

CS450: Numerical Analysis Optimization

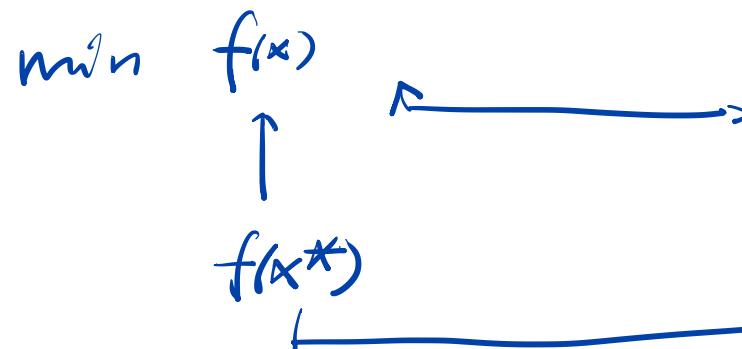
Howard Hao Yang

Assistant Professor, ZJU-UIUC Institute

15/12/2023

Today's Agenda

- Higher-order methods
 - Newton method
 - Quasi Newton method



GD:

first-order method

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

first-order derivative of f

$\nabla f(x_t)$

$\nabla^2 f(x_t)$

$x_t \rightarrow \nabla^2 f(x_t)$

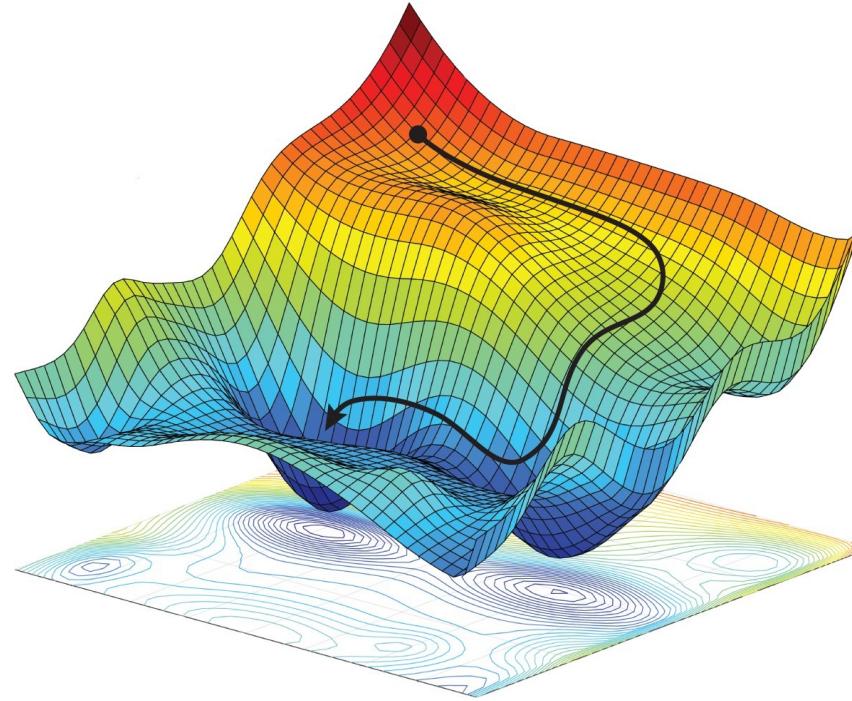
find $\nabla f(x) = 0$?

