

随机优化算法介绍

胡涛

北京大学信息科学技术学院

taohu@pku.edu.cn

备注：此笔记由北京大学文再文老师《大数据分析中的算法》课堂讲义整理而成

<http://bicmr.pku.edu.cn/~wenzw/bigdata/lect-sto.pdf>

Contents

1 概览	4
2 次梯度法	4
2.1 次梯度法(Subgradient methods)	4
2.2 投影次梯度法(Projected Subgradient methods)	6
2.3 随机次梯度法(Stochastic Subgradient Methods)	7
2.4 Azuma-Hoeffding Inequality	7
2.5 Adaptive stepsizes and metrics	7
2.5.1 Adaptive stepsize	7
2.5.2 Variable metric methods	7
2.5.3 Variable metric methods收敛性	8
2.5.4 Optimality Guarantees	8
3 梯度法	8
3.1 梯度法(GD)的收敛性	9
3.1.1 梯度法(GD)的另一种收敛性证明	11
3.1.2 backtracking line search	11
3.1.3 Gradient method for strongly convex function	12
3.1.4 利用二阶上界和二阶下界来证明GD收敛性	13
3.1.5 BB step	14
3.2 随机梯度法(SGD)的收敛性	15
4 Variance Reduction	17
4.1 回顾GD,SGD	17
4.2 SAG,SAGA,SVRG	18
4.3 SVRG算法理论证明	20
5 随机优化算法在深度学习中的应用	20
5.1 FG与SG	20
5.2 常见算法	20
5.2.1 SGD with momentum	20
5.2.2 Nesterov accelerated gradient (original version)	20
5.2.3 Nesterov accelerated gradient (momentum version)	21
5.2.4 Adaptive Subgradient Methods(Adagrad)	21
5.2.5 RMSprop	21
5.2.6 Adam	22
6 总结	23
7 附录	23
7.1 基本性质	23
7.2 Co-coercivity of gradient	27
7.2.1 co-coercivity的扩展	29
7.3 Projection Operator is non-expansive	29

1 概览

首先介绍Hoeffding Inequality:

martigale的定义可以参照[https://en.wikipedia.org/wiki/Martingale_\(probability_theory\)](https://en.wikipedia.org/wiki/Martingale_(probability_theory))

martingale difference sequence (MDS)的定义可以参照https://en.wikipedia.org/wiki/Martingale_difference_sequence

martigale和martingale difference sequence之间有一些联系。

总体需要优化的问题如下:

$$\min_{x \in R^n} f(x) \text{ 其中 } f_x \text{ 为需要优化的函数}$$

下面主要会介绍以下几种随机优化算法:

- 次梯度法

- 梯度法

- Variance Reduction

- 随机优化算法在深度学习中的应用

2 次梯度法

2.1 次梯度法(Subgradient methods)

次梯度法的流程如下:

$$x_{k+1} = x_k - \alpha_k g_k, g_k \in \partial f(x_k) \quad (2.1)$$

上述的式子等价于下面的式子:

$$x_{k+1} = \operatorname{argmin}_x f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 \quad (2.2)$$

公式2.1的具体推导如下:

泰勒公式二阶展开

$$f(x) \approx f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|_2^2$$

$$\text{则 } x_{k+1} = \operatorname{argmin}_x f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|_2^2$$

$$x_{k+1} = \operatorname{argmin}_x \langle g_k, x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|_2^2$$

$$x_{k+1} = \operatorname{argmin}_x 2\alpha_k \langle g_k, x - x_k \rangle + \|x - x_k\|_2^2$$

化简得到: $x_{k+1} = \operatorname{argmin}_x \langle x, x \rangle + \langle 2\alpha_k g_k - 2x_k, x \rangle$

上述问题有显式解:

$$x_{k+1} = x_k - \alpha_k g_k, \text{得证}$$

Theorem 1: Convergence of subgradient

Let $\alpha_k \geq 0$ be any non-negative sequence of stepsizes and the preceding assumptions hold. Let x_k be generated by the subgradient iteration. Then for all $K \geq 1$,

$$\sum_{k=1}^K \alpha_k [f(x_k) - f(x^*)] \leq \frac{1}{2} \|x_1 - x^*\|_2^2 + \frac{1}{2} \sum_{k=1}^K \alpha_k^2 M^2.$$

次梯度法的公式很简单，那么次梯度法的收敛性如何呢？下面给予证明。

首先证明一个引理：

值得注意的是上面的引理有两个假设：

- 最优解至少是bounded, 即存在 $x^* \in \operatorname{argmin}_x f(x)$ 并且 $f(x^*) > -\infty$
- 所有的次梯度都是bounded, 即 $\|g\|_2 \leq M \leq \infty$ 对所有的 x 和 $g \in \partial f(x)$ 都成立

下面给出具体证明：

由于 $f(x)$ 为凸函数，所以有：

$$\langle g_k, x^* - x_k \rangle \leq f(x^*) - f(x_k) \quad (2.3)$$

$$\begin{aligned} & \frac{1}{2} \|x_{k+1} - x^*\|_2^2 = \frac{1}{2} \|x_k - \alpha_k g_k - x^*\|_2^2 \\ & \text{拼凑, } = \frac{1}{2} \|x_k - x^*\|_2^2 - \alpha_k \langle g_k, x^* - x_k \rangle + \alpha_k^2 \|g_k\|_2^2 \\ & \text{利用凸函数性质 (2.3), } \leq \frac{1}{2} \|x_k - x^*\|_2^2 - \alpha_k (f(x_k) - f(x^*)) + \frac{\alpha_k^2}{2} M^2 \\ & \text{利用归纳法, 即得证.} \end{aligned}$$

引理证明完以后，下面接着证明次梯度法的收敛性。首先令 $\bar{x}_k = \frac{\sum_{k=1}^K \alpha_k x_k}{\sum_{k=1}^K \alpha_k}$ 。结合上面的引理很显然可以推导出：

$$f(\bar{x}_k) - f(x^*) \leq \frac{\sum_{k=1}^K \alpha_k x_k + \sum_{k=1}^K \alpha_k^2 M^2}{2 \sum_{k=1}^K \alpha_k}$$

可以得到以下几个结论：

- 根据实际应用中我们对步长的设置, $\sum_{k=1}^\infty \alpha_k = \infty$, 并且 $\frac{\sum_{k=1}^K \alpha_k^2 M^2}{2 \sum_{k=1}^K \alpha_k} \rightarrow 0$, 得知随着 K 增大, 式子左边会趋近于 0。
- 假设我们使用固定步长, $\alpha_k = \alpha$, $\|x_1 - x^*\| \leq R$, 那么可以得到：

$$f(\bar{x}_k) - f(x^*) \leq \frac{R^2}{2K\alpha} + \frac{\alpha M^2}{2}$$

- 如果使用固定步长, 上面的式子就不会趋近于 0 了, 因为有 $\frac{\alpha M^2}{2}$ 这一项。我们可以令步长 $\alpha_k = \frac{R}{M\sqrt{k}}$, 这样式子 $\frac{\alpha M^2}{2}$ 就会趋近于 0。

那么为什么 $f(\bar{x}_k) - f(x^*)$ 趋近于 0, 次梯度法就收敛呢?

2.2 投影次梯度法(Projected Subgradient methods)

考虑如下问题:

$$\min_{x \in C} f(x)$$

投影次梯度法步骤如下:

$$x_{k+1} = \pi_C(x_k - \alpha_k g_k), g_k \in \partial f(x_k)$$

其中的投影操作为:

$$\pi_C = \operatorname{argmin}_{y \in C} \|y - x\|_2^2$$

投影次梯度法也可以写成:

$$x_{k+1} = \operatorname{argmin}_{y \in C} f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{2\alpha} \|x - x_k\|_2^2$$

投影次梯度法的证明如下所示:

同样这里先证明一个引理:

Theorem 2: Convergence of projected subgradient method

Let $\alpha_k \geq 0$ be any non-negative sequence of stepsizes and the preceding assumptions hold. Let x_k be generated by the projected subgradient iteration. Then for all $K \geq 1$,

$$\sum_{k=1}^K \alpha_k [f(x_k) - f(x^*)] \leq \frac{1}{2} \|x_1 - x^*\|_2^2 + \frac{1}{2} \sum_{k=1}^K \alpha_k^2 M^2. \quad (20)$$

值得注意的是上面的引理同样有两个假设:

- 最优解至少是 bounded, 即存在 $x^* \in \operatorname{argmin}_x f(x)$ 并且 $f(x^*) > -\infty$
- 所有的次梯度都是 bounded, 即 $\|g\|_2 \leq M \leq \infty$ 对所有的 x 和 $g \in \partial f(x)$ 都成立

首先利用到了 convex set 上 projection 的 non-expansiveness (non-expansiveness 的证明见附录):

$$\begin{aligned} &\text{if } \pi_C = \operatorname{argmin}_{y \in C} \|y - x\|_2^2 \\ &\text{then, } \|y_1 - y_2\| \geq \|x_1 - x_2\| \end{aligned}$$

利用 π_C 的 non-expansiveness, 可以得到:

$$\|x_{k+1} - x^*\|_2^2 = \|\pi_C(x_k - \alpha_k g_k) - x^*\|_2^2 = \|\pi_C(x_k - \alpha_k g_k) - \pi_C(x^*)\|_2^2 \leq \|x_k - \alpha_k g_k - x^*\|_2^2$$

接下来的证明就和次梯度法类似了:

$$\begin{aligned} &\frac{1}{2} \|x_{k+1} - x^*\|_2^2 \leq \frac{1}{2} \|x_k - \alpha_k g_k - x^*\|_2^2 \\ &\text{拼凑, } = \frac{1}{2} \|x_k - x^*\|_2^2 - \alpha_k \langle g_k, x^* - x_k \rangle + \alpha_k^2 \|g_k\|_2^2 \\ &\text{利用凸函数性质(2.3), } \leq \frac{1}{2} \|x_k - x^*\|_2^2 - \alpha_k (f(x_k) - f(x^*)) + \frac{\alpha_k^2}{2} M^2 \\ &\text{利用归纳法, 即得证.} \end{aligned}$$

接下来的收敛性证明和次梯度法类似, 这里就不详细介绍了。

2.3 随机次梯度法(Stochastic Subgradient Methods)

2.4 Azuma-Hoeffding Inequality

2.5 Adaptive stepsizes and metrics

选择一个合适的度量方式(metrics)，或者一个更好的步长方案，往往能够实现更快的收敛。因此在梯度下降方法上一般都从以上两个方面来改进。

一个简单的方案是：

$$h(x) = \frac{1}{2}x^T Ax$$

其中A和数据项有关。

2.5.1 Adaptive stepsize

回顾上面我们介绍的随机次梯度法的bound:

$$E[f(\bar{x}_k) - f(x^*)] \leq E\left[\frac{R^2}{K\alpha_k} + \frac{1}{2K} \sum_{k=1}^K \alpha_k \|g_k\|_*^2\right]$$

很显然当k趋近于无穷大时， $\frac{1}{2K} \sum_{k=1}^K \alpha_k \|g_k\|_*^2$ 的收敛性无法保证，这里我们构造了一种步长规则，使得 $\frac{1}{2K} \sum_{k=1}^K \alpha_k \|g_k\|_*^2$ 收敛。

令 $\alpha_k = \frac{R}{\sqrt{\sum_{k=1}^K \|g_k\|_*^2}}$ ，则可以得到：

$$E[f(\bar{x}_k) - f(x^*)] \leq \frac{2R}{K} E[(\sum_{k=1}^K \|g_k\|_*^2)^{\frac{1}{2}}]$$

假设 $E[\|g_k\|_*^2] \leq M^2, \forall k$ ，上式可以进一步化简：

$$E[(\sum_{k=1}^K \|g_k\|_*^2)^{\frac{1}{2}}] \leq \sqrt{M^2 K} \leq M \sqrt{K}$$

可以看到，随着K增大，整个式子是收敛的。

2.5.2 Variable metric methods

Variable metric methods本质上就是选择不同的metric，即 H_k 矩阵的构造。

再来看看一下原问题：

$$x_{k+1} = \operatorname{argmin}_{x \in C} x_k + \langle g_k, x - x_k \rangle + \frac{1}{2} \langle x - x_k, H_k(x - x_k) \rangle$$

去掉无关变量：

$$x_{k+1} = \operatorname{argmin}_{x \in C} \langle g_k, x \rangle + \frac{1}{2} \langle x - x_k, H_k(x - x_k) \rangle$$

这就是梯度法的基础模型，下面举出几种常见的 H_k 的取法。

- 投影次梯度法。

$$H_k = \alpha_k I$$

- 牛顿法

$$H_k = \nabla^2 f(x_k)$$

- AdaGrad

$$H_k = \frac{1}{\alpha} \text{diag}(\sum_{k=1}^K g_k \cdot * g_k)^{\frac{1}{2}}$$

其中,点乘表示elementwise multiplication。diag表示将矢量按照对角线扩充为矩阵。

2.5.3 Variable metric methods收敛性

Theorem 9: Convergence of Variable metric methods

Let $H_k > 0$ be a sequence of positive define matrices, where H_k is a function of g_1, \dots, g_k . Let g_k be stochastic subgradient with $\mathbb{E}[g_k|x_k] \in \partial f(x_k)$. Then

$$\begin{aligned} \mathbb{E}\left[\sum_{k=1}^K (f(x_k) - f(x^*))\right] &\leq \frac{1}{2} \mathbb{E}\left[\sum_{k=2}^K (\|x_k - x^*\|_{H_k}^2 - \|x_k - x^*\|_{H_{k-1}}^2)\right] \\ &\quad + \frac{1}{2} \mathbb{E}\left[\|x_1 - x^*\|_{H_1}^2 + \sum_{k=1}^K \|g_k\|_{H_k}^2\right]. \end{aligned}$$

2.5.4 Optimality Guarantees

3 梯度法

梯度法的流程如下:

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

上述的式子等价于下面的式子:

$$x_{k+1} = \underset{x}{\operatorname{argmin}} f(x) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 \quad (2.4)$$

2.4式的具体推导可以参照次梯度法中的推导。

3.1 梯度法(GD)的收敛性

首先假设函数 f 是满足 L -smooth和 μ -strongly convex条件的。

Assumption

- f is L -smooth and μ -strongly convex.

lemma: Coercivity of gradients

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{L\mu}{L + \mu} \|x - y\|^2 + \frac{1}{L + \mu} \|\nabla f(x) - \nabla f(y)\|^2 \quad (42)$$

Theorem: Convergence rates of GD

Let $\alpha_k = \frac{2}{L+\mu}$ and let $\kappa = \frac{L}{\mu}$. Define $\Delta_k = \|x_k - x^*\|$. Then we get,

$$f(x_{T+1}) - f(x^*) \leq \frac{L\Delta_1^2}{2} \exp\left(-\frac{4T}{\kappa + 1}\right). \quad (43)$$

接下来证明收敛性，这里会用到coercivity of gradients的性质，大致的思路是先证明 $f(x_{T+1}) - f(x^*)$ 与 $f(x_{T+1}) - f(x^*)$ 之间的递推关系，然后利用 $\alpha_k = \frac{2}{L+\mu}$ 的特殊性得到 $f(x_{T+1}) - f(x^*) \leq \frac{L\Delta_1^2}{2} \exp\left(-\frac{4T}{\kappa + 1}\right)$

Proof of Theorem

•

$$\begin{aligned}
 \Delta_{k+1}^2 &= \|x_{k+1} - x^*\|_2^2 = \|x_k - \alpha_k \nabla f(x_k) - x^*\|_2^2 \\
 &= \|x_k - x^*\|_2^2 - 2\alpha_k \langle \nabla f(x_k), x_k - x^* \rangle + \alpha_k^2 \|\nabla f(x_k)\|_2^2 \\
 &= \Delta_k^2 - 2\alpha_k \boxed{\langle \nabla f(x_k), x_k - x^* \rangle} + \alpha_k^2 \|\nabla f(x_k)\|_2^2
 \end{aligned}$$

• By the lemma

$$\begin{aligned}
 \Delta_{k+1}^2 &\leq \Delta_k^2 - 2\alpha_k \left[\frac{L\mu}{L+\mu} \Delta_k^2 + \frac{1}{L+\mu} \|\nabla f(x_k)\|^2 \right] + \alpha_k^2 \|\nabla f(x_k)\|_2^2 \\
 &= (1 - 2\alpha_k \frac{L\mu}{L+\mu}) \Delta_k^2 + \left(-\frac{2\alpha_k}{L+\mu} + \alpha_k^2 \right) \|\nabla f(x_k)\|_2^2 \\
 &\leq (1 - 2\alpha_k \frac{L\mu}{L+\mu}) \Delta_k^2 + \left(-\frac{2\alpha_k}{L+\mu} + \alpha_k^2 \right) L^2 \Delta_k^2
 \end{aligned} \tag{44}$$

