# Stochastic Gradient Descent

§1 Classical Gradient Descent

Given $f : \mathbb{R}^d \to \mathbb{R}$ differentiable, convex
Aim: find $\min_w f(w)$, $w \in \mathbb{R}^d$
$\text{grad descent } \nabla f(w) = 0 \Leftrightarrow w \text{ is a min of } f(w)$
Idea:

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

Def (GD)
Input $\eta > 0$, $T \in \mathbb{N}$, $f : \mathbb{R}^d \to \mathbb{R}$
Initialize $w^0 = 0$
For $t = 1 \dots T$
$\quad w^{t} = w^{t-1} - \eta \nabla f(w^{t-1})$
Return $\bar{w} = \frac{1}{T}\sum_{t=1}^{T} w^t$ or best performing $w^t$
ie $\underset{t}{\arg\min} \, f(w^t)$, $t \leq T$

## §2 Subgradient

We want to generalize GD to non-diff functions

Recall $f : \mathbb{R}^d, f : S \to \mathbb{R}$ is convex iff for all $w \in S$, $f$ is differentiable iff $w \in S$

$\forall u \in S: \quad f(u) \geq f(w) + \langle \nabla u - w, \nabla f(w) \rangle$
$\quad\quad\quad$ tangent approx of $f$ at $w$

Lemma: Let $S$ be convex, open. Then a function $f : \mathbb{R}^d \to \mathbb{R}$ is convex
$\Leftrightarrow \forall w \in S \, \exists v \in \mathbb{R}^d \, \forall u \in S: \, f(u) \geq f(w) + \langle u - w, v \rangle$

Def: Let $f : S \to \mathbb{R}$ be a function at $w \in S$ satisfying $\forall u \in S: \, f(u) \geq f(w) + \langle u - w, v \rangle$
at a given point $w \in S$ is called subgradient of $f$ at $w$. The set of all subgradients is denoted $\partial f(w)$

Facts: i) $f$ is convex, then $\partial f(w) \neq \emptyset$ for $w \in S$
ii) $f$ convex and differentiable then $\partial f(w) = \{ \nabla f(w) \}$

Example $\partial f(x) = \begin{cases} \{ f'(x) \} & x \neq 0 \\ [f'_-, f'_+] & x = 0 \end{cases}, x \in \mathbb{R}$



## §3 Lemma 14.1

For later convergence statements we need

Lem 14.1 Let $v_1, \dots, v_T \in \mathbb{R}^d$, then any algorithm with initialization $w = 0$ and update rule
$\quad w^{t+1} = w^t - \eta v_t, \quad \eta > 0$
satisfies for all $w^* \in \mathbb{R}^d$

$$\frac{1}{T}\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2T}\sum_{t=1}^{T}\|v_t\|^2$$

In particular for every $B, \rho > 0$ if $\|v_1\|, \dots, \|v_T\| \leq \rho$ and $w^* \in B_B(0)$ then with $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ we have

$$\frac{1}{T}\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle \leq \frac{B\rho}{\sqrt{T}}$$

Proof: Using the polarization identities
$$\langle x, y \rangle = \frac{1}{2}\left( \|x\|^2 + \|y\|^2 - \|x - y\|^2 \right)$$

$$\Rightarrow \langle w^t - w^*, v_t \rangle = \frac{1}{\eta}\langle w^t - w^*, \eta v_t \rangle$$
$$= \frac{1}{2\eta}\left( -\|w^t - w^* - \eta v_t\|^2 + \|w^t - w^*\|^2 + \eta^2\|v_t\|^2 \right)$$
$$= \frac{1}{2\eta}\left( \|w^t - w^*\|^2 - \|w^{t+1} - w^*\|^2 \right) + \frac{\eta}{2}\|v_t\|^2 \quad (w^{t+1} = w^t - \eta v_t)$$

