

$$\min_{x \in \mathbb{R}^n} f(x)$$

Gradient Alg.

$$x \leftarrow x - t \nabla f(x)^T$$

of iterations $\sim \log \frac{1}{\epsilon}$

Newton's method

$$x \leftarrow x - t H(f(x))^{-1} \nabla f(x)^T$$

of iterations: $\log \log \frac{1}{\epsilon}$
if the initial value is close enough to a local solution.

- Which algorithm is better?

many iterations for Gradient alg.

Example: 10000

modest number of iterations for Newton's

Example: 20

cheap iterations \longleftrightarrow expensive iterations.

- Example: Imagine that $x \in \mathbb{R}^{1000}$

$\Rightarrow H(f(x))$ has 10^6 entries and we need to take its inverse at each iteration.

- we can compute $H(f(x))^{-1} \nabla f(x)^T$ without explicitly computing $H(f(x))^{-1}$, but it's still expensive

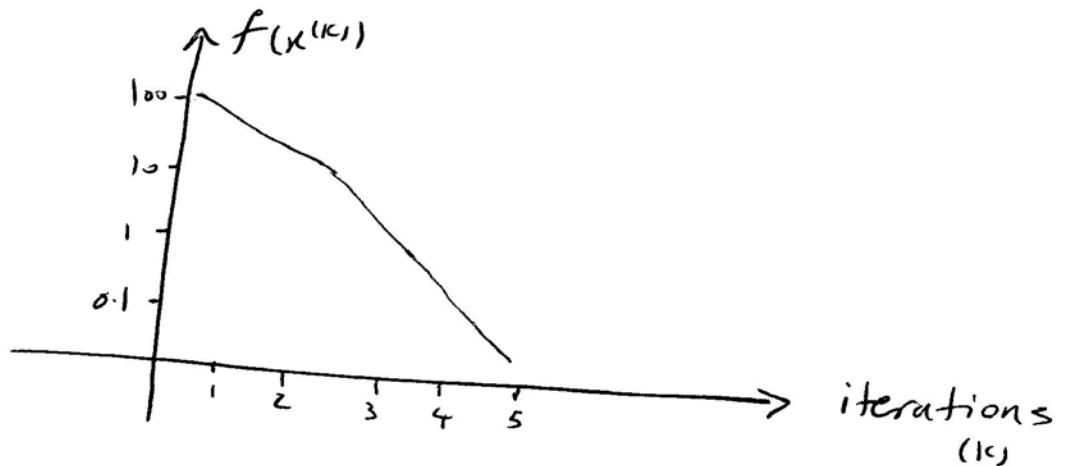
Summary:

(2)

Gradient algorithm needs many cheap iterations \longleftrightarrow to get a tolerance ϵ

Newton's method needs a few expensive iterations to get to a tolerance ϵ .

- How to evaluate the convergence? draw a linear-log plot where the horizontal axis shows iterations and the vertical axis shows $f(x)$.



- Can we design an algorithm with a convergence rate better than Gradient and Newton's methods?
- Recall that $\nabla f(x^{(k)}) \Delta x^{(k)}$ accounts for the reduction at every iteration and that's why we look for a descent direction: $\nabla f(x^{(k)}) \Delta x^{(k)} < 0$
- What if we optimize the direction Δx ?

$$\min_{\Delta x^{(k)}} \nabla f(x^{(k)}) \Delta x^{(k)}$$

→ Looks plausible but the minimum would be $-\infty$ corresponding to a $\Delta x^{(k)}$ with unbounded entries.

3

- Recall that $\Delta x^{(k)}$ accounts for direction and $t^{(k)}$ scales it down or up. So, why not normalize $\Delta x^{(k)}$.

- Consider an arbitrary norm $\|\cdot\|$.

- Example: $\| [1 \ -1 \ 3] \| = \begin{cases} \sqrt{1^2 + (-1)^2 + 3^2} \\ |1| + |-1| + |3| \\ \max(|1|, |-1|, |3|) \end{cases}$: different types of norms.

$$\Rightarrow \min_{\Delta x^{(k)}} \nabla f(x^{(k)}) \Delta x^{(k)} \quad \text{s.t.} \quad \|\Delta x^{(k)}\| = 1$$

- This is called steepest descent algorithm.

