

Numerical Optimization

Zhentao Shi

Jan 30, 2019

Numerical Optimization

Optimization is a key step to implement extremum estimators in econometrics. A general optimization problem is formulated as

$$\min_{\theta \in \Theta} f(\theta) \quad \text{s.t.} \quad g(\theta) = 0, h(\theta) \leq 0,$$

where $f(\cdot)$ is a criterion function, $g(\theta) = 0$ is an equality constraint, and $h(\theta) \leq 0$ is an inequality constraint.

Most established numerical optimization algorithms aim at finding a local minimum. However, there is no guarantee to locate the global minimum when multiple local minima exist.

Optimization without the equality and/or inequality constraints is called an *unconstrained* problem; otherwise it is called a *constrained* problem. The constraints can be incorporated into the criterion function via Lagrangian.

Methods

There are many optimization algorithms in the field of operational research; they are variants of a small handful of fundamental principles.

The essential idea for optimizing a twice-differentiable objective function is the Newton's method. A necessary condition for optimization is the first-order condition $s(\theta) = \partial f(\theta)/\partial \theta = 0$.

At an initial trial value θ_0 , if $s(\theta_0) \neq 0$, the search is updated by

$$\theta_{t+1} = \theta_t - (H(\theta_t))^{-1} s(\theta_t)$$

for $t = 0, 1, \dots$ where $H(\theta) = \frac{\partial s(\theta)}{\partial \theta}$ is the Hessian matrix. The algorithm iterates until $|\theta_{t+1} - \theta_t| < \epsilon$ (absolute criterion) and/or $|\theta_{t+1} - \theta_t|/|\theta_t| < \epsilon$ (relative criterion), where ϵ is a small positive number chosen as a tolerance level.

Newton's Method. Newton's method seeks the solution to $s(\theta) = 0$. Recall that the first-order condition is a necessary condition but not a sufficient condition. We still need to verify the second-order condition to verify whether a root to $s(\theta)$ is associated with a minimizer, a maximizer or a saddle point. If there are multiple minima, we compare the value at each to decide the global minimum.

It is clear that Newton's method requires computation of the gradient $s(\theta)$ and the Hessian $H(\theta)$. Newton's method converges at quadratic rate, which is fast.

Quasi-Newton Method. The most well-known quasi-Newton algorithm is **BFGS**. It avoids explicit calculation of the computationally expensive Hessian matrix. Instead,

starting from an initial (inverse) Hessian, it updates the Hessian by an explicit formula motivated from the idea of quadratic approximation.

Derivative-Free Method. [Nelder-Mead](#) is a simplex method. It searches a local minimum by reflection, expansion and contraction.

Implementation

R's optimization infrastructure has been constantly improving. [R Optimization Task View](#) gives a survey of the available CRAN packages.

For general-purpose nonlinear optimization, the package `optimx` (Nash 2014) effectively replaces R's default optimization commands. `optimx` provides a unified interface for various widely-used optimization algorithms. Moreover, it facilitates comparison among optimization algorithms.

Example

We use `optimx` to solve pseudo Poisson maximum likelihood estimation (PPML). If y_i is a continuous random variable, it obviously does not follow a Poisson distribution, whose support consists of non-negative integers. However, if the conditional mean model

$$E[y_i|x_i] = \exp(x_i'\beta),$$

is satisfied, we can still use the Poisson regression to obtain a consistent estimator of the parameter β even if y_i does not follow a conditional Poisson distribution.

If Z follows a Poisson distribution with mean λ , the probability mass function

$$\Pr(Z = k) = \frac{e^{-\lambda}\lambda^k}{k!}, \text{ for } k = 0, 1, 2, \dots,$$

so that

$$\log \Pr(Y = y|x) = -\exp(x'\beta) + y \cdot x'\beta - \log k!$$

Since the last term is irrelevant to the parameter, the log-likelihood function is

$$\ell(\beta) = \log \Pr(\mathbf{y}|\mathbf{x}; \beta) = -\sum_{i=1}^n \exp(x_i'\beta) + \sum_{i=1}^n y_i x_i'\beta.$$

In addition, it is easy to write the gradient

$$s(\beta) = \frac{\partial \ell(\beta)}{\partial \beta} = -\sum_{i=1}^n \exp(x_i'\beta) x_i + \sum_{i=1}^n y_i x_i.$$

and verify that the Hessian

$$H(\beta) = \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'} = -\sum_{i=1}^n \exp(x_i'\beta) x_i x_i'$$

is negative definite. Therefore, $\ell(\beta)$ is strictly concave in β .