$\nabla f(x^*) = 0$

x^*

$\nabla f(x^*) = 0$

Optimization



- The goal is to accomplish the following optimization problem

$$\min_x f(x) \quad \text{subject to} \quad x \in R^n$$

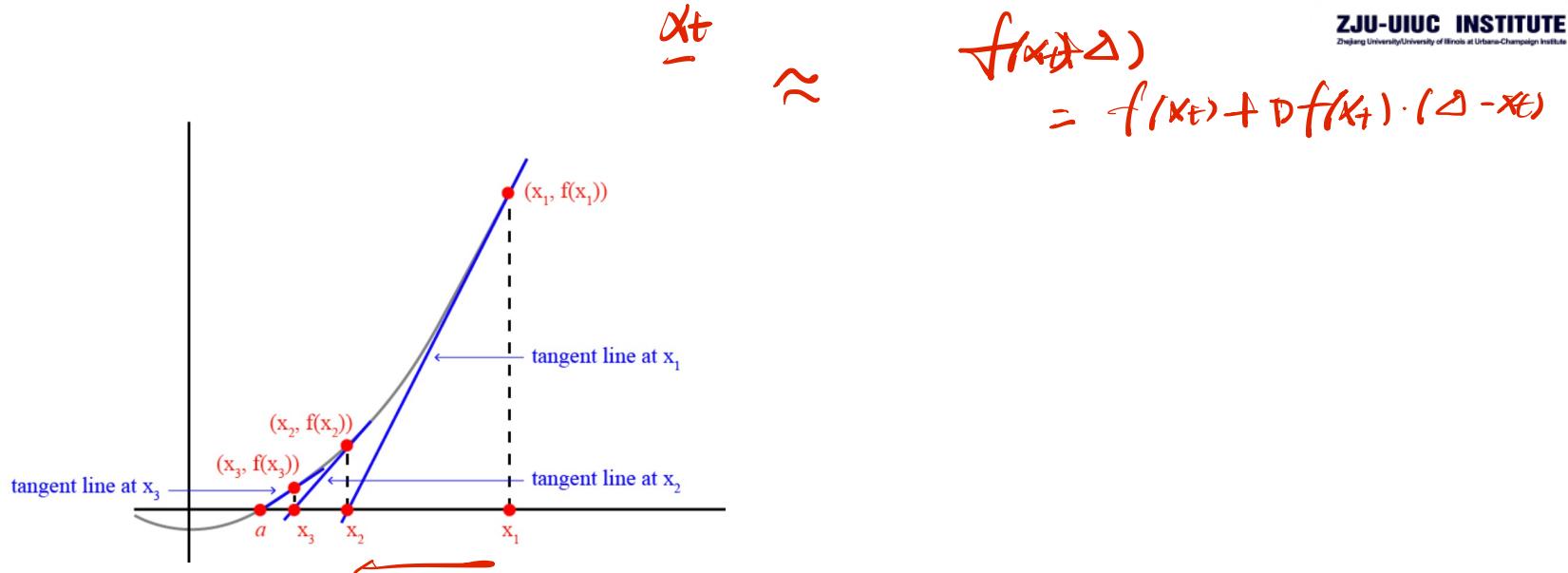
Optimality Conditions

- Suppose we found a candidate \mathbf{x}^* , how do we know it actually is (the) one?
- n -dimension case:
 - Necessary: $\nabla f(\mathbf{x}^*) = \mathbf{0}$ (i.e., \mathbf{x}^* is an extremal point)
 - Sufficient: $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $H_f(\mathbf{x}^*)$ is **positive semidefinite** (i.e., \mathbf{x}^* is a local minimum), where $H_f(\mathbf{x}^*)$ denotes the **Hessian matrix**

Optimality Conditions

- Suppose we found a candidate \mathbf{x}^* , how do we know it actually is (the) one?
- n -dimension case:
 - Necessary: $\nabla f(\mathbf{x}^*) = \mathbf{0}$ (i.e., \mathbf{x}^* is an extremal point)
 - Sufficient: $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $H_f(\mathbf{x}^*)$ is **positive semidefinite** (i.e., \mathbf{x}^* is a local minimum), where $H_f(\mathbf{x}^*)$ denotes the **Hessian matrix**
- Can we convert the optimization problem into solving nonlinear equation systems as $\nabla f(\mathbf{x}) = \mathbf{0}$, and then leverage techniques for root finding in nonlinear systems to solve for it?

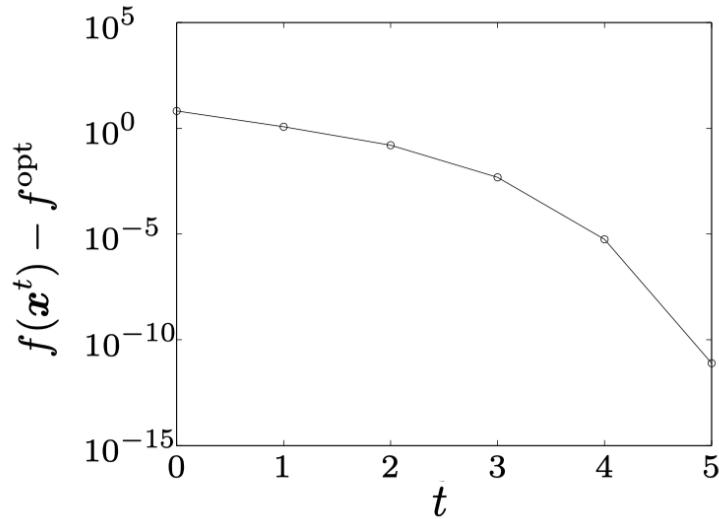
Recall: Newton's Method



- To find the root of $f(x) = 0$, algorithm details:

- 1) start with an initial guess x_0
- 2) keep iterating $x_{k+1} = g(x_k) = x_k - \frac{f(x_k)}{f'(x_k)}$

Newton's Method in Optimization



Steepest: fast

Bad:

①

$$\nabla^2 f(x_t)$$

$$\approx \mathcal{O}(n^2)$$

②

$$(\nabla^2 f(x_t))^{-1}$$

$$\approx \mathcal{O}(n^3)$$

- Given an unconstrained optimization problem,

$$\min_x f(x) \quad \text{subject to} \quad x \in R^n$$

- Newton's method runs as (solving $\nabla f(x) = 0$)

- Start with an initial point x_0 and update as follows

$$x_{t+1} = x_t - (\nabla^2 f(x_t))^{-1} \nabla f(x_t)$$

Newton's Method in Optimization (Cont'd)

- Quadratic convergence: attains ϵ accuracy within $O(\log \log \frac{1}{\epsilon})$ iterations
 - If the objective is quadratic, e.g., the following, Newton's method solves it in one iteration (why?)