- If $\|\cdot\| = \text{length of vector} \Rightarrow \Delta x^{(k)} = \frac{-\nabla f(x^{(k)})}{\|\nabla f(x^{(k)})\|}$
 ↓
 Normalized Gradient
 ↓
 Gradient method.

- If $\|y\| = y^T M y$ for some fixed matrix M ,

⇒ The optimal $\Delta x^{(k)}$ corresponds to Newton's method.

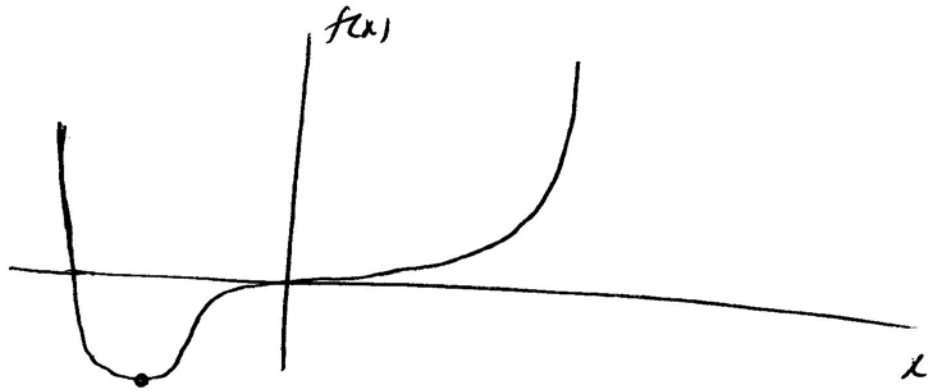
⇒ steepest descent includes Gradient and Newton's method.

- Assume we use one of the previous algorithms and generates a sequence:

$$x^{(0)} \rightarrow x^{(1)} \rightarrow x^{(2)} \dots \rightarrow x^{(k)} \rightarrow \dots$$

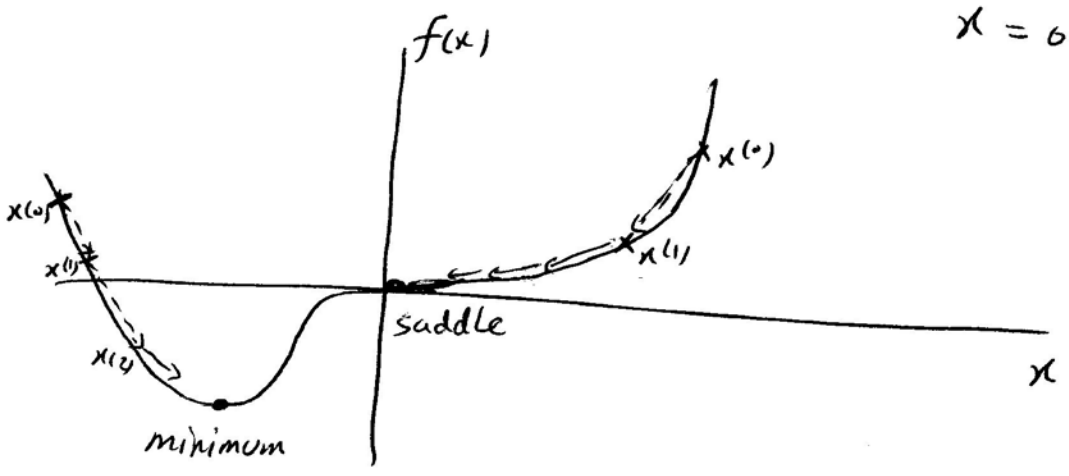
- Question: where does it converge to (if any)?
Local min or saddle point.

- Example:



If $x^{(0)} < 0 \Rightarrow$ it converges to a minimum

If $x^{(0)} \geq 0 \Rightarrow$ it converges to the saddle point $x = 0$



Theorem: Consider a local minimum x^* .

