

Lecture 14: Newton and Hessian Geometry

TTIC 31070 / CMSC 35470 / BUSF 36903 / STAT 31015
Convex Optimization

Prof. Zhiyuan Li

Spring 2026

From First Order to Second Order

First-order recap (Lectures 7–13). Local model = linear, with quadratic upper bound:

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

One global constant L ; same step size everywhere.

Lecture 14: Newton's method. Local model = the actual quadratic Taylor expansion:

$$m_x(d) = f(x) + \langle \nabla f(x), d \rangle + \frac{1}{2} \langle \nabla^2 f(x) d, d \rangle.$$

The Newton step is the *exact minimizer* of m_x .

Trade-offs.

- ▶ **Cost.** Each step needs the Hessian (or Hessian-vector products) and a linear solve.
- ▶ **Gain.** Affine-invariant; *quadratic* local convergence; one step exact on quadratics.

Hessian Geometry: A Metric That Moves

Definition 14.1 (Hessian local norms). Let f be twice differentiable at x with $\nabla^2 f(x) \succ 0$.

- ▶ **Primal local norm:** $\|u\|_x := \sqrt{\langle \nabla^2 f(x)u, u \rangle}$, $u \in E$.
- ▶ **Dual local norm:** $\|g\|_{x,*} := \sqrt{\langle g, (\nabla^2 f(x))^{-1}g \rangle}$, $g \in E^*$.

Same norm/dual-norm structure as Lectures 7–10, except the norm *moves with* x .

Convexity test (Lemma 14.1). f convex on open convex $U \Leftrightarrow \nabla^2 f(x) \succeq 0$ for all $x \in U$.

The *strict* version $\nabla^2 f(x) \succ 0$ makes the local quadratic model strictly convex — this is what makes the Newton direction uniquely defined.

Newton Direction and Newton Decrement

Definition 14.2 (Newton direction). At x with $\nabla^2 f(x) \succ 0$:

$$d_f(x) := -(\nabla^2 f(x))^{-1} \nabla f(x).$$

Newton decrement: $\lambda_f(x) := \|\nabla f(x)\|_{x,*} = \sqrt{\langle \nabla f(x), (\nabla^2 f(x))^{-1} \nabla f(x) \rangle}$.

Two key identities.

- ▶ $\|d_f(x)\|_x^2 = \lambda_f(x)^2$ (“Newton-step size in the local metric”)
- ▶ $\langle \nabla f(x), d_f(x) \rangle = -\lambda_f(x)^2$ (“descent rate per unit decrement”)

Damped Newton step. For $t \in (0, 1]$: $x^+ = x + t d_f(x)$. Full step: $t = 1$.

Stopping rule. $\lambda_f(x)^2/2 =$ predicted gap to the local quadratic minimizer.

Quadratic Model and Model Improvement

Lemma 14.2 (Newton minimizes the local quadratic model). With $m_x(d) := f(x) + \langle \nabla f(x), d \rangle + \frac{1}{2} \langle \nabla^2 f(x) d, d \rangle$:

$$d_f(x) = \operatorname{argmin}_d m_x(d) = \operatorname{argmin}_d \left\{ \langle \nabla f(x), d \rangle + \frac{1}{2} \|d\|_x^2 \right\},$$

and

$$m_x(d_f(x)) = f(x) - \frac{1}{2} \lambda_f(x)^2.$$

Three takeaways.

- ▶ Newton is *steepest descent in the Hessian metric* — with the metric chosen by the objective.
- ▶ $\frac{1}{2} \lambda_f(x)^2$ is the *model gap*: how much the local quadratic predicts f will drop.
- ▶ *Caveat*: this is a model gap, not a true gap to f^* — without extra structure (e.g., self-concordance, Lecture 16), it does not certify global suboptimality.

What the Decrement Measures

Lemma 14.3 (Decrement vs gradient norm). If $\mu I \preceq \nabla^2 f(x) \preceq MI$, then

$$\frac{1}{\sqrt{M}} \|\nabla f(x)\|_2 \leq \lambda_f(x) \leq \frac{1}{\sqrt{\mu}} \|\nabla f(x)\|_2.$$

The decrement is the gradient norm *measured in the inverse-Hessian geometry*.

