

Lecture 11: October 2

Lecturer: Geoff Gordon/Ryan Tibshirani

Scribes: Tongbo Huang, Shoou-I Yu

Note: *LaTeX template courtesy of UC Berkeley EECS dept.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

This lecture's notes illustrate some uses of various L^AT_EX macros. Take a look at this and imitate.

11.1 Matrix Differential Calculus

11.1.1 Review of Previous Class

Matrix differential is a solution to matrix calculation pain, it can be understood by either of the following:

- A compact way of writing Taylor expansion.
- Definition:
 $df = a(x; dx)[r(dx)]$ where r is the residual term
 $a(x, \cdot)$ is linear in the second argument
 $\frac{r(dx)}{\|dx\|} \rightarrow 0$ as $dx \rightarrow 0$

The derivative is linear, so it passes through addition and scalar multiplication.

It also generalizes Jacobian, Hessian, gradient and velocity.

Other topics covered: chain rule, product rule, bilinear functions, identities and identities theorems. Please refer to previous scribed notes for details.

11.1.2 Finding a maximum, minimum or saddle points

ID for $df(\mathbf{x})$	scalar x	vector \mathbf{x}	matrix X
scalar f	$df = a dx$	$df = \mathbf{a}^T d\mathbf{x}$	$df = \text{tr}(\mathbf{A}^T d\mathbf{X})$
vector \mathbf{f}	$d\mathbf{f} = \mathbf{a} dx$	$d\mathbf{f} = \mathbf{A} d\mathbf{x}$	
matrix F	$dF = \mathbf{A} dx$		

The principle: set coefficient of dX to 0 to find min, max, or saddle point:

- if $df = c(A; dX)[r(dX)]$ then $dX = tA$, $df = c(A; tA) = tc(A; A)$
- function is at min/max/saddle point iff $c(A; A) = 0$
- if c is any product, then $A = 0$

11.1.3 Infomax ICA

Suppose we have n training examples $x_i \in R^d$ and a scalar-valued, component-wise function g . We would like to find the $d \times d$ matrix W that maximizes the entropy of $y_i = g(Wx_i)$.

Detour: volume rule:

$$\text{vol}(AS) = |\det(A)|\text{vol}(S)$$

Interpretation: small determinant value means the existence of small eigenvalue, thus squash the volume flat, vice versa.

Back to infomax ICA. We have $y_i = g(Wx_i)$ where $dy_i = J(x; W)dx_i = J_i dx_i$. We want to maximize the entropy over the distribution of y :

$$\max_W \sum_i (-\ln(P(y_i))), \quad P(y_i) = \frac{P(x_i)}{|\det J(x_i; W)|}$$

And from

$$\max H(P(y)) = - \int P(y) \ln(P(y)) dy = -E(\ln(P(y)))$$

it is equivalent to maximizing

$$\max_W \sum_i \ln(|\det J(x_i; W)|)$$

since $P(x)$ is independent to W .

11.1.4 Solving ICA Gradient

Define $u_i = g'(Wx_i)$, $v_i = g''(Wx_i)$.

For gradient of $y_i = g(Wx_i)$:

$$\begin{aligned} dy_i &= g'(Wx_i) \circ d(Wx_i) \\ &= u_i \circ (Wdx_i) \\ &= \text{diag}(u_i)Wdx_i \end{aligned}$$

For gradient of $J_i = \text{diag}(u_i)W$:

$$\begin{aligned} dJ_i &= d(\text{diag}(u_i))W + \text{diag}(u_i)dW \\ &= \text{diag}(v_i \circ d(Wx_i))W + \text{diag}(u_i)dW \\ &= \text{diag}(u_i)dW + \text{diag}(v_i)\text{diag}(d(Wx_i))W \end{aligned}$$

Finally, define $\text{diag}(\alpha_i) = \text{diag}(u_i)^{-1} \text{diag}(v_i)$ solving the gradient of $\sum_i \ln(|\det J(x_i; W)|)$:

$$\begin{aligned}
 dL &= \sum_i d(\ln|\det J(x_i; W)|) \\
 &= \sum_i \text{tr}(J_i^{-1} dJ_i) \\
 &= \sum_i \text{tr}(W^{-1} dW + W^{-1} \text{diag}(u_i)^{-1} \text{diag}(v_i) \text{diag}(d(Wx_i))) \\
 &= \sum_i \text{tr}(W^{-1} dW) + \text{tr}(\text{diag}(\alpha_i) \text{diag}(d(Wx_i))) \\
 &= n \text{tr}(W^{-1} dW) + \sum_i \text{tr}(\alpha_i^T d(Wx_i)) \\
 &= n \text{tr}(W^{-1} dW) + \text{tr}(\sum_i x_i \alpha_i^T d(Wx_i)) \\
 &= nW^{-T} + \sum_i \alpha_i x_i^T \\
 &= nW^{-T} + C
 \end{aligned}$$