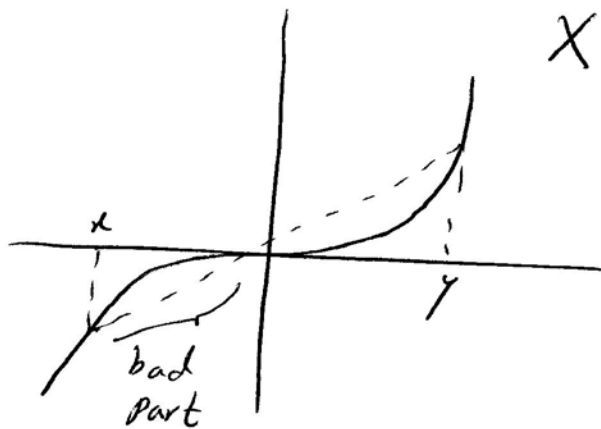
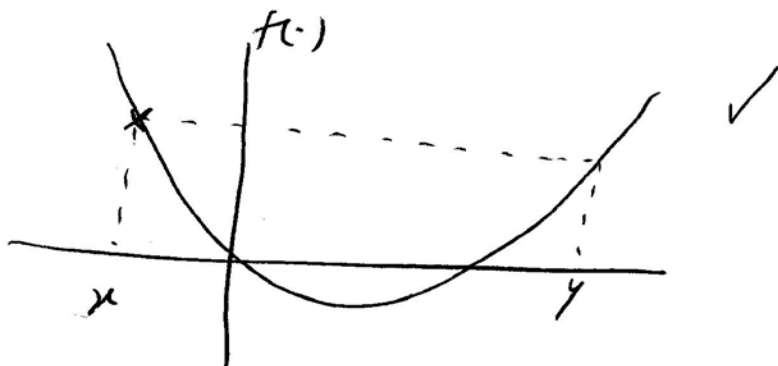
If $x^{(0)}$ is close enough to x^* , the sequence converges to x^* .

- I don't know the solution, how can I figure out my initial guess $x^{(0)}$ is good enough?

- Can we come up with a classes of functions for which $x^{(0)}$ doesn't matter?

- Yes, this is about convex optimization.

- A function $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for every two points x and y , the line connecting $(x, f(x))$ to $(y, f(y))$ is above the function. (segment)



Mathematical definition:

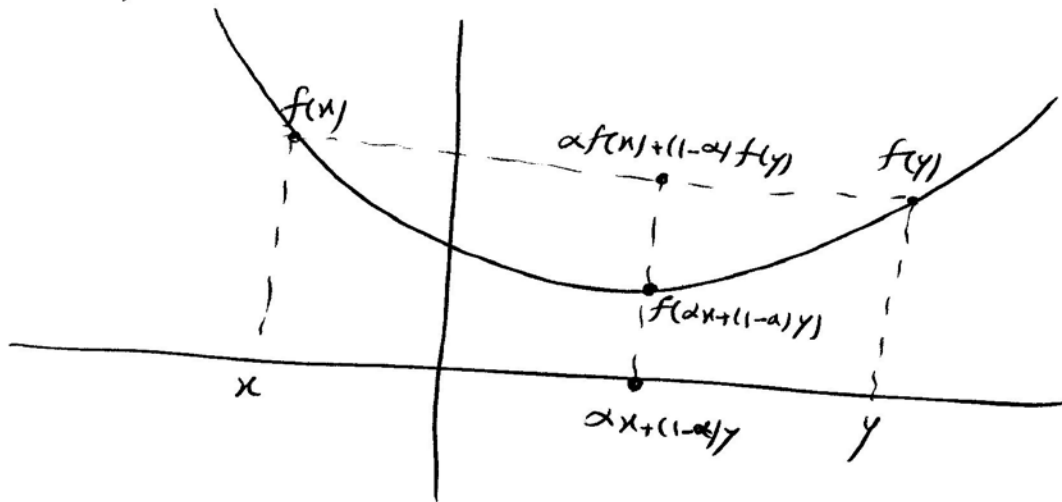
⑥

$f(\cdot): \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for every $x, y \in \mathbb{R}^n$

and $\alpha \in [0, 1]$, we have:

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

- Illustration:



- Examples of convex functions in \mathbb{R} :