Proof of Theorem

• $\alpha_k = \frac{2}{L+\mu}$

$$\begin{aligned}
 \Delta_{k+1}^2 &\leq \left(1 - \frac{4L\mu}{(L+\mu)^2}\right) \Delta_k^2 \\
 &= \left(\frac{L-\mu}{L+\mu}\right)^2 \Delta_k^2 = \left(\frac{\kappa-1}{\kappa+1}\right)^2 \Delta_k^2
 \end{aligned}$$

•

$$\begin{aligned}
 \Delta_{T+1}^2 &\leq \left(\frac{\kappa-1}{\kappa+1}\right)^{2T} \Delta_1^2 \\
 &= \Delta_1^2 \exp\left(2T \log\left(1 - \frac{2}{\kappa+1}\right)\right) \\
 &\leq \Delta_1^2 \exp\left(-\frac{4T}{\kappa+1}\right)
 \end{aligned}$$

•

$$f(x_{T+1}) - f(x^*) \leq \frac{L}{2} \Delta_{T+1}^2 \leq \frac{L\Delta_1^2}{2} \exp\left(-\frac{4T}{\kappa+1}\right)$$

3.1.1 梯度法(GD)的另一种收敛性证明

首先考虑固定步长的情况:

附录中介绍了quadratic upper bound:

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2}\|y - x\|_2^2$$

在这里令 $y = x - t\nabla f(x)$ 得到:

$$f(x - t\nabla f(x)) \leq f(x) - t(1 - \frac{Lt}{2})\|\nabla f(x)\|_2^2$$

x^+ 表示下一次迭代的 x , 这里假设步长满足如下:

$$x^+ = x - t\nabla f(x), 0 < t < \frac{1}{L}$$

所以有

$$\begin{aligned} f(x^+) &\leq f(x) - \frac{t}{2}\|\nabla f(x)\|_2^2????? \\ &\leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{t}{2}\|\nabla f(x)\|_2^2(\text{凸函数性质}) \\ &= f(x^*) + \frac{1}{2t}(\|x - x^*\|_2^2 - \|x - x^* - t\nabla f(x)\|_2^2)(\text{整理}) \\ &= f(x^*) + \frac{1}{2t}(\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2) \end{aligned}$$

所以有:

$$f(x^+) - f(x^*) = \frac{1}{2t}(\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2)$$

令 $x = x^{(i=1)}, x^+ = x^{(i)}$, 则
 $\sum_{i=1}^k(f(x^{(i)}) - f(x^*)) \leq \frac{1}{2t} \sum_{i=1}^k(\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2) \leq \frac{1}{2t}(\|x^{(0)} - x^*\|_2^2 - \|x^{(k)} - x^*\|_2^2) \leq \frac{1}{2t}\|x^{(0)} - x^*\|_2^2$
因为 $f(x^{(i)})$ 是递减的, 那么

$$f(x^{(k)}) - f(x^*) \leq \frac{1}{k} \sum_{i=1}^k(f(x^i) - f(x^*)) \leq \frac{1}{2kt}\|x^{(0)} - x^*\|_2^2$$

其中 k 在分母中, 因此梯度下降法达到 $f(x^{(k)}) - f(x^*) \leq \epsilon$ 的精度需要 $O(\frac{1}{\epsilon})$ 的时间复杂度。

3.1.2 backtracking line search

line search方法有很多, 这里只介绍backtracking line search:

initialize t_k at \hat{t} (for example $\hat{t} = 1$), take $t_k = \beta t_k$ until:
 $f(x - t_k \nabla f(x)) \leq f(x) - \alpha t_k \|\nabla f(x)\|_2^2$

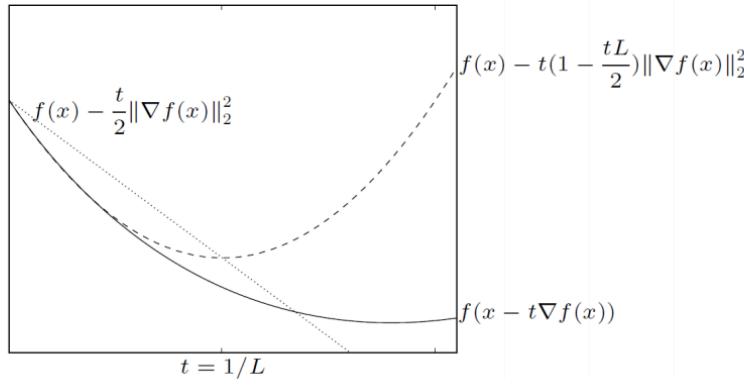
$0 < \beta < 1, \alpha = \frac{1}{2}$ mostly for simplify proofs

前面说到quadratic upper bound有:

$$f(x - t\nabla f(x)) \leq f(x) - t(1 - \frac{L}{2})\|\nabla f(x)\|_2^2$$

原函数,quadratic upper bound,line search对应的三条曲线如下所示(其中步长t为自变量,纵轴表示函数值,其中 $t = \frac{1}{L}$ 是quadratic upper bound的极小值点)。因此最好的步长t是line search方法和原函数的交点处的t值(记为 t_{opt}),t值太大也不可取,因为我们无法确定 t_{opt} 右侧的情况,但是能够保证在 t_{opt} 左侧:line search直线始终是原函数的linear upper bound。

line search with $\alpha = 1/2$ if f has a Lipschitz continuous gradient



3.1.3 Gradient method for strongly convex function

前面介绍的是L-Lipschitz Continuous Gradient的情况,利用了函数的quadratic upper bound来证明收敛性。这里主要介绍在 μ -strongly convex情况下GD收敛性的证明。

首先回顾co-coercivity性质之一:

co-coercivity of ∇g gives
 $(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{\mu L}{\mu + L}\|x - y\|_2^2 + \frac{1}{\mu + L}\|\nabla f(x) - \nabla f(y)\|_2^2$

令 $x^+ = x - t\nabla f(x)$, $0 < t \leq \frac{2}{\mu + L}$

则有:

$$\begin{aligned} \|x^+ - x^*\| &= \|x^+ - t\nabla f(x) - x^*\|_2^2 \text{ (拆开)} \\ &= \|x^+ - x^*\|_2^2 - 2t\nabla f(x)^T(x - x^*) + t^2\|\nabla f(x)\|_2^2 \\ &\leq \|x^+ - x^*\| - \frac{2t\mu L}{\mu + L}\|x - x^*\|_2^2 + (t^2 - \frac{2t}{t + \mu})\|\nabla f(x)\|_2^2 \text{ (此处有co-coercivity of gradient性质得到)} \end{aligned}$$

$$\begin{aligned} &\leq \left(1 - t \frac{2\mu L}{\mu+L}\right) \|x - x^*\|_2^2 + t \left(t - \frac{2}{\mu+L}\right) \|\nabla f(x)\|_2^2 \\ &\leq \left(1 - t \frac{2\mu L}{\mu+L}\right) \|x - x^*\|_2^2 \end{aligned}$$

综上，使用归纳法可以得到：

$$\|x^{(k)} - x^*\|_2^2 \leq c^k \|x^{(0)} - x^*\|_2^2, c = 1 - t \frac{2\mu L}{\mu+L}$$

所以有：

$$f(x^{(k)}) - f(x^*) \leq \frac{L}{2} \|x^{(k)} - x^*\| \leq \frac{c^k L}{2} \|x^{(0)} - x^*\|_2^2$$

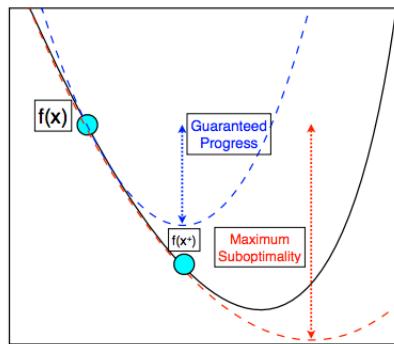
因此梯度下降法达到 $f(x^{(k)}) - f(x^*) \leq \epsilon$ 的精度需要 $O(\log(\frac{1}{\epsilon}))$ 的时间复杂度。//TODOTODO

3.1.4 利用二阶上界和二阶下界来证明GD收敛性

这里提供另外一种证明方法，比较直观理解GD。https://www.cs.ubc.ca/~schmidtm/Documents/2013_Notes_ConvexOptim.pdf

- We have bounds on x^+ and x^* :

$$f(x^+) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2, \quad f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2.$$



其中

$f(x^+) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$ 是有 L-Lipschitz Continuous gradient 的性质得到。(令 $y = x^+ = x - \alpha \nabla f(x)$, $\alpha = \frac{1}{L}$ ，带入即可。)

$f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$ 是由 μ -strongly convex 性质得到。

- We have bounds on x^+ and x^* :

$$f(x^+) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2, \quad f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2.$$

combine them to get

$$f(x^+) \leq f(x) - \frac{\mu}{L} [f(x) - f(x^*)]$$

$$f(x^+) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right) [f(x) - f(x^*)]$$

- This gives a linear convergence rate:

$$f(x^t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^t [f(x^0) - f(x^*)]$$

- Each iteration multiplies the error by a fixed amount.

