

MIDSUM MACHINE LEARNING

Q1 Given a standard observation model in additive noise

$$\underline{d(i)} = \underline{U_i^H w} + \underline{n(i)}$$

in matrix form for all observations

$$\underline{d} = \underline{Uw} + \underline{n}$$

\underline{d} is observed data vector

\underline{U} is matrix of known column vectors

w is unknown weight vector to be estimated

n is noise

Least Square Solution

find estimate \hat{w} that minimizes the squared distance (error) between d and $U\hat{w}$

$$\min_{\hat{w}} \|d - U\hat{w}\|^2$$

The solution is given by normal equation

$$U^H U \hat{w} = U^H d$$

Solving for \hat{w}

$$\boxed{\hat{w} = (U^H U)^{-1} \cdot U^H d}$$

\hat{w} is the weight vector that best fits the observed data in least squares sense.

U^H denotes conjugate transpose (Hermitian transpose)

For real-valued matrices \mathbf{U}^H is just transpose \mathbf{U}^T
 For Complex-valued matrices \mathbf{U}^H is transpose of \mathbf{U} with
 each element replaced by its Complex Conjugate.

(a) a) Show that the matrix $\phi(\phi^T\phi)^{-1}\phi^T$ takes any vector v and projects it onto the space spanned by the columns of ϕ . Use this result to show that the least square solution $w_m = \phi(\phi^T\phi)^{-1}\phi^T t$ corresponds to an orthogonal projection of the vector t onto the manifold S (space span by column of ϕ)

Soln

If we denote $v^* = (\phi^T\phi)^{-1}\phi^T v$, then

$$\phi(\phi^T\phi)^{-1}\phi^T v = \phi v^* \quad \text{--- (1)}$$

By definition ϕv^* is in the column space of ϕ .

In other words we have prove

$\phi(\phi^T\phi)^{-1}\phi^T$ can project a vector v into the column space of ϕ . Next we are required to prove the residue of the projection shown in below fig ($y - t$ is orthogonal to column space of ϕ)

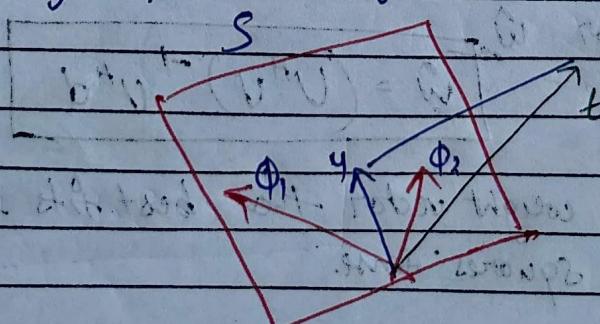


Figure - The least square regression function is obtained by

finding orthogonal projection of data vector t onto the subspace spanned by the basis function $\phi_j(x)$ in which each basis function is viewed as a vector. ϕ_j of length N with elements $\phi_j(x_n)$

Since we have

$$\begin{aligned}(y-t)^\top \phi &= (\phi w_{\text{MC}} - t)^\top \\&= (\phi(\phi^\top \phi)^{-1} \phi^\top t - t)^\top \phi \\&= t^\top (\phi(\phi^\top \phi)^{-1} \phi^\top - I)^\top \phi\end{aligned}$$

$$\begin{aligned}(1) \quad (y-t)^\top (\phi(\phi^\top \phi)^{-1} \phi^\top - I)^\top \phi &= 0 \\&= t^\top (\phi - \phi)\end{aligned}$$

which means $(y-t)$ is the left null space of ϕ and its orthogonal to the column space of ϕ .

Qb) Given joint probability for 3 observed data points

$$x_1 = 9 \quad x_2 = 9.5 \quad x_3 = 11$$

$$p(9, 9.5, 11; \mu, \sigma^2) = \prod_{i=1}^3 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

To maximize probability w.r.t μ and σ^2

Find MLE for normal distribution.

Step 1: \Rightarrow log likelihood

$$\ell(\mu, \sigma) = \text{Log}_e P = \sum_{i=1}^3 \left[-\log \left(3\sqrt{2\pi} + \frac{(x_i - \mu)^2}{2\sigma^2} \right) \right]$$

$$= -3 \log (3\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^3 (x_i - \mu)^2$$

Step 2:Maximize w.r.t. μ

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{2\sigma^2} \sum_{i=1}^3 (x_i - \mu) (-1)$$

$$= \frac{1}{2\sigma^2} \sum_{i=1}^3 (x_i - \mu)$$

Set to zero

$$\sum_{i=1}^3 (x_i - \mu) = 0 \Rightarrow \mu = \frac{1}{3} (x_1 + x_2 + x_3)$$

$$\mu = \frac{1}{3} (9 + 9.5 + 11) = \frac{29.5}{3} = 9.833$$

$$\boxed{\mu = 9.833}$$

Step 3: Maximize w.r.t \bar{z}

$$\frac{\partial l}{\partial \bar{z}} = -\frac{3}{8} + \frac{1}{8^3} \sum_{i=1}^3 (x_i - \bar{z})^2 = 0$$