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^T \mathbf{A} (\mathbf{x} - \mathbf{x}^*)$$

- Issues:
 - Need to **compute** a Hessian matrix $\nabla^2 f(\mathbf{x}_t)$ at every iteration
 - Need to **invert** the Hessian matrix $\nabla^2 f(\mathbf{x}_t)$ at every iteration
- Each of these two steps can be generally expensive

Quasi Newton Methods (1)

- Compare the following:

Gradient descent.

cheap ✓
slow

Quasi Newton method.

Newton method.

expensive
fast ✓

$$x_{t+1} = x_t - \eta_t \cdot I \cdot \nabla f(x_t)$$

$$x_{t+1} = x_t - \eta_t \cdot \tilde{H}_t \cdot \nabla f(x_t)$$

$$\left[\begin{array}{c} \\ \\ \end{array} \right]$$

$$x_{t+1} = x_t - \eta_t \cdot (\nabla^2 f(x_t))^{-1} \cdot \nabla f(x_t)$$

$$\left[\begin{array}{ccc} x & x & x \\ x & \vdots & x \\ x & x & x \end{array} \right] \in \mathbb{R}^{n \times n}$$

Quasi Newton Methods (2)

- Questions:

1) How are adaptive optimization methods (RMSprop, AdaGrad, Adam) linked to the previous form?

$$\left\{ \begin{array}{l} V_{t+1} = V_t + (\nabla f(x_t))^2 \\ x_{t+1} = x_t - \eta \frac{\nabla f(x_t)}{\sqrt{V_t}} \end{array} \right.$$

$$x_{t+1} = x_t - \eta_t \cdot \begin{bmatrix} \frac{1}{\sqrt{V_{1,t}}} & 0 & \dots & 0 \\ 0 & \frac{1}{\sqrt{V_{2,t}}} & & \\ & & \ddots & \\ & & & \frac{1}{\sqrt{V_{n,t}}} \end{bmatrix} \nabla f(x_t)$$

2) Can we do something similar?

Quasi Newton Methods (3)

$$f(x) : \mathbb{R} \rightarrow \mathbb{R}$$

$$x_{t+1} = x_t - \eta \cdot \frac{\frac{f'(x_t) \cdot}{(f'(x_t) - f'(x_{t-1}))}}{x_t - x_{t-1}} \approx f''(x_t)$$

ZJU-UIUC INSTITUTE
Zhejiang University/University of Illinois at Urbana-Champaign Institute

- Key idea: Approximate the Hessian matrix using only the gradient

information (without expensive computations):

$$\begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix}$$

$$x_{t+1} = x_t - \eta_t H_t \nabla f(x_t)$$

$$\begin{bmatrix} ? \\ ? \end{bmatrix}_n$$

where H_t is a surrogate for $(\nabla^2 f(x_t))^{-1}$

$$H_t \approx (\nabla^2 f(x_t))^{-1}$$

- Q: How to find a good approximation H_t for $(\nabla^2 f(x_t))^{-1}$ using
 - using only gradient information
 - using limited memory: $\underline{f''(x_t)} \leftarrow \text{poor estimate.}$ of $f(x_t), f'(x_t), f'(x_{t-1}), \dots$
 - achieving super-linear convergence

