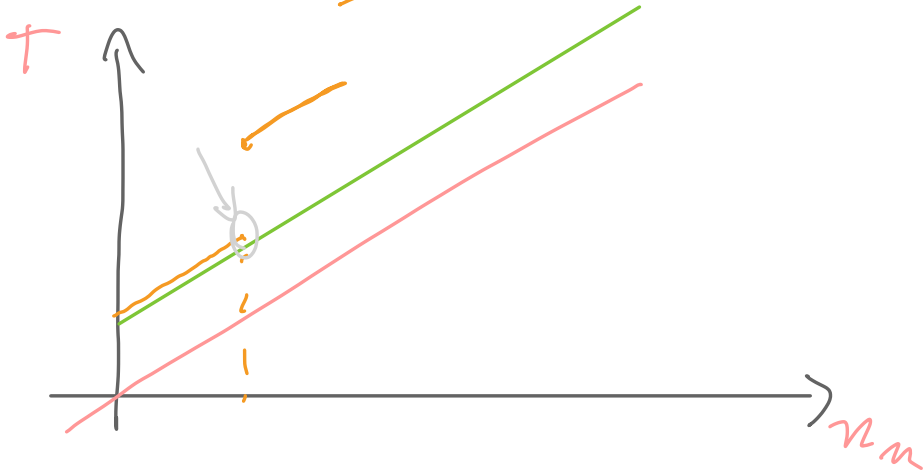
3) Probably a better computation model is

$$t \text{ of 1 iteration} = C + \alpha n_m$$

In our previous analysis
we considered $C=0$

$$\text{or } \alpha \gg C$$

IT GIVES LESS ADVANTAGE
TO SG



NOISE VARIANCE REDUCTION METHODS

• Basic idea

start with small batch size and progressively increase it

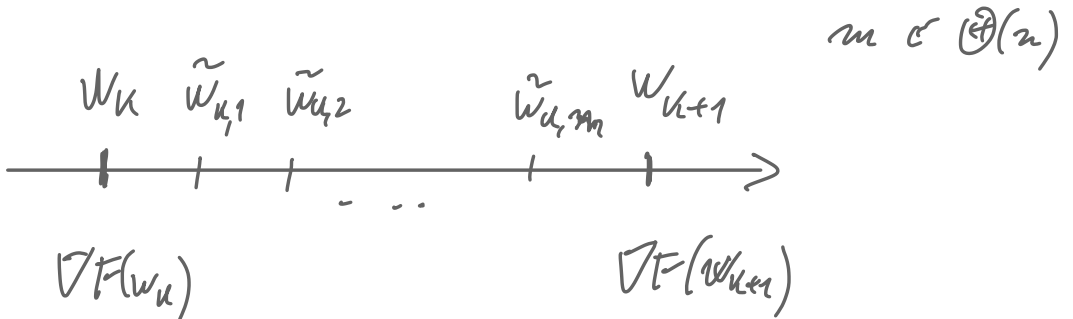
$$n_m(k) = \tau^k n_0 \quad \tau > 1$$

YOU CAN ACHIEVE THE SAME CONVERGENCE RATE (VS TIME) OF FB

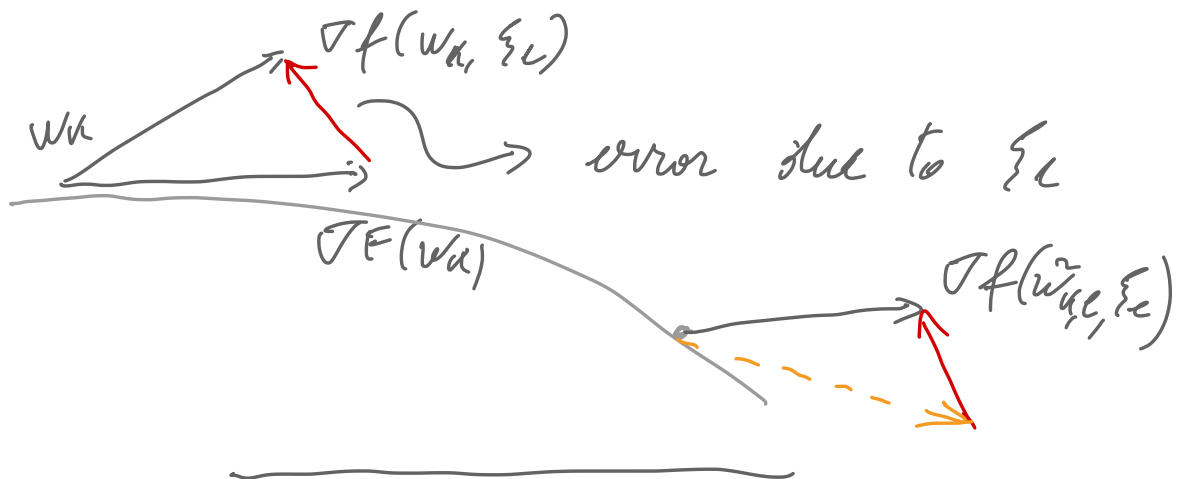
• GRADIENT AGGREGATION

SVRG

SAGA



$$\tilde{w}_{k,l+1} = \tilde{w}_{k,l} - \eta \left(\overbrace{\nabla f(\tilde{w}_{k,l}, \xi_l)}^{\text{current subgradient}} - \left(\nabla f(\underbrace{w_k}_{\text{reference}}, \xi_l) - \nabla F(w_k) \right) \right)$$