In operational reserach, the default optimization is minimization, although economics has utility maximization and cost minimization and statistics has maximum likelihood estimation and minimal least squared estimation. To follow this convention in operational research, here we formulate the *negative* log-likelihood.

```

# Poisson likelihood
poisson.loglik = function( b ) {
  b = as.matrix( b )
  lambda = exp( X %*% b )
  ell = -sum( -lambda + y * log(lambda) )
  return(ell)
}

```

To implement optimization in R, it is recommended to write the criterion as a function of the parameter. Data can be fed inside or outside of the function. If the data is provided as additional arguments, these arguments must be explicit. (In contrast, in Matlab the parameter must be the sole argument for the function to be optimized, and data can only be injected through a nested function.)

```

# implement both BFGS and Nelder-Mead for comparison.

```

```

library(AER)
library(numDeriv)
library(optimx)

## prepare the data
data("RecreationDemand")
y = RecreationDemand$trips
X = with(RecreationDemand, cbind(1, income))

## estimation
b.init = c(0, 1) # initial value
b.hat = optimx(b.init, poisson.loglik, method = c("BFGS", "Nelder-Mead"), control = list(reltol = 1e-07, abstol = 1e-07))
print(b.hat)

```

```

##           p1          p2    value fevals gevals niter convcode
## BFGS      1.177411 -0.09994222 261.1141     99     21     NA        0
## Nelder-Mead 1.167261 -0.09703975 261.1317     53     NA     NA        0
##           kkt1 kkt2 xt看imes
## BFGS          TRUE TRUE    0.01
## Nelder-Mead FALSE TRUE    0.00

```

Given the conditional mean model, nonlinear least squares (NLS) is also consistent in theory. NLS minimizes

$$\sum_{i=1}^n (y_i - \exp(x_i \beta))^2$$

A natural question is: why do we prefer PPML to NLS? My argument is that, PPML's optimization for the linear index is globally convex, while NLS is not. It implies that the numerical optimization of PPML is easier and more robust than that of NLS. I leave the derivation of the non-convexity of NLS as an exercise.

In practice no algorithm suits all problems. Simulation, where the true parameter is known, is helpful to check the accuracy of one's optimization routine before applying to an empirical problem,

where the true parameter is unknown. Contour plot is a useful tool to visualize the function surface/manifold in a low dimension.

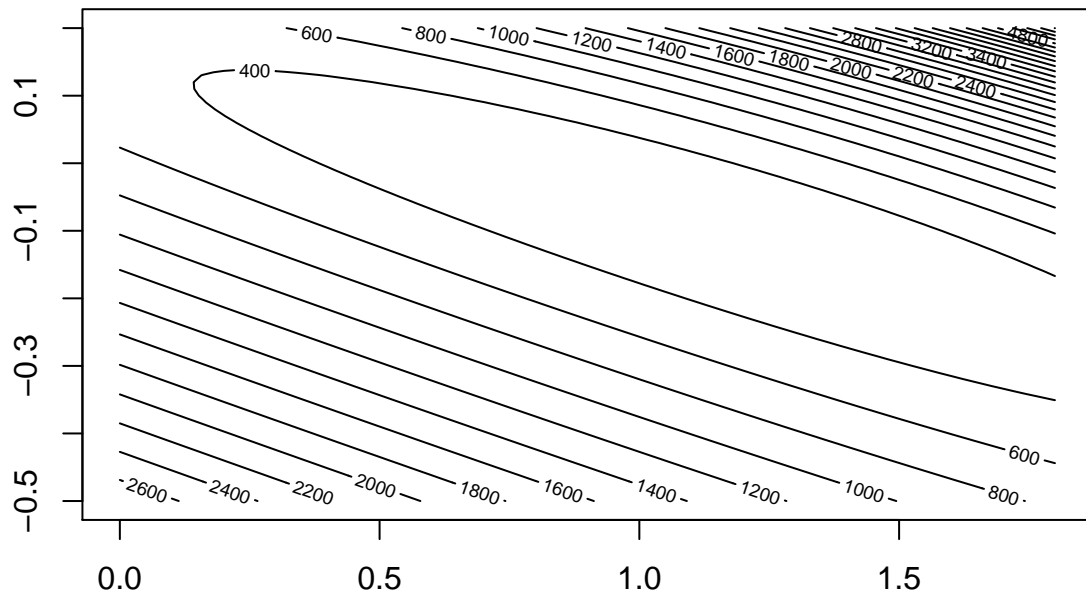
Example

```
x.grid = seq(0, 1.8, 0.02)
x.length = length(x.grid)
y.grid = seq(-0.5, 0.2, 0.01)
y.length = length(y.grid)

z.contour = matrix(0, nrow = x.length, ncol = y.length)

for (i in 1:x.length) {
  for (j in 1:y.length) {
    z.contour[i, j] = poisson.loglik(c(x.grid[i], y.grid[j]))
  }
}

contour(x.grid, y.grid, z.contour, 20)
```



For problems that demand more accuracy, third-party standalone solvers can be invoked via interfaces to R. For example, we can access [NLOpt](#) through the packages [nloptr](#). However, standalone solvers usually have to be compiled and configured. These steps are often not as straightforward as installing commercial Windows software.