$$f(x_t) \xrightarrow{?} f'(x_t) \xrightarrow{?} f''(x_t), f'(x_{t-1})$$

$$\xrightarrow{?} \frac{f'(x_t) - f'(x_{t-1})}{x_t - x_{t-1}}$$

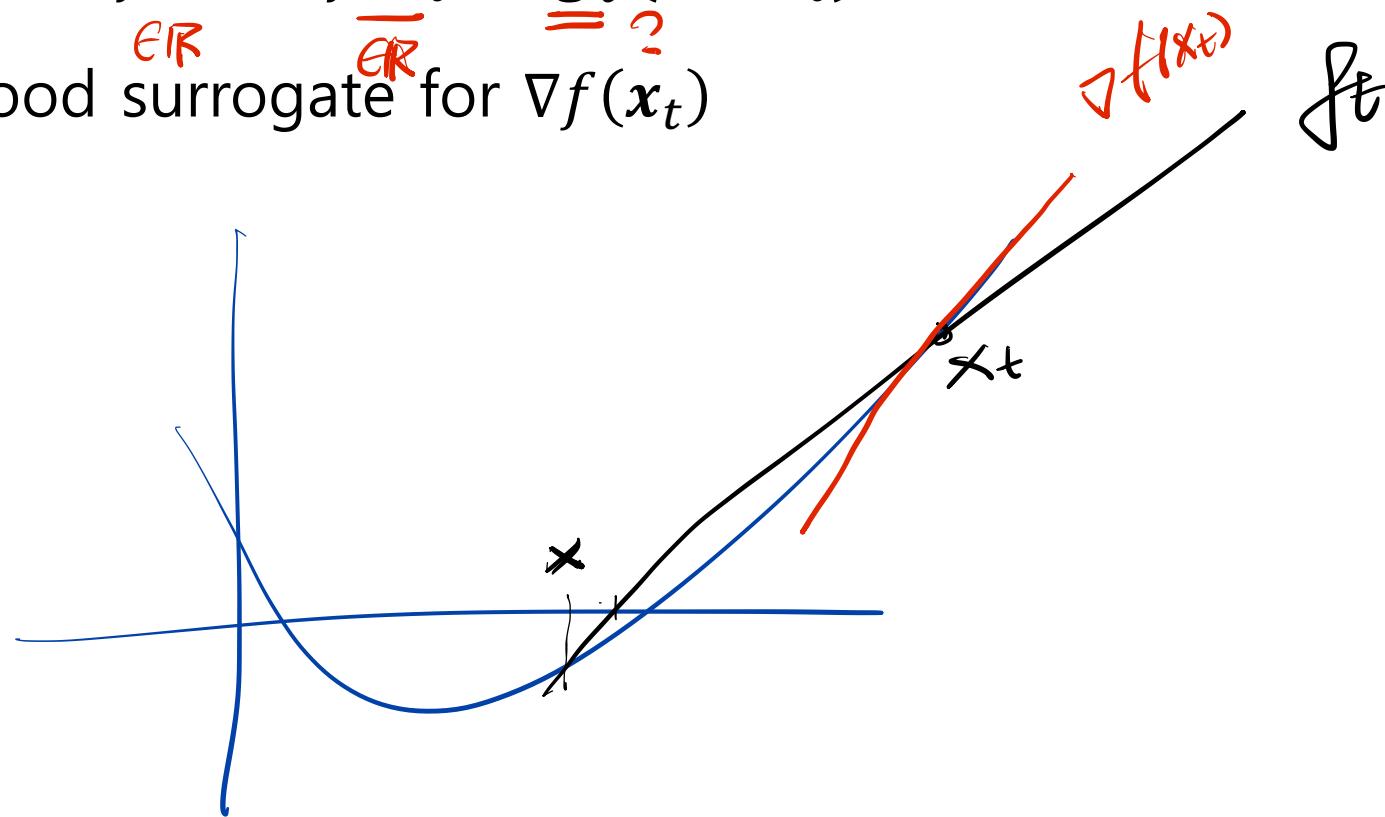
$$\text{of } f(x_t), f'(x_t), f'(x_{t-1}), f'(x_{t-2}), \dots$$

Approximating Hessian

- Approximating gradient: if x is close to x_t , some vector g_t satisfies

$$f(x) = f(x_t) + \underbrace{g_t^T(x - x_t)}_{\text{ER}} \stackrel{\text{ER}}{=} ?$$

then g_t can be a good surrogate for $\nabla f(x_t)$



Approximating Hessian (Cont'd)

$$\nabla f(\underline{x}) = \nabla f(\underline{x}_t) + \begin{bmatrix} ? & ? & ? \\ ? & ? & ? \\ ? & ? & ? \end{bmatrix}$$

- Likewise, if \underline{x} is close to \underline{x}_t , and some matrix G_{t+1} satisfies

$$\nabla f(\underline{x}) = \nabla f(\underline{x}_{t+1}) + G_{t+1}^T(\underline{x} - \underline{x}_{t+1})$$

then G_{t+1} can be a good surrogate for $\nabla^2 f(\underline{x}_{t+1})$

- By using information from previous iterations, the following holds

$$\nabla f(\underline{x}_t) = \nabla f(\underline{x}_{t+1}) + H_{t+1}^{-1}(\underline{x}_t - \underline{x}_{t+1}) \quad (*) \quad H_{t+1}^{-1} \approx \nabla^2 f(\underline{x}_{t+1})$$

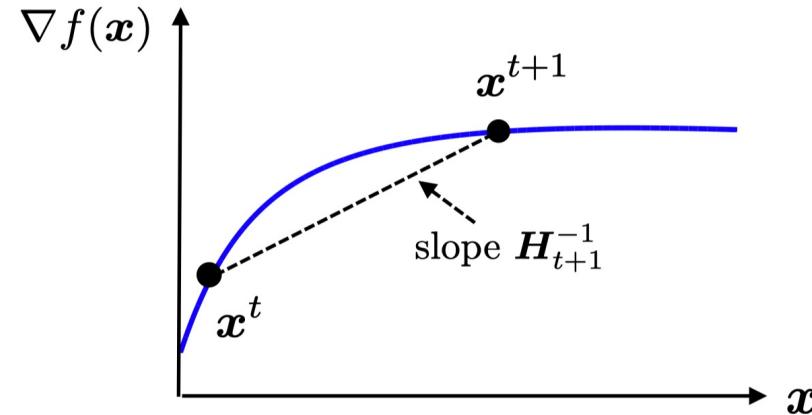
- H_{t+1}^{-1} can be regarded as G_{t+1}^{-1} evaluated at \underline{x}_t and \underline{x}_{t+1}

$$H_{t+1}^{-1} \approx \frac{\nabla f(\underline{x}_t) - \nabla f(\underline{x}_{t+1})}{\underline{x}_t - \underline{x}_{t+1}} \quad \left\{ \underline{x}_t, \nabla f(\underline{x}_t), \underline{x}_{t+1}, \nabla f(\underline{x}_{t+1}) \right\}$$

$$\in \mathbb{R}^n \quad \in \mathbb{R}^n \quad \in \mathbb{R}^n \quad \in \mathbb{R}^n$$

$$\approx \nabla^2 f(\underline{x}_{t+1})$$

Approximating Hessian (Cont'd)



