

Running Notes

Daniel Henderson

May 23, 2025

1 Preliminaries

1.1 Notation

We use the following notation

- $\|\cdot\|$ denotes the usual ℓ_2 norm for vectors \mathbf{x} in \mathbb{R}^n and $p = 2$ norm for matrices in $\mathbb{R}^{n \times m}$. i.e.,

$$\|\mathbf{x}\| := \left(\sum_i x_i^2 \right)^{1/2}$$
$$\|A\| := (\lambda_{\max}(A^\top A))^{1/2} = \max(\sigma(A))$$

- $\sigma(A) := \{\text{singular values of } A\}$.
- $A \in \mathbb{R}^{n \times n} \implies \sigma(A) = \{\text{eigenvalues of } A \text{ (i.e. spectrum)}\}$
- $\sigma_{\max}(A) := \max(\sigma(A))$ and $\sigma_{\min}(A) := \min(\sigma(A))$.
- $A \in \mathbb{R}^{n \times n} \implies \lambda_{\max}(A) := \sigma_{\max}(A)$ and $\lambda_{\min}(A) = \sigma_{\min}(A)$
- $\langle \cdot, \cdot \rangle$ denotes the usual inner product on \mathbb{R}^n , i.e.,

$$\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i = \|\mathbf{x}\| \|\mathbf{y}\| \cos(\theta)$$

where θ is the angle between \mathbf{x} and \mathbf{y} .

- $\mathcal{B}_r(\mathbf{x}) := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| < r\}$ is the open ball of radius r centered at $\mathbf{x} \in \mathbb{R}^n$.

1.2 Assumptions

We assume the following assumptions

- Let $\Omega \subset \mathbb{R}^n$ be a bounded open subset with a smooth boundary $\partial\Omega$, which we assume is a C^2 -manifold. Additionally, we often consider Ω and $\partial\Omega$ to be the $n - 1$ -dimensional hypersurface of a n -dimensional ball.
- Let be a function $f \in C^\infty(\Omega; \mathbb{R})$ is a continuous and infinitely-differentiable mapping from Ω to \mathbb{R} . In general, f is nonlinear and nonconvex on Ω .
- We assume C^2 regularity of the boundary $\partial\Omega$ so that the derivatives of f may be extended to $\partial\Omega$.
- The typefont \mathbf{x} will denote a vector in \mathbb{R}^n .

1.3 Definitions

Using the above notation, we assume our assumptions hold, and we define the following

Definition 1. f is **L -Lipschitz** if $\forall \mathbf{x}_1, \mathbf{x}_2$

$$\exists L \geq 0 : \|f(\mathbf{x}_1) - f(\mathbf{x}_2)\| \leq L\|\mathbf{x}_1 - \mathbf{x}_2\|$$

Definition 2. f has **ℓ -Lipschitz gradient**, or, f is **ℓ -smooth** if $\forall \mathbf{x}_1, \mathbf{x}_2$

$$\exists \ell \geq 0 : \|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq \ell\|\mathbf{x}_1 - \mathbf{x}_2\|$$

Definition 3. f has **ρ -Lipschitz Hessian** if $\forall \mathbf{x}_1, \mathbf{x}_2$

$$\exists \rho \geq 0 : \|\nabla^2 f(\mathbf{x}_1) - \nabla^2 f(\mathbf{x}_2)\| \leq \rho\|\mathbf{x}_1 - \mathbf{x}_2\|$$

Definition 4. f is **convex** if $\forall \mathbf{x}_1, \mathbf{x}_2$

$$\begin{aligned} f(\mathbf{x}_2) &\geq f(\mathbf{x}_1) + \langle \mathbf{x}_2 - \mathbf{x}_1, \nabla f(\mathbf{x}_1) \rangle \\ &= f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \end{aligned}$$

Definition 5. f is **strictly convex** if

$$\begin{aligned} \exists \mu > 0 : \nabla^2 f \succeq \mu I \\ \iff \lambda_{\min}(\nabla^2 f) \geq \mu > 0 \end{aligned}$$

Definition 6. f is α -**strongly convex** if $\forall \mathbf{x}_1, \mathbf{x}_2 \exists \alpha > 0$ s.t.

$$\begin{aligned} f(\mathbf{x}_2) &\geq f(\mathbf{x}_1) + \langle \nabla f(\mathbf{x}_1), \mathbf{x}_2 - \mathbf{x}_1 \rangle + \frac{\alpha}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2 \\ \iff \lambda_{\min}(\nabla^2 f(\mathbf{x})) &\geq -\alpha. \end{aligned}$$

Definition 7. \mathbf{x}^* is a **first-order stationary point** if $\|\nabla f(\mathbf{x}^*)\| = 0$.

