

Lecture 10. Convergence Rate of Gradient Descent

10.1 Gradient flow

At the end of the last lecture, we obtained the following lemma on the gradient descent for smooth functions.

Lemma (*Descent lemma*)

For an L -smooth differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (not necessarily convex), and $\eta \leq 1/L$, we have

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\eta}{2} \|\nabla f(\mathbf{x}_k)\|^2.$$

However, it is still not easy to show the convergence rate for the gradient descent. We now introduce a continuous version of the gradient descent instead, which is easier to analyse.

Definition (*Gradient flow*)

A *gradient flow* is a curve $\mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}^n$ following the direction of steepest descent of a function. Given a smooth convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a point $\hat{\mathbf{x}} \in \mathbb{R}^n$, the *gradient flow* of f with initial point $\hat{\mathbf{x}}$ is the solution to the following differential equation

$$\frac{d}{dt} \mathbf{x}_t = -\nabla f(\mathbf{x}_t), \quad \mathbf{x}_0 = \hat{\mathbf{x}}.$$

Here we use the notation $\mathbf{x}_t = \mathbf{x}(t)$ for convenience.

Applying the chain rule, $f(\mathbf{x}_t)$ is decreasing since

$$\frac{d}{dt} f(\mathbf{x}_t) = \left\langle \nabla f(\mathbf{x}_t), \frac{d}{dt} \mathbf{x}_t \right\rangle = -\langle \nabla f(\mathbf{x}_t), \nabla f(\mathbf{x}_t) \rangle \leq 0.$$

Now we can take the derivative

$$\begin{aligned} \frac{d}{dt} \|\mathbf{x}_t - \mathbf{x}^*\|^2 &= 2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \frac{d}{dt} \mathbf{x}_t \right\rangle \\ &= -2 \langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) \rangle \\ &= 2 \langle \mathbf{x}^* - \mathbf{x}_t, \nabla f(\mathbf{x}_t) \rangle \\ &\leq 2(f(\mathbf{x}^*) - f(\mathbf{x}_t)) \end{aligned}$$

by convexity. Then integrating both sides, we obtain that

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 - \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq 2Tf(\mathbf{x}^*) - 2 \int_0^T f(\mathbf{x}_t) dt \leq 2T(f(\mathbf{x}^*) - f(\mathbf{x}_T)),$$

which further gives that

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2}{2T} \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2T}.$$

10.2 Convergence of gradient descent with smoothness

We compare the gradient descent with the gradient flow. Assume the gradient descent iterates with a fixed step size η and an initial point $\hat{\mathbf{x}}$, i.e.,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k), \quad \mathbf{x}_0 = \hat{\mathbf{x}}.$$

For the gradient descent,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \int_{k\eta}^{(k+1)\eta} \nabla f(\mathbf{x}_k) dt.$$

For the gradient flow,

$$\mathbf{x}_{(k+1)\eta} = \mathbf{x}_{k\eta} - \int_{k\eta}^{(k+1)\eta} \nabla f(\mathbf{x}_t) dt.$$

Intuitively we know that, if ∇f does not change too fast, the gradient descent approximates the gradient flow.

Theorem

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex and L -smooth function. Choose $\eta \leq 1/L$, and let the gradient descent iterate with a fixed step size η . Then it holds that

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2T\eta}.$$

Proof

Analogously to the gradient flow, we calculate $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 - \|\mathbf{x}_k - \mathbf{x}^*\|^2$. Note that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \langle \mathbf{x}_{k+1} - \mathbf{x}^*, \mathbf{x}_{k+1} - \mathbf{x}^* \rangle \\ &= \langle \mathbf{x}_k - \mathbf{x}^* - \eta \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* - \eta \nabla f(\mathbf{x}_k) \rangle \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2. \end{aligned}$$

Similarly, it suffices to bound $-2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2$ by $f(\mathbf{x}_k) - f(\mathbf{x}^*)$. Since $\eta \leq 1/L$, we have

$$\eta^2 \|\nabla f(\mathbf{x}_k)\|^2 \leq 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}))$$

by the descent lemma. In addition, the convexity of f gives that

$$-2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle = 2\eta \langle \mathbf{x}^* - \mathbf{x}_k, \nabla f(\mathbf{x}_k) \rangle \leq 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)).$$