(very fast if μ/L is not too close to one)

3.1.5 BB step

Consider the problem

$$\min f(x)$$

- Steepest gradient descent method: $x^{k+1} := x^k - \alpha^k g^k$:

$$\alpha^k := \arg \min_{\alpha} f(x^k - \alpha g^k)$$

- Let $s^{k-1} := x^k - x^{k-1}$ and $y^{k-1} := g^k - g^{k-1}$.
- BB: choose α^k so that $D^k := \alpha^k I$ satisfies $D^k y^{k-1} = s^{k-1}$:

$$\alpha^k := \frac{(s^{k-1})^\top s^{k-1}}{(s^{k-1})^\top y^{k-1}} \text{ or } \alpha^k := \frac{(s^{k-1})^\top y^{k-1}}{(y^{k-1})^\top y^{k-1}}$$

3.2 随机梯度法(SGD)的收敛性

本节主要证明SGD的收敛性，先回顾一下基本公式：

- ERM(Empirical Risk Minimization) problem

$$\min_{x \in R^n} f(x) = \frac{1}{n} \sum_1^n f_i(x)$$

- 前面介绍的Gradient Descent如下：

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

- 本节的SGD如下：

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

其中的 s_k 表示从1,...n中采样。

和前面一样，在证明收敛性之前，给出四个基本假设方便后面证明。

- $f(x)$ is L-smooth.

$$\|\nabla f(x) - \nabla f(y)\|_2^2 \leq L\|x - y\|_2^2$$

- $f(x)$ is μ -strongly convex .

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \mu\|x - y\|_2^2$$

- $E_s[\nabla f_s(x)] = \nabla f(x)$

随机取的过程长远来看相当于按顺序取。

- $E_s\|\nabla f_s(x)\|^2 \leq M$

随机取的 x 的二阶导数不能无限大，有一个上界M。

证明会用到以下信息：

(1). α_k 为固定长度 α

(2). $E[X] = E[E[X|Y]]$

(3).strong monotonicity (coercivity) of ∇f : $(\nabla f(x) - \nabla f(y))^T(x - y) \geq \mu\|x - y\|_2^2, \forall x, y \in \text{dom } f$

在SGD收敛性证明中可以得到 $\mu\nabla_k^2 \leq \langle \nabla f(x_k), x_k - x^* \rangle$

- $$\begin{aligned}\Delta_{k+1}^2 &= \|x_{k+1} - x^*\|_2^2 = \|x_k - \alpha_k \nabla f_{s_k}(x_k) - x^*\|_2^2 \\ &= \|x_k - x^*\|_2^2 - 2\alpha_k \langle \nabla f_{s_k}(x_k), x_k - x^* \rangle + \alpha_k^2 \|\nabla f_{s_k}(x_k)\|_2^2 \\ &= \Delta_k^2 - 2\alpha_k \langle \nabla f_{s_k}(x_k), x_k - x^* \rangle + \alpha_k^2 \|\nabla f_{s_k}(x_k)\|_2^2\end{aligned}$$

- Using $E[X] = E[E[X|Y]]$:

$$\begin{aligned}\mathbb{E}_{s_1, \dots, s_k} [\langle \nabla f_{s_k}(x_k), x_k - x^* \rangle] &= \mathbb{E}_{s_1, \dots, s_{k-1}} [\mathbb{E}_{s_k} [\langle \nabla f_{s_k}(x_k), x_k - x^* \rangle]] \\ &= \mathbb{E}_{s_1, \dots, s_{k-1}} [\langle \mathbb{E}_{s_k} [\nabla f_{s_k}(x_k)], x_k - x^* \rangle] \\ &= \mathbb{E}_{s_1, \dots, s_{k-1}} [\langle \nabla f(x_k), x_k - x^* \rangle] \\ &= \mathbb{E}_{s_1, \dots, s_k} [\langle \nabla f(x_k), x_k - x^* \rangle]\end{aligned}$$

- By the strongly convexity

$$\mathbb{E}_{s_1, \dots, s_k} (\Delta_{k+1}^2) \leq (1 - 2\alpha\mu) \mathbb{E}_{s_1, \dots, s_k} (\Delta_k^2) + \alpha^2 M^2 \quad (49)$$

上面就得到了一个 ∇_{k+1}^2 的递归关系，并且利用 $0 \leq 2\alpha\mu \leq 1$ 的假设可以得到如下结果：

Proof of Theorem

- Taking induction from $k = 1$ to $k = T$, we have

$$\mathbb{E}_{s_1, \dots, s_T} (\Delta_{T+1}^2) \leq (1 - 2\alpha\mu)^T \Delta_1^2 + \sum_{i=0}^{T-1} (1 - 2\alpha\mu)^i \alpha^2 M^2 \quad (50)$$

- under the assumption that $0 \leq 2\alpha\mu \leq 1$, we have

$$\sum_{i=0}^{\infty} (1 - 2\alpha\mu)^i = \frac{1}{2\alpha\mu}$$

- Then

$$\mathbb{E}_{s_1, \dots, s_T} (\Delta_{T+1}^2) \leq (1 - 2\alpha\mu)^T \Delta_1^2 + \frac{\alpha M^2}{2\mu} \quad (51)$$

4 Variance Reduction

首先指出几个assumption, 只有在这些假设下, 下面的结论才会成立, 当然这些假设都是很合理的:

- $f(x)$ is L-smooth
- $f(x)$ is μ -strongly convex
- $E_s[\nabla f_s(x)] = \nabla f(x)$
- 其中 s 代表随机采样, 本条件的含义是随机采样的梯度的期望和不采用随机采样的梯度一样。
- $E_s\|\nabla f_s(x)\|^2 \leq M^2$
- GD 有线性(linear)收敛速度 $o(n \log(\frac{1}{\epsilon}))$
- GD 有次线性(sublinear)收敛速度 $o(\frac{1}{\epsilon})$

4.1 回顾GD,SGD

Gradient Descend:

$$\begin{aligned}
 \Delta_{k+1}^2 &= \|x_{k+1} - x^*\|_2^2 = \|x_k - \alpha \nabla f(x_k) - x^*\|_2^2 \\
 &= \Delta_k^2 - 2\alpha \langle \nabla f(x_k), x_k - x^* \rangle + \alpha^2 \|\nabla f(x_k)\|_2^2 \\
 &\leq (1 - 2\alpha\mu) \Delta_k^2 + \alpha^2 \|\nabla f(x_k)\|_2^2 (\mu - \text{strongly convex}) \\
 &\quad \text{TODO} \\
 &\leq (1 - 2\alpha\mu + \alpha^2 L^2) \Delta_k^2 (L - \text{smooth})
 \end{aligned}$$

Stochastic Gradient Descend:

$$\begin{aligned}
 E \Delta_{k+1}^2 &= E[\|x_k - \alpha \nabla f_{s_k}(x_k) - x^*\|_2^2] \\
 &= E \Delta_k^2 - 2\alpha E \langle \nabla f_{s_k}(x_k), x_k - x^* \rangle + \alpha^2 E \|\nabla f_{s_k}(x_k)\|_2^2 \\
 &= E \Delta_k^2 - 2\alpha E \langle \nabla f(x_k), x_k - x^* \rangle + \alpha^2 E \|\nabla f_{s_k}(x_k)\|_2^2
 \end{aligned}$$

$$\mathbb{E}\Delta_{k+1}^2 \leq \underbrace{(1 - 2\alpha\mu + 2\alpha^2 L^2)\mathbb{E}\Delta_k^2}_A + \underbrace{2\alpha^2 \mathbb{E} \|\nabla f_{s_k}(x_k) - \nabla f(x_k)\|_2^2}_B \quad (54)$$

