

Problem #1 → Given  $Q \in R^{N \times N}$ ,  $P \in R^{M \times M}$ , and  $B \in R^{M \times N}$

$$\text{Need to prove } (Q^{-1} + B^T P^{-1} B)^{-1} B^T P^{-1} = Q B^T (B Q B^T + P^{-1})^{-1}$$

Proof → Let's Multiply both sides by  $(Q^{-1} + B^T P^{-1} B)$

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

cancels out

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

Now, let's Multiply both sides by  $(B Q B^T + P)$

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

cancels out

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

$$\Rightarrow B^T P^{-1} B Q B^T + B^T P^{-1} P = \underbrace{Q^{-1} Q B^T}_{\text{cancels out}} + B^T P^{-1} B Q B^T$$

$$\Rightarrow B^T P^{-1} B Q B^T + B^T = B^T P^{-1} B Q B^T + B^T$$

$$\Rightarrow L.H.S = R.H.S, \text{ hence proved}$$

⊕ This equation is mentioned as a lemma in the paper.

Kalman filtering by Max Welling.

2nd part → Proof of Woodbury identity -

$$\text{Need to show } (A + B D^{-1} C)^{-1} = A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}$$

Proof → Let's Multiply both sides by  $(A + B D^{-1} C)$

$$\Rightarrow (A + B D^{-1} C) \cdot (A + B D^{-1} C)^{-1} = (A + B D^{-1} C) [A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}]$$

cancels out

$$\Rightarrow I = \underbrace{\{A A^{-1} - A A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}\}}_{\text{cancels out}} + \{B D^{-1} C A^{-1} + B D^{-1} C A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}\}$$

$$\Rightarrow I = \{I - B (D + C A^{-1} B)^{-1} C A^{-1}\} + \{ \dots \}$$

Date \_\_\_ / \_\_\_ / \_\_\_

$$\Rightarrow I = I - B(D + CA^{-1}B)^{-1}CA^{-1} + BD^{-1}CA^{-1} + BD^{-1}CA^{-1}B(D + CA^{-1}B)^{-1}CA^{-1}$$

Now, taking common  $(D + CA^{-1}B)^{-1}CA^{-1}$  from both the above terms

$$\Rightarrow I = I \cancel{+ BD^{-1}CA^{-1}} + [(D + CA^{-1}B)^{-1}CA^{-1} \cdot (B + BD^{-1}CA^{-1}B)]$$

Next, taking common  $BD^{-1}$  from above 2 terms

$$\Rightarrow I = I \cancel{+ BD^{-1}CA^{-1}} + [(D + CA^{-1}B)^{-1}CA^{-1} \cdot BD^{-1}(D + CA^{-1}B)]$$

Now, cancelling out  $(D + CA^{-1}B)^{-1}$  and  $(D + CA^{-1}B)$ , we get  $\rightarrow$

$$I = I \cancel{+ BD^{-1}CA^{-1}} + [CA^{-1} \cdot BD^{-1}]$$

$$\Rightarrow I = I \quad (\text{L.H.S} = \text{R.H.S}), \text{ hence proved}$$

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

$$y_1 = x_1^2 - x_2, y_2 = x_3^2 + 3x_2$$

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

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

$$\frac{\partial y_1}{\partial x_2} = \frac{\partial(x_1^2 - x_2)}{\partial x_2} = -1$$

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

$$\frac{\partial y_1}{\partial x_3} = \frac{\partial(x_1^2 - x_2)}{\partial x_3} = 0$$

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

$$\Rightarrow \frac{\partial y}{\partial x} = \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} \\ \frac{\partial y_1}{\partial x_3} & \frac{\partial y_2}{\partial x_3} \end{vmatrix} = \begin{vmatrix} 2x_1 & 0 \\ -1 & 3 \\ 0 & 2x_3 \end{vmatrix}$$

