Variance reduction for stochastic gradient descent

Даня Меркулов

Методы Оптимизации в Машинном Обучении. ФКН ВШЭ



We consider classic finite-sample average minimization:

$$\min_{x\in\mathbb{R}^p}f(x)=\min_{x\in\mathbb{R}^p}\frac{1}{n}\sum_{i=1}^nf_i(x)$$

The gradient descent acts like follows:

$$x_{k+1} = x_k - \frac{\alpha_k}{n}\sum_{i=1}^n \nabla f_i(x)$$

(GD)

• Iteration cost is linear in n.

We consider classic finite-sample average minimization:

$$\min_{x\in\mathbb{R}^p}f(x)=\min_{x\in\mathbb{R}^p}\frac{1}{n}\sum_{i=1}^n f_i(x)$$

The gradient descent acts like follows:

$$x_{k+1} = x_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(x) \tag{GD}$$

- Iteration cost is linear in n.
- Convergence with constant α or line search.



We consider classic finite-sample average minimization:

$$\min_{x\in\mathbb{R}^p}f(x)=\min_{x\in\mathbb{R}^p}\frac{1}{n}\sum_{i=1}^n f_i(x)$$

The gradient descent acts like follows:

$$x_{k+1} = x_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(x) \tag{GD}$$

- Iteration cost is linear in n.
- Convergence with constant α or line search.



We consider classic finite-sample average minimization:

$$\min_{x\in\mathbb{R}^p}f(x)=\min_{x\in\mathbb{R}^p}\frac{1}{n}\sum_{i=1}^n f_i(x)$$

The gradient descent acts like follows:

$$x_{k+1} = x_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(x) \tag{GD}$$

- Iteration cost is linear in n.
- Convergence with constant α or line search.

Let's/ switch from the full gradient calculation to its unbiased estimator, when we randomly choose i_k index of point at each iteration uniformly:

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k) \tag{SGD}$$

With $p(i_k = i) = \frac{1}{n}$, the stochastic gradient is an unbiased estimate of the gradient, given by:

$$\mathbb{E}[\nabla f_{i_k}(x)] = \sum_{i=1}^n p(i_k=i) \nabla f_i(x) = \sum_{i=1}^n \frac{1}{n} \nabla f_i(x) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x) = \nabla f(x)$$

This indicates that the expected value of the stochastic gradient is equal to the actual gradient of f(x).

Stochastic iterations are n times faster, but how many iterations are needed?

If ∇f is Lipschitz continuous then we have:

| Assumption | Deterministic Gradient Descent | Stochastic Gradient Descent |
|----------------------------|---------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| PL Convex Non-Convex | $O(\log(1/arepsilon)) \ O(1/arepsilon) \ O(1/arepsilon) \ O(1/arepsilon)$ | O(1/arepsilon) onumber on |

• Stochastic has low iteration cost but slow convergence rate.



Stochastic iterations are n times faster, but how many iterations are needed?

If ∇f is Lipschitz continuous then we have:

| Assumption | Deterministic Gradient Descent | Stochastic Gradient Descent |
|------------|--------------------------------|-----------------------------|
| PL | $O(\log(1/arepsilon))$ | O(1/arepsilon) |
| Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |
| Non-Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |

• Stochastic has low iteration cost but slow convergence rate.

• Sublinear rate even in strongly-convex case.



Stochastic iterations are n times faster, but how many iterations are needed?

If ∇f is Lipschitz continuous then we have:

| Assumption | Deterministic Gradient Descent | Stochastic Gradient Descent |
|------------|--------------------------------|-----------------------------|
| PL | $O(\log(1/arepsilon))$ | O(1/arepsilon) |
| Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |
| Non-Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |

• Stochastic has low iteration cost but slow convergence rate.

- Sublinear rate even in strongly-convex case.
- Bounds are unimprovable under standard assumptions.

Stochastic iterations are n times faster, but how many iterations are needed?

If ∇f is Lipschitz continuous then we have:

| Assumption | Deterministic Gradient Descent | Stochastic Gradient Descent |
|------------|--------------------------------|-----------------------------|
| PL | $O(\log(1/arepsilon))$ | O(1/arepsilon) |
| Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |
| Non-Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |

• Stochastic has low iteration cost but slow convergence rate.

- Sublinear rate even in strongly-convex case.
- Bounds are unimprovable under standard assumptions.
- Oracle returns an unbiased gradient approximation with bounded variance.



Stochastic iterations are n times faster, but how many iterations are needed?

If ∇f is Lipschitz continuous then we have:

| Assumption | Deterministic Gradient Descent | Stochastic Gradient Descent |
|------------|--------------------------------|-----------------------------|
| PL | $O(\log(1/arepsilon))$ | O(1/arepsilon) |
| Convex | O(1/arepsilon) | $O(1/arepsilon^2)$ |
| Non-Convex | O(1/arepsilon) | $O(1/\varepsilon^2)$ |

- Stochastic has low iteration cost but slow convergence rate.
 - Sublinear rate even in strongly-convex case.
 - Bounds are unimprovable under standard assumptions.
 - Oracle returns an unbiased gradient approximation with bounded variance.
- Momentum and Quasi-Newton-like methods do not improve rates in stochastic case. Can only improve constant factors (bottleneck is variance, not condition number).