- a worst case convergence rate of $\sim 1/T$ for SGD
- In practice, the actual convergence rate may be somewhat better than this bound.
- Initially, $B \ll A$ and we observe the linear rate regime, once $B > A$ we observe $1/T$ rate.
- How to reduce variance term B to speed up SGD?
 - SAG (Stochastic average gradient)
 - SAGA
 - SVRG (Stochastic variance reduced gradient)

4.2 SAG,SAGA,SVRG

SAG[2]与SAGA的算法如下：

SAG method

- SAG method (Le Roux, Schmidt, Bach 2012)

$$x_{k+1} = x_k - \frac{\alpha_k}{n} \sum_{i=1}^n g_k^i = x_k - \alpha_k \left(\frac{1}{n} (\nabla f_{s_k}(x_k) - g_{k-1}^{s_k}) + \frac{1}{n} \sum_{i=1}^n g_{k-1}^i \right) \quad (55)$$

where

$$g_k^i = \begin{cases} \nabla f_i(x_k) & \text{if } i = s_k, \\ \nabla g_{k-1}^i & \text{o.w.,} \end{cases} \quad (56)$$

and s_k is uniformly sampled from $\{1, \dots, n\}$

- complexity(# component gradient evaluations): $O(\max\{n, \frac{L}{\mu}\} \log(1/\epsilon))$
- need to store most recent gradient of each component.
- SAGA(Defazio, Bach,Julien 2014) is unbaised revision of SAG

$$x_{k+1} = x_k - \alpha_k (\nabla f_{i_k}(x_k) - g_{k-1}^{i_k} + \frac{1}{n} \sum_{i=1}^n g_{k-1}^i) \quad (57)$$

注意上面的n是指的全部样本，这样对于每个样本都需要保存整个系统梯度的副本来计算 $\frac{1}{n} \sum_{i=1}^n g_{k-1}^i$ ，会消耗巨大的空间，因此一般会采用mini-batch版本

的SAG算法。就算是这样也会需要保存每个batch的系统梯度的副本,这对于空间同样是一个很大的开销,最终占用的空间如下: $\frac{\text{total_simple_amount} * \text{system_gradient_amount}}{\text{batch_size}}$,SAGA算法同样会有这样的缺点。

下面来看变种SVRG[1]的算法:

```

Procedure SVRG
Parameters update frequency  $m$  and learning rate  $\eta$ 
Initialize  $\tilde{w}_0$ 
Iterate: for  $s = 1, 2, \dots$ 
     $\tilde{w} = \tilde{w}_{s-1}$ 
     $\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \nabla \psi_i(\tilde{w})$ 
     $w_0 = \tilde{w}$ 
    Iterate: for  $t = 1, 2, \dots, m$ 
        Randomly pick  $i_t \in \{1, \dots, n\}$  and update weight
         $w_t = w_{t-1} - \eta (\nabla \psi_{i_t}(w_{t-1}) - \nabla \psi_{i_t}(\tilde{w}) + \tilde{\mu})$ 
    end
    option I: set  $\tilde{w}_s = w_m$ 
    option II: set  $\tilde{w}_s = w_t$  for randomly chosen  $t \in \{0, \dots, m-1\}$ 
end

```

上面的n同样是指样本总数,而现实中一般都是使用的mini-batch的梯度下降算法,需要做的改动就是将 $\nabla \psi_i(\tilde{w})$ 由单个样本计算改为有一个batch的样本来计算。因此 $\tilde{\mu}$ 其实就是遍历了整个数据集。相当于对于每个s,都需要遍历整个数据集,因此速度很慢,但是一个s过后AVRG往往能够给出很好的结果。

SVRG的收敛速度确实很快,但是在现在主流的框架tensorflow,caffe,pytorch中都没有见到它的身影,本人也尝试将SVRG算法融入到这些框架中,但是发现有以下几点导致SVRG无法融入主流框架:

- 主流框架里面的梯度下降方法都是先把梯度算好以后,optimizer里面实现一些优化算法(ADAGrad,Momentum等)的逻辑。但是SVRG不一样,对于每一次求解一个新的梯度值,SVRG需要遍历整个数据集求解梯度,并且梯度的求解是一个不断迭代(迭代m步)的过程。SVRG算法的总体流程(求梯度和执行优化融为一体)和其他普通算法的流程(求梯度和执行优化完全隔绝)有很大区别。
- SVRG中的 $\nabla \psi_{i_t}(\omega_{t-1})$ 的求解比较复杂:需要计算系统在参数为 ω_{t-1} 的情况下新的梯度值,而 ω_{t-1} 对于每个s会有m个不同的取值,这是一个很复杂的流程。
- SVRG,SAG算法都假定数据集是有限数据集的。而在计算机视觉等领域,一般会有一个数据增强(image augmentation)的操作(random crop,flip,random resize etc.),这样会间接导致数据集变成一个无限数目的数据集,从而使得SVRG,SAG算法毫无用武之地。
- SAG算法及其变种的巨大内存开销也是一个主要原因。

SAG,SAGA,SVRG这类算法在实现过程中有一个需要注意的地方就是在每个epoch(每个s)开始的时候，不要重新shuffle数据集，否则就导致每个batch的梯度的对应关系就乱了。

4.3 SVRG算法理论证明

http://ranger.uta.edu/~heng/CSE6389_15_slides/SGD2.pdf

Theorem

Consider SVRG with option II. Assume that all $\psi_i(\omega)$ are convex and smooth, $P(\omega)$ is strongly convex. Let $\omega_* = \operatorname{argmin}_\omega P(\omega)$. Assume that m is sufficiently large so that

$$\alpha = \frac{1}{\gamma\eta(1-2L\eta)m} + \frac{2L\eta}{1-2L\eta} < 1$$

then we have geometric convergence in expectations for SVRG

$$\mathbb{E}P(\tilde{\omega}_s) \leq \mathbb{E}P(\tilde{\omega}_*) + \alpha^s[P(\tilde{\omega}_0) - P(\omega_*)]$$

Given any i, consider

$$g_i(\omega) = \psi_i(\omega) - \psi_i(\omega_*) - \nabla\psi_i(\omega_*)^T(\omega - \omega_*) \quad (13)$$

where $g_i(\omega_*) = \operatorname{argmin}_\omega g_i(\omega)$ and $\nabla g_i(\omega_*) = 0$

$$\begin{aligned} 0 &= g_i(\omega_*) \leq \min_\eta [g_i(\omega - \eta \nabla g_i(\omega))] \\ &\leq \min_\eta [g_i(\omega) - \eta \|\nabla g_i(\omega)\|^2 + 0.5L\eta^2 \|\nabla g_i(\omega)\|^2] \end{aligned} \quad (14)$$

Here it uses a well-known inequality for a function with $1/L$ -Lipschitz continuous gradient

$$f(x) - f(y) - \langle \nabla f(y), x - y \rangle \leq \frac{L}{2} \|x - y\|^2$$

From (14), we can get $\eta = 1/L$, then

$$0 = g_i(\omega_*) \leq g_i(\omega) - \frac{1}{2L} \|\nabla g_i(\omega)\|^2 \quad (15)$$

It can be rewrite as

$$\|\nabla g_i(\omega)\|^2 \leq 2Lg_i(\omega) \quad (16)$$

using the definition of $g_i(\omega)$ and $\nabla g_i(\omega) = \nabla \psi_i(\omega) - \nabla \psi_i(\omega_*)$, the (16) will be

$$\|\nabla \psi_i(\omega) - \nabla \psi_i(\omega_*)\|^2 \leq 2L[\psi_i(\omega) - \psi_i(\omega_*) - \nabla \psi_i(\omega_*)^T(\omega - \omega_*)] \quad (17)$$

By summing the inequality (17) over $i = \{1, \dots, n\}$, the fact that $P(\omega) = \frac{1}{n} \sum_{i=1}^n \nabla \psi_i(\omega)$ and $\nabla P(\omega_*) = 0$, we can get

$$n^{-1} \sum_{i=1}^n \|\nabla \psi_i(\omega) - \nabla \psi_i(\omega_*)\|^2 \leq 2L[P(\omega) - P(\omega_*)] \quad (18)$$

Use $\tilde{\mu} = \nabla P(\tilde{\omega})$ and let $v_t = \nabla \psi_{i_t}(\omega_{t-1}) - \nabla \psi_{i_t}(\tilde{\omega}) + \tilde{\mu}$, v_t is the approximate gradient of SVRG.