2nd part  $x = r \sin \theta \cos \phi, y = r \sin \theta \sin \phi, z = r \cos \theta$

here  $r > 0, 0 < \theta < \pi$  and  $0 \leq \phi < 2\pi$ .  $x = [x; y; z]$

and  $y = [r; \theta; \phi]$

Date / /

$$\text{Solution} \rightarrow \frac{\partial x}{\partial r} = \frac{\partial(r \sin \theta \cos \phi)}{\partial r} = \sin \theta \cos \phi$$

$$\frac{\partial x}{\partial \theta} = \frac{\partial(r \sin \theta \cos \phi)}{\partial \theta} = r \cos \theta \cos \phi$$

$$\frac{\partial x}{\partial \phi} = \frac{\partial(r \sin \theta \cos \phi)}{\partial \phi} = -r \sin \theta \sin \phi$$

$$\frac{\partial y}{\partial r} = \frac{\partial(r \sin \theta \sin \phi)}{\partial r} = \sin \theta \sin \phi$$

$$\frac{\partial y}{\partial \theta} = \frac{\partial(r \sin \theta \sin \phi)}{\partial \theta} = r \cos \theta \sin \phi \cdot \sin \phi$$

$$\frac{\partial y}{\partial \phi} = \frac{\partial(r \sin \theta \sin \phi)}{\partial \phi} = r \sin \theta \cos \phi$$

$$\frac{\partial z}{\partial r} = \frac{\partial(r \cos \theta)}{\partial r} = \cos \theta$$

$$\frac{\partial z}{\partial \theta} = \frac{\partial(r \cos \theta)}{\partial \theta} = -r \sin \theta$$

$$\frac{\partial z}{\partial \phi} = \frac{\partial(r \cos \theta)}{\partial \phi} = 0$$

$$\Rightarrow \begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} & \frac{\partial x}{\partial \phi} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} & \frac{\partial y}{\partial \phi} \\ \frac{\partial z}{\partial r} & \frac{\partial z}{\partial \theta} & \frac{\partial z}{\partial \phi} \end{bmatrix} = \begin{bmatrix} \sin \theta \cos \phi & \sin \theta \sin \phi & \cos \theta \\ r \cos \theta \cos \phi & r \cos \theta \sin \phi & -r \sin \theta \\ -r \sin \theta \sin \phi & r \sin \theta \cos \phi & 0 \end{bmatrix}$$

Problem 3 - Newton's Method to solve least squares loss in Linear regression

Solution → 1. The Hessian of the least squares loss function

$$L(w) = \frac{1}{2} \sum_{i=1}^n (x_i^T w - y_i)^2 \text{ w.r.t parameter vector } w$$

To find the Hessian, we will have to compute 2nd derivative

$$\nabla^2 L(w) \text{ w.r.t } w$$

so firstly, the 1st derivative (gradient) of  $L(w)$  =

$$\nabla L(w) = \frac{\partial L(w)}{\partial w} = \partial \left( \frac{1}{2} * \sum_{i=1}^n x_i^T w - y_i \right)^2 = \sum_{i=1}^n (x_i^T w - y_i) * x_i$$

Now, Hessian (2nd derivative) of  $L(w)$  w.r.t  $w$  =

$$\text{Hessian}(H) = \nabla^2 L(w) = \frac{\partial^2 L(w)}{\partial w^2} = \frac{\partial^2}{\partial w^2} \left( \sum_{i=1}^n (x_i^T w - y_i) x_i \right)$$

$$= \sum_{i=1}^n x_i^T x_i$$

So, the Hessian is the summation of outer products of the input data points  $x_i$ .

2: Show that first iteration of Newton's method gives us.

