

# Gradient descent

## Program for today :

- 1) Convex functions
- 2) Convex functions & Lipschitz continuity
- 3) GD : basic convergence theorem
- 4) Final remarks.

Next Time : Stochastic Gradient Descent.

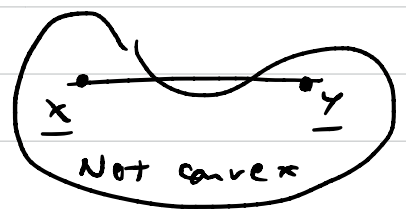
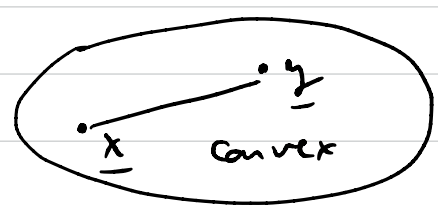
#.

## 1 Convex functions.

### Definition 1 : convex sets.

A set  $S$  is convex if for all  $\underline{x}, \underline{y} \in S$

then for all  $\lambda \in [0, 1]$  :  $\lambda \underline{x} + (1-\lambda)\underline{y} \in S$



(2)

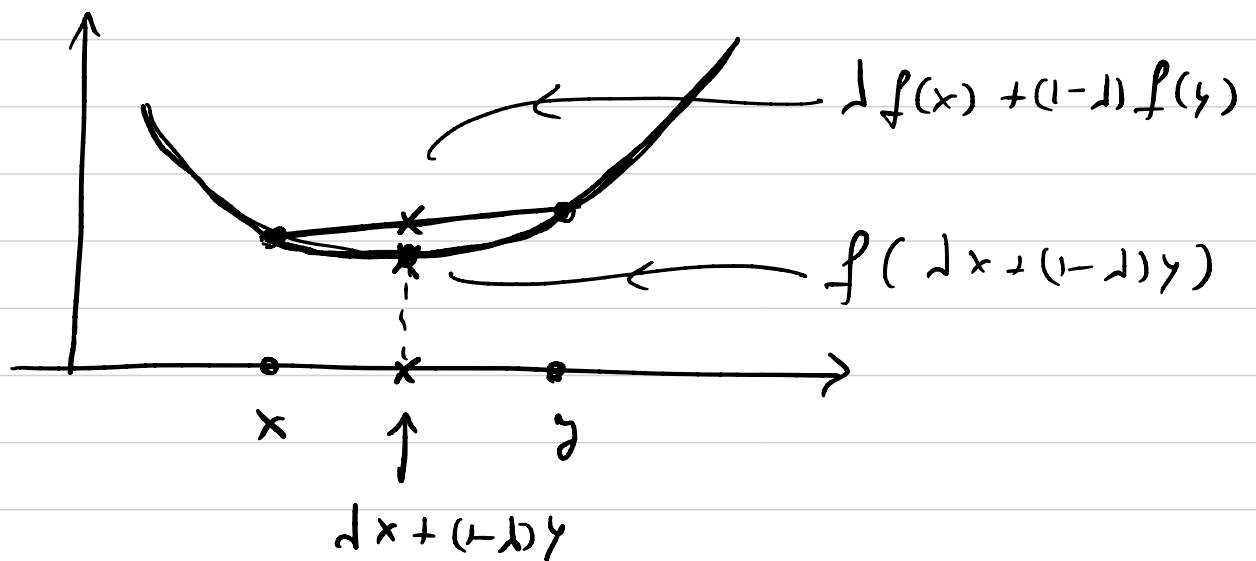
Definition 2; convex functions.

Let  $f: S \rightarrow \mathbb{R}$ ,  $S$  an open convex set.

We say that  $f$  is convex if for all  $\underline{x}, \underline{y} \in S$   
and all  $\lambda \in [0, 1]$ .

$$f(\lambda \underline{x} + (1-\lambda)\underline{y}) \leq \lambda f(\underline{x}) + (1-\lambda)f(\underline{y})$$

Picture for  $S = \mathbb{R}$  or  $S = ]a, b[$



The chord is above the function.