- ▶ Well-conditioned region $\mu \approx M$: $\lambda_f(x) \approx \|\nabla f(x)\|_2 / \sqrt{M}$.
- ▶ $\lambda_f(x) = 0 \Leftrightarrow \nabla f(x) = 0$. So λ_f is a valid stopping signal.

Why use λ_f instead of $\|\nabla f\|_2$?

- ▶ $\|\nabla f\|_2$ depends on the choice of coordinates.
- ▶ λ_f is *affine-invariant* (next slide).

Affine Invariance

Proposition 14.4 (Affine invariance of Newton). Let $T(z) = Az + b$ with A invertible, and $g(z) := f(Tz)$. Setting $x = Tz$:

$$d_g(z) = A^{-1}d_f(x), \quad \lambda_g(z) = \lambda_f(x).$$

Hence full and damped Newton trajectories are linked by T :

$$x_t = T(z_t) \quad \forall t.$$

Why this matters. Coordinate changes / preconditioning do not change Newton's behavior.

- ▶ GD: scaling axes changes step size, conditioning, convergence rate.
- ▶ Newton: identical iterate trajectories under any linear isomorphism.

Slogan. Newton chooses the “right” metric automatically, in every coordinate system.

Affine Invariance: Proof Sketch

Setup. $x = Tz = Az + b$, $g(z) := f(Tz)$. Pullback $A^* : E^* \rightarrow F^*$, $A^*\xi := \xi \circ A$.

Step 1 (chain rule). $\nabla g(z) = A^*\nabla f(x)$, $\nabla^2 g(z) = A^*\nabla^2 f(x)A$. The Hessian of g stays positive definite because A is an isomorphism.

Step 2 (Newton direction). Verify $A^{-1}d_f(x)$ solves the Newton equation for g :

$$\nabla^2 g(z)(A^{-1}d_f(x)) = A^*\nabla^2 f(x)d_f(x) = -A^*\nabla f(x) = -\nabla g(z).$$

By uniqueness of the Newton direction, $d_g(z) = A^{-1}d_f(x)$.

Step 3 (decrement). $(\nabla^2 g(z))^{-1}\nabla g(z) = A^{-1}(\nabla^2 f(x))^{-1}\nabla f(x)$, hence

$$\lambda_g(z)^2 = \langle \nabla g(z), (\nabla^2 g(z))^{-1}\nabla g(z) \rangle = \langle \nabla f(x), (\nabla^2 f(x))^{-1}\nabla f(x) \rangle = \lambda_f(x)^2.$$

Step 4 (trajectories). $T(z_t + \eta d_g(z_t)) = x_t + \eta d_f(x_t)$. The Armijo test is invariant because $\langle \nabla g, d_g \rangle = \langle \nabla f, d_f \rangle$ and $g(z + \eta d_g) = f(x + \eta d_f)$. Hence backtracking accepts identical step sizes. \square

Quadratic Exactness: One Step

Theorem 14.5 (Quadratic exactness). For $f(x) = \frac{1}{2}x^\top Qx + b^\top x + c$ with $Q \succ 0$, one full Newton step from any x lands at the unique minimizer:

$$x + d_f(x) = -Q^{-1}b = x^*.$$

Why. $\nabla f(x) = Qx + b$, $\nabla^2 f(x) = Q$, so $d_f(x) = -Q^{-1}(Qx + b) = -x - Q^{-1}b$, and $x + d_f(x) = -Q^{-1}b$.

What this tells us. On a quadratic, the local quadratic model *is* the global function. The Hessian-Lipschitz constant is 0. *Newton has nothing to approximate.*

For non-quadratic f : each step is exact on the local quadratic; Hessian-Lipschitz controls the discrepancy from one step to the next. This is the source of *quadratic* local convergence (next slides).

Example: Ridge Least Squares Solved in One Step

Problem. $f(w) = \frac{1}{2}\|Aw - b\|_2^2 + \frac{\sigma}{2}\|w\|_2^2, \quad \sigma > 0.$

Newton ingredients.

$$\nabla f(w) = A^\top(Aw - b) + \sigma w, \quad \nabla^2 f(w) = A^\top A + \sigma I.$$

One Newton step from any w .

$$w + d_f(w) = (A^\top A + \sigma I)^{-1} A^\top b.$$

This is exactly the ridge-regression closed form.