11.1.5 Natural Gradient

Define $L(W)$ as a function from $R^{d \times d}$ to R , then $dL = \text{tr}(G^T dW)$. So step S is:

$$S = \text{argmax}_S M(S), \quad M(S) = \text{tr}(G^T S) - \frac{\|SW^{-1}\|_F^2}{2}$$

which, in scalar case:

$$M = gS - \frac{S^2}{2W^2}$$

So, to find the max/min/saddle point:

$$M = \text{tr}(G^T S) - \frac{1}{2} \text{tr}(SW^{-1}W^{-T}S^T)$$

$$dM = \text{tr}(G^T dS) - \frac{1}{2} \text{tr}(dSW^{-1}W^{-T}S^T)$$

So, natural gradient becomes $G = W^{-1}W^{-T}$, and thus $GW^TW = S$. Using the gradient previously derived, $[W^{-T} + C]W^TW = W + CW^TW$.

11.1.6 More Info

- Minkas cheat sheet:
<http://research.microsoft.com/en-us/um/people/minka/papers/matrix/>
- Magnus & Neudecker. Matrix Differential Calculus. Wiley, 1999. 2nd ed.
<http://www.amazon.com/Differential-Calculus-Applications-Statistics-Econometrics/dp/047198633X>
- Bell & Sejnowski. An information-maximization approach to blind separation and blind deconvolution. Neural Computation, v7, 1995.

11.2 Newton's Method

Newton's method have two main applications: solving nonlinear equations and finding minima/maxima/saddles.

11.2.1 Solving Nonlinear Equations

For $x \in \mathbb{R}^d$ and $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ which is differentiable, we want to solve $f(x) = 0$. We perform first order Taylor approximation on $f(x)$,

$$f(y) \approx f(x) + J(x)(y - x) = \hat{f}(y)$$

where $J(x)$ is the Jacobian. We now try to solve for $\hat{f}(y) = 0$.

$$f(x) + J(x)(y - x) = 0$$

$$dx = y - x = -J(x)^{-1}f(x) \tag{11.1}$$

dx represents the update step for Newton's method.

We now work on the example of approximating the reciprocal of the Golden Number ϕ . The function and the derivative is as follows.

$$f(x) = \frac{1}{x} - \phi, \quad f'(x) = -\frac{1}{x^2}$$

And dx becomes the following.

$$dx = -J(x)^{-1}f(x) = x^2\left(\frac{1}{x} - \phi\right) = x - x^2\phi$$

The update rule is $x^+ = x + dx$. Figure 11.1 shows an iteration of Newton's method when $x = 1$.

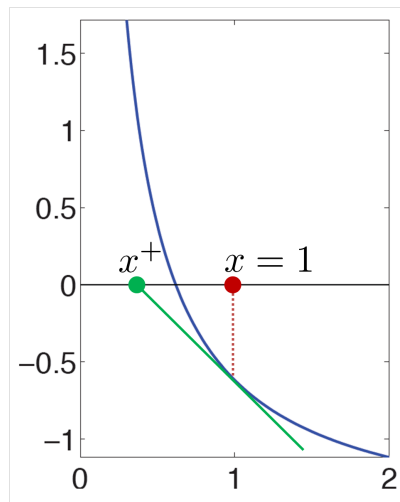


Figure 11.1: Example of one iteration of Newton's Method in solving non-linear equations.

We now perform error analysis of Newton's method. For value x , the error is $\epsilon = x\phi - 1$. For value x^+ , the error ϵ^+ is the following.

$$\begin{aligned}
\epsilon^+ &= x^+ \phi - 1 \\
&= (x + x - x^2 \phi) \phi - 1 \\
&= (x + x(1 - x\phi)) \phi - 1 \\
&= (x - x\epsilon) \phi - 1 \\
&= x\phi - 1 - x\epsilon\phi \\
&= \epsilon - x\epsilon\phi \\
&= \epsilon(1 - x\phi) \\
&= -\epsilon^2
\end{aligned}$$

This shows that if $\epsilon < 1$, then Newton's method has quadratic convergence. However, if $\epsilon > 1$, then Newton's method will diverge.

