

CS484. - Introduction to Machine Learning
Assignment - 1 (10 marks)

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

i) Given, $Q \in \mathbb{R}^{N \times N}$, $P \in \mathbb{R}^{M \times M}$, and $B \in \mathbb{R}^{M \times N}$
 we must verify the following identity:

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

Here, we use Sherman-Morrison-Woodbury formula to prove

$$(A + U C V)^{-1} = A^{-1} - A^{-1} U C C^{-1} + V A^{-1} (U)^{-1} V A^{-1}$$

Here,

A is invertible matrix

U and V are matrices

C is an invertible matrix

$$\text{so, } A = Q^{-1}$$

$$U = B^T$$

$$C = P^{-1}$$

$$V = B$$

Using these our LHS becomes :-

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

And the RHS becomes :-

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

Now to simplify let us multiply both sides with $(B^T P^{-1})$ to the equation. we obtain

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

so, As the matrix multiplication is associative we can rewrite the equation as:

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

Hence, the identity is proved.

2) Verify Woodbury Identity

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

Step 1 :- Multiply both sides with $(A + BD^{-1}C)$

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

Step 2 :- LHS becomes I , as we know

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

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

Step 3 :- Now let us solve the RHS, we get

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

Step 4 :- After simplifying RHS, we get

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

Step 5 :- Let us simplify further, we get

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

so, we get

$$I + A^{-1}BD^{-1}C - A^{-1}B^{-1}C$$

Step 6 :- Now combine LHS & RHS

$$I = I + \cancel{A^{-1}BD^2C} - \cancel{A^{-1}BD^{-1}C}$$

$$I = I + 0$$

$$I = I$$

Hence, we get LHS = RHS.

so, proved.

Problem - 2 :-

Given,

$$x = [x_1; x_2; x_3] \in \mathbb{R}^3$$

$$y = [y_1; y_2] \in \mathbb{R}^2$$

$$y_1 = x_1^2 - x_2$$

$$y_2 = x_3^2 + 3x_2$$

Now, we need to compute Jacobian matrix $\frac{\partial y}{\partial x}$

The Jacobian matrix is :

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

Partial Derivatives :-

$$\frac{\partial y_1}{\partial x_1} = \frac{\partial}{\partial x_1} (x_1^2 - x_2) = 2x_1 \quad \frac{\partial y_2}{\partial x_1} = \frac{\partial}{\partial x_1} (x_3^2 + 3x_2) = 0$$

$$\frac{\partial y_1}{\partial x_2} = \frac{\partial}{\partial x_2} (x_1^2 - x_2) = -1 \quad \frac{\partial y_2}{\partial x_2} = \frac{\partial}{\partial x_2} (x_3^2 + 3x_2) = 3$$

$$\frac{\partial y_1}{\partial x_3} = \frac{\partial}{\partial x_3} (x_1^2 - x_2) = 0 \quad \frac{\partial y_2}{\partial x_3} = \frac{\partial}{\partial x_3} (x_3^2 + 3x_2) = 2x_3$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} 2x_1 & -1 & 0 \\ 0 & 3 & 2x_3 \end{bmatrix}$$

2) Given, x, y and z are cylindrical coordinates.

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

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

$$z = r \cos \theta$$

where $r > 0$, $0 < \theta < \pi$, and $0 \leq \phi < 2\pi$

The Jacobian matrix is:

$$\frac{\partial \mathbf{x}}{\partial \mathbf{y}} = \frac{\partial [x; y; z]}{\partial [r; \theta; \phi]}$$

$$= \begin{bmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} & \frac{\partial x}{\partial \phi} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} & \frac{\partial y}{\partial \phi} \\ \frac{\partial z}{\partial r} & \frac{\partial z}{\partial \theta} & \frac{\partial z}{\partial \phi} \end{bmatrix}$$

Now, we need to calculate the partial derivatives of x, y, z .