- Rearranging terms, arrive at the following (which is termed the **secant equation**) $\Leftrightarrow (x_{t+1} - x_t) = \boxed{H_{t+1}} (\nabla f(x_{t+1}) - \nabla f(x_t))$
- Requires H_{t+1}^{-1} to map the displacement $x_{t+1} - x_t$ into the changes of $\nabla f(x_{t+1}) - \nabla f(x_t)$

Secant Equation (1)

$$x_1 + 9x_2 = 10$$

- The secant equation

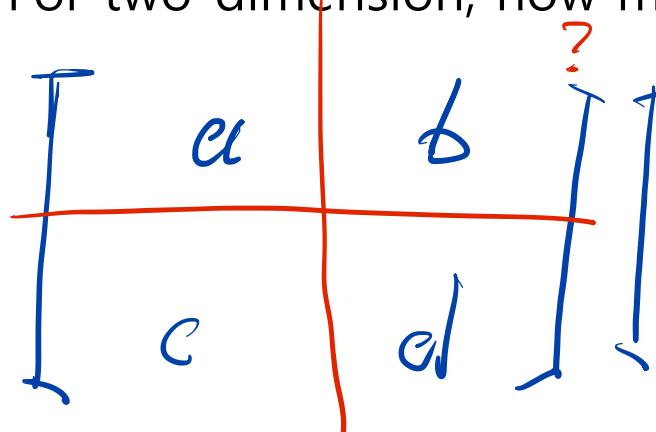
$$f(x) : \mathbb{R} \rightarrow \mathbb{R}$$

$$H_{t+1} = \frac{x_{t+1} - x_t}{f(x_{t+1}) - f(x_t)}$$

$$f(x_1, x_2) : \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$\underline{H_{t+1}} (\nabla f(\underline{x_{t+1}}) - \nabla f(\underline{x_t})) = \underline{x_{t+1}} - \underline{x_t}$$

- For one-dimension, how many possible solutions of H_{t+1} to the above equation?
- For two-dimension, how many possible solutions of H_{t+1} to the above equation?



$$\begin{bmatrix} \frac{\partial f}{\partial x_1}(x_{t+1}) - \frac{\partial f}{\partial x_1}(x_t) \\ \frac{\partial f}{\partial x_2}(x_{t+1}) - \frac{\partial f}{\partial x_2}(x_t) \end{bmatrix} = \begin{bmatrix} x_{1,t+1} - x_{1,t} \\ x_{2,t+1} - x_{2,t} \end{bmatrix}$$

- Admits infinitely many solutions, since degrees of freedom $O(n^2)$ in choosing H_{t+1} far exceeds number of constraints n

Secant Equation (2)

- The secant equation

$$\frac{\partial^2 f}{\partial x_1 \partial x_2}$$

$$\frac{\partial^2 f}{\partial x_2 \partial x_1}$$

$$\underline{H_t} (\nabla f(x_t) - \nabla f(x_{t-1})) = x_t - x_{t-1} \quad (*)'$$

$$\underline{H_{t+1}} (\nabla f(x_{t+1}) - \nabla f(x_t)) = x_{t+1} - x_t \quad (*)$$

- Admits infinitely many solutions, since degrees of freedom $O(n^2)$ in choosing H_{t+1} far exceeds number of constraints n

- Which H_{t+1} shall we choose? (solution to $(*)$)

① We want H_{t+1} to be symmetric (why?) : $H_{t+1} = H_{t+1}^T$

$$\downarrow \|H_{t+1} - H_t\|$$

② We want H_{t+1} to be "close" to H_t : $\min_H \|H - H_t\|$

③ If H_t is positive definite, we want H_{t+1} to also be positive definite

$$w^T \nabla^2 f(x_t) w \geq 0$$

Formalizing the Selection

$$\{x_t, x_{t+1}, \nabla f(x_t), \nabla f(x_{t+1})\}$$

- Chose an update by minimizing the closeness to H_t

$$H = H_t + uv^T$$

$$H^T = H_t + uu^T$$

$$+ b v v^T$$

$$\begin{aligned} & \min_H \|H - H_t\| \\ \text{s.t. } & H = H^T \\ & Hy_t = s_t \end{aligned}$$

(symmetric condition)

(secant equation)

where $y_t = \nabla f(x_{t+1}) - \nabla f(x_t)$ and $s_t = x_{t+1} - x_t$

- Exploit past information regarding H_t
- Each choice of norm gives a different H_{t+1}
- How to solve this?

