

TA SESSION 5: NUMERICAL OPTIMIZATION

Shuowen Chen¹

EC708: PhD Econometrics I (Spring 2020)

¹Parts of the materials are borrowed from teaching slides by Pierre Perron, Gianluca Violante, Jean-Jacques Forneron and textbook by Kenneth Train (2009). I thank Iván Fernández-Val, Hiroaki Kaido and Jean-Jacques Forneron for discussions of material coverage. I'm grateful to Yinuo Zhang for kindly sharing Violante's slides.

OUTLINE

- ▶ Optimization problem
- ▶ Full-Newton Method:
 - ▶ Newton-Ralphson
 - ▶ Gauss-Newton
- ▶ Quasi-Newton Method
 - ▶ Berndt-Hall-Hall-Hausman (BHHH) algorithm
 - ▶ Steepest Ascent
 - ▶ Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm
 - ▶ Davidon-Fletcher-Powell (DFP) algorithm
- ▶ Bonus Algorithms:
 - ▶ (Stochastic) Gradient Descent
 - ▶ Some comparison methods
 - ▶ Simplex method (Nelder-Meade/Polytope)
 - ▶ Simulated Annealing

OPTIMIZATION PROBLEM: GENERAL SETUP

We first consider the problem of **minimizing** a criterion function $f(\beta)$ with β being an $K \times 1$ vector of parameter.

Remark: for maximization, simply work with $-f(\beta)$.

Notations:

- ▶ Vector of first derivatives

$$g(\beta) = \frac{\partial f(\beta)}{\partial \beta}$$

- ▶ Matrix of second derivatives (Hessian)

$$G(\beta) = \frac{\partial^2 f(\beta)}{\partial \beta \partial \beta'}$$

GLOBAL OR LOCAL OPTIMUM?

Remark: A **sufficient condition** for a local minimum at $\hat{\beta}$ is

$$g(\hat{\beta}) = 0 \text{ and } G(\hat{\beta}) \text{ is positive definite}$$

Why not necessary? A counterexample: $y = \beta^4$ and $\hat{\beta} = 0$.

In practice no algorithm guarantees finding the global minimum.

One remedy: repeat the optimization with different starting values.

NUMERICAL EVALUATION OF DERIVATIVES

Recall the definition of first-order partial derivative of function f at β :

$$\frac{\partial f(\beta)}{\partial \beta_j} = \lim_{h \rightarrow 0} \frac{f(\beta_1, \dots, \beta_j + h, \dots, \beta_K) - f(\beta_1, \dots, \beta_j, \dots, \beta_K)}{h}$$

Numerically, we can approximate it by calculating

$$g_j(\hat{\beta}) = \frac{f(\hat{\beta} + h_j d_j) - f(\hat{\beta})}{h_j}, \quad j = 1, \dots, K$$

where

- ▶ $d_j = (0, \dots, 0, 1, 0, \dots, 0)$: unit vector with 1 in position j
- ▶ h_j : step length. Small for estimate to be close, larger enough than rounding error

Approximate second-order derivatives using similar notations:

$$G_{ji}(\hat{\beta}) = \frac{g_j(\hat{\beta} + h_i d_i) - g_j(\hat{\beta})}{h_i}, \quad i, j = 1, \dots, K$$

OPTIMIZATION PROBLEM: MLE

Consider **maximizing** the log likelihood function²

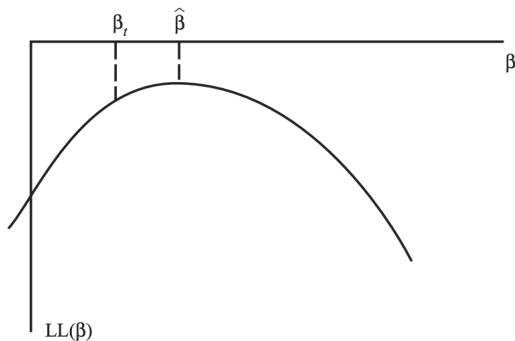
$$L(\beta) = \sum_{n=1}^N \log(P_n(\beta))/N,$$

- ▶ $P_n(\beta)$: prob of the observed outcome for decision maker n
- ▶ N : sample size
- ▶ β : $K \times 1$ vector of parameters

Goal: find $\hat{\beta} = \arg \max_{\beta} L(\beta)$

²To utilize the minimization packages in practice, usually work with $-L(\beta)$.