Summing over $t$ we get
$$\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle = \frac{1}{2\eta}\sum_{t=1}^{T}\left( \|w^t - w^*\|^2 - \|w^{t+1} - w^*\|^2 \right) + \frac{\eta}{2}\sum_{t=1}^{T}\|v_t\|^2$$
$$= \frac{1}{2\eta}\left( \|w^1 - w^*\|^2 - \|w^{T+1} - w^*\|^2 \right) + \frac{\eta}{2}\sum_{t=1}^{T}\|v_t\|^2$$
$$\leq \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2}\sum_{t=1}^{T}\|v_t\|^2$$

## §4 Stochastic Gradient Descent

Def (SGD)
Input $\eta > 0$, $T \in \mathbb{N}$, $f : S \to \mathbb{R}$
Initialize $w^1 = 0$
For $t = 1 \dots T$
$\quad$ Choose $v_t$ according to some prob.
$\quad$ distribution s.t. $\mathbb{E}[v_t | w^t] \in \partial f(w^t)$.
$\quad w^{t+1} = w^t - \eta v_t$
Return $\bar{w} = \frac{1}{T}\sum_{t=1}^{T} w^t$

Thm 14.8 Let $B, \rho > 0$, $f : S \to \mathbb{R}$ convex
Let $w^* \in \arg\min_{w: \|w\| \leq B} f(w)$. Assume that SGD runs over $T$ iterations and uses step size $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ and assume that
$\|v_1\|, \dots, \|v_T\| \leq \rho \quad a.s.$
Then $\mathbb{E}[f(\bar{w})] - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$ Therefore for given $\varepsilon > 0$ to achieve
$$\mathbb{E}[f(\bar{w})] - f(w^*) \leq \varepsilon$$
We need to run $T \geq 2\left(\frac{B\rho}{\varepsilon}\right)^2$ iterations of SGD

Proof: We use $v_{1:T} = v_1, \dots, v_T$. By Jensen's inequality
$$f(\bar{w}) - f(w^*) \leq \frac{1}{T}\sum_{t=1}^{T}\left( f(w^t) - f(w^*) \right)$$
$$\Rightarrow \mathbb{E}_{v_{1:T}}\left[ f(\bar{w}) - f(w^*) \right] \leq \mathbb{E}_{v_{1:T}}\left[ \frac{1}{T}\sum_{t=1}^{T} f(w^t) - f(w^*) \right]$$

Using Lemma 14.1
$$\mathbb{E}_{v_{1:T}}\left[ \frac{1}{T}\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle \right] \leq \frac{B\rho}{\sqrt{T}}$$

This means that it is enough to show
$$\mathbb{E}_{v_{1:T}}\left[ \frac{1}{T}\sum_{t=1}^{T}\left( f(w^t) - f(w^*) \right) \right] \leq \mathbb{E}_{v_{1:T}}\left[ \frac{1}{T}\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle \right]$$

Using linearity of expectation
$$\mathbb{E}_{v_{1:T}}\left[ \frac{1}{T}\sum_{t=1}^{T} \langle w^t - w^*, v_t \rangle \right] = \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{v_{1:T}}\left[ \langle w^t - w^*, v_t \rangle \right]$$

Next recall law of total expectation
for $\alpha, \beta$ random variables, $g$ some function
then $\mathbb{E}_\alpha[g(\alpha)] = \mathbb{E}_\beta[\mathbb{E}_\alpha[g(\alpha) | \beta]]$
Now put $\alpha = v_{1:t}$, $\beta = v_{1:t-1}$, we obtain
$$\mathbb{E}_{v_{1:T}}\left[ \langle w^t - w^*, v_t \rangle \right] = \mathbb{E}_{v_{1:t}}\left[ \langle w^t - w^*, v_t \rangle \right]$$
$$= \mathbb{E}_{v_{1:t-1}}\left[ \mathbb{E}_{v_t}\left[ \langle w^t - w^*, v_t \rangle \,\big|\, v_{1:t-1} \right] \right] = \mathbb{E}_{v_{1:t-1}}\left[ \langle w^t - w^*, \mathbb{E}[v_t | v_{1:t-1}] \rangle \right]$$
But since $w^t$ is determined $w^t = -\eta(v_1 + \dots + v_{t-1})$
it follows
$$\mathbb{E}_{v_t}[v_t | v_{1:t-1}] = \mathbb{E}_{v_t}[v_t | w^t] \in \partial f(w^t)$$
$$\Rightarrow \mathbb{E}_{v_{1:t-1}}\left[ \langle w^t - w^*, \mathbb{E}_{v_t}[v_t | v_{1:t-1}] \rangle \right]$$
$$\geq \mathbb{E}_{v_{1:t-1}}\left[ f(w^t) - f(w^*) \right]$$
$$\Rightarrow \mathbb{E}_{v_{1:T}}\left[ \langle w^t - w^*, v_t \rangle \right]$$
$$\geq \mathbb{E}_{v_{1:T}}\left[ f(w^t) - f(w^*) \right]$$

