

Handwritten Mathematical Analysis

① Linear Regression:

Dataset - 435742 rows, 13 columns

X (input matrix) - (435742, 4)

pollutant concentrations
of SO₂, NO₂, RSPM &
SPM

y (output label)
vector - (435742, 1)

AGI

Linear Regression equation -

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

where $X = [x_1, x_2, \dots, x_N]$, $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}$

$5 \times N$
 \uparrow
 $1+4$
 \uparrow
dummy
variable

$$y = X^T \theta \quad \text{--- (vectorized notation)}$$

\uparrow
(435742 x 5) (5 x 1)

$$P(y | X, \theta) = \prod_{n=1}^N P(y_n | x_n, \theta)$$

(N = total
no of
samples)

$$= \prod_{n=1}^N \mathcal{N}(y_n | x_n^T \theta, \sigma^2)$$

$x_n \in \mathbb{R}^4$

Taking negative log-likelihood,

$$-\log P(y | X, \theta) = -\log \prod_{n=1}^N \mathcal{N}(y_n | x_n^T \theta, \sigma^2)$$

$$= -\sum_{n=1}^N \log \mathcal{N}(y_n | x_n^T \theta, \sigma^2)$$

in LR model, $y = x^T \theta + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

Hence, $\log p(y_n | x_n, \theta) = -\frac{1}{2\sigma^2} (y_n - x_n^T \theta)^2 + \text{constant}$

Hence, the negative log likelihood turns out to be, :

$$\ell(\theta) = -\log(p(y|x, \theta)) = \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - x_n^T \theta)^2$$

$$L(\theta) = \frac{1}{2\sigma^2} (y - X\theta)^T (y - X\theta) = \frac{1}{2\sigma^2} \|y - X\theta\|^2$$

$$X \rightarrow [x_1, \dots, x_N]^T \in \mathbb{R}^{N \times D}$$

where $N = 435742$, $D = 4$ (features)

$$y \rightarrow [y_1, \dots, y_N]^T \in \mathbb{R}^{N \times 1}$$

Taking derivative of log likelihood & equating it to 0 to obtain OML (optimal parameter): \Rightarrow

$$\frac{dL}{d\theta} = \frac{d}{d\theta} \left(\frac{1}{2\sigma^2} (y - X\theta)^T (y - X\theta) \right)$$

$$= \frac{1}{2\sigma^2} \frac{d}{d\theta} (y^T y - 2y^T X\theta + \theta^T X^T X \theta)$$

$$= \frac{1}{\sigma^2} (-y^T X + \theta^T X^T X) \in \mathbb{R}^{1 \times D}$$

$$\frac{dL}{d\theta} = 0^T \Leftrightarrow$$

$$\theta_{ML}^T X^T X = y^T X$$

$$\Leftrightarrow$$

$$\theta_{ML}^T = y^T X (X^T X)^{-1}$$

MLE Equation \leadsto

$$\boxed{\theta_{ML} = (X^T X)^{-1} (X^T y)}$$

for non linear transformation (Polynomial Regression)

$$y = \phi(x)^T \theta$$

where $\phi(x) = [x^0 \ x^1 \ x^2 \ \dots \ x^K]^T$

on applying the Φ matrix,

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{K-1}(x_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \dots & \dots & \phi_{K-1}(x_N) \end{bmatrix} \in \mathbb{R}^{N \times K}$$

K = degree of polynomial

$x_n = D$ dimensional vector

$$\phi(x_n) : \mathbb{R}^D \rightarrow \mathbb{R} \text{ and } \Phi_{in} = \phi_i(x_n)$$

$$\cancel{\Phi(x_n) : \mathbb{R}^D \rightarrow \mathbb{R}^K} \quad \in (0, K) \quad \uparrow$$

~~Ques~~
In our case, after applying polynomial regression, we get,

$$y = \theta_0 + \theta_1(x_1 + x_2 + x_3 + x_4) + \theta_2(x_1^2 + x_2^2 + x_3^2 + x_4^2) + \dots$$

for each point (x_i, y)

Hence, MLE equation turns out to be

$$\boxed{\theta_{ML} = (\Phi^T \Phi)^{-1} (\Phi^T y)}$$

$$\Phi \in \mathbb{R}^{N \times K}$$

\uparrow
degree of polynomial

$$y \in \mathbb{R}$$

\uparrow
total no of points
 $N \times 1$

Maximum A Posteriori Estimation

to overcome the problem of overfitting in MLE estimation, we need prior distribution about what parameter values are plausible.