11.2.2 Finding Minima/Maxima/Saddles

For $x \in \mathbb{R}^d$ and $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ which is twice differentiable, we want to find $\min_x f(x)$. In the example, we only focus on minimizing f , but finding the maxima and saddle points are the same. We first define $g = f'$. Minimizing f is the same as finding x such that $g = f' = 0$. From Equation 11.1, the Newton's update is the following,

$$\begin{aligned}
d &= -J^{-1}g \\
&= -(g')^{-1}g \\
&= -(f'')^{-1}f' \\
&= -H^{-1}g
\end{aligned}$$

where H is the Hessian. We now show that Newton's method is a descent method if $H \succ 0$. We set $dx = td$ for $t > 0$. $r(dx)$ is the residual. Using first order Taylor expansion, we get the following.

$$\begin{aligned}
df &= g^T dx + r(dx) \\
&= g^T t \left(-(f'')^{-1} f' \right) + r(dx) \\
&= -t f'^T (f'')^{-1} f' + r(dx) \\
&= -t f'^T H^{-1} f' + r(dx)
\end{aligned}$$

If $H \succ 0$, then $H^{-1} \succ 0$, which makes the first term always negative, thus making Newton's method a descent method.

11.2.3 Newton's Method and Steepest Descent

Newton's method is a special case of steepest descent when the norm used is the Hessian norm. To find the step for steepest descent, we minimize the following.

$$\min_d g^T d + \frac{1}{2} \|d\|_H^2, \quad \|d\|_H^2 = \sqrt{d^T H d} \quad (11.2)$$

The solution to this minimization is $d = -H^{-1}g$. Steepest descent with a constraint or a penalty in the objective is equivalent. The equivalence will be covered when duality is covered in class. Figure 11.2 shows the difference between the direction of steps for gradient descent and Newton's method.

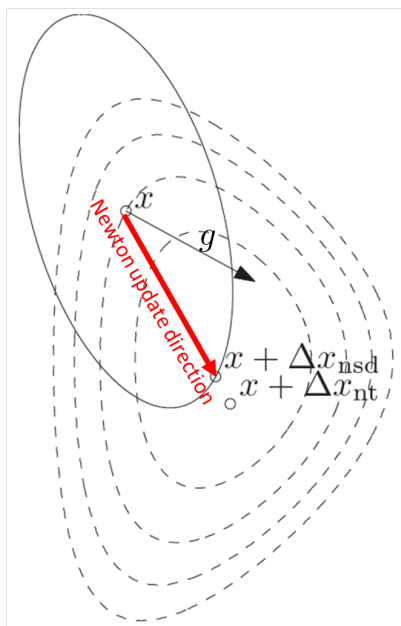


Figure 11.2: Direction of step for gradient descent and Newton's method.

11.2.4 Damped Newton

Damped Newton is combining the Newton's method with backtracking line search to make sure that the objective value does not increase.

```

Initialize  $x_1$ 
for  $k = 1, 2, \dots$ 
     $g_k = f'(x_k);$                                 gradient
     $H_k = f''(x_k);$                                 Hessian
     $d_k = -H_k^{-1}g_k;$                             Newton direction
     $t_k = 1;$                                        backtracking line search
    while  $f(x_k + t_k d_k) > f(x_k) + t_k g_k^T d_k / 3$     divide by 3 to make sure  $< \frac{1}{2}$  for future proofs
         $t_k = \beta t_k$                                  $\beta < 1$ 
     $x_{k+1} = x_k + t_k d_k$                             step
    
```

Damped Newton is affine invariant, meaning that suppose $g(x) = f(Ax + b)$, and we get Newton's updates x_1, x_2, \dots from $g(x)$ and y_1, y_2, \dots from $f(y)$, and if $y_1 = Ax_1 + b$, then $y_i = Ax_i + b \forall i$.

For damped Newton, if f is bounded below, then $f(x_k)$ converges. If f is strictly convex with bounded level sets, then x_k converges. Finally, damped Newton typically converges at quadratic rate in the neighborhood of x^* .