

Summary of MA 574 Fall 2024

- Part I : convergence of Gradient Descent
 - Nesterov's acceleration
 - Convergence of line search method
- Part II : subgradient method
 - proximal point method
 - (fast) proximal gradient descent
 - splitting methods for $\min_x f(x) + g(x)$
 - PDHG, ADMM, Douglas-Rachford
- Part III : randomized coordinate descent
 - Stochastic gradient descent
- Part IV : Riemannian gradient descent

All these methods are first order methods.
Using gradient

Usually we can prove at least $\lim \| \nabla f(x_k) \| = 0$

But a critical point may not be a global minimizer.

First-order methods almost always avoid strict saddle points

Full Length Paper | Series B | Published: 18 February 2019

Volume 176, pages 311–337, (2019) [Cite this article](#)

[Download PDF](#) ↓

✓ Access provided by Purdue University

[Jason D. Lee](#) ✉, [Ioannis Panageas](#), [Georgios Piliouras](#), [Max Simchowitz](#), [Michael I. Jordan](#) & [Benjamin Recht](#)

This paper covers first order methods
such as {
gradient descent
proximal point method
coordinate descent
Riemannian gradient descent

Today we go through the main ideas of
this paper

Step I: Dynamical System

Step II: Show that gradient descent with
only special initial guess can converge to saddle

Dynamical System for Optimization

Brief intro to stability analysis of dynamical system

Deterministic dynamical system

$$\mathcal{Q}: T \times X \rightarrow X \\ (t, x) \mapsto \mathcal{Q}(t, x)$$

Example: $T = \{0, 1, 2, 3, \dots\}$ $\min f(x)$

$$x_{t+h} = x_t - \eta \nabla f(x_t)$$

$$g = I - \eta \nabla f$$

$$\mathcal{Q}(t, x) = \underbrace{g \circ g \circ \dots \circ g}_{t\text{-fold}}(x)$$

Linearization near $x_* = 0$: $\nabla f(x) = \nabla f(x) - \nabla f(x_*) \approx \underbrace{H(x-x_*)}_{\text{Hessian}}$

Linearization gives $\Phi(t, x) = A^t x$

Assume $x_* = 0$

$$= (I - \alpha H)^t x \quad \leftarrow \begin{array}{l} \text{power } t \\ \text{is an integer} \end{array}$$

Equilibrium $\mathcal{Q}(t, \bar{x}) = \bar{x}, \forall t \in T$

Example: $\bar{x} - \eta \nabla f(\bar{x}) = \bar{x} \Leftrightarrow \nabla f(\bar{x}) = 0$

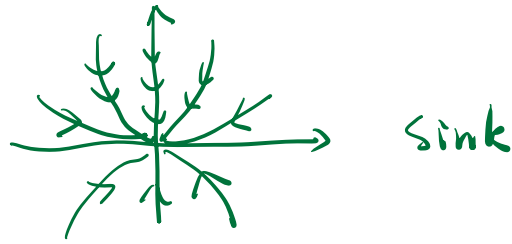
Stability near Equilibrium

Example: $x_t = A^t x_0, A \in \mathbb{R}^{2 \times 2}, x_0 \in \mathbb{R}^2$

① Let λ_1, λ_2 be eigenvalues of A

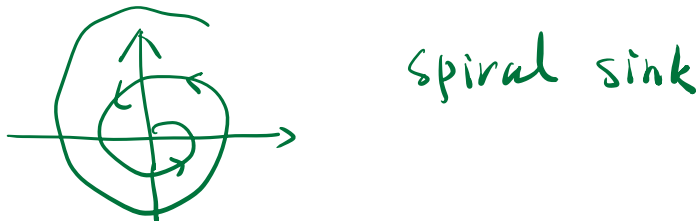
$\lambda_1, \lambda_2 \in \mathbb{R} \quad |\lambda_1| < 1, |\lambda_2| < 1$, then stable

meaning $x_t \rightarrow \bar{x} = 0$ for any $x_0 \neq 0$



② $\lambda_1, \lambda_2 \in \mathbb{R} \quad |\lambda_1| > 1, |\lambda_2| > 1$, source
unstable

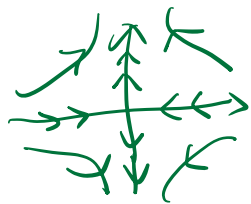
③ $\lambda_1, \lambda_2 = r e^{\pm i\theta}$, $r < 1$



④ $\lambda_1, \lambda_2 = r e^{\pm i\theta}$, $r > 1$



⑤ $\lambda_1, \lambda_2 \in \mathbb{R} \quad |\lambda_1| < 1, |\lambda_2| > 1$



saddle

$$\min_x f(x) = \frac{1}{2} x^T H x$$

$$\nabla^2 f(x) = H$$

$$(I - \eta \nabla f(x))^t = (I - \eta H)^t$$

Reference Michael Scharb 1987

Global stability of dynamical system

Theorem $g: X \rightarrow X$, $\mathcal{L}(t, x) = \underbrace{g \circ \dots \circ g}_t(x) = g^t(x)$
 \hookrightarrow finite dim linear space t -fold

If $\langle g(0) \rangle = 0$, $g'(0) = A$

② $X = X_1 \oplus X_2$ (A -invariant decomposition)

$$A(X_1) \subseteq X_1 \quad A(X_2) \subseteq X_2$$

$$\textcircled{3} \quad \|(A|_{X_1})^{-t}\| \leq c \lambda^t, \quad \lambda \in (0, 1)$$

$$\|(A|_{X_2})^t\| \leq c \lambda^t, \quad \forall t \in \mathbb{T} = \mathbb{N}$$