Definition 8. \mathbf{x}^* is an ϵ -**first-order stationary point** if $\|\nabla f(\mathbf{x}^*)\| \leq \epsilon$.

Definition 9. $\mathbf{x}^* \in \mathbb{R}^n$ is a **second-order stationary point** if $\|\nabla f(\mathbf{x}^*)\| = 0$ and $\nabla^2 f(\mathbf{x}^*) \succeq 0$.

Definition 10. if f has ρ -Lipschitz Hessian, $\mathbf{x}^* \in \mathbb{R}^n$ is a ϵ -**second-order stationary point** if

$$\|\nabla f(\mathbf{x}^*)\| \leq \epsilon \text{ and } \nabla^2 f(\mathbf{x}^*) \succeq -\sqrt{\rho}\epsilon$$

Remark. Note that the Hessian is not required to be positive definite, but it is required to have a small eigenvalue.

Definition 11. We consider the general form of our **unconstrained optimization problem** to be

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in K \subseteq \mathbb{R}^n \tag{1}$$

where K is a compact set in \mathbb{R}^n . We denote the **optimal solution** and **optimal value** of the optimization problem

$$\begin{aligned} \mathbf{x}^* &= \arg \min_{\mathbf{x} \in K} f(\mathbf{x}) \\ \mathbf{f}^* &= \min_{\mathbf{x} \in K} f(\mathbf{x}) \end{aligned}$$

where \mathbf{x}^* satisfies the first-order optimality condition, i.e., $\nabla f(\mathbf{x}^*) = 0$.

Definition 12. Let the **gradient flow** of f be a solution to the dynamical system defined as

$$\gamma'(t) = -\nabla f(\gamma(t))$$

where the evolution of our phase space is driven by the negative gradient of f . A **gradient flow line** of f is an integral curve $\gamma : [0, t_f] \rightarrow \Omega$ satisfying the above evolutionary system (ordinary-differential equation) subject to $\gamma(0) = \mathbf{x}_0$.

We aim to classify the phase space of the *gradient flow* of f on Ω . First we notice that for any critical point \mathbf{x}^* ,

$$\begin{aligned} \gamma(t) &= \mathbf{x}^* \quad \forall t \in [0, t_f] \\ \implies \gamma'(t) &= \mathbf{0} \quad \text{and} \quad -\nabla f(\gamma(t)) = -\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \because \mathbf{x}^* \text{ is a critical point} \\ \therefore \gamma(t) &= -\nabla f(\gamma(t)) \quad \forall t \in [0, t_f] \end{aligned}$$

Consequently, by the uniqueness of solutions for ordinary differential equations, if any flow line contains a *first-order* critical point \mathbf{x}^* , it must be a constant flow line.

Lemma A. *The function $f : \omega \rightarrow \mathbb{R}$ is nonincreasing along any flow-line $\gamma(t)$ and strictly decreasing along flow lines not containing a critical point \mathbf{x}^* .*

Proof. Let $\gamma : [0, t_f] \rightarrow \Omega$ be a flow line. Consider the composition $f \circ \gamma : [0, t_f] \rightarrow \mathbb{R}$, its derivative is

$$\begin{aligned} \frac{d}{dt}(f(\gamma(t))) &= \langle \nabla_{\gamma(t)}(f), \frac{d\gamma}{dt} \rangle \\ &= \langle \nabla_{\gamma(t)}(f), -\nabla_{\gamma(t)}(f) \rangle \\ &= -\langle \nabla_{\gamma(t)}(f), \nabla_{\gamma(t)}(f) \rangle \\ &= -|\nabla_{\gamma(t)}(f)|^2 \\ &\leq 0 \end{aligned}$$

Therefore, $f'(\gamma(t)) = 0$ iff $\gamma(t)$ is on a critical point of f . In particular, if $\gamma(t)$ does not contain in its image a critical point of f , then the above inequality implies that f is strictly decreasing along the integral curve $\gamma(t)$. \square

theorem. For all \mathbf{x} in the closed manifold $\overline{\Omega}$, there exists uniquely $\gamma_{\mathbf{x}}(t) : \mathbb{R} \rightarrow \overline{\Omega}$ such that $\gamma_{\mathbf{x}}(0) = \mathbf{x}$ and the limits

$$\lim_{t \rightarrow -\infty} \gamma_{\mathbf{x}}(t) \quad \text{and} \quad \lim_{t \rightarrow \infty} \gamma_{\mathbf{x}}(t)$$

exist and converge to *critical-points* of f .