SGD with constant stepsize does not converge



 $f \rightarrow \min_{x,y,z}$ Finite-sum problem

Main problem of SGD

$$f(x) = \frac{\mu}{2} \|x\|_2^2 + \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \langle a_i, x \rangle)) \to \min_{x \in \mathbb{R}^n}$$

Strongly convex binary logistic regression. m=200, n=10, mu=1.



Variance reduction methods



Key idea of variance reduction $V_{\alpha r}(x) = ? \quad \text{(F)} x = ?$

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation: $Z = \alpha(X - Y) + \mathbb{R}[Y]$

$$\mathbb{E}[Y] \qquad \mathbb{E}[Y] \qquad \mathbb{E}[Y] \qquad \mathbb{E}[Y] = \mathcal{L}(\mathbb{E}[X] - \mathbb{E}[Y]) + \mathbb{E}[Y]$$

• $\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

$$\begin{array}{l} \bullet \ \mathbb{E}[Z_\alpha] = \alpha \mathbb{E}[X] + (1-\alpha) \mathbb{E}[Y] \\ \bullet \ \mathrm{var}(Z_\alpha) = \alpha^2 \left(\mathrm{var}(X) + \mathrm{var}(Y) - 2 \mathrm{cov}(X,Y) \right) \end{array}$$

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

$$\begin{array}{l} \bullet \ \mathbb{E}[Z_\alpha] = \alpha \mathbb{E}[X] + (1-\alpha) \mathbb{E}[Y] \\ \bullet \ \mathrm{var}(Z_\alpha) = \alpha^2 \left(\mathrm{var}(X) + \mathrm{var}(Y) - 2 \mathrm{cov}(X,Y) \right) \\ \bullet \ \mathrm{If} \ \alpha = 1 : \ \mathrm{no} \ \mathrm{bias} \end{array}$$

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^{2} (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

If α < 1: potential bias (but reduced variance).



Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^{2} (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

- If $\alpha < 1$: potential bias (but reduced variance).
- Useful if Y is positively correlated with X.

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^{2} (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

- If $\alpha < 1$: potential bias (but reduced variance).
- Useful if Y is positively correlated with X.

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^{2} (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

- If $\alpha < 1$: potential bias (but reduced variance).
- Useful if Y is positively correlated with X.

Application to gradient estimation ?

• SVRG: Let $X=\nabla f_{i_k}(x^{(k-1)})$ and $Y=\nabla f_{i_k}(\tilde{x}),$ with $\alpha=1$ and \tilde{x} stored.

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation:

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^{2} (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

- If $\alpha < 1$: potential bias (but reduced variance).
- Useful if Y is positively correlated with X.

Application to gradient estimation ?

• SVRG: Let
$$X = \nabla f_{i_{\star}}(x^{(k-1)})$$
 and $Y = \nabla f_{i_{\star}}(\tilde{x})$, with $\alpha = 1$ and \tilde{x} stored

•
$$\mathbb{E}[Y] = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\tilde{x})$$
 full gradient at \tilde{x} ;

$$f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) \longrightarrow \min_{x \in \mathbb{R}^n} f_i(x)$$

$$\nabla f(x) = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(x)$$

Principle: reducing variance of a sample of X by using a sample from another random variable Y with known expectation: $\begin{aligned} & \begin{array}{c} & & \\ & & \\ & (seO) \end{array} \\ & X_{K+1} = X_{K-} \quad d_{K} \cdot \nabla f_{i_{K}}(X_{K}) \end{aligned}$

$$Z_\alpha = \alpha(X-Y) + \mathbb{E}[Y]$$

•
$$\mathbb{E}[Z_{\alpha}] = \alpha \mathbb{E}[X] + (1 - \alpha) \mathbb{E}[Y]$$

• $\operatorname{var}(Z_{\alpha}) = \alpha^2 (\operatorname{var}(X) + \operatorname{var}(Y) - 2\operatorname{cov}(X, Y))$
• If $\alpha = 1$: no bias

- If $\alpha < 1$: potential bias (but reduced variance).
- Useful if Y is positively correlated with X.

Application to gradient estimation ?

• SVRG: Let
$$X = \nabla f_{i_k}(x^{(k-1)})$$
 and $Y = \nabla f_{i_k}(\tilde{x})$, with $\alpha = 1$ and \tilde{x} stored.

•
$$\mathbb{E}[Y] = \frac{1}{n} \sum_{i=1}^{n} \nabla f_i(\tilde{x})$$
 full gradient at \tilde{x} ;
• $X - Y = \nabla f_{i_k}(x^{(k-1)}) - \nabla f_{i_k}(\tilde{x})$
CTANO:
(SVRG) $X_{k+1} - X_k - d_k \left[\nabla f_{i_k}(X_k) - \nabla f_{i_k}(\tilde{x}) + \nabla f(\tilde{x}) \right]$

+ nogc nët
volto is
rpaguant
& rotane

rpaguers. & HOZOAR

2×SGD+

• Maintain table, containing gradient g_i of f_i , i = 1, ..., n



SAG (Stochastic average gradient, Schmidt, Le Roux, and Bach 2013) • Maintain table, containing gradient g_i of f_i , i = 1, ..., nXPAHIM • Initialize $x^{(0)},$ and $g_i^{(0)} = \nabla f_i(x^{(0)}), \ i = 1$ $\int_{i=1}^{\infty} f_{L}(x) \rightarrow \min$ n n

