

G.D. & Back tracking line search.

To find an inexact approximation for the stepsize α , we use the wolf condition.

Wolf condition is a combination of 2 different conditions.

① Armijo Condition : (Sufficient decrease condition)

$$F(x_k + \alpha p_k) \leq F(x_k) + c_1 \alpha \nabla f_k^T p_k$$

$$[x_k \Rightarrow \beta_k] \quad = \quad L(\beta_k) + \alpha \nabla(\beta_k) \leq c_1 \alpha \|\nabla \beta_k\|^2$$

L = loss function, β_k = coefficient vector for k^{th} step.

c_1 = chosen $\forall c_1 \in (0, 1)$ [usually close to 10^{-4}]

$\nabla(\beta_k)$ = diff. of β_k = direction.

② Curvature condition: It ensures that stepsizes are not too short.

$$\nabla F(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f_k^T p_k \quad [x_k = \beta_k \text{ \& } p_k = \nabla(\beta_k)]$$

$$\rightarrow \nabla(\beta_k + \alpha_k \nabla(\beta_k))^T \nabla(\beta_k) \geq c_2 \nabla(\beta_k)^T \nabla(\beta_k) \quad \forall c_2 \in (0, 1)$$

It is so that, if slope is very negative, the direction will have negative slope later as well.

Sudocode:

- ① choose initial step size, α_0 . Set $\alpha = \alpha_0$ [1 in this case]
- ② choose $\rho \in (0, 1)$.
- ③ choose a constant, $c \in (0, 1)$
- ④ Repeat till $L(\beta_k - \alpha \nabla(L(\beta_k))) \leq L(\beta_k) + c\alpha \|\nabla(L(\beta_k))\|^2$
- ⑤ Return $\alpha_k = \alpha = \text{Stepsize}$.

PROBLEM 2: Quasi-Newton METHOD :

MOTIVATION: Calculating Hessian is expensive, we want something cheaper.

SOLUTION: SECANT CONDITION?

Hessian is the second derivative of the loss function, to make an appropriate approximation, we can use the definition of derivative.

$$J'(x) = \lim_{h \rightarrow 0} \frac{J(x+h) - J(x)}{h}, \text{ how small is } h??$$

-- (1)

How accurate we want our $J'(x)$ to be.

$$\begin{aligned} B = \text{Hessian} &= \text{Second derivative} = 1^{\text{st}} \text{ derivative of gradient} \\ &= \nabla(L(\beta_{k+1}))' = \frac{\nabla L(\beta_{k+1}) - \nabla L(\beta_k)}{\beta_{k+1} - \beta_k} \end{aligned}$$

-- (2)

How $\beta_{k+1} - \beta_k = h = \text{stepsize}$, equation (1) & (2) are equivalent -

$$\Rightarrow \boxed{H_{k+1} = \frac{\nabla(L(\beta_{k+1})) - \nabla(L(\beta_k))}{\beta_{k+1} - \beta_k}}$$

It is a good Approximation for the Hessian, found from the secant condition

$$\Rightarrow \begin{cases} H_{k+1} S_k = z_k \\ S_k = \beta_{k+1} - \beta_k \\ z_k = \nabla(L(\beta_{k+1})) - \nabla(L(\beta_k)) \end{cases}$$

PSEUDO-CODE : (using limited memory BFGS)

I (a) Initialize, β_i (truncated normal), $H_i = [\text{Identity matrix}]$

(b) calculate loss (-ve log likelihood)

(c) calculate gradient (g_i)

II

(1) compute $p = -(H_{i-1}) g_{i-1}$

(2) $\alpha = \text{step size}$ from line search

(3) update $\beta_i = \beta_{i-1} + \alpha * p$

(4) calculate loss

(5) update gradient

(6) update hessian $H_i = (1 - \rho_i) (\text{diag } \beta_{i-1}) (1 - \rho_i^T) + \rho_i s_i s_i^T$

$$\rho_i = \rho_i * s_i * z_i^T$$

$$\rho_i = \frac{1}{z_i^T s_i}$$

$$z_i = g_i - g_{i-1}$$

$$s = \alpha * p$$

$g_i = \text{gradient at } i^{\text{th}} \text{ step}$

(7) check for convergence, if not,
go to step (1) of (II) part.

if yes, return β_i