With respect to i_t , expectation can be obtained as

$$\begin{aligned} \mathbb{E}\|v_t\|^2 &= \mathbb{E}\|\nabla \psi_{i_t}(\omega_{t-1}) - \nabla \psi_{i_t}(\tilde{\omega}) + \tilde{\mu}\|^2 \\ &\leq 2\mathbb{E}\|\nabla \psi_{i_t}(\omega_{t-1}) - \nabla \psi_{i_t}(\omega_*)\|^2 + 2\mathbb{E}\|\nabla \psi_{i_t}(\tilde{\omega}) - \nabla \psi_{i_t}(\omega_*)\|^2 \\ &= 2\mathbb{E}\|\nabla \psi_{i_t}(\omega_{t-1}) - \nabla \psi_{i_t}(\omega_*)\|^2 + 2\mathbb{E}\|\nabla \psi_{i_t}(\tilde{\omega}) - \nabla \psi_{i_t}(\omega_*)\|^2 \\ &\quad - \mathbb{E}\|\nabla \psi_{i_t}(\tilde{\omega}) - \nabla \psi_{i_t}(\omega_*)\|^2 \\ &\leq 2\mathbb{E}\|\nabla \psi_{i_t}(\omega_{t-1}) - \nabla \psi_{i_t}(\omega_*)\|^2 + 2\mathbb{E}\|\nabla \psi_{i_t}(\tilde{\omega}) - \nabla \psi_{i_t}(\omega_*)\|^2 \\ &\leq 4L[P(\omega_{t-1}) - P(\omega_*) + P(\tilde{\omega}) - P(\omega_*)] \end{aligned} \quad (19)$$

The first inequality uses $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$

The second inequality uses $\mathbb{E}\|X - \mathbb{E}X\|^2 \leq \mathbb{E}\|X\|^2$

The third one uses (18)

The update form of SVRG is $\omega_t = \omega_{t-1} - \eta v_t$, conditioned on ω_{t-1}

$$\begin{aligned}\mathbb{E}\|\omega_t - \omega_*\|^2 &= \mathbb{E}\|\omega_{t-1} - \omega_* - \eta v_t\|^2 \\ &= \|\omega_{t-1} - \omega_*\|^2 - 2\eta(\omega_{t-1} - \omega_*)^\top \mathbb{E}v_t + \eta^2 \mathbb{E}\|v_t\|^2\end{aligned}$$

Here $\mathbb{E}v_t = \mathbb{E}[\nabla\psi_{i_t}(\omega_{t-1}) - \nabla\psi_{i_t}(\tilde{\omega}) + \tilde{\mu}] = \nabla P(\omega_{t-1})$
Using (19) then we can get

$$\begin{aligned}\mathbb{E}\|\omega_t - \omega_*\|^2 &\leq \|\omega_{t-1} - \omega_*\|^2 - 2\eta(\omega_{t-1} - \omega_*)^\top \nabla P(\omega_{t-1}) \\ &\quad + 4L\eta^2[P(\omega_{t-1}) - P(\omega_*) + P(\tilde{\omega}) - P(\omega_*)]\end{aligned}\tag{20}$$

By convexity of $P(\omega)$ that

$$-(\omega_{t-1} - \omega_*)^\top \nabla P(\omega_{t-1}) \leq P(\omega_*) - P(\omega_{t-1})\tag{21}$$

$$\begin{aligned}\mathbb{E}\|\omega_t - \omega_*\|^2 &\leq \|\omega_{t-1} - \omega_*\|^2 - 2\eta[P(\omega_*) - P(\omega_{t-1})] \\ &\quad + 4L\eta^2[P(\omega_{t-1}) - P(\omega_*) + P(\tilde{\omega}) - P(\omega_*)] \\ &= \|\omega_{t-1} - \omega_*\|^2 - 2\eta(1 - 2L\eta)[P(\omega_{t-1}) - P(\omega_*)] + 4L\eta^2[P(\tilde{\omega}) - P(\omega_*)]\end{aligned}\tag{22}$$

In each fixed stage s , $\tilde{\omega} = \tilde{\omega}_{s-1}$ and $\tilde{\omega}_s$ is selected after all updates have completed. By summing the inequality over $t = 1, \dots, m$, taking expectation with all the history

$$\begin{aligned}\mathbb{E}\|\omega_m - \omega_*\|^2 + 2\eta(1 - 2L\eta)m\mathbb{E}[P(\tilde{\omega}_s) - P(\omega_*)] &\leq \mathbb{E}\|\tilde{\omega} - \omega_*\|^2 + 4Lm\eta^2\mathbb{E}[P(\tilde{\omega}) - P(\omega_*)] \\ &\leq \frac{2}{\gamma}\mathbb{E}[P(\tilde{\omega}) - P(\omega_*)] + 4Lm\eta^2\mathbb{E}[P(\tilde{\omega}) - P(\omega_*)]\end{aligned}\tag{23}$$

From above inequality, we can have

$$\begin{aligned} & 2\eta(1 - 2L\eta)m\mathbb{E}[P(\tilde{\omega}_s) - P(\omega_*)] \\ & \leq \frac{2}{\gamma}\mathbb{E}[P(\tilde{\omega}) - P(\omega_*)] + 4Lm\eta^2\mathbb{E}[P(\tilde{\omega}) - P(\omega_*)] \end{aligned} \quad (24)$$

which can be also rewrite as

$$\mathbb{E}[P(\tilde{\omega}_s) - P(\tilde{\omega}_*)] \leq \alpha\mathbb{E}[P(\tilde{\omega}_{s-1}) - P(\omega_*)] \quad (25)$$

where

$$\alpha = \frac{1}{\gamma\eta(1 - 2L\eta)m} + \frac{2L\eta}{1 - 2L\eta} \quad (26)$$

5 随机优化算法在深度学习中的应用

5.1 FG与SG

FG:Full Gradient.即遍历整个数据集，显然不够实用。

$$x^{k+1} = x^k - \alpha_k g'(x^k) = x^k - \frac{\alpha_k}{n} \sum_{i=1}^n f_i'(x^k)$$

SG:Stochastic Gradient,随机遍历某些数据集，比较实用。

$$x^{k+1} = x^k - \alpha_k f_{i_k}'(x^k)$$

上面是标准的SGD,即每一次针对一个样本，这样由导致结果比较随机，因此提出mini-batch SGD(batch size为m):

$$x^{k+1} = x^k - \alpha_k g'(x^k) = x^k - \frac{\alpha_k}{m} \sum_{i=1}^m f_i'(x^k)$$

5.2 常见算法

5.2.1 SGD with momentum

$$\begin{aligned} v^{(t+1)} &= \mu^{(t)}v^{(t)} - \alpha^{(t)}\nabla f_i(x^{(t)}) \\ x^{(t+1)} &= x^{(t)} + v^{(t+1)} \end{aligned}$$

综合起来就是:

$$x^{(t+1)} = x^{(t)} - \alpha^{(t)}\nabla f_i(x^{(t)}) + \mu^{(t)}v^{(t)}$$

对于每一个变量需要同时也开辟一个空间保存动量 $v^{(t)}$

5.2.2 Nesterov accelerated gradient (original version)

$$v^{(t+1)} = (1 + \mu^{(t)})x^{(t)} - \mu^{(t)}x^{(t-1)}$$

$$x^{(t+1)} = v^{(t+1)} - \alpha^{(t)}\nabla f_i(v^{(t+1)})$$

其中 $\mu^{(t)} = \frac{t+2}{t+5}\alpha^{(t)}$ 可以通过 line search 得到。
对于每一个变量需要同时开辟一个空间保存 $x^{(t-1)}$, 然后
在 $x^{(t)}$ 上 in-place 计算得到 $x^{(t+1)}$

5.2.3 Nesterov accelerated gradient (momentum version)

$$v^{(t+1)} = \mu^{(t)}v^{(t)} - \alpha^{(t)}\nabla f_i(x^{(t)} + \mu^t v^{(t)})$$

$$x^{(t+1)} = x^{(t)} + v^{(t+1)}$$

对于每一个变量需要同时开辟一个空间保存动量 $v^{(t)}$

5.2.4 Adaptive Subgradient Methods(Adagrad)

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ
Require: Initial parameter θ
Require: Small constant δ , perhaps 10^{-7} , for numerical stability
 Initialize gradient accumulation variable $r = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$
 Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta\theta$
end while

5.2.5 RMSprop

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ .
Require: Initial parameter θ
Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers.
 Initialize accumulation variables $r = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$
 Compute parameter update: $\Delta\theta = -\frac{\epsilon}{\sqrt{\delta+r}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta+r}}$ applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta\theta$
end while