- Maintain table, containing gradient g_i of $f_i,\,i=1,\ldots,n$
- Initialize $x^{(0)}$, and $g_i^{(0)} = \nabla f_i(x^{(0)})$, $i=1,\ldots,n$
- At steps k=1,2,3,..., pick random $i_k\in\{1,\ldots,n\},$ then let

 $g_{i_k}^{(k)} = \nabla f_{i_k}(x^{(k-1)}) \quad (\text{most recent gradient of } f_{i_k})$



- Maintain table, containing gradient g_i of $f_i,\,i=1,\ldots,n$
- Initialize $x^{(0)}$, and $g_i^{(0)} = \nabla f_i(x^{(0)})$, $i=1,\ldots,n$
- At steps k=1,2,3,..., pick random $i_k\in\{1,\ldots,n\}$, then let

$$g_{i_k}^{(k)} = \nabla f_{i_k}(x^{(k-1)}) \quad (\text{most recent gradient of } f_{i_k})$$

Set all other $g_i^{(k)}=g_i^{(k-1)}$, $i
eq i_k$, i.e., these stay the same

Update

$$x^{(k)} = x^{(k-1)} - \alpha_k \frac{1}{n} \sum_{i=1}^n g_i^{(k)}$$

- Maintain table, containing gradient g_i of $f_i,\,i=1,\ldots,n$
- Initialize $x^{(0)}$, and $g_i^{(0)} = \nabla f_i(x^{(0)})$, $i=1,\ldots,n$
- At steps k=1,2,3,..., pick random $i_k\in\{1,\ldots,n\},$ then let

$$g_{i_k}^{(k)} = \nabla f_{i_k}(x^{(k-1)}) \quad (\text{most recent gradient of } f_{i_k})$$

Set all other $g_i^{(k)} = g_i^{(k-1)}$, $i \neq i_k$, i.e., these stay the same

Update

$$x^{(k)} = x^{(k-1)} - \alpha_k \frac{1}{n} \sum_{i=1}^n g_i^{(k)}$$

• SAG gradient estimates are no longer unbiased, but they have greatly reduced variance

- Maintain table, containing gradient g_i of $f_i,\,i=1,\ldots,n$
- Initialize $x^{(0)}$, and $g_i^{(0)} = \nabla f_i(x^{(0)}), \ i=1,\ldots,n$
- At steps k=1,2,3,..., pick random $i_k\in\{1,\ldots,n\},$ then let

$$g_{i_k}^{(k)} = \nabla f_{i_k}(x^{(k-1)}) \quad (\text{most recent gradient of } f_{i_k})$$

Set all other $g_i^{(k)} = g_i^{(k-1)} \text{, } i \neq i_k \text{, i.e., these stay the same}$

Update

$$x^{(k)} = x^{(k-1)} - \alpha_k \frac{1}{n} \sum_{i=1}^n g_i^{(k)}$$

- SAG gradient estimates are no longer unbiased, but they have greatly reduced variance
- Isn't it expensive to average all these gradients? Basically just as efficient as SGD, as long we're clever:

$$x^{(k)} = x^{(k-1)} - \alpha_k \underbrace{\left(\frac{1}{n}g_i^{(k)} - \frac{1}{n}g_i^{(k-1)} + \underbrace{\frac{1}{n}\sum_{i=1}^{n}g_i^{(k-1)}}_{\text{old table average}}\right)}_{\text{new table average}} + \underbrace{\text{CTO UMOCTG}}_{\text{TAKAS XC}} + \underbrace{\text{CTO U$$

Assume that $f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$, where each f_i is differentiable, and ∇f_i is Lipschitz with constant L. Denote $\bar{x}^{(k)} = \frac{1}{k} \sum_{l=0}^{k-1} x^{(l)}$, the average iterate after k-1 steps.





• Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.

- Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.
- This is $\mathcal{O}\left(\frac{1}{k}\right)$ convergence rate for SAG. Compare to $\mathcal{O}\left(\frac{1}{k}\right)$ rate for GD, and $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ rate for SGD.



- Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.
- This is $\mathcal{O}\left(\frac{1}{k}\right)$ convergence rate for SAG. Compare to $\mathcal{O}\left(\frac{1}{k}\right)$ rate for GD, and $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ rate for SGD.
- But, the constants are different! Bounds after k steps:



- Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.
- This is $\mathcal{O}\left(\frac{1}{k}\right)$ convergence rate for SAG. Compare to $\mathcal{O}\left(\frac{1}{k}\right)$ rate for GD, and $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ rate for SGD.
- But, the constants are different! Bounds after k steps:
 - GD: $\frac{L \|x^{(0)} x^{\star}\|^2}{2k}$



- Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.
- This is $\mathcal{O}\left(\frac{1}{k}\right)$ convergence rate for SAG. Compare to $\mathcal{O}\left(\frac{1}{k}\right)$ rate for GD, and $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ rate for SGD.
- But, the constants are different! Bounds after k steps:

• GD:
$$\frac{L\|x^{(0)}-x^{\star}\|^2}{2k}$$

• SAG: $\frac{48n[f(x^{(0)})-f^{\star}]+128L\|x^{(0)}-x^{\star}\|^2}{k}$

- Result stated in terms of the average iterate $\bar{x}^{(k)}$, but also can be shown to hold for the best iterate $x_{best}^{(k)}$ seen so far.
- This is $\mathcal{O}\left(\frac{1}{k}\right)$ convergence rate for SAG. Compare to $\mathcal{O}\left(\frac{1}{k}\right)$ rate for GD, and $\mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$ rate for SGD.
- But, the constants are different! Bounds after k steps:
 - GD: $\frac{L \|x^{(0)} x^{\star}\|^2}{2h}$
 - SAG: $\frac{48n[f(x^{(0)})-f^{\star}]+128L\|x^{(0)}-x^{\star}\|^2}{L}$
- So the first term in SAG bound suffers from a factor of n; authors suggest smarter initialization to make $f(x^{(0)}) f^*$ small (e.g., they suggest using the result of n SGD steps).
SAG convergence Assume further that each f_i is strongly convex with parameter μ .



+ numerineri exception KAK Y GD

SAG convergence

Assume further that each f_i is strongly convex with parameter μ .

i Theorem

SAG, with a step size $\alpha = \frac{1}{16L}$ and the same initialization as before, satisfies

$$\mathbb{E}[f(x^{(k)})] - f^{\star} \leq \left(1 - \min\left(\frac{\mu}{16L}, \frac{1}{8n}\right)\right)^k \left(\frac{3}{2}\left(f(x^{(0)}) - f^{\star}\right) + \frac{4L}{n}\|x^{(0)} - x^{\star}\|^2\right) \leq \frac{1}{2} \left(\frac{1}{2}\left(\frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k\right)^k - \frac{1}{2}\left(\frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k\right)^k - \frac{1}{2}\left(\frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}$$

- This is linear convergence rate $\mathcal{O}(\gamma^k)$ for SAG. Compare this to $\mathcal{O}(\gamma^k)$ for GD, and only $\mathcal{O}\left(\frac{1}{k}\right)$ for SGD.
- Like GD, we say SAG is adaptive to strong convexity.

SAG convergence

Assume further that each f_i is strongly convex with parameter μ .

i Theorem

SAG, with a step size $\alpha = \frac{1}{16L}$ and the same initialization as before, satisfies

$$\mathbb{E}[f(x^{(k)})] - f^{\star} \leq \left(1 - \min\left(\frac{\mu}{16L}, \frac{1}{8n}\right)\right)^k \left(\frac{3}{2}\left(f(x^{(0)}) - f^{\star}\right) + \frac{4L}{n}\|x^{(0)} - x^{\star}\|^2\right) \leq \frac{1}{2} \left(\frac{1}{2}\left(\frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k\right)^k - \frac{1}{2}\left(\frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k\right)^k - \frac{1}{2}\left(\frac{1}{2}\right)^k - \frac{1}{2}\left(\frac{$$

- This is linear convergence rate $\mathcal{O}(\gamma^k)$ for SAG. Compare this to $\mathcal{O}(\gamma^k)$ for GD, and only $\mathcal{O}\left(\frac{1}{k}\right)$ for SGD.
- Like GD, we say SAG is adaptive to strong convexity.
- Proofs of these results not easy: 15 pages, computed-aided!

• Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.



- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.



- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.
 - Choose initial L_0



- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.
 - Choose initial L_0
 - Increase *L*, until the following satisfies

$$f_{i_k}(x^{k+1}) \leq f_{i_k}(x^k) + \nabla f_{i_k}(x^k)(x^{k+1}-x^k) + \frac{L}{2}\|x^{k+1}-x^k\|_2^2$$

- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.
 - Choose initial L_0
 - Increase *L*, until the following satisfies

$$f_{i_k}(x^{k+1}) \leq f_{i_k}(x^k) + \nabla f_{i_k}(x^k)(x^{k+1}-x^k) + \frac{L}{2}\|x^{k+1}-x^k\|_2^2$$

• Decrease *L* between iterations



- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.
 - Choose initial L_0
 - Increase *L*, until the following satisfies

$$f_{i_k}(x^{k+1}) \leq f_{i_k}(x^k) + \nabla f_{i_k}(x^k)(x^{k+1}-x^k) + \frac{L}{2}\|x^{k+1}-x^k\|_2^2$$

• Decrease *L* between iterations

• Since stochastic gradient $g(x^k) \to \nabla f(x^k)$ you can use its norm to track convergence (which is not true for SGD!)



- Note, that the method in vanilla formulation is not applicable to the large neural networks training, due to the memory requirements.
- In practice you can use backtracking strategy to estimate Lipschitz constant.
 - Choose initial L₀
 - Increase *L*, until the following satisfies

$$f_{i_k}(x^{k+1}) \leq f_{i_k}(x^k) + \nabla f_{i_k}(x^k)(x^{k+1}-x^k) + \frac{L}{2}\|x^{k+1}-x^k\|_2^2$$

- Decrease L between iterations
- Since stochastic gradient $g(x^k) \to \nabla f(x^k)$ you can use its norm to track convergence (which is not true for SGD!)
- For the generalized linear models (this includes LogReg, LLS) you need to store much less memory $\mathcal{O}(n)$ instead of $\mathcal{O}(pn)$.