NLopt offers an [extensive list of algorithms](#).

Example

We first carry out the Nelder-Mead algorithm in NLOPT.

```
library(nloptr)
## optimization with Nloptr

opts = list(algorithm = "NLOPT_LN_NELDERMEAD", xtol_rel = 1e-07, maxeval = 500)

res_NM = nloptr(x0 = b.init, eval_f = poisson.loglik, opts = opts)
print(res_NM)

##
## Call:
## nloptr(x0 = b.init, eval_f = poisson.loglik, opts = opts)
##
##
## Minimization using NLopt version 2.4.2
##
## NLopt solver status: 4 ( NLOPT_XTOL_REACHED: Optimization stopped because
## xtol_rel or xtol_abs (above) was reached. )
##
## Number of Iterations.....: 118
## Termination conditions:  xtol_rel: 1e-07 maxeval: 500
## Number of inequality constraints:  0
## Number of equality constraints:    0
## Optimal value of objective function: 261.114078295329
## Optimal value of controls: 1.177397 -0.09993984

# 'SLSQP' is indeed the BFGS algorithm in NLopt, though 'BFGS' doesn't
# appear in the name
opts = list(algorithm = "NLOPT_LD_SLSQP", xtol_rel = 1e-07)
```

To invoke BFGS in NLOPT, we must code up the gradient $s(\beta)$, as in the function `poisson.log.grad()` below.

```
poisson.loglik.grad = function(b) {
  b = as.matrix(b)
  lambda = exp(X %*% b)
  ell = -colSums(-as.vector(lambda) * X + y * X)
  return(ell)
}
```

We compare the analytical gradient with the numerical gradient to make sure the function is correct.

```
# check the numerical gradient and the analytical gradient
b = c(0, 0.5)
grad(poisson.loglik, b)

## [1] 6542.46 45825.40
```

```
poisson.loglik.grad(b)
```

```
##           income  
## 6542.46 45825.40
```

With the function of gradient, we are ready for BFGS.

```
res_BFGS = nloptr(x0 = b.init, eval_f = poisson.loglik, eval_grad_f = poisson.loglik.grad,  
                 opts = opts)  
print(res_BFGS)
```

```
##  
## Call:  
##  
## nloptr(x0 = b.init, eval_f = poisson.loglik, eval_grad_f = poisson.loglik.grad,  
##       opts = opts)  
##  
##  
## Minimization using Nlopt version 2.4.2  
##  
## Nlopt solver status: 4 ( NLOPT_XTOL_REACHED: Optimization stopped because  
## xtol_rel or xtol_abs (above) was reached. )  
##  
## Number of Iterations.....: 38  
## Termination conditions: xtol_rel: 1e-07  
## Number of inequality constraints: 0  
## Number of equality constraints: 0  
## Optimal value of objective function: 261.114078295329  
## Optimal value of controls: 1.177397 -0.09993984
```

Constrained Optimization

- `optimx` can handle simple box-constrained problems.
- `constrOptim` can handle linear constrained problems.
- Some algorithms in `nloptr`, for example, `NLOPT_LD_SLSQP`, can handle nonlinear constrained problems.

Convex Optimization

If a function is convex in its argument, then a local minimum is a global minimum. Convex optimization is particularly important in high-dimensional problems. The readers are referred to Boyd and Vandenberghe (2004) for an accessible comprehensive treatment. They claim that “convex optimization is technology; all other optimizations are arts.” This is true to some extent.

Example

- linear regression model MLE
- Lasso (Su, Shi, and Phillips 2016)

- Relaxed empirical likelihood (Shi 2016).