Algorithm 8.6 RMSProp algorithm with Nesterov momentum

Require: Global learning rate ϵ , decay rate ρ , momentum coefficient α .
Require: Initial parameter θ , initial velocity v .
 Initialize accumulation variable $r = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.
 Compute interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$
 Accumulate gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$
 Compute velocity update: $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{r}}$ applied element-wise)
 Apply update: $\theta \leftarrow \theta + v$
end while

5.2.6 Adam

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)
Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1]$.
 (Suggested defaults: 0.9 and 0.999 respectively)
Require: Small constant δ used for numerical stabilization. (Suggested default:
 10^{-8})
Require: Initial parameters θ
 Initialize 1st and 2nd moment variables $s = \mathbf{0}$, $r = \mathbf{0}$
 Initialize time step $t = 0$
 while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with
 corresponding targets $\mathbf{y}^{(i)}$.
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$
 $t \leftarrow t + 1$
 Update biased first moment estimate: $\hat{s} \leftarrow \rho_1 s + (1 - \rho_1) \mathbf{g}$
 Update biased second moment estimate: $\hat{r} \leftarrow \rho_2 r + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$
 Correct bias in first moment: $\hat{s} \leftarrow \frac{\hat{s}}{1 - \rho_1^t}$
 Correct bias in second moment: $\hat{r} \leftarrow \frac{\hat{r}}{1 - \rho_2^t}$
 Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$ (operations applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta\theta$
 end while

6 总结

7 附录

7.1 基本性质

一些有用的链接:

https://www.cs.ubc.ca/~schmidtm/Documents/2013_Notes_ConvexOptim.pdf
<http://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>
<http://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>
<http://bicmr.pku.edu.cn/~wenzw/opt2015/lect-gm.pdf>
<http://bicmr.pku.edu.cn/~wenzw/bigdata/lect-sto.pdf>

- convex function

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \forall \lambda \in [0, 1], x, y$$

上面的式子也称为jensen不等式，
对于凸函数一阶可微的情况下等价于：

$$f(y) \geq f(x) + \nabla f(x)(y - x)$$

将 $f(y)$ 在 x 处二阶展开,可以得到如下结果:

$$f(y) = f(x) + \nabla f(x)(y - x) + \frac{\nabla^2 f(x)^2}{2\beta^2}(y - x)^2$$

很显然有:

$$f(y) \geq f(x) + \nabla f(x)(y - x)$$

对于凸函数二阶可微的情况下等价于：

$$\nabla^2 f(x) \geq 0$$

凸函数梯度存在单调性(monotonicity):

a differentiable function f is convex if and only if $\text{dom } f$ is convex and
 $(f(x) - f(y))^T(x - y) \geq 0$,for all $x,y \in \text{dom } f, x \neq y$

注意“单调性,定义域集合为凸”和“函数 f 为凸”是等价的。

如果 $f(x)$ 为凸函数,那么

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T(y - x) \\ f(x) &\geq f(y) + \nabla f(y)^T(x - y) \end{aligned}$$

结合两者即得证。

- M-Lipschitz Continuous

$$\|f(x) - f(y)\| \leq M\|x - y\|$$

M-Lipschitz Continuous有如下性质(利用梯度的定义即可得证):

$$\|\nabla f(x)\| \leq M$$

- L-Lipschitz Smoothness 也称为L-Lipschitz continuous gradient, 意思是 f 的梯度满足L-Lipchitz Continuous

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

此处不要求f为凸函数

L-Lipschitz Smoothness有以下性质：

- (1). $\nabla^2 f(x) \leq LI$
- (2). $\frac{L}{2}x^T x - f(x)$ 为凸函数。
- (3). $f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2}\|x - y\|_2^2$
- (4). $\frac{1}{2L}\|\nabla f(x)\|_2^2 \leq f(x) - f(x^*) \leq \frac{L}{2}\|x - x^*\|_2^2$

(1).

利用导数的定义很容易证明。

(2).

利用 ∇f 的 Lipschitz Continuity:

$$\|\nabla f(x) - \nabla f(y)\| \|x - y\| \leq L\|x - y\|_2^2$$

利用柯西不等式可得：

$$(\nabla f(x) - \nabla f(y))^T (x - y) \leq \|\nabla f(x) - \nabla f(y)\| \|x - y\| \leq L\|x - y\|_2^2$$

移到同一边，整理得：

$$(\nabla f(x) - Lx - \nabla f(y) + Ly)^T (x - y) \leq 0$$

换号：

$$(Lx - \nabla f(x) - (Ly - \nabla f(y)))^T (x - y) \leq 0$$

其中 $Lx - \nabla f(x)$ 是 $\frac{L}{2}x^T x - f(x)$ 的梯度。

且 $\frac{L}{2}x^T x - f(x)$ 的 dom 是 convex set，由于 gradient monotonicity 的等价性知：

$\frac{L}{2}x^T x - f(x)$ 是凸函数。

(3).

将 $f(y)$ 在 x 出泰勒展开：

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(z)(y - x)$$

并且由于 $\nabla^2 f(z) \leq LI$:

$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{L}{2}\|y - x\|_2^2$$

此处的 $f(x) + \nabla f(x)^T (y - x) + \frac{L}{2}\|y - x\|_2^2$ 是 $f(y)$ 的二阶上界 (quadratic upper bound)，是关于 y 的二次函数。

(4). 首先证明 $f(x) - f(x^*) \leq \frac{L}{2}\|x - x^*\|_2^2$:

由于 x^* 为 f 的极值点，所以 $f(x^*) = 0$ ，并且：

$$f(x) \leq f(x^*) + \nabla f(x^*)^T (x - x^*) + \frac{L}{2}\|x - x^*\|_2^2$$

化简得证，接下来证明 $\frac{1}{2L}\|\nabla f(x)\|_2^2 \leq f(x) - f(x^*)$

因为 x^* 为极值点，所以 $f(x^*) \leq f(y)$:

$$f(x^*) \leq \inf_{y \in \text{dom } f} (f(x) + \nabla f(x)^T (y - x) + \frac{L}{2}\|y - x\|_2^2)$$

因为 f 的定义域为 R^n ，这里不妨令 $y = x - \frac{1}{L}\nabla f(x)$

化简整理即可证明左边不等式。

L-Lipschitz Continuous 相当于确定了 f 的二阶上界。

- μ -strongly convex characterization

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\mu}{2}\|y - x\|_2^2, \text{ for } \forall y, x$$

Properties of Lipschitz-Continuous Gradient

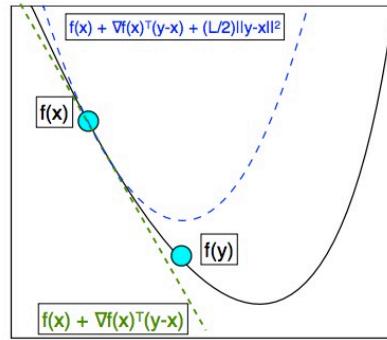
- From Taylor's theorem, for some z we have:

$$f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(z)(y - x)$$

- Use that $\nabla^2 f(z) \preceq L I$.

$$f(y) \leq f(x) + \nabla f(x)^T(y - x) + \frac{L}{2} \|y - x\|^2$$

- Global quadratic upper bound on function value.



同时也等价于：

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y) - \frac{\mu}{2}\lambda(1-\lambda)\|x-y\|_2^2, \forall \lambda \in [0, 1], x, y$$

μ -strongly convex characterization有以下性质：

(1). $\nabla^2 f(x) \geq \mu I$

(2). strong monotonicity (coercivity) of ∇f : $(\nabla f(x) - \nabla f(y))^T(x - y) \geq \mu\|x - y\|_2^2, \forall x, y \in \text{dom } f$

(3). $f(x) - \frac{\mu}{2}x^T x$ 为凸函数.

(4). $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2}\|x - y\|_2^2$

(5). $\frac{\mu}{2}\|x - x^*\|_2^2 \leq f(x) - f(x^*) \leq \frac{1}{2\mu}\|\nabla f(x)\|_2^2$

(1). $f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{\mu}{2}\|y - x\|_2^2$, for $\forall y, x$
同时又有在x点处泰勒展开:

$$f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{\nabla^2 f(z)}{2}\|y - x\|_2^2, \text{for } \forall y, x$$

所以:

$$f(x) + \nabla f(x)^T(y - x) + \frac{\nabla^2 f(z)}{2}\|y - x\|_2^2 \geq f(x) + \nabla f(x)^T(y - x) + \frac{\mu}{2}\|y - x\|_2^2$$

整理即得证。

(2). 由二阶导数定义知:

(3).