🔊 🖸 🛛 13

• The step size α_k and the convergence rate of the method are determined by the constant L for f(x), where $L = \max_{1 \le i \le n} L_i$, L_i is the Lipschitz constant for the function f_i

- The step size α_k and the convergence rate of the method are determined by the constant L for f(x), where $L = \max_{1 \le i \le n} L_i$, L_i is the Lipschitz constant for the function f_i
- When selecting components with a probability proportional to L_i , the constant L can be reduced from $\max_i L_i$ to $\overline{L} = \sum_i L_i / N$:

$$\begin{split} (x) &= \frac{1}{n} \sum_{i=1}^{n} f_i(x) \\ &= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \frac{f_i(x)}{L_i} \\ &= \frac{1}{\sum_k L_k} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \left(\sum_k \frac{L_k}{n} \frac{f_i(x)}{L_i} \right) \end{split}$$

With this approach, the component with a larger value of L_i is selected more often.

g

- The step size α_k and the convergence rate of the method are determined by the constant L for f(x), where $L = \max_{1 \le i \le n} L_i$, L_i is the Lipschitz constant for the function f_i
- When selecting components with a probability proportional to L_i , the constant L can be reduced from $\max_i L_i$ to $\overline{L} = \sum_i L_i / N$:

$$\begin{aligned} &(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) \\ &= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \frac{f_i(x)}{L_i} \\ &= \frac{1}{\sum_k L_k} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \left(\sum_k \frac{L_k}{n} \frac{f_i(x)}{L_i} \right) \end{aligned}$$

With this approach, the component with a larger value of L_i is selected more often.

g

To ensure convergence, component selection should be carried out according to the rule: with probability 0.5, select from a uniform distribution, with probability 0.5, select with probabilities L_i / ∑_i L_j.

- The step size α_k and the convergence rate of the method are determined by the constant L for f(x), where $L = \max_{1 \le i \le n} L_i$, L_i is the Lipschitz constant for the function f_i
- When selecting components with a probability proportional to L_i , the constant L can be reduced from $\max_i L_i$ to $\overline{L} = \sum_i L_i / N$:

$$\begin{aligned} &(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x) \\ &= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \frac{f_i(x)}{L_i} \\ &= \frac{1}{\sum_k L_k} \sum_{i=1}^{n} \sum_{j=1}^{L_i} \left(\sum_k \frac{L_k}{n} \frac{f_i(x)}{L_i} \right) \end{aligned}$$

With this approach, the component with a larger value of L_i is selected more often.

g

- To ensure convergence, component selection should be carried out according to the rule: with probability 0.5, select from a uniform distribution, with probability 0.5, select with probabilities L_i / ∑_i L_j.
- To generate with probabilities $L_i / \sum_j L_j$, there is an algorithm with complexity $O(\log N)$.

• Initialize: $\tilde{x} \in \mathbb{R}^d$



- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1 \ {\rm to} \ {\rm \#}$ of epochs



- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$



- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - $\bullet \ \ {\rm Initialize} \ x_0 = \tilde{x}$



- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - $\bullet \ \ {\rm Initialize} \ x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)
 - Pick $i_t \in \{1, ..., n\}$ uniformly at random

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)
 - Pick $i_t \in \{1, \dots, n\}$ uniformly at random
 - $\bullet \ x_t = x_{t-1} \alpha \left[\nabla f_{i_t}(x_{t-1}) \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)
 - Pick $i_t \in \{1, ..., n\}$ uniformly at random
 - $\bullet \ x_t = x_{t-1} \alpha \left[\nabla f_{i_t}(x_{t-1}) \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$
 - Update $\tilde{x} = x_m$

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)
 - Pick $i_t \in \{1, ..., n\}$ uniformly at random
 - $\bullet \ x_t = x_{t-1} \alpha \left[\nabla f_{i_t}(x_{t-1}) \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$
 - Update $\tilde{x} = x_m$

SVRG ~ GD

B noral

JUOXA



- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m) • Pick $i_t \in \{1, \dots, n\}$ uniformly at random • $x_t = x_{t-1} - \alpha \left[\nabla f_{i_t}(x_{t-1}) - \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$ • Update $\tilde{x} = x_m$

Notes:

• Two gradient evaluations per inner step.

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m) • Pick $i_t \in \{1, \dots, n\}$ uniformly at random • $x_t = x_{t-1} - \alpha \left[\nabla f_{i_t}(x_{t-1}) - \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$ • Update $\tilde{x} = x_m$

- Two gradient evaluations per inner step.
- Two parameters: length of epochs + step-size α .

- Initialize: $\tilde{x} \in \mathbb{R}^d$
- For $i_{epoch} = 1$ to # of epochs
 - Compute all gradients $\nabla f_i(\tilde{x})$; store $\nabla f(\tilde{x}) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{x})$
 - Initialize $x_0 = \tilde{x}$
 - For t = 1 to length of epochs (m)
 - Pick $i_t \in \{1, \dots, n\}$ uniformly at random
 - $\bullet \ x_t = x_{t-1} \alpha \left[\nabla f_{i_t}(x_{t-1}) \nabla f_{i_t}(\tilde{x}) + \nabla f(\tilde{x}) \right]$
 - Update $\tilde{x} = x_m$

- Two gradient evaluations per inner step.
- Two parameters: length of epochs + step-size α .
- Linear convergence rate, simple proof.