Now it may be shown that the flow map operator $T : \bar{\Omega} \times \mathbb{R} \rightarrow \bar{\Omega}$ defined as $T(\mathbf{x}, t) := \gamma_{\mathbf{x}}(t)$. This implies theoretically that we may perform analysis of our phase space as constructed by a smooth union of integral curves. Consequently, by our previous lemma, if the flow map T contains a critical point \mathbf{x}^* then it ought to be that $T(\mathbf{x}^*, t) \equiv \text{const} \quad \forall t$, otherwise, T is descending for all of time.

Definition 13. A **scheme** for solving a general form *unconstrained optimization problem* is a one-parameter family of iteration operators:

$$T_h : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{where} \quad \mathbf{x}_{k+1} = T_h(\mathbf{x}_k), \quad h \in (0, h_0]$$

where h_0 is a constant and h is the step size. The scheme is well-defined such that the triplet (\mathbf{x}_0, h, T_h) satisfy

1. *Consistency*: $\forall \mathbf{x} \in K$

$$(T_h(\mathbf{x}) - \mathbf{x}) h^{-1} + \nabla f(\mathbf{x}) \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

implying that a single step approximates the continuous gradient flow w/ local error $\mathcal{O}(h^{p+1})$ where p is the global order of the scheme.

2. *Stability*: $\exists c > 0, h_0 > 0 : \forall h \in (0, h_0]$, and for all $\mathbf{x}_1, \mathbf{x}_2$ in a *neighborhood* $N \subset K$ around an *optimal solution* \mathbf{x}^*

$$\|T_h(\mathbf{x}_1) - T_h(\mathbf{x}_2)\| \leq (1 - ch)\|\mathbf{x}_1 - \mathbf{x}_2\|$$

where c is a constant that depends on the scheme and h is the step size. Or, equivalently, the scheme is *stable* if $\exists c > 0, h_0 > 0 : \forall h \in (0, h_0]$, and for all \mathbf{x} in a *neighborhood* about \mathbf{x}^* , each step results in a strict decrease of by atleast a factor of $1 - ch$, i.e.,

$$\|J(T_h(\mathbf{x}))\| \leq (1 - ch)$$

where $J(T_h(\mathbf{x}))$ is the Jacobian of the scheme.

3. *Convergence*: $\forall \mathbf{x}_0 \in N \subset K$, s.t. N is some neighborhood around a strict minimizer \mathbf{x}^* . and $\forall \epsilon > 0$

$$\exists K \in \mathbb{N} : \forall k > K, \mathbf{x}_k \in N \text{ and } d(\mathbf{x}_k, \mathbf{x}^*) \leq \epsilon.$$

Definition 14. \mathbf{x}^* is **non-degenerate** if $\nabla^2 f(\mathbf{x}^*)$ is non-singular.

Definition 15. The **level set** of f at c is the set of points $\mathbf{x} \in \Omega$ such that $f(\mathbf{x}) = c$, i.e.,

$$L_c = \{\mathbf{x} \in \Omega : f(\mathbf{x}) = c\}$$

The level set L_c is a smooth manifold with boundary ∂L_c . The **sublevel** and **superlevel** sets of

f at c are the sets of points $\mathbf{x} \in \Omega$ such that $f(\mathbf{x}) \leq c$ and $f(\mathbf{x}) \geq c$, respectively, i.e.,

$$L_c^- = \{\mathbf{x} \in \Omega : f(\mathbf{x}) \leq c\}$$

$$L_c^+ = \{\mathbf{x} \in \Omega : f(\mathbf{x}) \geq c\}.$$

The sublevel set L_c^- is a smooth manifold with boundary ∂L_c^- and the superlevel set L_c^+ is a smooth manifold with boundary ∂L_c^+ .

theorem. For a Morse function f on Ω , the gradient of f is either zero or orthogonal to the tangent space of the level set L_c at $\mathbf{x} \in L_c$.