Thus, we obtain that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 - \|\mathbf{x}_k - \mathbf{x}^*\|^2 &= -2\eta \langle \mathbf{x}_k - \mathbf{x}^*, \nabla f(\mathbf{x}_k) \rangle + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 \\ &\leq 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)) + 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})) \\ &= 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_{k+1})). \end{aligned}$$

Summing over both sides from 0 to $T - 1$, it implies that

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 - \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \sum_{k=1}^T 2\eta (f(\mathbf{x}^*) - f(\mathbf{x}_k)) \leq 2T\eta (f(\mathbf{x}^*) - f(\mathbf{x}_T)),$$

which is equivalent to

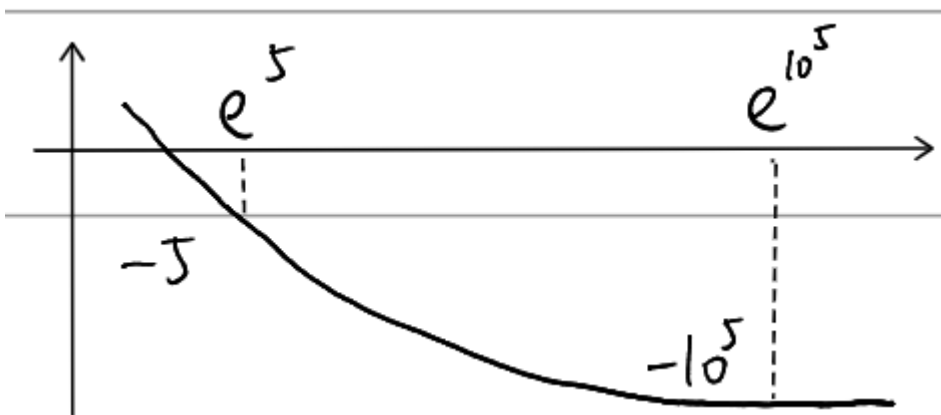
$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2}{2T\eta}.$$

If we hope $|f(\mathbf{x}_T) - f^*| < \varepsilon$, we need to run the gradient descent

$T = \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\varepsilon\eta} = O(1/\varepsilon)$ steps. If the initial point \mathbf{x}_0 is far from \mathbf{x}^* , and ε is sufficiently small, the gradient descent is slow. Unfortunately, consider the

following function

$$f(x) = \begin{cases} -\log x, & x < e^{10^5} \\ -10^5, & x \geq e^{10^5} \end{cases}$$



This function is convex and 1-smooth. Hence $x_{k+1} = x_k + \eta/x_k \leq x_k + 1/x_k$ and the convergence rate of the gradient descent will be very small if $x_0 = 1$.

Question

Under which assumptions the gradient descent converges rapidly?

10.3 Strongly convex functions

Recall that, if we run the gradient descent for a quadratic function $f(x) = ax^2$ where $a \in \mathbb{R}_{>0}$, it gives that $x_{k+1} = (1 - 2a\eta)x_k$ and thus $f(x_k) = a(1 - 2a\eta)^{2k}x_0^2$. Clearly $f(x_k)$ converges to the optimal value 0 at an exponential rate.

We now introduce the following definition, which requires the function is a bit "better" than some quadratic function.

Definition (Strong convexity)

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is *strongly convex* with $\mu > 0$ if $f(x) - \frac{\mu}{2}\|x\|^2$ is convex.

There are some other forms of quadratic functions. Why don't we choose other functions such as $x^T Q x$ or $\|x - y_0\|^2$ for some given y_0 ? In fact, these functions mentioned can just achieve a similar effect to $\frac{\mu}{2}\|x\|^2$. For example, $x^T Q x$ is almost equivalent to $\lambda_{\max}(Q)x^T x = \lambda_{\max}(Q)\|x\|^2$. In addition,

$$\|x - y_0\|^2 = \|x\|^2 + \underbrace{\|y_0\|^2}_{\text{constant}} - \underbrace{2\langle x, y_0 \rangle}_{\text{don't affect convexity}}.$$

Hence, all quadratic functions achieve similar effects to $\|x\|^2$.