Now Summing over $t \leq T$ and dividing by $T$ gives the desired inequality.
To achieve
$$\mathbb{E}[f(\bar{w})] - f(w^*) \leq \frac{B\rho}{\sqrt{T}} \leq \varepsilon$$
$$\Rightarrow \quad T \geq 2\left(\frac{B\rho}{\varepsilon}\right)^2 \qquad \square$$

## §5 Stochastic Gradient Descent for Risk Minimization

In Learning Theory we want to minimize
$$L_\mathcal{D}(w) = \mathbb{E}_{z \sim \mathcal{D}}[\ell(w, z)] \quad \leftarrow \text{ loss function}$$

SGD allows us to directly minimize $L_\mathcal{D}$
For simplicity we assume that $\ell(\cdot, z)$ is differentiable for all $z \in Z$
We construct the random $v_t$ as follows
Sample $z \sim \mathcal{D}$
$\Rightarrow \quad v_t = \nabla \ell(w^t, z)$
where the gradient is taken w.r.t. $w$
Interchanging expectation and gradient
$$\mathbb{E}[v_t | w^t] = \mathbb{E}_z[\nabla \ell(w^t, z)]$$
$$= \nabla \mathbb{E}_z[\ell(w^t, z)] = \nabla L_\mathcal{D}(w^t) \in \partial L_\mathcal{D}(w^t)$$
This same argument can be applied to the subgradient case
Let $v_t \in \partial \ell(w^t, z)$ for sample $z \sim \mathcal{D}$
Then by def
for $\quad \ell(u, z) - \ell(w^t, z) \geq \langle u - w^t, v_t \rangle$
$$\Rightarrow \quad L_\mathcal{D}(u) - L_\mathcal{D}(w^t) \geq \mathbb{E}[\langle u - w^t, v_t \rangle | w^t]$$
$$= \langle u - w^t, \mathbb{E}[v_t | w^t] \rangle$$
$$\Rightarrow \mathbb{E}[v_t | w^t] \in \partial L_\mathcal{D}(w^t)$$

Def (SGD for minimizing $L_\mathcal{D}$)
Input $\eta > 0$, $T \in \mathbb{N}$
Initialize $w^1 = 0$
For $t = 1 \dots T$
$\quad$ Sample $z \sim \mathcal{D}$
$\quad$ Pick $v_t \in \partial \ell(w^t, z)$
$\quad w^{t+1} = w^t - \eta v_t$
Return $\bar{w} = \frac{1}{T}\sum_{t=1}^{T} w^t$

By Theorem 14.8

Cor: Let $B, \rho > 0$, $L_\mathcal{D}$ convex such that
$$\|\partial \ell(w^1, z)\|, \dots, \|\partial \ell(w^T, z)\| \leq \rho \quad a.s.$$
Let $w^*$ be a minimizer of $L_\mathcal{D}(w)$. Then to achieve
$$\mathbb{E}[L_\mathcal{D}(\bar{w})] - L_\mathcal{D}(w^*) \leq \varepsilon$$
for given $\varepsilon > 0$ we need to run SGD with stepsize
$\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ and $T \geq 2\left(\frac{B\rho}{\varepsilon}\right)^2$ iterations