A class of common convex optimization can be reliably implemented in R. `Rmosek` is an interface in R to access Mosek. Mosek is a high-quality commercial solver dedicated to convex optimization. It offers free academic licenses. (`Rtools` is a prerequisite to install `Rmosek`.)

```
lo1 <- list()
lo1$sense <- "max"
lo1$c <- c(3, 1, 5, 1)
lo1$A <- Matrix::Matrix(c(3, 1, 2, 0, 2, 1, 3, 1, 0, 2, 0, 3), nrow = 3, byrow = TRUE,
  sparse = TRUE)
lo1$bc <- rbind(blc = c(30, 15, -Inf), buc = c(30, Inf, 25))
lo1$bx <- rbind(blx = c(0, 0, 0, 0), bux = c(Inf, 10, Inf, Inf))
r <- Rmosek::mosek(lo1) # as today (Feb 10, 2019), the latest version of
# Rmosek is not functional. I can make it work on R3.4, but not the latest
# R3.5
```

A survey and demonstration paper can be found here (Gao and Shi 2018).

Stochastic Gradient Descent (SGD)

When the sample size is huge and the number of parameters is also big, the evaluation of the gradient can be prohibitively expensive. Stochastic Gradient Descent (SGD) uses a small batch of the sample to evaluate the gradient in each iteration. It can significantly save computational time. In complex optimization problems such as training a neural network, SGD is the de facto optimization procedure.

However, SGD involves tuning parameters (say, the batch size and the learning rate) that can dramatically affect the outcome, in particular in nonlinear problems. Careful experiments must be carried out before serious implementation.

Below is an example of SGD in the PPMLE with sample size 100,000 and the number of parameters 100. SGD is usually much faster.

The new functions are defined with the data explicitly as arguments. Because in SGD each time the log-likelihood function and the gradient are evaluated at a different subsample.

```
poisson.loglik = function( b, y, X ) {
  b = as.matrix( b )
  lambda = exp( X %*% b )
  ell = -mean( -lambda + y * log(lambda) )
  return(ell)
}

poisson.loglik.grad = function( b, y, X ) {
  b = as.matrix( b )
  lambda = as.vector( exp( X %*% b ) )
  ell = -colMeans( -lambda * X + y * X )
  ell_eta = ell
```

```

    return(ell_eta)
}

3
##### generate the artificial data
set.seed(898)
nn = 1e+05
K = 100
X = cbind(1, matrix(runif(nn * (K - 1)), ncol = K - 1))
b0 = rep(1, K)/K
y = rpois(nn, exp(X %*% b0))

b.init = runif(K)
b.init = 2 * b.init/sum(b.init)

# and these tuning parameters are related to N and K

n = length(y)
test_ind = sample(1:n, round(0.2 * n))

y_test = y[test_ind]
X_test = X[test_ind, ]

y_train = y[-test_ind]
X_train = X[-test_ind, ]

# optimization parameters

# sgd depends on * eta: the learning rate * epoch: the averaging small batch
# * the initial value

max_iter = 5000
min_iter = 20
eta = 0.01
epoch = round(100 * sqrt(K))

b_old = b.init

pts0 = Sys.time()
# the iteration of gradient
for (i in 1:max_iter) {

    loglik_old = poisson.loglik(b_old, y_train, X_train)
    i_sample = sample(1:length(y_train), epoch, replace = TRUE)
    b_new = b_old - eta * poisson.loglik.grad(b_old, y_train[i_sample], X_train[i_sample,

```



```

    ])
    loglik_new = poisson.loglik(b_new, y_test, X_test)
    b_old = b_new # update

    criterion = loglik_old - loglik_new
    cat("the ", i, "-th criterion = ", criterion, "\n")

    if (criterion < 1e-04 & i >= min_iter)
        break
}

```

```

## the 1 -th criterion = 0.1876909
## the 2 -th criterion = 0.02972239
## the 3 -th criterion = -0.001035201
## the 4 -th criterion = -0.004724442
## the 5 -th criterion = -0.00698445
## the 6 -th criterion = -0.007083379
## the 7 -th criterion = -0.007508545
## the 8 -th criterion = -0.007070038
## the 9 -th criterion = -0.007254321
## the 10 -th criterion = -0.007142421
## the 11 -th criterion = -0.007302179
## the 12 -th criterion = -0.008023021
## the 13 -th criterion = -0.006883816
## the 14 -th criterion = -0.007415384
## the 15 -th criterion = -0.007686572
## the 16 -th criterion = -0.007324292
## the 17 -th criterion = -0.00751597
## the 18 -th criterion = -0.007777976
## the 19 -th criterion = -0.007587205
## the 20 -th criterion = -0.007214732

```

```

cat("point estimate =", b_new, ", log_lik = ", loglik_new, "\n")

```

```

## point estimate = -0.00693619 0.0130523 0.007441082 0.005964518 0.003117607 -0.00894338 -0.007214732

```

```

pts1 = Sys.time() - pts0
print(pts1)

