

1. Problem.

1. verify the following identity.

$$(Q^{-1} + B^T P^{-1} B)^{-1} B^T P^{-1} = Q B^T (B Q B^T + P)^{-1}$$

So, the given identity is.

$$(Q^{-1} + B^T P^{-1} B)^{-1} B^T P^{-1} = Q B^T (B Q B^T + P)^{-1}$$

we can solve left side of the derivation first

$$(Q^{-1} + B^T P^{-1} B)^{-1} B^T P^{-1}$$

now multiply both sides by  $(B Q B^T + P)$  on right

$$(Q^{-1} + B^T P^{-1} B)^{-1} B^T P^{-1} (B Q B^T + P) = Q B^T$$

we got the derivation after multiplying on right.

now multiply both sides on the left by

$$(Q^{-1} + B^T P^{-1} B).$$

$$B^T P^{-1} (B Q B^T + P) = (Q^{-1} + B^T P^{-1} B) Q B^T$$

simplify the left side

$$B^T P^{-1} (B Q B^T + P) = Q^{-1} Q B^T + B^T P^{-1} B Q B^T$$

Simplify the right side  $QB^T$

We have  $B^T P^{-1} (BQB^T B^T + P)$  on right side  
and  $Q^{-1} QB^T$  on side.

after simplification.

$$B^T P^{-1} (BQB^T B^T + P) = Q^{-1} QB^T$$

$$B^T P^{-1} BQB^T B^T + B^T P^{-1} P = QB^T$$

Since matrix multiplication is associative  
arrange terms as you follow

$$B^T P^{-1} (BQB^T B^T) + B^T P^{-1} P = QB^T$$

$P^{-1}$  in the  $P$  is identity matrix

because any matrix multiply by its inverse  
it is identity matrix

$$B^T P^{-1} (BQB^T B^T) + B^T = QB^T$$

subtract  $B^T$  on both side

$$B^T P^{-1} (BQB^T B^T) = QB^T - B^T$$

To isolate  $B^T P^{-1} (BQB^T B^T)$ , multiply

both sides on left by  $(BQ^T B^T)^{-1}$

$$(BQ^T B^T)^{-1} B^T P^{-1} (BQ^T B^T) = (BQ^T B^T)(QB^T - B^T)$$

$\therefore$  left side simplifies to  $P^{-1}$

$$P^{-1} = (BQ^T B^T)^{-1} (QB^T - B^T)$$

finally to isolate  $P$ , take the inverse of both sides.

$$P = [(BQ^T B^T)^{-1} (QB^T - B^T)]^{-1}$$

this completes the derivation of given identity.



2. Verify the following Woodberg identity.

$$(A + BD^{-1}C)^{-1} = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

given identity is.

$$(A + BD^{-1}C)^{-1} = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

Multiply both sides on the right  $D + CA^{-1}B$

$$(A + BD^{-1}C)^{-1}(D + CA^{-1}B) = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}(D + CA^{-1}B)$$

$\therefore$  now multiply both sides on the left by  $A + BD^{-1}C$

$$(A + BD^{-1}C)^{-1}(D + CA^{-1}B) = (A + BD^{-1}C)A^{-1} -$$

$$(A + BD^{-1}C)A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}(D + CA^{-1}B)$$

simplify the left side.

$$(D + CA^{-1}B) = (A + BD^{-1}C)A^{-1} - (A + BD^{-1}C)$$

$$A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}(D + CA^{-1}B)$$

cancel terms on both sides.

$$(D + CA^{-1}B) = (A + BD^{-1}C)A^{-1} - (A + BD^{-1}C)A^{-1}B$$

since it is identity matrix we can  
arrange terms

$$(D + CA^{-1}B) = A^{-1}(A + BD^{-1}C) - A^{-1}(A + BD^{-1}C)B$$

factor out  $A^{-1}$

$$(D + CA^{-1}B) = A^{-1}(A + BD^{-1}C - (A + BD^{-1}C)B)$$

$$(D + CA^{-1}B) = A^{-1}(A - AB + BD^{-1}CB)$$

$\therefore$  now isolate  $(D + CA^{-1}B)$  by multiplying  
both sides by  $A$ .

$$A(D + CA^{-1}B) = A - AB + BD^{-1}CB$$

$\therefore$  now isolate  $(D + CA^{-1}B)$  by dividing both  
sides by  $A$ .

$$D + CA^{-1}B = I - AB + BD^{-1}CB$$

Subtract  $I - AB$  from both sides.

$$D + CA^{-1}B - (I - AB) = BD^{-1}CB$$

simplify

$$D - I + CA^{-1}B + AB = BD^{-1}CB$$

Rearrange terms

$$(A + BD^{-1}C) - I = B(A - I + (A^{-1}B)D^{-1}C)$$

Isolate the left side

$$(A + BD^{-1}C) - I = B(A - I + (A^{-1}B)D^{-1}C)$$

$\therefore$  finally add  $I$  to both sides.

$$(A + BD^{-1}C) = I + B(A - I + (A^{-1}B)D^{-1}C)$$

finally after the verification the equation/derivative is.

$$(A + BD^{-1}C)^{-1} = A^{-1} - A^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

Hence Proved.