SAGA

given i let $\nabla f(w_{[i]}, i)$
be the previous gradient
computed on point i

$$w_{k+1} = w_k - \gamma \left(\nabla f(w_k, i) - \left(\nabla f(w_{[i]}, i) - \frac{1}{n} \sum_{j=1}^n \nabla f(w_{[j]}, j) \right) \right)$$

Convergence in terms of time

FB $n\hat{K}$ $\ln \frac{1}{\epsilon}$

SG $\hat{K}^2 \frac{1}{\epsilon}$

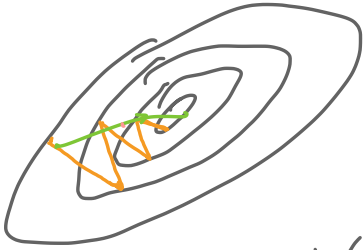
\hat{K} condition number
 $\hat{K} = \frac{L}{c}$

SAGA/
SVRG $(n + \hat{K})$ $\ln \frac{1}{\epsilon}$

SAGA, SVRG win on FB

POPULAR METHODS FOR OPTIMIZATION

• MOMENTUM



$$w_{k+1} = w_k - \alpha_k \nabla f(w_k, \xi_k) +$$

$$\beta_k (w_k - w_{k-1})$$

MOMENTUM TERM

$$\alpha_k = \alpha$$

$$\beta_k = \beta$$

SNOWBALL METHOD

$$FB \quad \text{error}(k) \sim \rho^k \quad \rho = \frac{\hat{k} - 1}{\hat{k} + 1}$$

\hat{k} condition number

$$\text{Snow Ball} \quad \text{error}(k) \sim (\rho')^k \quad \rho' = \frac{\sqrt{\hat{k}} - 1}{\sqrt{\hat{k}} + 1}$$

• NESTEROV METHOD

$$w_{k+1} = w_k - \alpha \nabla f(\underbrace{w_k + \beta(w_k - w_{k-1})}_{\text{Nesterov's look-ahead}}, \xi_k) \\ + \beta(w_k - w_{k-1})$$

UNDER THE SAME HYPOTHESES
FOR WHICH SG HAS $\frac{1}{K}$ CONV.

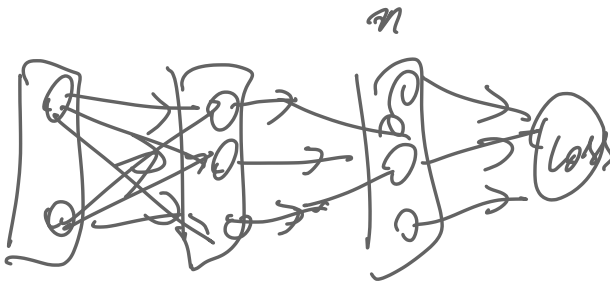
NESTEROV METHOD HAS $\frac{1}{K^2}$
(IMPOSSIBLE TO DO BETTER)

• COORDINATE DESCENT METHOD

$$(w_{k+1})_i = (w_k)_i - \alpha \frac{\partial f(w_k, \xi)}{\partial (w_k)_i}$$

OPTIMIZATION FOR NEURAL NETWORKS

→ ADAM, ADAGRAD, RMS Prop



$$L(\underline{W}) = L(\underline{Q}^n) = L(\theta_1^n, \dots, \theta_{m_n}^n)$$

output of node 1 of layer n → $\theta_1^n = a_1^n \left(\sum_{j=1}^{m_{n-1}} w_{j,1}^{n-1} o_j^{n-1} \right)$

$m_n = \#$ of neurons of layer n

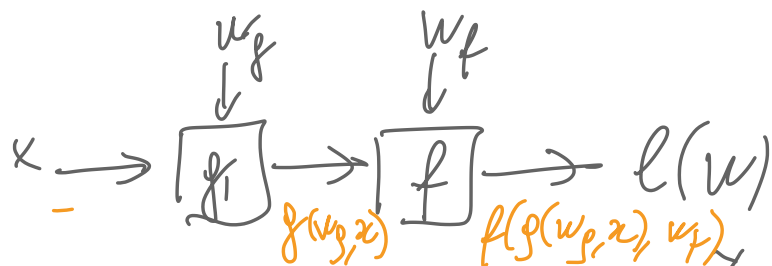
$$Q_1^n(x) = Q^n(x, W_1^n)$$

↑
activation function at the
1st neuron of layer n

$$l(w) = f(g(w))$$

$$l'(w) = f'(g(w)) \times g'(w)$$

$$f'(g(w)) = f'(x) \Big|_{x=g(w)}$$



$$l = f(w_f, g(w_g, x))$$

$$\frac{\partial l}{\partial w_g} = \frac{\partial f(w_f, g(w_g, x))}{\partial y} \times \frac{\partial g(w_g, x)}{\partial w_g}$$

$$= \left. \frac{\partial f(w_f, y)}{\partial y} \right|_{y=g(w_g, x)} \times \frac{\partial g(w_g, x)}{\partial w_g}$$