Contrast with gradient descent.

$$w^+ = w - \eta(A^\top(Aw - b) + \sigma w).$$

Same residual; *scalar* η vs. *matrix* $(A^\top A + \sigma I)^{-1}$. Newton replaces the step size with the curvature inverse.

Local Quadratic Convergence

Theorem 14.6 (Local quadratic convergence). Suppose $\nabla f(x^*) = 0$, $\nabla^2 f(x^*) \succ 0$, and $\nabla^2 f$ is ρ -Lipschitz on a ball $B(x^*, r)$. Then for x_0 close enough,

$$\|x_{t+1} - x^*\|_2 \leq C \|x_t - x^*\|_2^2.$$

Doubling-precision-per-step. Once in the basin of attraction, the error *squares* every step.

- ▶ One step: $10^{-2} \rightarrow 10^{-4} \rightarrow 10^{-8} \rightarrow 10^{-16}$. **Doubling-precision-per-step.**
- ▶ Compare GD: $\|x_{t+1} - x^*\|_2 \leq \rho \|x_t - x^*\|_2$ for some $\rho < 1$ — linear, not quadratic.

Caveat. This is a *local* statement. Outside the basin, Newton can diverge or take unboundedly large steps. Globalization (next part) handles this.

Local Convergence: Proof Sketch

Setup: $H_* := \nabla^2 f(x^*) \succ 0$. Near x^* , $\nabla^2 f(x) \succeq (\mu/2)I$, so $\|\nabla^2 f(x)^{-1}\|_2 \leq 2/\mu$.

Newton residual. Write $e := x - x^*$, $x^+ = x + d_f(x)$:

$$x^+ - x^* = \nabla^2 f(x)^{-1}(\nabla^2 f(x)e - (\nabla f(x) - \nabla f(x^*))).$$

Hessian-Lipschitz: replace $\nabla^2 f(x^* + se)$ with $\nabla^2 f(x)$.

$$\nabla^2 f(x)e - (\nabla f(x) - \nabla f(x^*)) = \int_0^1 (\nabla^2 f(x) - \nabla^2 f(x^* + se))e \, ds.$$

ρ -Lipschitz \Rightarrow integrand has norm $\leq \rho(1-s)\|e\|_2^2$. Integrating: $\leq \frac{\rho}{2}\|e\|_2^2$.

Combine. $\|x^+ - x^*\|_2 \leq \frac{\rho}{\mu}\|e\|_2^2$. Set $C = \rho/\mu$. □

Heart. The integral remainder of Taylor's theorem; Hessian-Lipschitz controls it; $\|H^{-1}\| \leq 2/\mu$ converts norm-of-residual into norm-of-error.

Logistic Regression: Newton as IRLS

Problem. $f(w) = \sum_{i=1}^m \log(1 + \exp(-b_i a_i^\top w)) + \frac{\sigma}{2} \|w\|_2^2$.

Hessian. Let $q_i(w) := 1/(1 + \exp(b_i a_i^\top w))$. Then

$$\nabla^2 f(w) = A^\top D(w)A + \sigma I, \quad D(w) = \text{diag}(q_i(w)(1 - q_i(w))).$$

Newton step. Solve

$$(A^\top D(w)A + \sigma I)d = -\nabla f(w).$$

This is **iteratively reweighted least squares (IRLS)**. Weights $q_i(1 - q_i)$:

- ▶ Large near uncertain examples ($q_i \approx 1/2$).
- ▶ Small for confident examples ($q_i \approx 0$ or ≈ 1).

Curvature analysis. $\nabla^2 f \succeq \sigma I$ (regularization gives strong convexity). $\nabla^2 f$ is Hessian-Lipschitz with constant $\propto \sum_i \|a_i\|_2^3$.

Cost Model: Iterations vs Wall Clock

Compare per iteration.

- ▶ Gradient step: $O(n)$ + first-order oracle.
- ▶ Newton step: Hessian (or H-vector products) + linear solve $Hd = -g$.