A^{-1} on X_1 } is a contraction
 A on X_2 }

A on X_1 is an expansion

$$\left(\begin{array}{l} \text{Example: } A = \begin{pmatrix} 2 & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \quad X_1 = \text{Span}\left\{\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right\} \\ X_2 = \text{Span}\left\{\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right\} \end{array} \right)$$

then locally near 0, there exists local manifolds M_1, M_2

which are tangent to X_1, X_2 s.t.

$$\textcircled{1} \quad X \in M_1 \cap B \Leftrightarrow \exists X_t \rightarrow 0, \text{ s.t. } g^t(X_t) = X$$

\downarrow
 local ball in X

$$X_t = (g^{-1})^t(X) \rightarrow 0$$

$$\textcircled{2} \quad X \in M_2 \cap B \Leftrightarrow X_t = g^t(X) \rightarrow 0 \leftarrow \text{exponential rate}$$

$$\textcircled{3} \quad \text{If } X \notin M_2 \cap B, \exists t \text{ s.t. } g^t(X) \notin B$$

Remark: ① Only a local result.

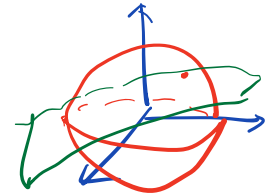
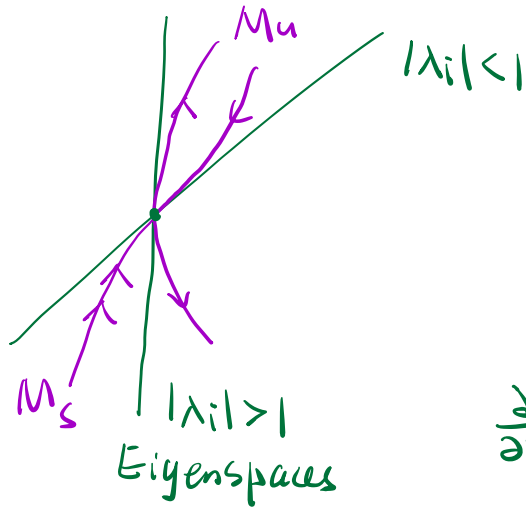
$$\textcircled{2} \quad \text{If } X = X_s \oplus X_c \oplus X_u$$

| | | |
|-------------------|-------------------|-------------------|
| $ \lambda < 1$ | $ \lambda = 1$ | $ \lambda > 1$ |
| M_s | M_c | M_u |

For center manifold M_c , then

- 1) $g^t(x)$ escapes B (subexponential rate)
- 2) $g^t(x) \rightarrow 0$ (subexponential rate)
- 3) $g^t(x) \in B$, but not converging to 0.

Now just focus on a saddle point x_*



strict saddle points

$$\nabla f(x_*) = 0 \quad \lambda_{\min}(\nabla^2 f(x_*)) < 0$$

$$\frac{\partial}{\partial x} (I - \eta \nabla f(x_*)) = I - \eta \nabla^2 f(x_*)$$

Gradient Descent

$$x_{t+1} = x_t - \alpha \nabla f(x_t) = g(x_t)$$

$$g = I - \alpha \nabla f$$

Linearization : $A = I - \alpha H$

$$H = \nabla^2 f(x_*)$$

$$\lambda_i(A) = 1 - \alpha \lambda_i(H)$$

If $\lambda_i(H) < 0$, then $\lambda_i(A) > 1$

\Rightarrow the unstable manifold M_u is non-trivial
 $\dim(M_u) \geq 1$

Assume $\lambda_i(H) \neq 0, \forall i \Rightarrow |\lambda_i(A)| \neq 1, \forall i$

\Rightarrow no center manifold

$\Rightarrow \dim(M_s) = d - \dim(M_u) \leq d-1 < d$

$$\Rightarrow \text{Measure}(M_S \cap B) = 0$$

If $x_t \rightarrow x^*$, ^{saddle point} then $\exists t_0$ s.t.

$$x_{t_0} \in M_S \cap B$$

← Theorem from last time

$$x_{t_0-1} = g^{-1}(x_{t_0}) \in g^{-1}(M_S \cap B)$$

$$\Rightarrow x_0 \in g^{-t_0}(M_S \cap B)$$

$$\Rightarrow x_0 \in \bigcup_{t=0}^{\infty} g^{-t}(M_S \cap B)$$

$$\Rightarrow \text{measure} \left(\bigcup_{t=0}^{\infty} g^{-t}(M_S \cap B) \right)$$

$$\leq \sum_{t=0}^{\infty} \text{measure} \left(g^{-t}(M_S \cap B) \right)$$

(Lemma: If $\lambda_i(Dg(x)) \neq 0, \forall x \in \mathbb{R}^d$
then $\text{meas}(B) = 0 \Rightarrow \text{meas}(g^{-1}(B)) = 0$)

$$\stackrel{\downarrow}{=} 0$$

Countable Union of Measure Zero Sets has measure 0

Theorem If x^* is a strict saddle,

then $\{x_0 \mid x_t \rightarrow x^*\}$ has measure 0.

So if we take a random initial guess, then the probability of this initial guess lying in the set $\{x_0 \mid x_t \rightarrow x^*\}$ is zero.