The above theorem implies that at a stationary point \mathbf{x}^* , a level set $L_{\mathbf{x}^*}$ is reduced to a single point when \mathbf{x}^* is a local minimum or maximum. Otherwise, the level set may have a singularity such as a self-intersection or a cusp.

2 Overview of Optimization Schemes

2.1 Gradient Descent (GD)

Continuous Model

TODO: Explain the gradient descent phase space and the level set equations for the continuous model.

Scheme

The **gradient descent** line search scheme is

$$T_h(\mathbf{x}) = \mathbf{x} - h\nabla f(\mathbf{x})$$

first-order ($p = 1$) and contractive when $\nabla^2 f \succeq \mu I \succeq 0$.

2.2 Newton's Method (NM)

Continuous Model

TODO: Explain the Newton's method phase space and the level set equations for the continuous model.

Scheme

The **Newton's method** line search scheme is

$$T_h(\mathbf{x}) = \mathbf{x} - h\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$$

second-order ($p = 2$) and contractive when $\nabla^2 f \succeq \mu I \succeq 0$.

2.3 Trust Region Methods (TR)

Continuous Model

TODO: Explain the trust region method phase space and the level set equations for the continuous model.

Scheme

The **trust region method** line search scheme is

$$T_h(\mathbf{x}) = \mathbf{x} + \arg \min_{\boldsymbol{\tau}} m_{\mathbf{x}}(\boldsymbol{\tau})$$

where $m_{\mathbf{x}}(\boldsymbol{\tau}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\tau} \rangle + \frac{1}{2} \langle \boldsymbol{\tau}, \nabla^2 f(\mathbf{x}) \boldsymbol{\tau} \rangle$ is the quadratic approximation of f at \mathbf{x} and $\|\boldsymbol{\tau}\| \leq \Delta$ is the trust region constraint.

2.4 Quasi-Newton Methods (QN)

Continuous Model

TODO : Explain the quasi-newton method phase space and the level set equations for the continuous model.

Scheme

The **quasi-newton method** line search scheme is

$$T_h(\mathbf{x}) = \mathbf{x} - hB\nabla f(\mathbf{x})$$

where $B \approx \nabla^2 f^{-1}(\mathbf{x})$ is a positive-definite approximation of the Hessian. The quasi-newton method is a first-order ($p = 1$) scheme and contractive when $\nabla^2 f \succeq \mu I \succeq 0$.

3 Theory

3.1 Morse Theory

Note that Ω is a bounded subset of \mathbb{R}^n , so its closure $\bar{\Omega} = \Omega \cup \partial\Omega$ is a compact subset of \mathbb{R}^n , by the Heine-Borel Theorem. Also, the boundary $\partial\Omega$ is sufficiently smooth, so we can apply the theory of smooth manifolds. The closure $\bar{\Omega}$ is a compact subset of \mathbb{R}^n and is a smooth manifold with boundary $\partial\Omega$. The interior Ω is an open subset of \mathbb{R}^n and is a smooth manifold.

Definition 16. A point $\mathbf{x}^* \in \Omega$ is a **critical point** of f if the differentiable map $df_p : T_p\Omega \rightarrow \mathbb{R}$ is zero. (Here $T_p\Omega$ is a tangent space of the Manifold M at p .) The set of critical points of f is denoted by $\text{crit}(f)$.

Definition 17. A point $\mathbf{x}^* \in \Omega$ is a **non-degenerate critical point** of f if the Hessian $H_p f$ is non-singular.

Definition 18. The **index** of a *non-degenerate critical point* \mathbf{x}^* is defined to be the dimension of the negative eigenspace of the Hessian $H_p f$.

- local minima at \mathbf{x}^* have index 0.
- local maxima at \mathbf{x}^* have index n .
- saddle points at \mathbf{x}^* have index k where $0 < k < n$.

We reserve the integers $c_0, c_1, \dots, c_i, \dots, c_n$ to denote the number of critical points of index i .