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

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

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

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

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

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

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

$\phi \text{ and } \theta \text{ mix } \Rightarrow \alpha \cos \theta \sin \phi$

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

$\theta \text{ and } \phi \text{ mix } \Rightarrow \alpha \sin \theta \cos \phi$

$$z = \alpha \cos \theta \quad : \text{all Cartesian coordinates are}$$

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

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

$$\frac{\partial z}{\partial \phi} = \frac{\partial}{\partial \phi} (\alpha \cos \theta) = \begin{pmatrix} 0 & \frac{x}{\alpha} & \frac{y}{\alpha} \\ \frac{-x}{\alpha} & 0 & \frac{y}{\alpha} \\ \frac{-y}{\alpha} & \frac{-x}{\alpha} & 0 \end{pmatrix}$$

The Jacobian matrix is:

$$\frac{\partial x}{\partial y} = \begin{pmatrix} \sin \theta \cos \phi & \alpha \cos \theta \cos \phi & \alpha \sin \theta \sin \phi \\ \sin \theta \sin \phi & \alpha \cos \theta \sin \phi & \alpha \sin \theta \cos \phi \\ \cos \theta & -\alpha \sin \theta \cos \theta \sin \phi & \alpha \cos \theta \cos \phi \\ \cos \theta \sin \phi & \alpha \cos \theta \sin \phi & \alpha \sin \theta \cos \phi \end{pmatrix}$$

Problem - 3 :-

Given,

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

Here, we need to find Hessian of Least Square Loss

Step 1:- express the loss function in matrix form
Loss function can be written as:

$$L(w) = \frac{1}{2} \|xw - y\|^2$$

so, we can re-write the above loss function as

$$L(w) = \frac{1}{2} (xw - y)^T (xw - y)$$

Step 2:- compute the gradient of $L(w)$ with respect to w .

we get;

$$\frac{\partial L(w)}{\partial w} = x^T (xw - y)$$

Step 3:- Now, we compute Hessian; which is the derivative of gradient

$$\frac{\partial^2}{\partial w^2} L(w) = \frac{\partial}{\partial w} (x^T (xw - y)) = x^T x$$

So, Hessian is $x^T x$

2) Given,

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

Now, we need to show that the first iteration of Newton's method gives us the above equation.

Newton's method is an iterative optimization algorithm. It updates parameters w according to the rule:

$$w_{\text{new}} = w - H(w)^{-1} \nabla_w L(w)$$

where,

- * w_{new} is the updated parameter vector
- * $H(w)$ is the Hessian matrix
- * $\nabla_w L(w)$ is the gradient of loss function

Step 1 :- Substitute the values of $H(w)$ and $\nabla_w L(w)$. From previous problem we know that,

$$H(w) = X^T X$$

$$\nabla_w L(w) = X^T (X w - y)$$

$$w_{\text{new}} = w - (X^T X)^{-1} X^T (X w - y)$$

Step 2 :- Now, we need to simplify the expression above

$$\omega_{\text{new}} = \omega - (X^T X)^{-1} X^T X \omega + (X^T X)^{-1} X^T y$$

here, $(X^T X)^{-1} X^T X \omega = \omega$ because $(X^T X)^{-1}$
 $X^T X$ is an identity matrix
so,

$$\omega_{\text{new}} = \omega - \omega + (X^T X)^{-1} X^T y$$

$$\omega_{\text{new}} = (X^T X)^{-1} X^T y$$

The ω_{new} is the solution to the normal equations for the least squares problem:

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

This is the optimal solution to the least squares problem, which means that Newton's method converges to the optimal solution just after one iteration.

Problem - 4 :-

1) Here, we are asked to show that minimizing the least square loss subject to an L_p -norm constraint is equivalent to minimizing the regularized least square loss with an L_p -norm regularization term.

Given,

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

Here,

- * w is the model parameters (weights)
- * $f(x_n; w)$ is the model prediction for data point x_n
- * t_n is the target value for data point x_n .
- * $\|w\|_p$ is the L_p -norm of the vector w , which can represent L_1 -norm and L_2 -norm depending on the choice of P .

we need to show that this constrained optimization is equivalent to the below unconstrained optimization problem:

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

where λ is regulation parameter.

Step 1: Form the Lagrangian.
To solve the constrained optimization problem, we use Lagrange multipliers. The idea is to transform the constrained into unconstrained one by introducing a Lagrange multiplier λ .