(4). 定义

(5). 证法和上面类似。相当于二阶下界:

μ -strongly convex相当于确定了f的二阶下界。

Properties of Strong-Convexity

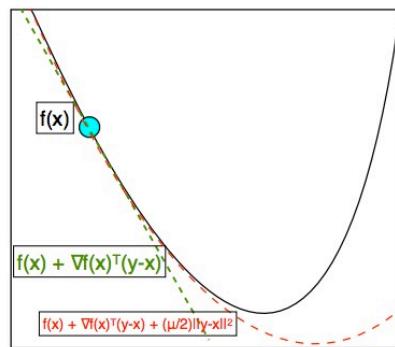
- From Taylor's theorem, for some z we have:

$$f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(z)(y - x)$$

- Use that $\nabla^2 f(z) \succeq \mu I$.

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{\mu}{2}\|y - x\|^2$$

- Global quadratic upper bound on function value.



L-Lipschitz-Smooth 和 μ -strongly convex 共同对应着f的二阶上界和二阶下界，两者之间的性质有遥相呼应的关系。

7.2 Co-coercivity of gradient

<http://www.seas.ucla.edu/~vandenbe/236C/lectures/gradient.pdf>

if f is convex with $\text{dom } f = R^n$, and $\frac{L}{2}x^T x - f(x)$ is convex, then
 $(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{1}{L}\|\nabla f(x) - \nabla f(y)\|_2^2, \forall x, y$
this property is known as co-coercivity of ∇f with parameter $\frac{1}{L}$

co-coercivity 的证明如下:

构造以下两个函数:

$$f_x(z) = f(z) - \nabla f(x)^T z, f_y(z) = f(z) - \nabla f(y)^T z,$$

那么可以得到以下两个结论:

- (1). 通过求导数可以知道: $f_x(z)$ 在 $z=x$ 处取极值, $f_y(z)$ 在 $z=y$ 处取极值.
- (2). $\frac{L}{2}z^T z - f_x(z)$ 为凸函数。 (因为其二阶导数等于 L , 大于 0)

因为 $f_x(z)$ 满足 L-Lipschitz continuity, 且在 x 处取极值, 所以有:
 $f_x(y) - f_x(x) \geq \frac{1}{2L}\|\nabla f_x(y)\|_2^2$

考虑 $f_x(y)$:

$$f(y) - f(x) - \nabla f(x)^T(y - x) = f_x(y) - f_x(x) \geq \frac{1}{2L}\|\nabla f_x(y)\|_2^2 = \frac{1}{2L}\|\nabla f(y) - \nabla f(x)\|_2^2$$

考虑 $f_y(x)$:

$$f(x) - f(y) - \nabla f(y)^T(x - y) = f_y(x) - f_y(y) \geq \frac{1}{2L}\|\nabla f_y(x)\|_2^2 = \frac{1}{2L}\|\nabla f(x) - \nabla f(y)\|_2^2$$

上面两个式子相加即得证。

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{1}{L}\|\nabla f(x) - \nabla f(y)\|_2^2, \forall x, y$$

co-coercivity 之间的等价性:

Lipschitz continuity of $\nabla f(x) \Rightarrow$
convexity of $\frac{L}{2}x^T x - f(x) \Rightarrow$
co-coercivity of $\nabla f(x) \Rightarrow$
Lipschitz continuity of $\nabla f(x)$

下面给出证明:

(1). Lipschitz continuity of $\nabla f(x)$ 到 convexity of $\frac{L}{2}x^T x - f(x)$ 的充分性:

TODO

(2). convexity of $\frac{L}{2}x^T x - f(x)$ 到 co-coercivity of $\nabla f(x)$ 的充分性 上面已经证明。

(3). co-coercivity of $\nabla f(x)$ 到 Lipschitz continuity of $\nabla f(x)$ 的充分性 证明如下:

$$\|\nabla f(x) - \nabla f(y)\|_2^2 \leq L(\nabla f(x) - \nabla f(y))^T(x - y)$$

利用柯西不等式:

$$\|\nabla f(x) - \nabla f(y)\|_2^2 \leq L(\nabla f(x) - \nabla f(y))^T(x - y) \leq L\|\nabla f(x) - \nabla f(y)\|_2\|(x - y)\|_2$$

化简即得证。

上面的等价性刻画co-coercivity of $\nabla f(x)$,Lipschitz continuity of $\nabla f(x)$,convexity of $\frac{L}{2}x^T x - f(x)$ 三者之间的关系。

7.2.1 co-coercivity的扩展

if f is strongly convex and ∇f is Lipschitz Continuous, then

- (1). $g(x) = f(x) - \frac{\mu}{2}x^T x$ is convex.
- (2). ∇g is Lipschitz Continuous with parameter $L - \mu$
- (3). co-coercivity of ∇g gives

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{\mu L}{\mu + L} \|x - y\|_2^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

下面依次给出证明:

(1).

直接对 $g(x)$ 求二阶导数得到 $\nabla^2 g(x) = \nabla^2 f(x) - \mu I$,因此为凸函数。

(2).这里看到网上的一个答案证明的是小于等于 $L + \mu$,

<https://math.stackexchange.com/questions/1645272/extension-of-co-coercivity-in-strongly-convex-functions>
这个上界太宽泛,需要进一步压缩。证明如下:

由 ∇g 的co-coercivity的性质知道:

$$(\nabla g(x) - \nabla g(y))^T(x - y) \geq \frac{1}{L} \|\nabla g(x) - \nabla g(y)\|_2^2, \forall x, y$$

将 $g(x)$ 定义带入不等式,整理得到:

$$(\nabla g(x) - \nabla g(y))^T(x - y) \leq \frac{L\mu + \mu^2}{L+2\mu} \|x - y\|_2^2 + \frac{1}{L+2\mu} \|\nabla f(x) - \nabla f(y)\|_2^2$$

可以发现这是一个比 $\frac{\mu L}{\mu + L} \|x - y\|_2^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|_2^2$ 更紧的上界,所以(3)得证。

7.3 Projection Operator is non-expansive

$$\begin{aligned} & \text{if } \pi_C = \operatorname{argmin}_{y \in C} \|y - x\|_2^2, \\ & \text{then, } \|y_1 - y_2\| \geq \|x_1 - x_2\| \end{aligned}$$

这里的集合C必须是非空凸集。

可以参考:<https://math.stackexchange.com/questions/1426343/prove-that-projection-operator-is-non-expansive>

下面给出具体证明:

首先需要知道projection operator的variational characterization(又叫做Bourbaki-Cheney-Goldstein inequality,具体可以参考Proposition 1.1.9 in the book Convex Optimization Theory by Dimitri Bertsekas):

$$\begin{aligned} & \text{if } \pi_C = \operatorname{argmin}_{y \in C} \|y - x\|_2^2, \\ & \text{then, } \langle x_1 - y_1, x - y_1 \rangle \leq 0 \end{aligned}$$

下面利用projection operator的variational characterization来证明non-expansive:
variational characterization :

$$\langle x_1 - y_1, x - y_1 \rangle \leq 0, \forall x$$

所以可得:

$$\langle x_1 - y_1, y_2 - y_1 \rangle \leq 0$$

同理可得:

$$\langle x_2 - y_2, y_1 - y_2 \rangle \leq 0$$

将上面两个不等式相加, 整理, 并且运用柯西不等式:

$$\begin{aligned} & \langle x_1 - x_2, y_2 - y_1 \rangle + \langle y_2 - y_1, y_2 - y_1 \rangle \leq 0 \\ & \text{then, } \langle y_2 - y_1, y_2 - y_1 \rangle \leq \langle x_2 - x_1, y_2 - y_1 \rangle \\ & \text{then, } \langle y_2 - y_1, y_2 - y_1 \rangle \leq \langle x_2 - x_1, y_2 - y_1 \rangle \leq \|x_2 - x_1\| \|y_2 - y_1\| \\ & \text{then, } \|y_2 - y_1\|^2 \leq \|x_2 - x_1\| \|y_2 - y_1\| \end{aligned}$$

所以有:

$$\|y_2 - y_1\| \leq \|x_2 - x_1\|$$

8 Useful Links

Machine Learning Summer School: http://learning.mpi-sws.org/mlss2016/slides/mlss_2016_cadiz_slides.pdf

References

- [1] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger, editors, *NIPS*, pages 315–323, 2013.
- [2] Nicolas L. Roux, Mark Schmidt, and Francis R. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2663–2671. Curran Associates, Inc., 2012.