GRAPHICAL ILLUSTRATION OF MLE



Finding $\hat{\beta}$ is a hill-climbing process:

1. Given a starting point, decide what direction and how far to climb
2. Update point of starting and keep climbing
3. Stop when some criteria are met

We formalize this process via math, starting from [full-Newton method](#)

NEWTON-RAPHSON METHOD

Take a second-order Taylor expansion of $f(\beta)$ around the optimum $\hat{\beta}$:

$$f(\beta) \approx f(\hat{\beta}) + (\beta - \hat{\beta})' g(\hat{\beta}) + \frac{1}{2} (\beta - \hat{\beta})' G(\hat{\beta}) (\beta - \hat{\beta})$$

Differentiate w.r.t β yields

$$g(\beta) \approx g(\hat{\beta}) + G(\hat{\beta})(\beta - \hat{\beta})$$

Since $g(\hat{\beta}) = 0$ by definition, we have

$$\hat{\beta} \approx \beta - G^{-1}(\hat{\beta})g(\beta)$$

The derivation suggests the following iteration procedure:

1. Pick an initial guess β_0 and compute $\beta_1 = \beta_0 - G^{-1}(\beta_0)g(\beta_0)$
2. Compute $\beta_2 = \beta_1 - G^{-1}(\beta_1)g(\beta_1)$
3. Continue step 2 until convergence

Remark: $G^{-1}(\hat{\beta})$ is infeasible, so use iterated estimate instead

NEWTON–RAPHSON TO LOG LIKELIHOOD FUNCTION

Define the gradient in iteration k , β_k :

$$g_k = \left. \frac{\partial L(\beta)}{\partial \beta} \right|_{\beta=\beta_k}$$

and Hessian:

$$H_k = \left(\frac{\partial^2 L(\beta)}{\partial \beta \partial \beta'} \right) \Big|_{\beta=\beta_k}$$

The Newton–Raphson update in iteration $k + 1$ is:

$$\beta_{k+1} = \beta_k - H_k^{-1} g_k$$

CONVERGENCE CRITERIA

If f is **exactly quadratic**, then Newton–Raphson finds the optimum in **one–step**. Consider $f(\beta) = a + b\beta + c\beta^2$, whose optimum is $\hat{\beta} = -\frac{b}{2c}$.

- ▶ Its first and second order derivatives, evaluated at iteration k , are $b + 2c\beta_k$ and $2c$ respectively
- ▶ Hence $\beta_{k+1} = \beta_k - \frac{1}{2c}(b + 2c\beta_k) = -\frac{b}{2c} = \hat{\beta}$

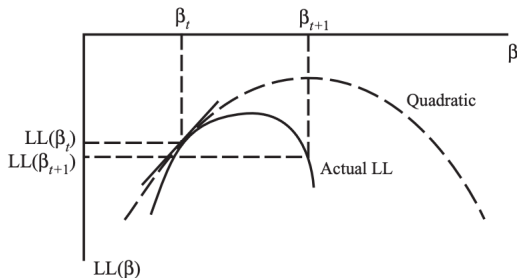
In practice, convergence can be defined in many ways:

- ▶ $f(\beta_{k+1})$ close to $f(\beta_k)$
- ▶ β_{k+1} close to β_k
- ▶ $g(\beta_{k+1})$ close to $g(\beta_k)$

Convergence can be hard because

- ▶ Local versus global optimums
- ▶ Objective functions being flat (weakly identified IV and GMM)

OVERSHOOTING AND STEPSIZE



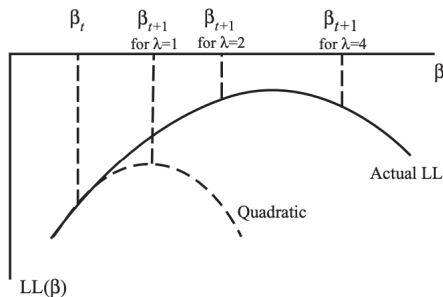
- ▶ The Newton–Raphson method updates from β_k to β_{k+1} , but $L(\beta_{k+1}) < L(\beta_k)$. The difference between β_{k+1} and β_k too large
- ▶ To ensure each iteration provides an increase in $L(\beta)$, we introduce stepsize λ_k :

$$\beta_{k+1} = \beta_k - \lambda_k G^{-1}(\beta_k) g(\beta_k)$$

DETERMINING THE STEPSIZE: BACKTRACKING LINE SEARCH

For each iteration, start with a large λ and gradually shrink it to guarantee an increase in $L(\beta)$. One such procedure for iteration k :

1. Start with $\lambda_k = 1$, if $L(\beta_{k+1}) > L(\beta_k)$, move to β_{k+1} and start iteration $k + 1$
2. If $L(\beta_{k+1}) < L(\beta_k)$, set $\lambda_k = \frac{1}{2}$ and try again
3. Continue step 2 until we find the first λ_k that yields $L(\beta_{k+1}) > L(\beta_k)$. Now move on to iteration $k + 1$.