Formalizing the Selection (Cont'd)

- The secant equation stipulates the following

$$\mathbf{H}_{t+1}^{-1} \mathbf{s}_t = \mathbf{y}_t$$

- Suppose we already have \mathbf{H}_t^{-1} and update \mathbf{H}_{t+1}^{-1} as follows (why?)

$$\mathbf{H}_{t+1}^{-1} = \underbrace{\mathbf{H}_t^{-1}}_{\times} + \underbrace{a\mathbf{u}\mathbf{u}^T + b\mathbf{v}\mathbf{v}^T}_{\times}$$

- Substitute into the equation above

$$(\mathbf{H}_t^{-1} + a\mathbf{u}\mathbf{u}^T + b\mathbf{v}\mathbf{v}^T) \mathbf{s}_t = \mathbf{y}_t$$

Formalizing the Selection (Cont'd)

$$\Delta f(x_t) - \Delta f(x_{t+1})$$

- Substitute into the equation above gives

$$y_t - H_t^{-1} s_t = a u (u^T s_t) + b v (v^T s_t)$$

setting $u = y_t$ and $v = H_t^{-1} s_t$ allows us to solve for a and b

$$a \cdot (y_t^T s_t) = 1 \Rightarrow a = \frac{1}{y_t^T s_t}$$

$$b (s_t^T H_t^{-1} s_t) = -1 \Rightarrow b = -\frac{1}{s_t^T H_t^{-1} s_t}$$

Formalizing the Selection (Cont'd)

- The secant equation stipulates the following

$$\mathbf{H}_{t+1}^{-1} \mathbf{s}_t = \mathbf{y}_t$$

- Suppose we already have \mathbf{H}_t^{-1} and update \mathbf{H}_{t+1}^{-1} as follows (why?)

$$\mathbf{H}_{t+1}^{-1} = \mathbf{H}_t^{-1} + a\mathbf{u}\mathbf{u}^T + b\mathbf{v}\mathbf{v}^T$$

- Substitute into the equation above gives

$$\mathbf{y}_t - \mathbf{H}_t^{-1} \mathbf{s}_t = a\mathbf{u}(\mathbf{u}^T \mathbf{s}_t) + b\mathbf{v}(\mathbf{v}^T \mathbf{s}_t)$$

setting $\mathbf{u} = \mathbf{y}_t$ and $\mathbf{v} = \mathbf{H}_t^{-1} \mathbf{s}_t$ allows us to solve for a and b , given as

$$a = \frac{1}{\mathbf{y}_t^T \mathbf{s}_t}, \quad b = -\frac{1}{\mathbf{s}_t^T \mathbf{H}_t^{-1} \mathbf{s}_t}$$

Formalizing the Selection (Cont'd)

- The update of \mathbf{H}_{t+1}^{-1} , based on \mathbf{H}_t^{-1} can then be formally written as

$$\mathbf{H}_{t+1}^{-1} = \mathbf{H}_t^{-1} + \frac{\mathbf{y}_t \mathbf{y}_t^T}{\mathbf{y}_t^T \mathbf{s}_t} - \frac{\mathbf{H}_t^{-1} \mathbf{s}_t \mathbf{s}_t^T \mathbf{H}_t^{-1}}{\mathbf{s}_t^T \mathbf{H}_t^{-1} \mathbf{s}_t} \quad (*)$$

- The updating is **rank-2** and hence can be executed efficiently
- This updating is known as the **Broyden-Fletcher-Goldfarb-Shanno (BFGS)** method

BFGS

Broyden, Fletcher, Goldfarb, Shanno



Formalizing the Selection (Cont'd)

- We can now update \mathbf{H}_{t+1}^{-1} based on \mathbf{H}_t^{-1}

$$\mathbf{H}_{t+1}^{-1} = \mathbf{H}_t^{-1} + \frac{\mathbf{y}_t \mathbf{y}_t^T}{\mathbf{y}_t^T \mathbf{s}_t} - \frac{\mathbf{H}_t^{-1} \mathbf{s}_t \mathbf{s}_t^T \mathbf{H}_t^{-1}}{\mathbf{s}_t^T \mathbf{H}_t^{-1} \mathbf{s}_t} \quad (*)$$

- But the algorithm updates using the following form

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{H}_t \nabla f(\mathbf{x}_t)$$

- We need \mathbf{H}_t rather than \mathbf{H}_t^{-1}
 - Q: can we construct a similar update formula?

Formalizing the Selection (Cont'd)

- Sherman-Morrison-Woodbury formula

$$(A + UDV)^{-1} = A^{-1} - A^{-1}U(D^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

- Homework: By setting

$$\begin{aligned} \mathbf{U} &= \mathbf{V}^T = [\mathbf{H}_t^{-1} \mathbf{s}_t \ \mathbf{y}_t] \\ \mathbf{D} &= \begin{bmatrix} -1/(\mathbf{s}_t^T \mathbf{H}_t^{-1} \mathbf{s}_t) & \mathbf{0} \\ \mathbf{0} & 1/(\mathbf{y}_t^T \mathbf{s}_t) \end{bmatrix} \end{aligned}$$