Recall that, a function is convex iff its hessian matrix is positive semidefinite. The hessian matrix of $f(x) - \frac{\mu}{2}\|x\|^2$ is

$$\nabla^2 f(x) - \frac{\mu}{2}\nabla^2(\|x\|^2) = \nabla^2 f(x) - \frac{\mu}{2}\nabla^2(x^\top x) = \nabla^2 f(x) - \mu I.$$

Lemma

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable function. Then f is μ -strongly convex iff $\nabla^2 f(x) \succeq \mu I_n$. Namely, for all $x \in \mathbb{R}^n$, $\lambda_{\min}(\nabla^2 f(x)) \geq \mu$.

We also have the following lemma similar to the first order condition for convexity and smoothness.

Lemma

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function. Then f is μ -strongly convex iff for all $x, y \in \mathbb{R}^n$,

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

Proof

Let $g(x) = f(x) - \frac{\mu}{2}\|x\|^2$. By the first order condition for convexity, $g(x)$ is convex iff for all $x, y \in \mathbb{R}^n$,

$$g(y) \geq g(x) + \langle \nabla g(x), y - x \rangle.$$

Note that $\nabla g(x) = \nabla f(x) - \mu x$. So it gives that $g(x)$ is convex iff for all $x, y \in \mathbb{R}^n$,

$$g(y) \geq g(x) + \langle \nabla f(x), y - x \rangle - \mu \langle x, y - x \rangle.$$

The last inequality is equivalent to

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

Rearranging it, we obtain that

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \mu \langle \mathbf{x}, \mathbf{y} \rangle + \frac{\mu}{2} \|\mathbf{x}\|^2 + \frac{\mu}{2} \|\mathbf{y}\|^2 \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2. \end{aligned}$$

As a corollary, above lemma implies that $f(\mathbf{y}) > f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$ for any $\mathbf{x} \neq \mathbf{y}$. Hence, f is strictly convex.

Example

1. An affine functions $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ can not be strongly convex since it is not strictly convex.
2. $f(x) = -\log x$ can not be strongly convex since $f'(x) = -\frac{1}{x}$ and $f''(x) = \frac{1}{x^2}$ and we can not find out such μ when $x \rightarrow 0$.
3. $f(x) = ax^2, a > 0$ is $2a$ -strongly convex.
4. $f(x) = x^4$ is not strongly convex since $f'(x) = 4x^3, f''(x) = 12x^2$ and we can not find such $\mu > 0$.
5. $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Q}\mathbf{x}$ where $\mathbf{Q} \succ 0$ is strongly convex. Because $\nabla^2 f(\mathbf{x}) = 2\mathbf{Q}$, f is $2\lambda_{\min}(\mathbf{Q})$ -strongly convex.

Recall the property of monotone gradient for convex functions. We have a similar corollary.

Corollary

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function. Then f is μ -strongly convex iff for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

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

10.4 Convergence of gradient descent with strong convexity

We now establish the convergence of gradient descent with strong convexity. First consider the gradient flow again. By strong convexity, we can bound the derivative

as follows

$$\begin{aligned}\frac{d}{dt} \|\mathbf{x}_t - \mathbf{x}^*\|^2 &= 2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \frac{d}{dt} \mathbf{x}_t \right\rangle \\ &= -2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) \right\rangle \\ &= -2 \left\langle \mathbf{x}_t - \mathbf{x}^*, \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*) \right\rangle \\ &\leq -2\mu \|\mathbf{x}_t - \mathbf{x}^*\|^2.\end{aligned}$$

For a time-continuous non-negative process $u_t = u(t)$, if $\frac{d}{dt} u_t = -\alpha u_t$, then we have $u_t = u_0 \exp(-\alpha t)$. The same result holds if we replace the equality by an inequality.

Theorem (Gronwall's lemma)

For a time-continuous non-negative process $u_t = u(t)$, if $\frac{d}{dt} u_t \leq -\alpha_t u_t$, then we have

$$u_T \leq u_0 \exp\left(-\int_0^T \alpha_t dt\right).$$

Applying the Gronwall's lemma, we conclude $\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \exp(-2\mu T)$ immediately, which gives an exponential decay rate. Intuitively, as the discretization version of the gradient flow, the gradient descent for strongly convex functions should also follow the exponential decay.