Problem 2:

1. Given  $x = [x_1 : x_2 : x_3] \in \mathbb{R}^3$  and  $y = [y_1 : y_2] \in \mathbb{R}^2$

where  $y_1 = x_1^2 - x_2$  and  $y_2 = x_3^2 + 3x_2$  compute  $dy/dx$ .

differentiation in Partial ways.

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \end{bmatrix}$$

$$\frac{\partial y_1}{\partial x_1} = \frac{\partial}{\partial x_1} (x_1^2 - x_2)$$

$$= 2x_1$$

$$\frac{\partial y_2}{\partial x_1} = \frac{\partial}{\partial x_1} (x_3^2 + 3x_2) = 0$$

$$\frac{\partial y_1}{\partial x_2} = \frac{\partial}{\partial x_2} (x_1^2 - x_2)$$

$$\frac{\partial y_2}{\partial x_2} = \frac{\partial}{\partial x_2} (x_3^2 + 3x_2)$$

$$\frac{\partial y_1}{\partial x_3} = \frac{\partial}{\partial x_3} (x_1^2 - x_2)$$

$$= 0$$

$$\frac{\partial y_2}{\partial x_3} = \frac{\partial}{\partial x_3} (x_3^2 + 3x_2)$$

$$= 2x_3 + 0$$

$$= 2x_3$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} \\ \frac{\partial y_1}{\partial x_3} & \frac{\partial y_2}{\partial x_3} \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} 2x_1 & 0 \\ -1 & 6x_2 \\ 0 & 2x_3 \end{bmatrix}$$

$\therefore$  The  $2 \times 3$  derivative matrix  $\partial y / \partial x$  is as shown above.

2.

$$x = r \sin \theta \cos \phi$$

$$y = r \sin \theta \sin \phi$$

$$z = r \cos \theta$$

$$x = [x; y; z] \Rightarrow [x_1; x_2; x_3]$$

$$y = [r; \theta; \phi] \Rightarrow [y_1; y_2; y_3]$$

$$\frac{\partial x}{\partial y} = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_2}{\partial y_1} & \frac{\partial x_3}{\partial y_1} \\ \frac{\partial x_1}{\partial y_2} & \frac{\partial x_2}{\partial y_2} & \frac{\partial x_3}{\partial y_2} \\ \frac{\partial x_1}{\partial y_3} & \frac{\partial x_2}{\partial y_3} & \frac{\partial x_3}{\partial y_3} \end{bmatrix}$$



to compute jacobian  $\partial x/\partial y$

$$\partial x/\partial r = \sin(\theta) * \cos(\phi)$$

$$\partial x/\partial \theta = r * \cos(\theta) * \cos(\phi)$$

$$\partial x/\partial \phi = -r * \sin(\theta) * \sin(\phi)$$

$$\partial y/\partial r = \sin(\theta) * \sin(\phi)$$

$$\partial y/\partial \theta = r * \cos(\theta) * \sin(\phi)$$

$$\partial y/\partial \phi = r * \sin(\theta) * \cos(\phi)$$

$$\partial z/\partial r = \cos(\theta)$$

$$\partial z/\partial \theta = -r * \sin(\theta)$$

$$\partial z/\partial \phi = 0$$

After Arranging this into  $3 \times 3$  jacobian,

$$\begin{bmatrix} \sin(\theta) * \cos(\phi) & r * \cos(\theta) * \cos(\phi) & -r * \sin(\theta) * \sin(\phi) \\ \sin(\theta) * \sin(\phi) & r * \cos(\theta) * \sin(\phi) & r * \sin(\theta) * \cos(\phi) \\ \cos(\theta) & -r * \sin(\theta) & 0 \end{bmatrix}$$

$\therefore$  the jacobian  $\partial x/\partial y$  is the  $3 \times 3$  matrix

Problem 3:

⇒

1. The Hessian of the least square loss is

$$L(w) = \left(\frac{1}{2}\right) \sum_{i=1}^n (x_{-i}^T w - y_{-i})^2$$

where  $x_{-i}$  is input

$y_{-i}$  is the  $i$ th target and  $w$  is

Parameter vector

we should take 2nd Partial derivative

$$\partial^2 L / \partial w_j \partial w_k = \sum_{i=1}^n x_{-i,j} x_{-i,k}$$

forms the Hessian matrix  $H$

$$H = X^T X$$

$X$  is design matrix, each row is a feature vector  $x_{-i}$ .