Cost of the linear solve.

- ▶ Dense direct (Cholesky): $O(n^3)$ flops, $O(n^2)$ memory.
- ▶ Iterative (CG): $O(n^2) \times$ (effective condition number).
- ▶ Structured (sparse, low-rank, Kronecker): can be near-linear.

When Newton wins.

- ▶ Small/medium n (linear solve cheap).
- ▶ Hessian has structure (sparse, low-rank, $A^\top DA + \sigma I$).
- ▶ Need high accuracy (log log phase pays off).
- ▶ Ill-conditioned Hessian where GD struggles.

When GD/AGD wins. Large-scale learning where n^2 memory is impossible and modest accuracy suffices.

Globalization: Damped Newton with Backtracking

Algorithm (Backtracking Damped Newton). Parameters $\alpha \in (0, 1/2)$, $\beta \in (0, 1)$.

1. Compute $d_t = d_f(x_t)$, $\lambda_t = \lambda_f(x_t)$.
2. If $\lambda_t = 0$: stop.
3. Set $\eta = 1$.
4. **While** $x_t + \eta d_t \notin \text{dom } f$ **or** $f(x_t + \eta d_t) > f(x_t) + \alpha \eta \langle \nabla f(x_t), d_t \rangle$: $\eta \leftarrow \beta \eta$.
5. $x_{t+1} = x_t + \eta d_t$.

Two checks.

- ▶ *Domain test*: keep x_t inside $\text{dom } f$. Critical for log barriers (Lecture 16).
- ▶ *Armijo test*: actual descent matches α -fraction of the predicted descent.

Slogan. Backtracking does *not* change the Newton direction. It only checks whether the local quadratic model is reliable at the current scale.

Why Backtracking Terminates

Proposition 14.7 (Descent + finite Armijo). If f is differentiable and $\nabla^2 f(x) \succ 0$, $\nabla f(x) \neq 0$:

- ▶ Newton direction is strict descent: $\langle \nabla f(x), d_f(x) \rangle = -\lambda_f(x)^2 < 0$.
- ▶ For every Armijo $\alpha \in (0, 1)$, there is $\bar{\eta} > 0$ such that every $\eta \in (0, \bar{\eta}]$ satisfies the domain and Armijo conditions.

Backtracking eventually accepts: tested values $\eta = 1, \beta, \beta^2, \dots$ reach $\bar{\eta}$ in finite steps.

Caveat: positive curvature is essential. If $\nabla^2 f(x)$ has a negative direction, $d_f(x)$ may point to a *saddle* or *maximum* of the quadratic model. The descent identity uses positive definiteness. **Today's lecture is the convex regime.**

Nonconvex variants modify H , use trust regions, or regularize: see cubic regularization below.

Two-Phase Damped Newton Theorem

Theorem 14.8 (Globalized convergence). Assume $\mu I \preceq \nabla^2 f$ and $\nabla^2 f$ is ρ -Lipschitz globally; $f^* > -\infty$. Define $\lambda_N := \min\{1, 3(1 - 2\alpha)\} \mu^{3/2}/\rho$ and $\gamma_N := \alpha\beta\lambda_N^2$. Then:

1. **Damped phase** ($\lambda_t \geq \lambda_N$): $f(x_{t+1}) \leq f(x_t) - \gamma_N$.
2. **Quadratic phase** ($\lambda_t < \lambda_N$): backtracking accepts $\eta = 1$, and $\lambda_{t+1} \leq \frac{\rho}{2\mu^{3/2}} \lambda_t^2 \leq \frac{1}{2} \lambda_t$.

Total complexity.

- ▶ Damped phase iterations: $\leq \lceil (f(x_0) - f^*)/\gamma_N \rceil$ — each guarantees a constant drop.
- ▶ Quadratic phase: $\log \log(1/\varepsilon)$ iterations to reach decrement $\leq \varepsilon$.

Slogan. Constant-drop damped phase + doubling-precision quadratic phase.