(3)

## Alternative characterisation of a convex function:

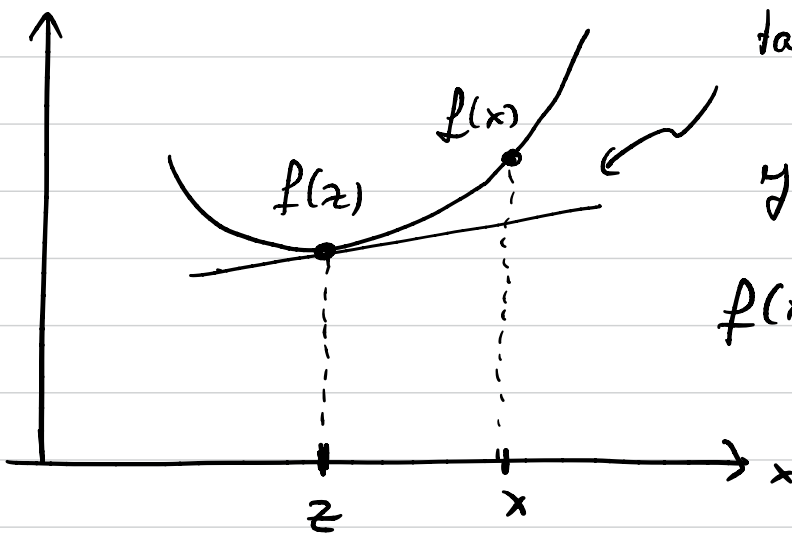
Lemma: Let  $S$  an open convex set and

let  $f: S \rightarrow \mathbb{R}$ . Then  $f$  is convex if

$\forall z \in S$  there exist  $\underline{n}$  such that

$$f(x) \geq f(z) + \langle \underline{n}, x - z \rangle, \quad \forall x \in S$$

Intuition:



tangent line at  $p$ :

$$y = f(z) + f'(z)(x - z)$$

$$f(x) \geq f(z) + f'(z)(x - z)$$

i.e

$$f(x) - f(z) \geq \underbrace{f'(z)}_{\underline{n}} (x - z)$$

If  $f$  is differentiable this intuition can be made rigorous and  $\underline{n} = \nabla f(z)$ , Gradient at  $z$ .

(4)

Proof of Lemma.

Fix  $z \in S$ . Assume  $\exists \underline{v}$  s.t.  $\forall x, y \in S$

$$f(\underline{x}) \geq f(z) + \langle \underline{v}, \underline{x} - z \rangle$$

$$f(\underline{y}) \geq f(z) + \langle \underline{v}, \underline{y} - z \rangle$$

Then  $\forall \lambda \in [0, 1]$  multiply first eqn by  $\lambda$  and second equation by  $1-\lambda$ , and finally sum them:

$$\lambda f(\underline{x}) + (1-\lambda)f(\underline{y}) \geq f(z) + \langle \underline{v}, \lambda \underline{x} + (1-\lambda)\underline{y} - z \rangle$$

Now set  $z = \lambda \underline{x} + (1-\lambda)\underline{y}$ . We get:

$$\lambda f(\underline{x}) + (1-\lambda)f(\underline{y}) \geq f(\lambda \underline{x} + (1-\lambda)\underline{y})$$

which means that  $f$  is convex 

Remark: in the proof  $\underline{v}$  depends on  $z$  only, eventually,

we choose  $z = \lambda \underline{x} + (1-\lambda)\underline{y}$ .

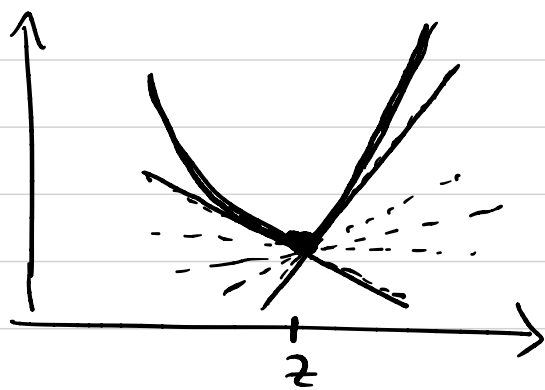
(5)

Converse Lemma: Let  $f: S \rightarrow \mathbb{R}$  and assume it is convex. Then given any  $z \in S$  there exist  $\underline{v}$  such that for all  $x \in S$

$$f(x) \geq f(z) + \langle \underline{v}, (x-z) \rangle.$$

If  $f$  is differentiable at  $z$ , then  $\underline{v}$  is unique (for this  $z$ ) and  $\underline{v} = \nabla f(z)$ .