multiply by 8^3

$$\Rightarrow -3\bar{z}^2 + \sum_{i=1}^3 (x_i - \bar{z})^2 = 0$$

$$\bar{z}^2 = \frac{1}{3} \sum_{i=1}^3 (x_i - \bar{z})^2$$

put value of x_1, x_2, x_3 and \bar{z}

$$\begin{aligned} \bar{z}^2 &= \frac{1}{3} [(9 - (9 - 9.833))^2 + (9.5 - 9.833)^2 \\ &\quad + (11 - 9.833)^2] \end{aligned}$$

$$= (-0.833)^2 + (-0.333)^2 + (1.167)^2$$

$$\bar{z}^2 = \frac{2.166}{3} = 0.722$$

$$\bar{z} = \sqrt{0.722} = 0.85$$

(Q) Determine the condition under when the least square solution is equivalent to maximum likelihood estimates.

Soln Suppose we have linear model

$$y = Xw + \epsilon$$

y is vector of observed output

X is design matrix

w is the vector of parameters to estimate

ϵ is error (noise) vector

Least squares

The least square method find w that minimizes the sum of squared errors.

$$\min_w \|y - Xw\|^2$$

Maximum likelihood

Under Gaussian noise assumption the likelihood of observing y given X and w is

$$p(y|X, w) = N(y|Xw, \sigma^2 I)$$

log likelihood is

$$\log p(y|X, w) = -\frac{1}{2\sigma^2} \|y - Xw\|^2 + \text{const}$$

maximizing log likelihood w.r.t w is equivalent to minimizing $\|y - xw\|^2$

- \Rightarrow if the noise is i.i.d Gaussian with zero mean and constant variance, the least square estimator and the M.L.E for w are same ~~✓~~
- \Rightarrow if the noise is not Gaussian least square is not generally the same as MLE ~~✓~~

- Q2) code block is attached
 Q3) code block is attached

Question 2

Computes the gradient of the quadratic function of x given the starting points $x=[1,2,3]$ and then uses the result of the gradient to feed the next iterations, with new points. Prints out the result of the function at each iteration till 3rd run. Use Python script to print the results.

```
import numpy as np

# Define the quadratic function: f(x) = x^T x
def f(x):
    return np.dot(x, x)

# Its gradient: grad_f(x) = 2x
def grad_f(x):
    return 2 * x

# Starting point
x = np.array([1.0, 2.0, 3.0])

# Learning rate
alpha = 0.1

print("Iteration results:")
for i in range(3):
    fx = f(x)
    print(f"Iteration {i+1}: f(x) = {fx:.4f}, x = {x}")
    x = x - alpha * grad_f(x)

Iteration results:
Iteration 1: f(x) = 14.0000, x = [1. 2. 3.]
Iteration 2: f(x) = 8.9600, x = [0.8 1.6 2.4]
Iteration 3: f(x) = 5.7344, x = [0.64 1.28 1.92]
```

Question 3

With given input and output relation $x = [-1, -0.8, -0.6, -0.4, -0.2, 0, 0.2, 0.4, 0.6, 0.8, 1]$ $t = [-4.9, -3.5, -2.8, 0.8, 0.3, -1.6, -1.3, 0.5, 2.1, 2.9, 5.6]$ Please fit a curve with $M=4$ Gaussian basis functions having unity variance.

```
import numpy as np
import matplotlib.pyplot as plt

# Given data
x = np.array([-1, -0.8, -0.6, -0.4, -0.2, 0, 0.2, 0.4, 0.6, 0.8, 1])
t = np.array([-4.9, -3.5, -2.8, 0.8, 0.3, -1.6, -1.3, 0.5, 2.1, 2.9, 5.6])

# Number of Gaussian basis functions
```

```

M = 4

# Centers of the Gaussians (evenly spaced in x range)
centers = np.linspace(np.min(x), np.max(x), M)
variance = 1.0 # Unity variance

# Design matrix using Gaussian basis functions
def design_matrix(x, centers, variance):
    Phi = np.zeros((len(x), len(centers)))
    for i, c in enumerate(centers):
        Phi[:, i] = np.exp(-0.5 * ((x - c) ** 2) / variance)
    return Phi

Phi = design_matrix(x, centers, variance)

# Fit weights using least squares
w = np.linalg.pinv(Phi) @ t

# Predict function
def predict(x_new):
    Phi_new = design_matrix(x_new, centers, variance)
    return Phi_new @ w

# Plotting
x_plot = np.linspace(-1.2, 1.2, 200)
t_pred = predict(x_plot)

plt.figure(figsize=(8, 5))
plt.scatter(x, t, color='red', label='Data')
plt.plot(x_plot, t_pred, label='Gaussian Basis Fit (M=4)')
plt.xlabel('x')
plt.ylabel('t')
plt.title('Curve Fitting with 4 Gaussian Basis Functions (Unity Variance)')
plt.legend()
plt.grid(True)
plt.show()

```

Curve Fitting with 4 Gaussian Basis Functions (Unity Variance)