Adaptivity or scaling



Very popular adaptive method. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)})$, and update for $j = 1, \dots, p$:

$$\begin{split} v_j^{(k)} &= v_j^{k-1} + (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

Notes:

• AdaGrad does not require tuning the learning rate: $\alpha > 0$ is a fixed constant, and the learning rate decreases naturally over iterations.

Very popular adaptive method. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)}),$ and update for $j=1,\ldots,p:$

$$\begin{split} v_j^{(k)} &= v_j^{(k-1)} + (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- AdaGrad does not require tuning the learning rate: $\alpha > 0$ is a fixed constant, and the learning rate decreases naturally over iterations.
- The learning rate of rare informative features diminishes slowly.

Very popular adaptive method. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)}),$ and update for $j=1,\ldots,p:$

$$\begin{split} v_j^{(k)} &= v_j^{k-1} + (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- AdaGrad does not require tuning the learning rate: $\alpha > 0$ is a fixed constant, and the learning rate decreases naturally over iterations.
- The learning rate of rare informative features diminishes slowly.
- Can drastically improve over SGD in sparse problems.

Very popular adaptive method. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)}),$ and update for $j=1,\ldots,p:$

$$\begin{split} v_j^{(k)} &= v_j^{k-1} + (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- AdaGrad does not require tuning the learning rate: $\alpha > 0$ is a fixed constant, and the learning rate decreases naturally over iterations.
- The learning rate of rare informative features diminishes slowly.
- Can drastically improve over SGD in sparse problems.
- Main weakness is the monotonic accumulation of gradients in the denominator. AdaDelta, Adam, AMSGrad, etc. improve on this, popular in training deep neural networks.

Very popular adaptive method. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)}),$ and update for $j=1,\ldots,p:$

$$\begin{split} v_j^{(k)} &= v_j^{k-1} + (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- AdaGrad does not require tuning the learning rate: $\alpha > 0$ is a fixed constant, and the learning rate decreases naturally over iterations.
- The learning rate of rare informative features diminishes slowly.
- Can drastically improve over SGD in sparse problems.
- Main weakness is the monotonic accumulation of gradients in the denominator. AdaDelta, Adam, AMSGrad, etc. improve on this, popular in training deep neural networks.
- The constant ϵ is typically set to 10^{-6} to ensure that we do not suffer from division by zero or overly large step sizes.

RMSProp (Tieleman and Hinton, 2012)

An enhancement of AdaGrad that addresses its aggressive, monotonically decreasing learning rate. Uses a moving average of squared gradients to adjust the learning rate for each weight. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)})$ and update rule for $j = 1, \dots, p$:

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

Notes:

• RMSProp divides the learning rate for a weight by a running average of the magnitudes of recent gradients for that weight.



RMSProp (Tieleman and Hinton, 2012)

An enhancement of AdaGrad that addresses its aggressive, monotonically decreasing learning rate. Uses a moving average of squared gradients to adjust the learning rate for each weight. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)})$ and update rule for $j = 1, \dots, p$:

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- RMSProp divides the learning rate for a weight by a running average of the magnitudes of recent gradients for that weight.
- Allows for a more nuanced adjustment of learning rates than AdaGrad, making it suitable for non-stationary problems.



RMSProp (Tieleman and Hinton, 2012)

An enhancement of AdaGrad that addresses its aggressive, monotonically decreasing learning rate. Uses a moving average of squared gradients to adjust the learning rate for each weight. Let $g^{(k)} = \nabla f_{i_k}(x^{(k-1)})$ and update rule for $j = 1, \dots, p$:

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \frac{g_j^{(k)}}{\sqrt{v_j^{(k)} + \epsilon}} \end{split}$$

- RMSProp divides the learning rate for a weight by a running average of the magnitudes of recent gradients for that weight.
- Allows for a more nuanced adjustment of learning rates than AdaGrad, making it suitable for non-stationary problems.
- Commonly used in training neural networks, particularly in recurrent neural networks.

Adadelta (Zeiler, 2012)

An extension of RMSProp that seeks to reduce its dependence on a manually set global learning rate. Instead of accumulating all past squared gradients, Adadelta limits the window of accumulated past gradients to some fixed size w. Update mechanism does not require learning rate α :

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ \tilde{g}_j^{(k)} &= \frac{\sqrt{\Delta x_j^{(k-1)} + \epsilon}}{\sqrt{v_j^{(k)} + \epsilon}} g_j^{(k)} \\ x_j^{(k)} &= x_j^{(k-1)} - \tilde{g}_j^{(k)} \\ \Delta x_j^{(k)} &= \rho \Delta x_j^{(k-1)} + (1-\rho) (\tilde{g}_j^{(k)})^2 \end{split}$$

Notes:

• Adadelta adapts learning rates based on a moving window of gradient updates, rather than accumulating all past gradients. This way, learning rates adjusted are more robust to changes in model's dynamics.
Adadelta (Zeiler, 2012)

An extension of RMSProp that seeks to reduce its dependence on a manually set global learning rate. Instead of accumulating all past squared gradients, Adadelta limits the window of accumulated past gradients to some fixed size w. Update mechanism does not require learning rate α :

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ \tilde{g}_j^{(k)} &= \frac{\sqrt{\Delta x_j^{(k-1)} + \epsilon}}{\sqrt{v_j^{(k)} + \epsilon}} g_j^{(k)} \\ x_j^{(k)} &= x_j^{(k-1)} - \tilde{g}_j^{(k)} \\ \Delta x_j^{(k)} &= \rho \Delta x_j^{(k-1)} + (1-\rho) (\tilde{g}_j^{(k)})^2 \end{split}$$

- Adadelta adapts learning rates based on a moving window of gradient updates, rather than accumulating all past gradients. This way, learning rates adjusted are more robust to changes in model's dynamics.
- The method does not require an initial learning rate setting, making it easier to configure.