Remark: When  $f$  is not differentiable at  $z$   $\underline{v}$  still exists as long as  $f$  is convex, however it is non-unique.



All slopes between the left and right derivatives can serve as  $\underline{v}$  and we have

$$f(x) \geq f(z) + \underline{v}(x-z).$$

⑥

This leads us to the definition:

### Definition 3. Subgradient

Any  $\underline{v}_z$  that fulfills the condition

$$f(\underline{x}) \geq f(\underline{z}) + \langle \underline{v}_z, \underline{x} - \underline{z} \rangle$$

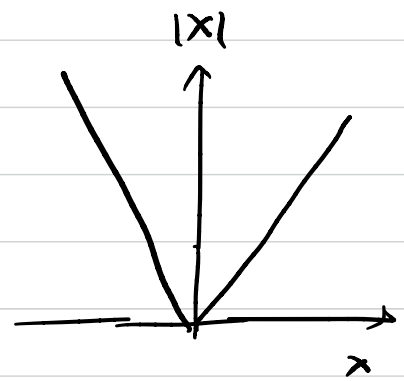
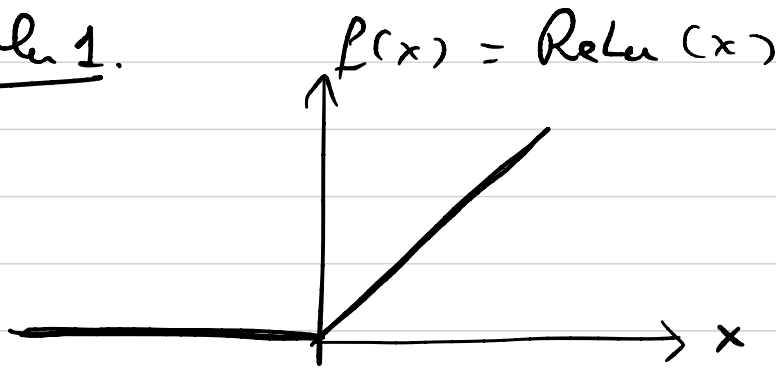
is called a subgradient. The set of all such  $\underline{v}_z$ 's is called the differential set. We denote

$$\text{this by } \underline{v}_z \in \underbrace{\partial f(z)}.$$

notation for differential set

When  $f$  is differentiable at  $z$  we have

$$\partial f(z) = \{ \nabla f(z) \}.$$

Example 1.

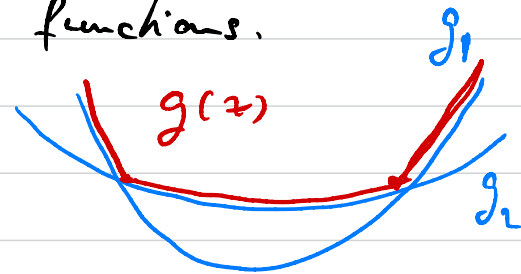
$$\partial f = \begin{cases} \{0\} & \text{if } x < 0 \\ \{+1\} & \text{if } x > 0 \\ [0, 1] & \text{if } x = 0 \end{cases}$$

$$\partial f = \begin{cases} -1, & x < 0 \\ +1, & x > 0 \\ [-1, +1] & x = 0 \end{cases}$$

Example 2 (important and classical).

$g_1, \dots, g_r$  convex differentiable functions.

Let  $g(z) = \max_{i=1, \dots, r} g_i(z)$



Then  $g(z)$  is a convex function

Proof:  $g_j$  are differentiable and convex so:

$$g(x) = \max_i g_i(x) \geq g_j(x) \geq g_j(z) + \langle \nabla g_j(z), x - z \rangle$$

Now take  $j = \arg \max_i g_i(z)$ . Thus  $g_j(z) = g(z)$

$$\Rightarrow g(x) \geq g(z) + \underbrace{\langle \nabla g_j(z), x - z \rangle}_{\substack{\text{with } j = \arg \max_i g_i(z) \\ \text{not unique.}}}$$

8

Proof of concave lemma in differentiable case.