FUTHER ISSUES WITH NEWTON–RAPHSON METHOD

Objective function not necessarily **globally convex**³

- ▶ Hessian may not be **positive** definite⁴
- ▶ Newton–Raphson update actually moves to opposite direction
- ▶ One Remedy: **regularization**, instead of using $G(\beta_k)^{-1}$ directly, use

$$\left[G(\beta_k) + \mu_k I_K \right]^{-1},$$

where $\mu_k > 0$: guarantees positive definiteness

Computing Hessian is costly

- ▶ Too many function evaluations are required
- ▶ Numerically calculated Hessian might be ill–behaved (singular)
- ▶ Consider algorithms that **approximate** the Hessian.

Let's start with **Gauss–Newton** method, which is designed for calculation of nonlinear least–squared estimates

³For maximization problem objective function is not **globally concave**.

⁴For maximization problem replace with **negative** definite

GAUSS-NEWTON METHOD

Consider the following general nonlinear model:

$$y_t = f(x_t; \beta) + u_t, \quad t = 1, \dots, T$$

- ▶ $u_t \sim i.i.d. \mathcal{N}(0, \sigma^2)$. Assume σ^2 is known here.
- ▶ x_t : $M \times 1$ exogenous regressors
- ▶ β : $K \times 1$ vector of parameters
- ▶ $f(\cdot)$: some function satisfying some regularity conditions

GAUSS-NEWTON METHOD CONT.D

Due to normality assumption, we have the log likelihood function

$$L(\beta) = -\frac{T}{2} \log 2\pi - \frac{T}{2} \log \sigma^2 - \frac{1}{2\sigma^2} S(\beta),$$

where $S(\beta) = \sum_{t=1}^T [y_t - f(x_t; \beta)]^2 = \sum_{t=1}^T u_t^2$. It suffices to work with $S(\beta)$. Its first and second order derivatives w.r.t β are

$$g(\beta) = 2 \sum_{t=1}^T \frac{\partial u_t}{\partial \beta} u_t, \quad G(\beta) = 2 \sum_{t=1}^T \left[\frac{\partial u_t}{\partial \beta} \frac{\partial u_t}{\partial \beta'} + \frac{\partial^2 u_t}{\partial \beta \partial \beta'} u_t \right]$$

In $G(\beta)$, the blue term is usually small relative to the red term, so we **neglect** it and use the following:

$$\beta_{k+1} = \beta_k - \left[\sum_{t=1}^T \frac{\partial u_t}{\partial \beta} \frac{\partial u_t}{\partial \beta'} \right]^{-1} \Big|_{\beta=\beta_k} \sum_{t=1}^T \frac{\partial u_t}{\partial \beta} u_t \Big|_{\beta=\beta_k}$$

REMARKS ON THE GAUSS-NEWTON METHOD

- ▶ The method doesn't compute Hessian
- ▶ Has an OLS interpretation. Let $z_t = -\partial u_t / \partial \beta$, then

$$\beta_{k+1} = \beta_k + \left(\sum_{t=1}^T z_t z_t' \right)^{-1} \Big|_{\beta=\beta_k} \sum_{t=1}^T z_t u_t \Big|_{\beta=\beta_k}$$

- ▶ Similar modification can be incorporated as in the Newton-Raphson: Marquart quadratic hill climbing

$$\beta_{k+1} = \beta_k + \left(\sum_{t=1}^T z_t z_t' + \mu I \right)^{-1} \Big|_{\beta=\beta_k} \sum_{t=1}^T z_t u_t \Big|_{\beta=\beta_k}$$

QUASI-NEWTON METHOD

Full-Newton method numerically evaluates the Hessian.

Quasi-Newton methods avoid it by approximating the Hessian, with differences in terms of how to approximate.

General procedures:

1. Specify initial β_0 and G_0 . In each iteration k
2. Compute Quasi-Newton direction $\Delta\beta_k = -G(\beta_k)^{-1}g(\beta_k)$
3. Determine the stepsize λ_k (by backtracking line search)
4. Compute $\beta_{k+1} = \beta_k + \lambda_k\Delta\beta_k$
5. Compute $G(\beta_{k+1})$ from $G(\beta_k)$

I will use G_k and g_k to denote $G(\beta_k)$ and $g(\beta_k)$ now.