Theorem

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is n L -smooth and μ -strongly convex function. Choose $\eta \leq 1/L$, and let the gradient descent iterate with a fixed step size η . Then it holds that

$$\|\mathbf{x}_T - \mathbf{x}^*\|^2 \leq (1 - \mu\eta)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Proof

By strong convex,

$$\begin{aligned}
\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \langle (\mathbf{x}_k - \mathbf{x}^*) - \eta \nabla f(\mathbf{x}_k), (\mathbf{x}_k - \mathbf{x}^*) - \eta \nabla f(\mathbf{x}_k) \rangle \\
&= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle \\
&\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta \left(f(\mathbf{x}_k) - f(\mathbf{x}^*) + \frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \right) \\
&= (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \eta^2 \|\nabla f(\mathbf{x}_k)\|^2 - 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \\
&\leq (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 + 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})) - 2\eta (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \\
&= (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\eta (f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)) \\
&\leq (1 - \mu\eta) \|\mathbf{x}_k - \mathbf{x}^*\|^2
\end{aligned}$$

where the second inequality is due to the *descent lemma*

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{\eta}{2} \|\nabla f(\mathbf{x}_k)\|^2.$$

The function value also has an exponential decay. Since f is L -smooth, we have

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}^*), \mathbf{x}_T - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_T - \mathbf{x}^*\|^2,$$

which gives the following corollary.

Corollary

$$f(\mathbf{x}_T) - f^* \leq \frac{L}{2} (1 - \mu\eta)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

10.5 Condition number

For a quadratic function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top Q \mathbf{x}$ where $Q \succeq 0$, we already have the following facts:

- $f(\mathbf{x})$ is $\lambda_{\min}(Q)$ -strongly convex;
- $f(\mathbf{x})$ is $\lambda_{\max}(Q)$ -smooth.

Applying the above theorem, if we take the step size $\eta = 1/\lambda_{\max}(Q)$, $\{\mathbf{x}_n\}$ will converge at an exponential rate of $1 - \frac{\lambda_{\min}(Q)}{\lambda_{\max}(Q)}$ (since $\mathbf{x}^* = \mathbf{0}$). Recall the in last lecture, we have shown $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ as long as $\eta < 2/\lambda_{\max}(Q)$, which means that $\eta = 1/\lambda_{\max}(Q)$ is not necessary for convergence. Since we hope $1 - \mu\eta$ is as small as possible, we may choose a greater value of η to obtain a better rate.

Now let us calculate the optimal convergence rate of $f(\mathbf{x})$. Since $\nabla f(\mathbf{x}) = Q\mathbf{x}$, we have

$$\mathbf{x}_{k+1} = (I - \eta Q) \mathbf{x}_k.$$

Applying the eigen-decomposition $Q = U\Lambda U^\top$, where $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ and $UU^\top = I$.

Then $I - \eta Q = U\Lambda'U^\top$ where $\Lambda' = \text{diag}\{1 - \eta\lambda_1, 1 - \eta\lambda_2, \dots, 1 - \eta\lambda_n\}$. So

$$\mathbf{x}_k = (U\Lambda'U^\top)^k \mathbf{x}_0 = U(\Lambda')^k U^\top \mathbf{x}_0.$$

Let $\mathbf{y}_k = U^\top \mathbf{x}_k$. We have

$$\mathbf{y}_k = (\Lambda')^k \mathbf{y}_0 = \text{diag}\{(1 - \eta\lambda_1)^k, (1 - \eta\lambda_2)^k, \dots, (1 - \eta\lambda_n)^k\} \mathbf{y}_0.$$

Note that $\mathbf{x}^* = \arg \min f(\mathbf{x}) = \mathbf{0}$. Thus,

$$\begin{aligned} \|\mathbf{x}_k - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k\|^2 = \mathbf{x}_k^\top \mathbf{x}_k = \mathbf{x}_k^\top U U^\top \mathbf{x}_k \\ &= \mathbf{y}_k^\top \mathbf{y}_k = \|\mathbf{y}_k\|^2 = \sum_{i=1}^n (1 - \eta\lambda_i)^{2k} y_{0,i}^2 \\ &\leq \max_{1 \leq i \leq n} \{(1 - \eta\lambda_i)^{2k}\} \|\mathbf{y}_0\|^2 \\ &= \left(\max_{1 \leq i \leq n} \{1 - \eta\lambda_i\} \right)^{2k} \|\mathbf{y}_0\|^2 \\ &= \max \{|1 - \eta\lambda_{\min}|, |1 - \eta\lambda_{\max}|\}^{2k} \|\mathbf{y}_0\|^2 \\ &= \max \{|1 - \eta\lambda_{\min}|, |1 - \eta\lambda_{\max}|\}^{2k} \|\mathbf{x}_0\|^2 \end{aligned}$$