We assume  $f$  is concave. Then

$$f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$$

$$\Rightarrow f(x) \geq \frac{1}{\lambda} [f(\lambda x + (1-\lambda)y) - (1-\lambda)f(y)]$$

$$f(x) \geq f(y) + \frac{1}{\lambda} [f(\lambda x + (1-\lambda)y) - f(y)]$$

$$\forall \lambda \in [0, 1].$$

$$\Rightarrow f(x) \geq f(y) + \lim_{\lambda \rightarrow 0^+} \frac{f(y + \lambda(x-y)) - f(y)}{\lambda}$$

if  $f$  is differentiable the limit is  $\langle \underline{\nabla} f(y), \underline{x-y} \rangle$

(by definition of gradient or by Taylor expansion, ...)

Thus  $\underline{\omega}_y = \underline{\nabla} f(y)$ .





## 2 Lipshitz functions (and convexity)

The Lipshitz condition provides a strong form of continuity. When combined with convexity nice properties emerge.

### Definition 4. $\rho$ -Lipshitz functions

$f$  is  $\rho$ -Lipshitz if  $\forall \underline{x}, \underline{y} \in S$  an open set we have

$$|f(\underline{x}) - f(\underline{y})| \leq \rho \|\underline{x} - \underline{y}\|.$$

Lemma Let  $S$  an open convex set and

$f: S \rightarrow \mathbb{R}$  a convex function,

$f$  is  $\rho$ -Lipshitz if and only if for all  $z \in S$

and  $\underline{v}_z \in \partial f(z)$  we have  $\|\underline{v}_z\| \leq \rho$ .

For differentiable  $f$  in particular  $\|\nabla f(z)\| \leq \rho$ .

Proof

First direction: assume  $\|v_z\| \leq \rho$ .

Since  $f$  is convex and  $v_z$  is a subgradient we

have

$$f(x) \geq f(z) + \langle v_z, x - z \rangle$$

$$\Rightarrow f(z) - f(x) \leq \langle v_z, z - x \rangle$$

$$\leq \|v_z\| \|z - x\| \text{ by Cauchy-Schwarz.}$$

$$\leq \rho \|z - x\|$$

by exchanging  $z$  &  $x$ :  $f(x) - f(z) \leq \rho \|x - z\|$ .

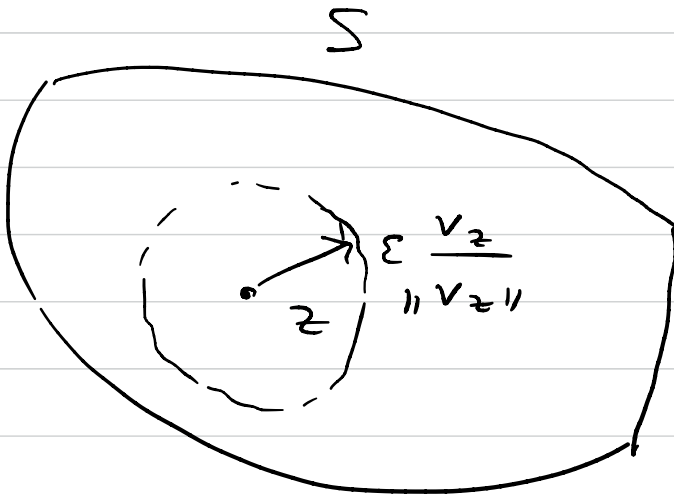
Thus  $|f(x) - f(z)| \leq \rho \|x - z\|$  and  $f$  is

$\rho$ -Lipschitz.

Converse direction: assume  $f$  is  $\rho$ -Lipschitz.

Then  $\forall x, z \in S$   $|f(x) - f(z)| \leq \rho \|x - z\|$

Take  $x = z + \frac{\rho v_z}{\|v_z\|}$ ,  $v_z \in \partial f(z)$   
 (Recall  $f$  is convex so subgradients exist)



By convexity of  $f$ :

$$f(\underline{x}) \geq f(\underline{z}) + \langle \underline{v}_z, \underline{x} - \underline{z} \rangle$$

$$f(\underline{x}) - f(\underline{z}) \geq \underbrace{\langle \underline{v}_z, \frac{\epsilon \underline{v}_z}{\|\underline{v}_z\|} \rangle}_{\epsilon \|\underline{v}_z\|}$$

By  $\epsilon$ -Lipschitzness  $f(\underline{x}) - f(\underline{z}) \leq \epsilon \|\underline{x} - \underline{z}\|$

Thus  $\underbrace{\epsilon \|\underline{x} - \underline{z}\|}_{\epsilon} \geq \epsilon \|\underline{v}_z\|$

$\Rightarrow \|\underline{v}_z\| \leq \epsilon$



### 3 Gradient Descent.

GD is a way to "walk" efficiently through  $S$  in discrete time steps in order to reach minimum of  $f$  in  $S$ .