Adadelta (Zeiler, 2012)

An extension of RMSProp that seeks to reduce its dependence on a manually set global learning rate. Instead of accumulating all past squared gradients, Adadelta limits the window of accumulated past gradients to some fixed size w. Update mechanism does not require learning rate α :

$$\begin{split} v_j^{(k)} &= \gamma v_j^{(k-1)} + (1-\gamma) (g_j^{(k)})^2 \\ \tilde{g}_j^{(k)} &= \frac{\sqrt{\Delta x_j^{(k-1)} + \epsilon}}{\sqrt{v_j^{(k)} + \epsilon}} g_j^{(k)} \\ x_j^{(k)} &= x_j^{(k-1)} - \tilde{g}_j^{(k)} \\ \Delta x_j^{(k)} &= \rho \Delta x_j^{(k-1)} + (1-\rho) (\tilde{g}_j^{(k)})^2 \end{split}$$

- Adadelta adapts learning rates based on a moving window of gradient updates, rather than accumulating all past gradients. This way, learning rates adjusted are more robust to changes in model's dynamics.
- The method does not require an initial learning rate setting, making it easier to configure.
- Often used in deep learning where parameter scales differ significantly across layers.

Adam (Kingma and Ba, 2014)¹²

 $\langle 1 \rangle$

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

(1)

EMA:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)}\right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

(1 1)

Notes:

• It corrects the bias towards zero in the initial moments seen in other methods like RMSProp, making the estimates more accurate.

Update:

Adam (Kingma and Ba, 2014)¹²

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

EMA:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)} \right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

Notes:

- It corrects the bias towards zero in the initial moments seen in other methods like RMSProp, making the estimates more accurate.
- Одна из самых цитируемых научных работ в мире

Update:

Adam (Kingma and Ba. 2014)^{1 2}

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

EMA:

Update:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)}\right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

- It corrects the bias towards zero in the initial moments seen in other methods like RMSProp, making the estimates more accurate.
- Одна из самых цитируемых научных работ в мире
- В 2018-2019 годах вышли статьи, указывающие на ошибку в оригинальной статье

Adam (Kingma and Ba. 2014)^{1 2}

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

EMA:

Update:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)}\right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

- It corrects the bias towards zero in the initial moments seen in other methods like RMSProp. making the estimates more accurate.
- Одна из самых цитируемых научных работ в мире
- В 2018-2019 годах вышли статьи, указывающие на ошибку в оригинальной статье
- Не сходится для некоторых простых задач (даже выпуклых)

Adam (Kingma and Ba. 2014) ^{1 2}

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

EMA:

Update:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)} \right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

- It corrects the bias towards zero in the initial moments seen in other methods like RMSProp. making the estimates more accurate.
- Одна из самых цитируемых научных работ в мире
- В 2018-2019 годах вышли статьи, указывающие на ошибку в оригинальной статье
- Не сходится для некоторых простых задач (даже выпуклых)
- Почему-то очень хорошо работает для некоторых сложных задач

Adam (Kingma and Ba. 2014)¹²

Combines elements from both AdaGrad and RMSProp. It considers an exponentially decaying average of past gradients and squared gradients.

EMA:

Update:

$$\begin{split} \text{EMA:} & m_{j}^{(k)} = \beta_{1} m_{j}^{(k-1)} + (1-\beta_{1}) g_{j}^{(k)} \\ & v_{j}^{(k)} = \beta_{2} v_{j}^{(k-1)} + (1-\beta_{2}) \left(g_{j}^{(k)}\right)^{2} \\ \text{Bias correction:} & \hat{m}_{j} = \frac{m_{j}^{(k)}}{1-\beta_{1}^{k}} \\ & \hat{v}_{j} = \frac{v_{j}^{(k)}}{1-\beta_{2}^{k}} \\ \text{Update:} & x_{j}^{(k)} = x_{j}^{(k-1)} - \alpha \; \frac{\hat{m}_{j}}{\sqrt{\hat{v}_{j}} + \epsilon} \end{split}$$

- It corrects the bias towards zero in the initial moments seen in other methods like RMSProp. making the estimates more accurate.
- Одна из самых цитируемых научных работ в мире
- В 2018-2019 годах вышли статьи, указывающие на ошибку в оригинальной статье
- Не сходится для некоторых простых задач (даже выпуклых)
- Почему-то очень хорошо работает для некоторых сложных задач
- Гораздо лучше работает для языковых моделей, чем для задач компьютерного зрения - почему?

¹Adam: A Method for Stochastic Optimization

²On the Convergence of Adam and Beyond

AdamW (Loshchilov & Hutter, 2017)

Addresses a common issue with ℓ_2 regularization in adaptive optimizers like Adam. Standard ℓ_2 regularization adds $\lambda \|x\|^2$ to the loss, resulting in a gradient term λx . In Adam, this term gets scaled by the adaptive learning rate $\left(\sqrt{\hat{v}_j} + \epsilon\right)$, coupling the weight decay to the gradient magnitudes.