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

solution → We know that  $L(w) = \frac{1}{2} \sum_{i=1}^n (x_i^T w - y_i)^2$

$$\text{Gradient of } L(w) = \sum_{i=1}^n (x_i^T w - y_i) x_i = X^T X w - X^T y$$

$$\text{Hessian of } L(w) = \sum_{i=1}^n x_i^T x_i = X^T X$$

Now, Newton's method sets the parameters to an initial guess  $w_0$ , then iteratively them. Let  $w_t$  be the current parameters on iteration  $t$ . Then, the updated parameter  $w_{t+1}$  as Newton's method gives :

$$w_{t+1} = w_t - H_L(w_t)^{-1} \cdot \nabla L(w_t) \quad \begin{matrix} \text{(Product of Hessian inverse)} \\ \text{(and gradient)} \end{matrix}$$

$$= w_t - (X^T X)^{-1} \cdot (X^T X w_t - X^T y)$$

$$= w_t - (X^T X)^{-1} \cdot (X^T X) w_t + (X^T X)^{-1} X^T y$$

$$\Rightarrow w_{t+1} = (X^T X)^{-1} X^T y \quad , \text{ hence, Newton's method converges in first iteration.}$$

5. Need to prove the following theorem:-

Suppose the function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex & differentiable & that is gradient  $\nabla f$  is Lipschitz continuous with constant  $L > 0$ , then if we run gradient descent for  $K$  iterations with a fixed Learning rate  $0 < \alpha \leq \frac{1}{L}$ , it will satisfies:-

$$f(x^K) - f(x^*) \leq \frac{1}{2\alpha K} \|x^0 - x^*\|_2^2, \quad (\text{where } f(x^*) \text{ is the optimal value})$$

We will prove the given hints step by step in order to prove this theorem finally.

$$1) \text{ Hints 1: } f(x^{K+1}) \leq f(x^K) - (1 - \frac{1}{2} L \alpha) \cdot \alpha \|\nabla f(x^K)\|_2^2$$

**Proof :-** Given Lipschitz continuity of the gradient ( $\nabla f$ ), we know that:-  $\|f(y) - f(x)\|_2 \leq L \|y - x\|_2$

Now, applying this gradient update equation of gradient descent:-

$$f(x^{K+1}) \leq f(x^K) + \nabla f(x^K)^T (x^{K+1} - x^K) + \frac{L}{2}$$

$$f(x^{K+1}) \leq f(x^K) + \nabla f(x^K)^T (x^{K+1} - x^K) + \frac{L}{2} \|x^{K+1} - x^K\|_2^2$$

Using gradient descent update rule:-

$$x^{K+1} = x^K - \alpha \nabla f(x^K)$$

putting this above value in our equation:-

$$f(x^{K+1}) \leq f(x^K) - \alpha \nabla f(x^K)^T \nabla f(x^K) + \frac{L\alpha^2}{2} \|\nabla f(x^K)\|_2^2$$

$$\Rightarrow f(x^{K+1}) \leq f(x^K) - \alpha \|\nabla f(x^K)\|_2^2 + \frac{L\alpha^2}{2} \|\nabla f(x^K)\|_2^2$$

$$\Rightarrow f(x^{K+1}) \leq f(x^K) - \left(1 - \frac{1}{2} L \alpha\right) \alpha \|\nabla f(x^K)\|_2^2$$

$$\text{Now, since } \alpha \leq \frac{1}{L}, \text{ we have } 1 - \frac{1}{2} L \alpha \geq \frac{1}{2}$$

Hence, inequality will be -

$$f(x^{K+1}) \leq f(x^K) - \frac{\alpha}{2} \|\nabla f(x^K)\|_2^2$$

$$2) \text{ Hint 2: } f(x^{k+1}) \leq f(x^*) - \frac{\alpha}{2}$$

$$\text{Hint 2: } f(x^{k+1}) - f(x^*) \leq \frac{1}{2\alpha} \left( \|x^k - x^*\|_2^2 - \|x^k - \alpha \nabla f(x^k) - x^*\|_2^2 \right)$$

Now, using Hint 1 inequality we got:-

$$f(x^{k+1}) \leq f(x^k) - \frac{\alpha}{2} \|\nabla f(x^k)\|_2^2$$

$$\Rightarrow f(x^{k+1}) - f(x^*) \leq f(x^k) - f(x^*) - \frac{\alpha}{2} \|\nabla f(x^k)\|_2^2$$