GD algorithm. Let  $S \ni 0$  open convex set and  $f$  convex let  $f: S \rightarrow \mathbb{R}$ .

a) start at  $w^1 = 0$

b) at steps  $t = 1 \dots T$  do

$$\underline{w}^{t+1} = \underline{w}^t - \gamma \underbrace{\nabla f(\underline{w}^t)}$$

if  $f$  is differentiable and otherwise pick any subgradient. Not here we denote  $\nabla f(\underline{w}^t)$  by slight abuse of notation.

c) Output:  $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^t$ .

Here  $\gamma$  is called the "step size" or "rate".

Theorem: "GD approach to optimal values"

$$\text{Let } \underline{w}^* = \underset{\|\underline{w}\| \leq B}{\operatorname{argmin}} f(\underline{w})$$



the minimizer of  $f$  in a ball  $B \ni 0$ .

Then after  $T$  steps and step size  $\gamma = \sqrt{\frac{B^2}{\rho^2 T}}$

we have:

$$0 \leq f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

Remark:

If we ask for  $0 \leq f(\bar{w}) - f(w^*) \leq \epsilon$

it is enough to take  $\frac{B\rho}{\sqrt{T}} \leq \epsilon \Rightarrow T \geq \frac{B^2\rho^2}{\epsilon^2}$

$$\text{and } \gamma = \sqrt{\frac{B^2}{\rho^2 T}} \leq \sqrt{\frac{B^2}{\rho^2 \frac{B^2\rho^2}{\epsilon^2}}} = \frac{\epsilon}{\rho^2}$$

$$T \geq \frac{B^2\rho^2}{\epsilon^2} \quad \& \quad \gamma \leq \frac{\epsilon}{\rho^2} \quad \text{to get } \epsilon \text{ close to Min.}$$

Proof of Theorem.

$$f(\bar{w}) - f(w^*) = f\left(\frac{1}{T} \sum_{t=1}^T w^t\right) - f(w^*)$$

$$\stackrel{\text{convexity}}{\leq} \frac{1}{T} \sum_{t=1}^T f(w^t) - f(w^*)$$

$$= \frac{1}{T} \sum_{t=1}^T (f(w^t) - f(w^*))$$

$$\stackrel{\text{convexity}}{\leq} \frac{1}{T} \sum_{t=1}^T \langle \nabla f(w^t), w^t - w^* \rangle$$

$$= \frac{1}{2\gamma T} \sum_{t=1}^T \langle \gamma \nabla f(w^t), w^t - w^* \rangle$$

$$= \frac{1}{2\gamma T} \sum_{t=1}^T \left\{ -\|w^t - w^* - \gamma \nabla f(w^t)\|^2 + \|w^t - w^*\|^2 + \|\gamma \nabla f(w^t)\|^2 \right\}$$

$$\stackrel{\text{use GD step}}{=} \frac{1}{2\gamma T} \sum_{t=1}^T \left\{ -\|w^{t+1} - w^*\|^2 + \|w^t - w^*\|^2 + \gamma^2 \|\nabla f(w^t)\|^2 \right\}$$

$$= \frac{1}{2\gamma T} \sum_{t=1}^T \left\{ -\|\underline{w}^{t+1} - \underline{w}^*\|^2 + \|\underline{w}^t - \underline{w}^*\|^2 \right\} \\ + \frac{\gamma}{2T} \sum_{t=1}^T \|\nabla f(\underline{w}^t)\|^2$$

The first sum is a "telescopic sum" and all successive terms cancel except for first and last one:

$$= \frac{1}{2\gamma T} \left\{ \|\underline{w}^1 - \underline{w}^*\|^2 - \|\underline{w}^{T+1} - \underline{w}^*\|^2 \right\} \\ + \frac{\gamma}{2T} \sum_{t=1}^T \|\nabla f(\underline{w}^t)\|^2$$