We would like to choose η to minimize $\max \{|1 - \eta\lambda_{\min}|, |1 - \eta\lambda_{\max}|\}$, which means

$\eta = \frac{2}{\lambda_{\min} + \lambda_{\max}} = \frac{2}{\mu + L}$. In this case,

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left| \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \right|^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 = \left| \frac{L - \mu}{L + \mu} \right|^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Definition (Condition number)

Given an $n \times n$ positive definite matrix $Q \succ 0$, its *condition number* is defined by

$$\kappa(Q) = \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} \geq 1.$$

The argument above reveals that for quadratic functions, the convergence rate of gradient descent depends on $\left(\frac{\kappa-1}{\kappa+1}\right)^2$.

Example

- For $Q = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix}$, its condition number is 2, so the convergence rate is $\frac{1}{3}$.
- For $Q = \begin{pmatrix} \frac{1}{100} & 0 \\ 0 & 1 \end{pmatrix}$, its condition number is 100, so the convergence rate is $\frac{99}{101}$.

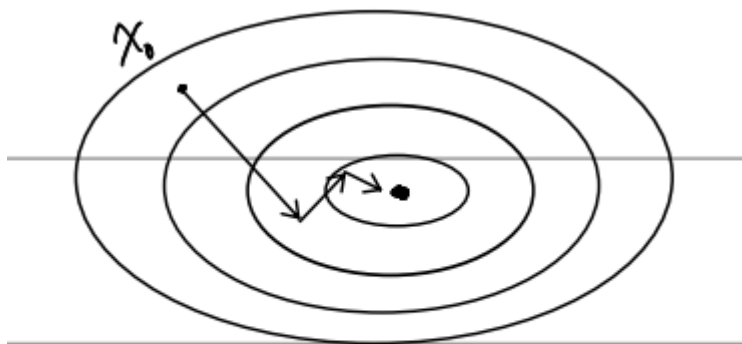
These two examples shows that the gradient descent may converge very slowly when the coefficient matrix has a large condition number.

For nonquadratic functions, we can approximate them locally (near the minimum point \mathbf{x}^*) by the following Taylor series

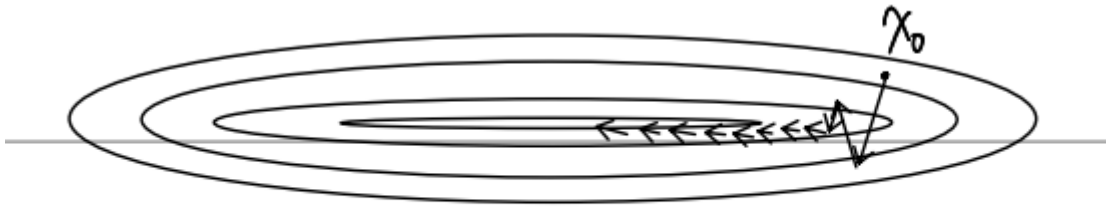
$$\begin{aligned} f(\mathbf{x}) &\approx f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 f(\mathbf{x}^*) (\mathbf{x} - \mathbf{x}^*) \\ &= f(\mathbf{x}^*) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \nabla^2 f(\mathbf{x}^*) (\mathbf{x} - \mathbf{x}^*) \end{aligned}$$

Hence, in the neighborhood of the minimum point \mathbf{x}^* , the convergence rate depends on $\left(\frac{\kappa(\nabla^2 f(\mathbf{x}^*)) - 1}{\kappa(\nabla^2 f(\mathbf{x}^*)) + 1} \right)^2$. If the condition number of $\nabla^2 f(\mathbf{x}^*)$ is large, the results given by gradient descent with fixed step size cannot converge rapidly.

Here are some well-conditioned and ill-conditioned examples:



$$Q = \text{diag} \{ 1/2, 1 \}$$



$$Q = \text{diag} \{ 0.01, 1 \}$$