Helper Lemma: Hessian-Lipschitz Taylor Estimate

Lemma 14.4 (Hessian-Lipschitz Taylor). If $\nabla^2 f$ is ρ -Lipschitz on \mathbb{R}^n , then for all x, s :

$$f(x + s) \leq f(x) + \langle \nabla f(x), s \rangle + \frac{1}{2} s^\top \nabla^2 f(x) s + \frac{\rho}{6} \|s\|_2^3.$$

Proof. Let $\phi(\theta) := f(x + \theta s)$, $\phi''(\theta) = s^\top \nabla^2 f(x + \theta s) s$. Integral Taylor:

$$f(x + s) = f(x) + \phi'(0) + \int_0^1 (1 - \theta) \phi''(\theta) d\theta.$$

Add/subtract $s^\top \nabla^2 f(x) s$. Hessian-Lipschitz: $s^\top (\nabla^2 f(x + \theta s) - \nabla^2 f(x)) s \leq \rho \theta \|s\|_2^3$.

Then $\int_0^1 \theta(1 - \theta) d\theta = \frac{1}{6}$. \square

What this says. The L -smoothness upper bound ($\frac{L}{2} \|s\|_2^2$) is replaced by the *true* quadratic at x + a cubic correction — sharper whenever the Hessian varies slowly.

Two-Phase Proof (1/3): Apply Lemma 14.4 to the Newton Direction

Setup. Write $g := \nabla f(x)$, $d := d_f(x)$, $\lambda := \lambda_f(x)$. From Lemma 14.2:

$$g^\top d = -\lambda^2, \quad d^\top \nabla^2 f(x) d = \lambda^2.$$

Bound on $\|d\|_2$. $\nabla^2 f \succeq \mu I \Rightarrow \mu \|d\|_2^2 \leq \lambda^2 \Rightarrow \|d\|_2 \leq \lambda/\sqrt{\mu}$.

Apply Lemma 14.4 with $s = td$, $t \in [0, 1]$.

$$f(x + td) \leq f(x) + t \underbrace{g^\top d}_{=-\lambda^2} + \frac{t^2}{2} \underbrace{d^\top \nabla^2 f(x) d}_{=\lambda^2} + \frac{\rho t^3}{6} \underbrace{\|d\|_2^3}_{\leq \lambda^3/\mu^{3/2}}.$$

Boxed Taylor estimate (along the Newton direction).

$$f(x + td) \leq f(x) + t\lambda^2 \left(-1 + \frac{t}{2} + \frac{\rho t \lambda}{6\mu^{3/2}} \right), \quad 0 \leq t \leq 1.$$

The bracket controls whether Armijo accepts step t .

Two-Phase Proof (2/3): Damped Phase ($\lambda \geq \lambda_N$)

Goal. Show backtracking accepts a step that gives a uniform objective drop $\gamma_N = \alpha\beta\lambda_N^2$.

Pick the trial step $\hat{t} := \frac{\lambda_N}{\lambda} \in (0, 1]$. For any $t \in (0, \hat{t}]$, $t\lambda \leq \lambda_N$, so the boxed Taylor estimate gives

$$f(x + td) \leq f(x) - t\lambda^2 \left(\underbrace{1 - \frac{t}{2}}_{\leq 1/2} - \underbrace{\frac{\rho t \lambda}{6\mu^{3/2}}}_{\leq \rho\lambda_N/(6\mu^{3/2})} \right).$$

Tune λ_N . Recall $\lambda_N \leq 3(1 - 2\alpha)\mu^{3/2}/\rho$, so $\frac{\rho\lambda_N}{6\mu^{3/2}} \leq \frac{1-2\alpha}{2}$. Hence

$$\frac{t}{2} + \frac{\rho t \lambda}{6\mu^{3/2}} \leq \frac{1}{2} + \frac{1-2\alpha}{2} = 1 - \alpha,$$

giving $f(x + td) \leq f(x) - \alpha t \lambda^2 = f(x) + \alpha t g^\top d$. (**Armijo at every $t \leq \hat{t}$.**)

Backtracking accepts $t_{\text{acc}} \geq \beta\hat{t}$. Hence $f(x) - f(x^+) \geq \alpha\beta\lambda_N\lambda \geq \alpha\beta\lambda_N^2 = \gamma_N$.