2.

Given,

The first iteration of Newton's method gives us  $w^* = (X X^T)^{-1} X y$

new rule is

$$w \leftarrow w - H^{-1} \nabla L(w)$$

$$= w - (X^T X)^{-1} X^T (Xw - y)$$

as of the definition

$$w \leftarrow (X^T X)^{-1} X^T y$$

$\therefore$  Newton method converges immediately in one iteration for linear regression.

Hence  $w^* = (X^T X)^{-1} X^T y$



#### 4. Problem

1. The constrained optimization Problem

$$\min L(w) = \sum (f(x_n; w) - t_n)^2$$

$$\text{Subject to } \|w\|_P \leq \gamma$$

To convert Lagrangian form

$$L(w, \lambda) = \sum (f(x_n; w) - t_n)^2 + \lambda (\|w\|_P - \gamma)$$

Setting the derivative w.r.t  $w$  to 0.

$\therefore$  It gives optimality condition:

$$\nabla L(w, \lambda) = \nabla L(w) + \lambda \nabla \|w\|_P = 0$$

which is same as the optimality condition for

$$\min L(w) = \sum (f(x_n; w) - t_n)^2 + \lambda \|w\|_P$$

So the two problems are equivalent

At the optimal  $w$ , the constraint is satisfied as equality.

$$\|w\|_P = \gamma$$

plugging this into the Lagrange dual function

$$L(w, \lambda) = L(w) + \lambda(Y - \hat{Y}) = L(w)$$

2.  $\lambda = \gamma$  makes the solutions equivalent for the hyperparameter  $\lambda$  and  $\gamma$ .

1.  $\gamma$  controls the constraint boundary formulation

2.  $\lambda$  controls the regularization strength in the regularized formulation.

3.  $\lambda$  is directly optimized  $\gamma$  is imposed as a constraint.

4.  $\lambda$  and  $\gamma$  play similar roles in controlling model complexity, with  $\lambda$  being the optimized parameter.



## 5. Problem

1. By gradient update rule and continuity of  $\nabla f$ .

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

using this with  $x = x(k)$ ,  $y = x(k+1)$ ;

$$f(x(k+1)) \geq f(x(k)) + \nabla f(x(k))^T (x(k+1) - x(k))$$

$$\text{GD update } x(k+1) = x(k) - a \nabla f(x(k))$$

$$\begin{aligned} f(x(k+1)) &\geq f(x(k)) - a \|\nabla f(x(k))\|_2^2 \\ &\leq f(x(k)) - (1 - 2a/2) a \|\nabla f(x(k))\|_2^2 \end{aligned}$$

So  $f$  decreases by less  $(a/2) \|\nabla f(x(k))\|_2^2$

Per iteration

2. By convexity of  $f$ :

By Lipschitz condition

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

Applying this to  $x = x(k)$  and  $y = x(k+1)$ ,

$$\begin{aligned} f(x(k+1)) &\leq f(x(k)) - a \|\nabla f(x(k))\|_2^2 + \\ &\quad (2a^2/2) \|\nabla f(x(k))\|_2^2 \end{aligned}$$



using  $\alpha \leq 1/2$ , we get

$$\begin{aligned} f(x(k+1)) &\leq f(x(k)) - (1 - L\alpha/2)\alpha \|\nabla f(x(k))\|^2 \\ &\leq f(x(k)) - (1/2)\alpha \|\nabla f(x(k))\|^2 \end{aligned}$$

3. Summing over  $k$  iterations and using  $f(x) \geq f(x^*)$ :

$$\begin{aligned} \sum_{k=0}^{K-1} [f(x(k)) - f(x^*)] &\leq (1/2)\alpha \\ &\quad \|\nabla f(x(0))\|^2 \end{aligned}$$

$$\text{So } f(x(K)) - f(x^*) \leq (1/2\alpha K) \|\nabla f(x(0))\|^2$$

Therefore, gradient descent ~~converges~~ converges  
at rate  $O(1/K)$ . This also proved that  
convergence rate for gradient descent on  
convex differentiable functions.