The Lagrangian for the constrained is:-

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

Step 2: Deriving the Unconstrained form
Now, we can transform the constrained problem into an unconstrained problem by minimizing the Lagrangian.

By re-arranging Lagrangian, we get:

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

Since γ is a constant it doesn't affect optimization over w . we can drop $\rightarrow \gamma$

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

Thus, we get the transformed constrained problem into the following unconstrained problem:

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

The above is the form of regularized least squares problem, where λ controls the strength of the regularization and $||w||_p^p$ is the regularization term.

2) Relationship between Hyperparameters
and λ and γ :

The hyperparameter λ in the constrained formulation controls the strength of the regularization, while γ is a threshold on the norm of the parameter vector w in the constrained formulation.

- * λ and γ are inversely related: Specifically as λ increases the regularization term $\lambda \|w\|_p$ becomes more dominant, forcing the norm of w to decrease. Hence, a large λ implies a small γ , i.e., a smaller upper bound on the l_p -norm of w .
- * Conversely, if γ is large (the constraint is loose), this implies that the norm of w can be larger, so λ should be small to allow for less penalization.

Problem 5:- we are proving the convergence result for gradient descent on convex functions with Lipschitz continuous gradients. The problem asks us to show that after K iterations of gradient descent with a fixed learning rate $\alpha \leq \frac{1}{L}$, the function satisfies the inequality:

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

where $f(x^*)$ is the optimal value of the convex function f .

Step 1:- Prove $f(x^{(K+1)}) \leq f(x^{(K)})$
The iterative update rule for gradient descent is:

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

This says that we move in the direction of negative gradient scaled by a step size α .

We use the Lipschitz continuity of the gradient to derive a bound for $f(x^{(K+1)})$. From the Lipschitz property:-

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

Substituted $x^{(k+1)} - x^{(k)}$ into inequality:

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

This simplifies to:

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

for $\alpha \leq \frac{1}{L}$, we have $\alpha - \frac{L\alpha^2}{2} \geq \frac{\alpha}{2}$.

So,

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

This shows that the function value decreases after each iteration.

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

We want to bound the difference $f(x^{(k)}) - f(x^*)$, where x^* is the optimal solution. From the convexity of f , we have:

$$f(x^{(k)}) - f(x^*) \leq \nabla f(x^{(k)})^T (x^{(k)} - x^*)$$

Now use the gradient descent update

$x^{(k+1)} = x^{(k)} - \alpha \nabla f(x^{(k)})$. The difference between $x^{(k+1)}$ and x^* is:

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

Taking the squared norm of both sides

$$\|x^{(k+1)} - x^*\|_2^2 = \|x^{(k)} - x^*\|_2^2 - 2\alpha \nabla f(x^{(k)})^\top (x^{(k)} - x^*)$$
$$\|x^{(k)} - x^*\|_2^2 + \alpha^2 \| \nabla f(x^{(k)}) \|_2^2$$

Using the convexity bound $\nabla f(x^{(k)})^\top (x^{(k)} - x^*) \geq f(x^{(k)}) - f(x^*)$, we get:

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

Rearranging the terms:-

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

Dividing by 2α , we obtain:

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

Step 3: Prove $\sum_{k=1}^K (f(x^{(k)}) - f(x^*)) \leq \frac{1}{2\alpha} (\|x^{(0)} - x^*\|_2^2 - \|x^{(K)} - x^*\|_2^2)$

Now, sum the inequality from step 2 over all iterations (from $k=1$ to K):

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

since $\|x^{(k)} - x^*\|_2^2 \geq 0$, we conclude

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

Dividing both sides by K , we get the average decrease in the function value:

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

since $f(x^{(k)}) \leq \frac{1}{K} \sum_{k=1}^K f(x^{(k)})$, this implies:

$$f(x^{(k)}) - f(x^*) \leq \frac{1}{2dK} \|x^{(0)} - x^*\|_2^2$$

we have proven that gradient descent with a step size $\alpha \leq \frac{1}{L}$ converges to the optimal solution at a rate of $O(1/K)$, with a final bound:

$$f(x^{(k)}) - f(x^*) \leq \frac{1}{2dK} \|x^{(0)} - x^*\|_2^2$$