AdamW decouples weight decay from the gradient adaptation step.

Update rule:

$$\begin{split} m_j^{(k)} &= \beta_1 m_j^{(k-1)} + (1-\beta_1) g_j^{(k)} \\ v_j^{(k)} &= \beta_2 v_j^{(k-1)} + (1-\beta_2) (g_j^{(k)})^2 \\ \hat{m}_j &= \frac{m_j^{(k)}}{1-\beta_1^k}, \quad \hat{v}_j = \frac{v_j^{(k)}}{1-\beta_2^k} \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \left(\frac{\hat{m}_j}{\sqrt{\hat{v}_j} + \epsilon} + \lambda x_j^{(k-1)}\right) \end{split}$$

Notes:

• The weight decay term $\lambda x_j^{(k-1)}$ is added after the adaptive gradient step.

AdamW (Loshchilov & Hutter, 2017)

Addresses a common issue with ℓ_2 regularization in adaptive optimizers like Adam. Standard ℓ_2 regularization adds $\lambda \|x\|^2$ to the loss, resulting in a gradient term λx . In Adam, this term gets scaled by the adaptive learning rate $\left(\sqrt{\hat{v}_j} + \epsilon\right)$, coupling the weight decay to the gradient magnitudes.

AdamW decouples weight decay from the gradient adaptation step.

Update rule:

$$\begin{split} m_j^{(k)} &= \beta_1 m_j^{(k-1)} + (1 - \beta_1) g_j^{(k)} \\ v_j^{(k)} &= \beta_2 v_j^{(k-1)} + (1 - \beta_2) (g_j^{(k)})^2 \\ \hat{m}_j &= \frac{m_j^{(k)}}{1 - \beta_1^k}, \quad \hat{v}_j = \frac{v_j^{(k)}}{1 - \beta_2^k} \\ x_j^{(k)} &= x_j^{(k-1)} - \alpha \left(\frac{\hat{m}_j}{\sqrt{\hat{v}_j} + \epsilon} + \lambda x_j^{(k-1)} \right) \end{split}$$

Notes:

- The weight decay term $\lambda x_{j}^{(k-1)}$ is added after the adaptive gradient step.
- Widely adopted in training transformers and other large models. Default choice for huggingface trainer.

 $f \rightarrow \min_{x,y,z}$ Adaptivity or scaling

A lot of them



♥ O Ø 22

How to compare them? AlgoPerf benchmark



NanoGPT speedrun



24

Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

1. Compute gradient G_k .



Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.



Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)



Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)
- 4. Update: $W_{k+1} = W_k \alpha P_L G_k P_R$.



Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)
- 4. Update: $W_{k+1} = W_k \alpha P_L G_k P_R$.

Notes:

• Aims to capture curvature information more effectively than first-order methods.

Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)
- 4. Update: $W_{k+1} = W_k \alpha P_L G_k P_R$.

- Aims to capture curvature information more effectively than first-order methods.
- Computationally more expensive than Adam but can converge faster or to better solutions in terms of steps.

Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)
- 4. Update: $W_{k+1} = W_k \alpha P_L G_k P_R$.

- Aims to capture curvature information more effectively than first-order methods.
- Computationally more expensive than Adam but can converge faster or to better solutions in terms of steps.
- Requires careful implementation for efficiency (e.g., efficient computation of inverse matrix roots, handling large matrices).

Stands for **S**tochastic Hessian-Approximation Matrix Preconditioning for Optimization Of deep networks. It's a method inspired by second-order optimization designed for large-scale deep learning.

Core Idea: Approximates the full-matrix AdaGrad pre conditioner using efficient matrix structures, specifically Kronecker products.

For a weight matrix $W \in \mathbb{R}^{m \times n}$, the update involves preconditioning using approximations of the statistics matrices $L \approx \sum_k G_k G_k^T$ and $R \approx \sum_k G_k^T G_k$, where G_k are the gradients.

Simplified concept:

- 1. Compute gradient G_k .
- 2. Update statistics $L_k = \beta L_{k-1} + (1-\beta)G_kG_k^T$ and $R_k = \beta R_{k-1} + (1-\beta)G_k^TG_k$.
- 3. Compute preconditioners $P_L = L_k^{-1/4}$ and $P_R = R_k^{-1/4}$. (Inverse matrix root)
- 4. Update: $W_{k+1} = W_k \alpha P_L G_k P_R$.

- Aims to capture curvature information more effectively than first-order methods.
- Computationally more expensive than Adam but can converge faster or to better solutions in terms of steps.
- Requires careful implementation for efficiency (e.g., efficient computation of inverse matrix roots, handling large matrices).
- Variants exist for different tensor shapes (e.g., convolutional layers).

Muon³

$$\begin{split} W_{t+1} &= W_t - \eta (G_t G_t^\top)^{-1/4} G_t (G_t^\top G_t)^{-1/4} \\ &= W_t - \eta (US^2 U^\top)^{-1/4} (USV^\top) (VS^2 V^\top)^{-1/4} \\ &= W_t - \eta (US^{-1/2} U^\top) (USV^\top) (VS^{-1/2} V^\top) \\ &= W_t - \eta US^{-1/2} SS^{-1/2} V^\top \\ &= W_t - \eta UV^\top \end{split}$$