Putting Lipschitz property i.e.  $f(x^k) - f(x^*) \leq \nabla f(x^k)^T (x^k - x^*)$

$$\Rightarrow f(x^{k+1}) - f(x^*) \leq \nabla f(x^k)^T (x^k - x^*) - \frac{\alpha}{2} \|\nabla f(x^k)\|_2^2$$

$$\Rightarrow f(x^{k+1}) - f(x^*) \leq \frac{1}{2\alpha} \|x^k - x^*\|_2^2 - \left[ \left( \|\nabla f(x^k)\|_2^2 \right) - \nabla f(x^k)^T (x^k - x^*) \right]$$

By applying gradient descent update rule., we get -

$$\Rightarrow f(x^{k+1}) - f(x^*) \leq \frac{1}{2\alpha} \|x^k - x^*\|_2^2 - \|\nabla f(x^k) - (x^k - x^*)\|_2^2$$

$$\text{Hint 3: } \sum_{k=1}^K [f(x^k) - f(x^*)] \leq \frac{1}{2\alpha} \|x_0 - x^*\|_2^2$$

Summation of Hint 2 formulae/inequality from  $k=1$  to  $K \rightarrow$

$$\sum_{k=1}^K (f(x^{k+1}) - f(x^*)) \leq \sum_{k=1}^K \left( \frac{1}{2\alpha} \|x^k - x^*\|_2^2 - \|\nabla f(x^k) - (x^k - x^*)\|_2^2 \right)$$

$$\leq \sum_{k=1}^K (f(x^{k+1}) -$$

$$\leq \frac{1}{2\alpha} \sum_{k=1}^K \|x^k - x^*\|_2^2 - \sum_{k=1}^K \|\nabla f(x^k) - (x^k - x^*)\|_2^2$$

Since,  $\sum_{k=1}^K (f(x^{k+1}) - f(x^*)) = f(x^{K+1}) - f(x^*)$

on applying property,  $f$  is decreasing on every iteration

$$f(x^{K+1}) \leq f(x^K) \leq \dots \leq f(x^2) \leq f(x^1)$$

Date \_\_\_\_\_

Solutions

Re-evaluating L.H.S as  $K(f(x') - f(x^*))$ 

$$\Rightarrow K(f(x') - f(x^*)) \leq \frac{1}{2\alpha} \sum_{k=1}^K \|x^k - x^*\|_2^2 - \sum_{k=1}^K \| \nabla f(x^k) - (x^k - x^*) \|_2^2$$

Dividing both sides by  $K$ 

$$\Rightarrow f(x') - f(x^*) \leq \frac{1}{2\alpha K} \sum_{k=1}^K \|x^k - x^*\|_2^2 - \frac{1}{K} \sum_{k=1}^K \| \nabla f(x^k) - (x^k - x^*) \|_2^2$$

Taking Limit as  $K \rightarrow \infty$ ,  $\frac{1}{K} \sum_{k=1}^K \|x^k - x^*\|_2^2$ approaches  $\frac{1}{K} \|x^0 - x^*\|_2^2$ and  $\frac{1}{K} \sum_{k=1}^K \| \nabla f(x^k) - (x^k - x^*) \|_2^2$  approaches 0.As  $k$  gets larger,  $\| \nabla f(x^k) - (x^k - x^*) \|$  becomes smaller, since gradient descent is converging.Hence, in limit  $f(x') - f(x^*) \leq \frac{1}{2\alpha} \|x^0 - x^*\|_2^2$ 

Therefore, we can say here -

$$\sum_{k=1}^K f(x^{k+1}) - f(x^*) \leq \frac{1}{2\alpha} \|x^0 - x^*\|_2^2$$

$$\Rightarrow f(x^k) - f(x^*) \leq \frac{1}{2\alpha k} \|x^0 - x^*\|_2^2 \quad \boxed{\text{Theorem proved.}}$$

Thus, this completes the proof and depicts that gradient descent with the given conditions converges at a specified rate.