Remark. For each objective function f we are interested in determining the critical points of f

Remark. The **Morse function** is a smooth function $f : \Omega \rightarrow \mathbb{R}$ such that all critical points of f are non-degenerate.

theorem. Let f be a Morse function on Ω , then the Euler characteristic of Ω is given by

$$\chi(\Omega) = \sum_{i=0}^n (-1)^i c_i$$

where c_i is the number of critical points of index i .

Remark. The Euler characteristic $\chi(\Omega)$ is a topological invariant of the manifold Ω and is independent of the choice of Morse function f . The Euler characteristic is a measure of the "shape" of the manifold and can be used to distinguish between different topological spaces. The Euler characteristic may be defined by the alternating sum of the Betti numbers b_i of the manifold Ω

$$\chi(\Omega) = \sum_{i=0}^n (-1)^i b_i$$

where b_i is the i -th Betti number of the manifold Ω .

theorem. (Sard's theorem) Let f be a Morse function on Ω , then the image $f(\text{crit}(f))$ has Lebesgue measure zero in \mathbb{R} .

Remark. We state a particular instance of Sard's theorem for continuous scalar-valued functions f , which was first proved by Anothony P. Morse in 1939. The theorem asserts that the image of the critical points of a Morse function is a set of measure zero in \mathbb{R} . This means that the critical points of a Morse function are "rare" in the sense that they do not form a dense subset of the manifold Ω . Consequently, selecting $\mathbf{x} \in \Omega$ at random will almost never yeild a critical point of f .

Remark. The property that $\mathbf{x}^* \in \Omega$ being a *critical point* of a Morse function f is not dependent of the metric of $\Omega \subset \mathbb{R}^n$ (and consequently, the norm induced by the metric)

3.2 Analysis of Gradient Descent (GD)

theorem. Assume f is ℓ -smooth and α -strongly convex and that $\epsilon > 0$. If we iterate the gradient descent *scheme* with $h = h_0 = \frac{1}{\ell}$ held fixed, i.e.,

$$T_h(\mathbf{x}_k) = \mathbf{x}_k - \frac{1}{\ell} \nabla f(\mathbf{x}_k),$$

then $d(\mathbf{x}_k, \mathbf{x}^*) \leq \epsilon$ for all $k > K$ where K is chosen to satisfy

$$\frac{2\ell}{\alpha} \cdot \log \left(\frac{d(\mathbf{x}_0, \mathbf{x}^*)}{\epsilon} \right) \leq K.$$

Remark. Under ℓ -smoothness and α -strong convexity assumptions in a neighborhood Ω about \mathbf{x}^* , it may be shown directly from the above theorem above that the *GD scheme* converges linearly to the optimal solution \mathbf{x}^* at a rate of

$$\frac{d(\mathbf{x}_k, \mathbf{x}^*)}{d(\mathbf{x}_{k-1}, \mathbf{x}^*)} \leq 1 - \frac{\alpha}{\ell}$$

where $d(\mathbf{x}_k, \mathbf{x}^*)$ is the distance between the current iterate \mathbf{x}_k and the optimal solution \mathbf{x}^* . The convergence rate is linear in the sense that the distance between the current iterate and the optimal solution decreases by a factor of $1 - \frac{\alpha}{\ell}$ at each iteration. (Ref: TODO)

Remark. Convergence to a first-order stationary point trivially implies convergence to a ϵ -first-order stationary point. Similarly, convergence to a second-order stationary point trivially implies convergence to a ϵ -second-order stationary point.

theorem. Assume f is ℓ -smooth, then for any $\epsilon > 0$, if we iterate the GD scheme with $h = h_0 = \frac{1}{\ell}$ held fixed starting from $\mathbf{x}_0 \in \Omega$ where Ω is a neighborhood of \mathbf{x}^* , then the number of iterations K required to achieve the stopping condition $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$ is at most

$$\left\lceil \frac{\ell}{\epsilon^2} (f(\mathbf{x}_0) - f(\mathbf{x}^*)) \right\rceil$$

Remark. **TODO and Questions**

- State how we use theorems in when performing analysis from the results of our experiments.
- What is the relationship between ℓ and α ?
- In practice do we know how to compute ℓ and α ?
- What is the relationship between ℓ and ρ ?
- In practice do we know how to compute ℓ and ρ ?

Remark. JinSaddle.pdf Section 3.1 - Strict Saddle Property