QUASI-NEWTON ALGORITHM: BHHH

Very suitable for **log likelihood function maximization** as it uses **score** to approximate Hessian:

$$s_n(\beta_k) = \left. \frac{\partial \log P_n(\beta)}{\partial \beta} \right|_{\beta=\beta_k},$$

For a log likelihood function, **gradient is average scores**:

$$g_k = \sum_{n=1}^N \frac{s_n(\beta_k)}{N}$$

Outer product of observation n 's score is the $K \times K$ matrix:

$$s_n(\beta_k)s_n(\beta_k)' = \begin{pmatrix} s_n^1 s_n^1 & s_n^1 s_n^2 & \cdots & s_n^1 s_n^K \\ s_n^2 s_n^1 & s_n^2 s_n^2 & \cdots & s_n^2 s_n^K \\ \vdots & \vdots & \ddots & \vdots \\ s_n^K s_n^1 & s_n^K s_n^2 & \cdots & s_n^K s_n^K \end{pmatrix}$$

where s_n^j is the j -th element of $s_n(\beta_k)$

BHHH UPDATE

Average outer product:

$$G_k = \frac{1}{N} \sum_{n=1}^N s_n(\beta_k) s_n(\beta_k)'$$

BHHH update:

$$\beta_{k+1} = \beta_k + \lambda_k G_k^{-1} g_k$$

Why does this work?

- ▶ At maximum, the average score is zero and thus average outer product becomes the variance of scores in the sample
- ▶ This variance, like Hessian, provides a measure of the log-likelihood function's curvature, **A higher variance implies a greater curvature**
- ▶ Formally speaking, this is the **information matrix equality**:

$$\mathbb{E}\left[\frac{\partial L(\beta)}{\partial \beta} \frac{\partial L(\beta)}{\partial \beta'}\right] = -\mathbb{E}\left[\frac{\partial^2 L(\beta)}{\partial \beta \partial \beta'}\right]$$

QUASI-NEWTON ALGORITHM: STEEPEST ASCENT

BHHH can suffer from singularity or bad scaling in practice. One alternative formula is the following:

$$\beta_{k+1} = \beta_k + \lambda_k g_k$$

This is using the identity matrix to approximate the Hessian⁵.

- ▶ Identity matrix is pd, guarantees an increase in each iteration
- ▶ Steepest as it provides the greatest possible increase in $L(\beta)$ for the distance between β_k and β_{k+1} , but choice of λ_k is critical. Need to be small, which slows convergence.
- ▶ For minimization problem, this is called the gradient descent. We will briefly talk about stochastic gradient descent, which is popular in deep learning.

⁵Again we are considering maximization here, for minimization we consider descent and negative identity matrix

BFGS AND DFP

BFGS is the algorithm behind matlab's `fminunc`. The update of G_{k+1} from G_k is the following:

$$G_{k+1} = G_k + \frac{y_k y_k'}{y_k' y_k} - \frac{G_k \gamma_k \gamma_k' G_k}{\gamma_k' G_k \gamma_k},$$

where $\gamma_k = \beta_{k+1} - \beta_k$ and $y_k = g_{k+1} - g_k$. Using Sherman–Morrison formula, we have

$$G_{k+1}^{-1} = \left(I - \frac{\gamma_k y_k'}{y_k' \gamma_k}\right) G_k^{-1} \left(I - \frac{y_k \gamma_k'}{\gamma_k' y_k}\right) + \frac{\gamma_k \gamma_k'}{\gamma_k' \gamma_k}$$

DFP precedes BFGS, but gets superseded by the latter:

$$G_{k+1}^{DFP} = G_k^{DFP} + \frac{\gamma_k \gamma_k'}{\gamma_k' y_k} - \frac{G_k^{DFP} y_k y_k' G_k^{DFP}}{y_k' G_k^{DFP} y_k}$$

The two update swap the role of γ_k and y_k .

GRADIENT DESCENT

We've seen steepest ascent in maximization problem. Now consider the minimizing an unconstrained, smooth convex function

$$\min_x f(x),$$

where f is convex and twice differentiable. The Newton–Raphson method updates x as follows:

$$x_{k+1} = x_k - G_k^{-1} g_k$$

while gradient descent updates as follows:

$$x_{k+1} = x_k - \lambda_k g_k$$

Reminder: G_k and g_k respectively denote Hessian and gradient evaluated at x_k . λ_k denotes the step size.