4. Given, ① Minimizing ordinal least square loss subject to  
 $L_p$  norm - Minimize  $L(w) = \sum f(x_n, w) - t_n)^2 ; \|w\|_p^p \leq \gamma$

② Minimizing regularized least square loss with  $L_p$  regularization.  
 Minimize  $L(w) = \sum (f(x_n, w) - t_n)^2 + \lambda \|w\|_p^p$

Starting with OLS equation, applying Lagrange multiplier:

$$L_1(w, \alpha) = \sum (f(x_n, w) - t_n)^2 + \alpha (\|w\|_p^p - \gamma)$$

on taking partial derivative of  $L_1$  w.r.t.  $w$  and  $\alpha$

$$\frac{\partial L_1}{\partial w} = 2 \sum f(x_n, w) - t_n \frac{\partial f}{\partial w} - \alpha p \|w\|_p^{p-1} w / \|w\|_p = 0$$

$$\frac{\partial L_1}{\partial \alpha} = \|w\|_p^p - \gamma = 0$$

$\Rightarrow \|w\|_p = \gamma$ , putting this in above equation -

$$2 \sum f(x_n, w) - t_n \frac{\partial f}{\partial w} - \alpha p \gamma^{(p-1)} w / \gamma = 0$$

$$\text{or } 2 \sum f(x_n, w) - t_n \frac{\partial f}{\partial w} - \alpha p \|w\|_p^{p-1} w / \|w\|_p = 0$$

Next, applying Lagrangian in Regularized Least square with Lagrange Multiplier  $\beta$  -

$$L_2(w, \beta) = \sum (f(x_n, w) - t_n)^2 + \beta \|w\|_p^{p-1}$$

on taking partial derivative. of  $L_2$  w.r.t  $w$  and  $\beta$

$$\frac{\partial L_2}{\partial w} = 2 \sum (f(x_n, w) - t_n) \frac{\partial f}{\partial w} + \beta (p-1) \|w\|_p^{p-2} w = 0$$

$$\frac{\partial L_2}{\partial \beta} = \|w\|_p^{p-1} = 0$$

$\hookrightarrow$  we can write  $\|w\|_p^p = 0$ , as  $p$  is constant

This means, we get  $\|w\|_p^p = \gamma$  and also  $\|w\|_p^p = 0$

$$\Rightarrow \gamma = 0 \quad (\text{under given conditions})$$

$$\text{Hence, } \min_w L(w) = \sum_{n=1}^N (f(x_n; w) - t_n)^2, \text{ s.t. } \|w\|_p^p \leq \gamma$$

$$\Leftrightarrow \min_w L(w) \sum_{n=1}^N (f(x_n; w) - t_n)^2 + \lambda \|w\|_p^p$$

In summary, minimizing OLS with  $L_p$  norm is equivalent to minimizing  $L_p$  regularised least square loss when  $\gamma = 0$ , where  $\lambda$  controls the strength of regularization.

## Part 2 Relationship between $\lambda$ and $\gamma$ .

In the regularised least square loss problem, the hyperparameter  $\lambda$  controls the strength of regularization.

The larger the value of  $\lambda$ , the stronger the regularization and more will be penalizing of weight vector  $w$ .

On the other hand,  $L_p$  norm constraint,  $\gamma$  sets an upper bound on the  $L_p$  norm of the weight vector  $w$ . If  $\gamma$  is small, it means  $L_p$  norm constraint is strict, and  $w$  must be close to the origin. If  $\gamma$  is large, it allows  $w$  to have a larger  $L_p$  norm.

The relationship between  $\lambda$  and  $\gamma$  depends on the particular optimization problem and value of other parameters, generally  $\lambda$  controls trade-off between fitting the data and having smaller  $L_p$  norm  $w$ .

The accurate relationship would require other details of dataset, model, optimization problem etc. In general, these hyperparameters are tuned by cross-validation to find best balance between fitting and regularization for a given task.