Two-Phase Proof (3/3): Quadratic Phase ($\lambda < \lambda_N$)

Full step accepted. The boxed Taylor estimate at $t = 1$ with $1 \cdot \lambda < \lambda_N$ gives the same Armijo conclusion as in the damped phase, so backtracking accepts $\eta = 1$. Set $x^+ := x + d$.

Newton residual at x^+ . Since $g + \nabla^2 f(x)d = 0$, the fundamental theorem of calculus gives

$$\nabla f(x^+) = \int_0^1 (\nabla^2 f(x + sd) - \nabla^2 f(x)) d ds.$$

Bound the residual. Hessian-Lipschitz gives $\|\nabla^2 f(x + sd) - \nabla^2 f(x)\|_2 \leq \rho s \|d\|_2$. So

$$\|\nabla f(x^+)\|_2 \leq \int_0^1 \rho s \|d\|_2^2 ds = \frac{\rho}{2} \|d\|_2^2 \leq \frac{\rho}{2\mu} \lambda^2.$$

Convert to decrement. $\nabla^2 f(x^+) \succeq \mu I \Rightarrow \lambda_f(x^+) \leq \frac{1}{\sqrt{\mu}} \|\nabla f(x^+)\|_2$. Combine:

$$\lambda_f(x^+) \leq \frac{\rho}{2\mu^{3/2}} \lambda^2 \leq \frac{1}{2} \lambda \quad (\text{since } \lambda < \lambda_N).$$



Aside: Cubic Regularization

Yu. Nesterov & B. T. Polyak, “Cubic regularization of Newton method and its global performance,” *Math. Program.* 108 (2006).



Yurii Nesterov



Boris T. Polyak

Idea. If f has ρ -Lipschitz Hessian, Taylor gives the cubic upper model

$$f(x + s) \leq f(x) + \langle \nabla f(x), s \rangle + \frac{1}{2} s^\top \nabla^2 f(x) s + \frac{\rho}{6} \|s\|_2^3.$$

Method. Pick $M \geq \rho$. At each step, choose s_t minimizing the cubic upper model:

$$s_t \in \operatorname{argmin}_s \left\{ \langle \nabla f(x_t), s \rangle + \frac{1}{2} s^\top \nabla^2 f(x_t) s + \frac{M}{6} \|s\|_2^3 \right\}.$$

Rate (convex case). If $\|x - x^*\|_2 \leq R$ on the level set $\{f \leq f(x_0)\}$:

$$f(x_T) - f^* = O(MR^3/T^2).$$

Bonuses. Cubic term makes the model coercive even with indefinite Hessian — handles nonconvex Newton too.

Summary

Newton = exact minimizer of the local quadratic Taylor model.

- ▶ Hessian local norm $\|u\|_x$ + dual local norm $\|g\|_{x,*}$ define a metric that moves with x .
- ▶ Newton direction $d_f = -(\nabla^2 f)^{-1} \nabla f$; decrement $\lambda_f = \|\nabla f\|_{x,*}$.
- ▶ Affine-invariant under invertible affine changes of variables.

Behavior.

- ▶ Quadratic objective: *exact* in one step.
- ▶ Locally near a nondegenerate min: *quadratic* convergence (ρ/μ rate).
- ▶ Globally with $\mu I \preceq \nabla^2 f$ and $\nabla^2 f$ ρ -Lipschitz: damped phase + quadratic phase (Theorem 14.8).

Cost. Linear solve per iteration. Worth it for high accuracy and structured Hessians.

Next. Lecture 15: self-concordant Newton calculus — replace the Euclidean Hessian-Lipschitz constant by an affine-invariant local condition that survives near boundaries (e.g., $-\log t$). This is the right setting for interior-point methods in Lecture 16.

Bibliographic Notes

- ▶ **S. Boyd & L. Vandenberghe, Convex Optimization, Cambridge (2004)**. Chapter 9 covers Newton's method, the damped variant, and the two-phase analysis used today.
- ▶ **Yu. Nesterov, Introductory Lectures on Convex Optimization, Kluwer (2004)**. Chapter 1.2: classical Newton; Chapter 4: the self-concordant generalization that Lecture 16 will pick up.

Cubic regularization (Nesterov & Polyak 2006) is cited inline on the cubic-regularization slide.