Now use  $\underline{w}^1 = 0$  (initialization of GD) :

$$\leq \frac{1}{2\gamma T} \|\underline{w}^1\|^2 + \frac{\gamma}{2T} \sum_{t=1}^T \|\nabla f(\underline{w}^t)\|^2 \\ \leq e^2 \text{ by Lipschitz condition.}$$

$$\leq \frac{B^2}{2\gamma T} + \frac{\gamma e^2}{2}$$

Now we choose  $\gamma$  to balance the two terms

(best possible  $\gamma$ ) :

$$\frac{B^2}{2\gamma T} = \frac{\gamma P^2}{2} \Rightarrow \gamma = \sqrt{\frac{B^2}{e^2 T}} = \frac{B}{e\sqrt{T}}$$

Moreover with this choice the upper bound is

$$\begin{aligned} \frac{B^2}{2\gamma T} + \frac{\gamma P^2}{2} &= \frac{B^2 \sqrt{T}}{2 \frac{B}{P} T} + \frac{B}{e} \frac{e^2}{2\sqrt{T}} \\ &= \frac{BP}{2\sqrt{T}} + \frac{Be}{2\sqrt{T}} = \frac{3P}{\sqrt{T}} \end{aligned}$$





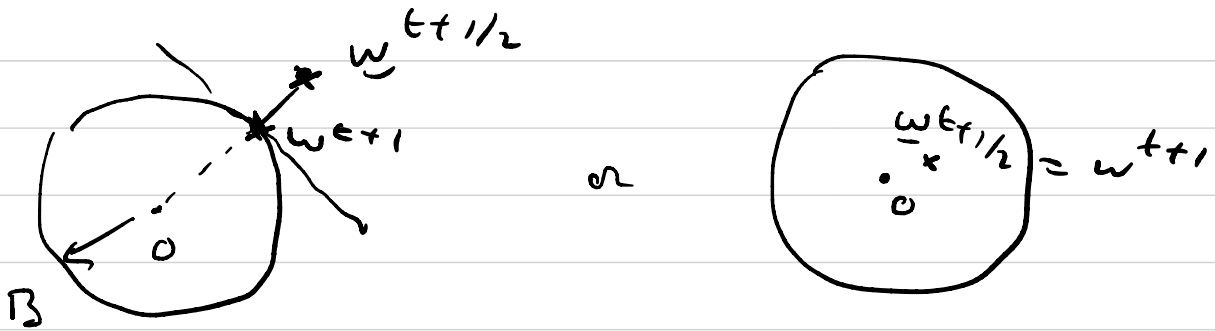
Remarks.

1) Instead of taking constant step size we can choose  $\gamma_t = \frac{B}{\sqrt{t}}$ . Steps are bigger at beginning and then basically unchanged.

2) In the theorem  $w^* \in \text{Ball}(0, B)$ . But the result does not guarantee that  $\bar{w} = \frac{1}{T} \sum_{t=1}^T w^t$  is in this ball. If we want to make sure that it is we can modify algo by adding a projection step:

$$\left\{ \begin{array}{l} \underline{w}^{t+1/2} = \underline{w}^t - \gamma \underbrace{\nabla f(\underline{w}^t)}_{\text{or } \underline{v}^t \in \partial f(\underline{w}^t)} \\ \underline{w}^{t+1} = \underset{\underline{w} \in \text{Ball}(0, B)}{\text{argmin}} \|\underline{w} - \underline{w}^{t+1/2}\| \end{array} \right.$$

Because  $B(0, \beta)$  is convex the minimizer  $\underline{w}^{t+1}$  in "projection step" is unique.



The proof of them is almost same (see ULM).

3) One could consider other averages instead of  $\frac{1}{T} \sum_{t=1}^T w^t$ , for example average over last  $\Delta T$  steps.

4) Notion of strong convexity:

$$f(\underline{x}) \geq f(\underline{z}) + \langle \nabla f(\underline{z}), \underline{z} - \underline{x} \rangle + \frac{\delta}{2} \|\underline{x} - \underline{z}\|^2$$

(excludes linear pieces and requires curvature.)

$$\Leftrightarrow f(\lambda \underline{x} + (1-\lambda)\underline{y}) \leq \lambda f(\underline{x}) + (1-\lambda)f(\underline{y}) - \frac{\delta}{2} \lambda(1-\lambda) \|\underline{x} - \underline{y}\|^2$$

yields better GD estimates.