- The BFGS is equivalent to the following

$$\mathbf{H}_{t+1} = (\mathbf{I} - \rho_t \mathbf{s}_t \mathbf{y}_t^T) \mathbf{H}_t (\mathbf{I} - \rho_t \mathbf{y}_t \mathbf{s}_t^T) + \rho_t \mathbf{s}_t \mathbf{s}_t^T$$

where $\rho_t = \frac{1}{\mathbf{y}_t^T \mathbf{s}_t}$

An Alternative Interpretation

- The update of \mathbf{H}_{t+1} based on \mathbf{H}_t is also equivalent to the following

$$\min_{\mathbf{H}} \langle \mathbf{H}_t, \mathbf{H}^{-1} \rangle - \log \det (\mathbf{H}_t \mathbf{H}^{-1})$$

$$\text{s.t. } \mathbf{H} = \mathbf{H}^T \quad \dots \quad (\text{symmetric condition})$$

$$\mathbf{H} \mathbf{y}_t = \mathbf{s}_t \quad \dots \quad (\text{secant equation})$$

- In essence, minimizing KL divergence with constraints on the secant equation

Quasi Newton Method

- Suppose we have an unconstrained optimization problem,

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{x} \in R^n$$

- Start with an initial point \mathbf{x}_0 and update as follows

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{H}_t \nabla f(\mathbf{x}_t)$$

where η_t is determined by **line search**

- Update \mathbf{H}_{t+1} by the following

$$\mathbf{H}_{t+1} = (\mathbf{I} - \rho_t \mathbf{s}_t \mathbf{y}_t^T) \mathbf{H}_t (\mathbf{I} - \rho_t \mathbf{y}_t \mathbf{s}_t^T) + \rho_t \mathbf{s}_t \mathbf{s}_t^T$$

where $\mathbf{y}_t = \nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}_t)$, $\mathbf{s}_t = \mathbf{x}_{t+1} - \mathbf{x}_t$, and $\rho_t = \frac{1}{\mathbf{y}_t^T \mathbf{s}_t}$

Convergence Rate

- Assume f is strongly convex and has Lipschitz-continuous Hessian.
- **Theorem:** BFGS method converges as

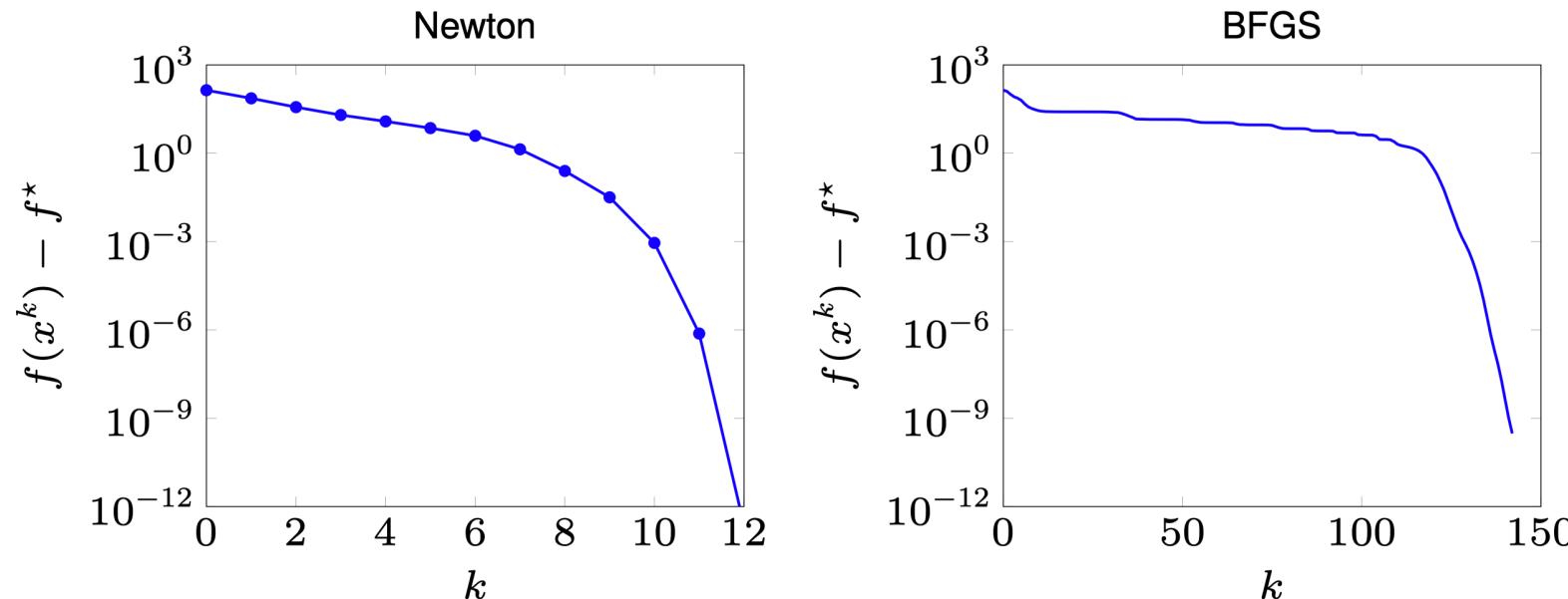
$$\lim_{t \rightarrow \infty} \frac{\|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*\|_2}{\|\boldsymbol{x}_t - \boldsymbol{x}^*\|_2} = 0$$

- Superlinear convergence rate (why): number of required iterations to achieve a certain accuracy is larger than Newton methods but smaller than Gradient descent