STOCHASTIC GRADIENT DESCENT

- ▶ Historically gradient descent not very popular in nonlinear/nonconvex optimization problem because it gets stuck at local minima
- ▶ For deep network, local minima are close to global minimum in terms of prediction, so not problematic
- ▶ Computing gradient can be expensive in deep network. One way to be more efficient is the stochastic gradient descent
- ▶ Instead of evaluating the gradient at each observation (with sample size N), in each iteration we consider a subsample $(x_1^\star, \dots, x_m^\star)$, $m < N$ and compute⁶

$$\beta_{k+1} = \beta_k - \lambda_k g_k^\star,$$

where g_k^\star denotes gradient of subsample evaluated at β_k

- ▶ In practice, prefer small m (like 1): cheaper to compute and avoids overfitting (Goodfellow et al, 2016). Good finite and large sample behavior (Toulis and Airolidi, 2017).

⁶On how to tune λ_k , refer to [this](#).

COMPARISON METHODS

- ▶ Most of the methods we discussed are **gradient-based**: use info on the slope and possibly on curvature.
- ▶ Gradient-based methods get stuck if objective function is non-smooth; started at different initial points, more likely to lead to a different local optimum
- ▶ An alternative type of method is **comparison-based**: compute objective function at several points and pick the one yielding the optimum value
- ▶ Comparison method better behaved with non-smooth objective functions; stochastic comparison method more likely to find global optimum (in theory)

NONLINER SIMPLEX METHOD: NELDER–MEADE/POLYTOPE

This is `fminsearch` in Matlab

1. Choose initial simplex⁷ $\{x_1, x_2, \dots, x_{n+1}\} \in \mathbb{R}^n$
2. Sort simplex vertices in descending order

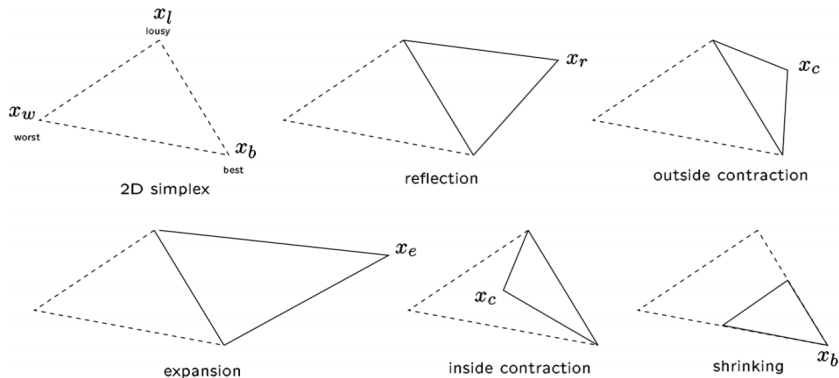
$$f(x_i) \geq f(x_{i+1}) \geq \dots, \forall i$$

3. Find the smallest i such that $f(x_i^R) < f(x_i)$, where x_i^R is the reflection of x_i . If x_i^R exists, replace x_i with x_i^R and go back to step 2. Otherwise, go to step 4.
4. If width of the current simplex smaller than tolerance, stop. Otherwise go to step 5
5. For $i = 1, \dots, n$, set $x_i^S = \frac{x_i + x_{i+1}}{2}$ to shrink the simplex. Go back to step 1.

Doesn't require information on derivatives, hence suitable for non-smooth objective functions.

⁷Think of it as an n -dimensional version of a triangle.

SOME SIMPLEX MOVES



Intuition: Nelder–Mead starts with a simplex and modifies it at each iteration using one of the moves. The sequence of moves to be performed is chosen based on the relative values of the objective functions at each of the points.

SIMULATED ANNEALING

Suitable for finding a global minimum among many local ones.

1. Draw z from $\mathcal{N}(0, 1)$ and perturb initial guess x_0 : $x_1 = x_0 + \lambda z$
2. If $f(x_1) < f(x_0)$, move to step 3. Otherwise **accept stochastically**: accept if

$$\frac{f(x_1) - f(x_0)}{|f(x_0)|} < \tau c,$$

where $c \sim U[0, 1]$ and $\tau > 0$ is **temperature paramter**⁸. If not, go back to step 1 and redraw.

3. Compute $|x_1 - x_0|$. Stop if less than tol, otherwise, update initial guess ($x_0 = x_1$) and back to step 1.

Fun fact: I used this once in an IO project on moment inequalities, the results suck...big time.

⁸Along the iterations decrease $\tau \rightarrow 0$ to cool down (reduce randomness).

CENTRAL DIFFERENCE

For numerical differentiation, it is recommended to use

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x-h)}{2h},$$

which is more accurate in practice.

The multivariate version should be fairly straightforward. [◀ Back](#)