```

```

## Time difference of 0.6563041 secs

```

```

# optimx is too slow for this dataset. Nelder-Mead method is too slow for
# this dataset

```

```

# thus we only sgd with NLOptr

```

```

opts = list(algorithm = "NLOPT_LD_SLSQP", xtol_rel = 1e-07, maxeval = 5000)

```

```
pts0 = Sys.time()
res_BFGS = nloptr::nloptr(x0 = b.init, eval_f = poisson.loglik, eval_grad_f = poisson.loglik.grad,
  opts = opts, y = y_train, X = X_train)
print(res_BFGS)
```

```
##
## Call:
##
## nloptr::nloptr(x0 = b.init, eval_f = poisson.loglik, eval_grad_f = poisson.loglik.grad,
##   opts = opts, y = y_train, X = X_train)
##
##
## Minimization using NLOpt version 2.4.2
##
## NLOpt solver status: 4 ( NLOPT_XTOL_REACHED: Optimization stopped because
## xt看ol_rel or xt看ol_abs (above) was reached. )
##
## Number of Iterations.....: 50
## Termination conditions:  xt看ol_rel: 1e-07 maxeval: 5000
## Number of inequality constraints:  0
## Number of equality constraints:    0
## Optimal value of objective function:  0.817239347882672
## Optimal value of controls: 0.008377464 0.00223468 0.006802827 0.009527097 -0.009179097 0.00
## -0.001861563 0.0250962 0.007914374 0.01430903 0.0185239 0.0009192213
## -0.003079815 0.02091482 0.01838836 0.02317564 0.007484981 0.02760165
## 0.02395294 0.02581524 0.0008703182 0.003018651 -0.0008777799 0.005505568
## -0.006333724 0.007045993 0.009869454 0.01635177 0.02444754 0.01859429
## 0.002282465 -0.006204774 0.01787497 0.00871074 0.01634583 0.01452785
## 0.01535949 0.02114798 0.02882988 0.003990018 0.006488755 0.01821596
## 0.01224189 -0.001992572 0.009518121 -0.004068174 -0.00955518 -0.008724836
## 0.01379455 0.01182887 0.02293765 0.01606992 0.01559747 0.02075959
## -0.002392985 -0.003346411 0.006898647 0.0120779 0.02070444 0.01023695
## 0.01356451 0.01390674 0.01992954 0.00914424 0.01823229 0.009599003
## 0.005063411 0.001639757 -0.003360399 0.004146783 0.0004618288 0.01613158
## 0.00628123 -0.007416353 0.01045353 0.02667767 0.01690433 0.009531548
## 0.01430839 0.005220393 0.02352873 0.01921597 0.01270081 0.00527763
## 0.008045643 0.002385956 0.02469504 0.002604459 0.02668898 -0.01424429
## 0.01709918 0.008055604 0.007352139 0.01098655 0.002287008 0.0190659
## 0.0004742784 0.01790932 -0.0007725258 0.02297663
```

```
pts1 = Sys.time() - pts0
print(pts1)
```

```
## Time difference of 6.020936 secs
```

```
b_hat_nlopt = res_BFGS$solution
```

```
#### evaluation in the test sample
cat("\n\n\n\n\n\n\n\n")

cat("log lik in test data by sgd = ", poisson.loglik(b_new, y = y_test, X_test),
    "\n")

## log lik in test data by sgd = 0.8266132

cat("log lik in test data by nlopt = ", poisson.loglik(b_hat_nlopt, y = y_test,
    X_test), "\n")

## log lik in test data by nlopt = 0.8263522

cat("log lik in test data by oracle = ", poisson.loglik(b0, y = y_test, X_test),
    "\n")

## log lik in test data by oracle = 0.8254042
```

References

- Boyd, Stephen, and Lieven Vandenberghe. 2004. *Convex Optimization*. Cambridge university press.
- Gao, Zhan, and Zhentao Shi. 2018. “Two Examples of Convex-Programming-Based High-Dimensional Econometric Estimators.” *arXiv Preprint arXiv:1806.10423*.
- Nash, John C. 2014. “On Best Practice Optimization Methods in R.” *Journal of Statistical Software* 60 (2): 1–14.
- Shi, Zhentao. 2016. “Econometric Estimation with High-Dimensional Moment Equalities.” *Journal of Econometrics* 195 (1). Elsevier: 104–19.
- Su, Liangjun, Zhentao Shi, and Peter CB Phillips. 2016. “Identifying Latent Structures in Panel Data.” *Econometrica* 84 (6). Wiley Online Library: 2215–64.