Numerical Example

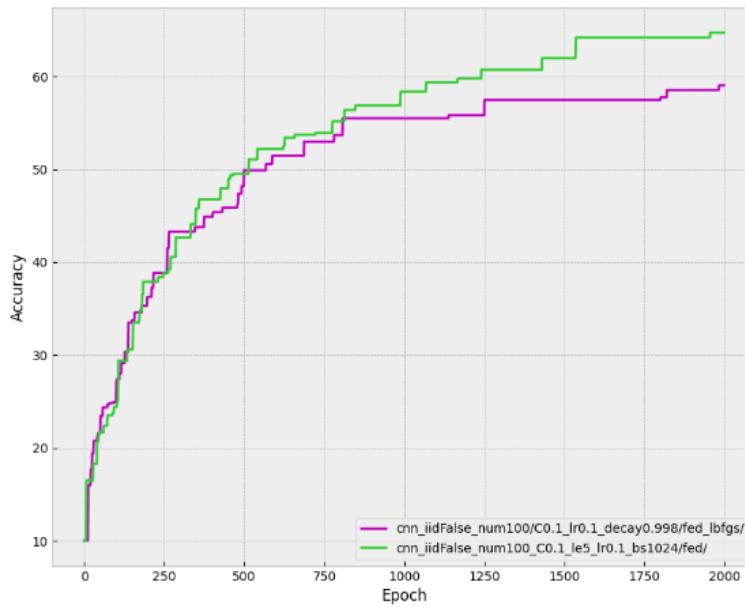
- Example from Vandenberghe's lecture notes: Newton versus BFGS on LP barrier problem, for n=100, m=500

$$\min_x \quad c^T x - \sum_{i=1}^m \log(b_i - a_i^T x)$$

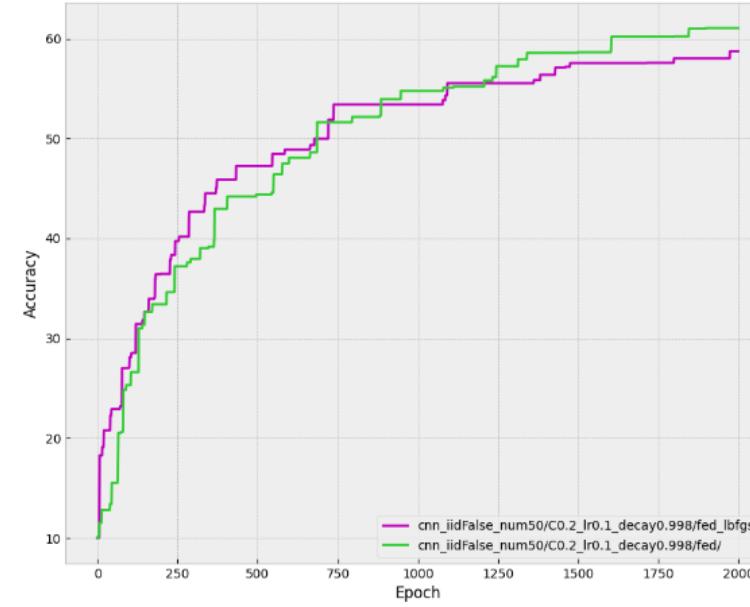


Stochastic Quasi Newton Method

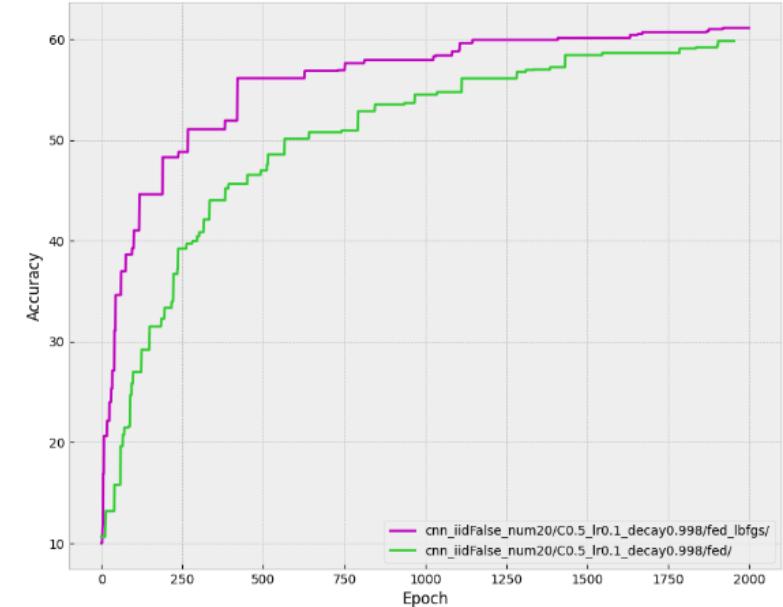
- Can we use the stochastic gradient in BFGS?



Batch size = 500



Batch size = 1000



Batch size = 2500

Learning Objectives

- Gradient descent and accelerating methods
- Variants of (stochastic) gradient descent—adaptive optimization methods
- Quasi Newton method