According to Bayes Theorem,

$$p(\theta | x, y) = \frac{p(y | x, \theta) p(\theta)}{p(y | x)}$$

to find the MAP estimate,

$$\log(p(\theta | x, y)) = \log p(y | x, \theta) + \log p(\theta) + \text{constant}$$

the gradient of the -ve log posterior with respect to θ is

$$\frac{-d(\log p(\theta | x, y))}{d\theta} = -\frac{d}{d\theta}(\log p(y | x, \theta)) - \frac{d}{d\theta} \log p(\theta)$$

after solving, we get,

$$-\log p(\theta | x, y) = \frac{1}{2\sigma^2} (y - \Phi\theta)^T (y - \Phi\theta) + \frac{1}{2b^2} \theta^T \theta + \text{constant}$$

considering gaussian prior $p(\theta) = \mathcal{N}(0, b^2 \mathbf{I})$ on the parameters θ .

$$-\frac{d}{d\theta}(\log p(\theta | x, y)) = \frac{1}{\sigma^2} (\theta^T \Phi^T \Phi - y^T \Phi) + \frac{1}{b^2} \theta^T$$

Support vector machines

on simplification, we get

$$\frac{1}{\sigma^2} (\Theta^T \Phi^T \Phi - y^T \Phi) + \frac{1}{b^2} \Theta^T = 0^T$$

$$\Theta^T \left(\frac{1}{\sigma^2} \Phi^T \Phi + \frac{1}{b^2} \mathbf{I} \right) - \frac{1}{b^2} y^T \Phi = 0^T$$

$$\Theta^T = y^T \Phi \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1}$$

$$\# \quad \Theta_{\text{MAP}} = \left(\Phi^T \Phi + \frac{\sigma^2}{b^2} \mathbf{I} \right)^{-1} \Phi^T y \quad \#$$

↑
MAP estimation
optimal parameter

$\Phi \in \mathbb{R}^{N \times K}$, $y \in \mathbb{R}^{N \times 1}$
 $N \rightarrow$ total no of datapoints
 $K \rightarrow$ degree of polynomial.

The extra term ensures that
 it is symmetric & strictly
 positive definite &
~~hence~~ (also works as an effect of
Regularizer)

Principal Component Analysis

$X^{N \times D}$ = original dataset (feature matrix)

$N = 434752$ datapoints, $D = 4$ features

aim \rightarrow to calculate $\tilde{X}^{N \times D}$ that has similar dimensions as that of original dataset (feature matrix) but which have a significantly lower intrinsic dimensionality.

\rightarrow First we normalize the dataset,

$$X_{\text{norm}} = \frac{X - \mu}{\sigma}, \text{ where}$$

$$\mu = \frac{1}{N} \sum_{i=1}^N X_i$$

σ = standard dev of points in X

dataset \rightarrow

$$X_{\text{norm}}^T = \{x_1, \dots, x_N\};$$

$x_n \in \mathbb{R}^D$ & with mean 0.

\rightarrow Covariance matrix \rightarrow

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T = \frac{1}{N} X_{\text{norm}}^T X_{\text{norm}}$$

dimension of covariance matrix $\rightarrow (4, 4)$

Now for principal components, we need to find eigenvalues & eigenvectors of the covariance matrix S .

$$S \vec{v} = \lambda \vec{v} \quad \text{--- (1)}$$

where $\vec{v} \rightarrow$ eigenvector &

$\lambda \rightarrow$ corresponding eigenvalue of S .

$$(S \vec{v} - \lambda \vec{v}) = 0$$

$$\vec{v} (S - \lambda I) = 0 \quad \text{--- (2)}$$

$(S - \lambda I)$ has to be non invertible,
hence, $\det(S - \lambda I) = 0$. — (2)

since, S is a 4×4 matrix,

on solving eqn (3), we get 4 different values of λ .

Using $S\vec{v} = \lambda\vec{v}$,

for each λ , we get a column vector \vec{v} .

~~Higher the value~~

hence, we get 4 eigenvalues of S and 4 corresponding eigen vectors of S .

Higher the value of λ , higher will be the variance contributed by that particular eigenvector,

hence, we sort the eigenvectors corresponding to decreasing order of eigenvalues

$$\text{eigenvec} = [\vec{v}_1 \ \vec{v}_2 \ \vec{v}_3 \ \vec{v}_4]$$

where $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4$.
(decreasing order)

$B =$ matrix that takes ^{for} M eigen vectors
where $M \rightarrow$ lower dimensional
space = $\frac{2}{\text{(in our case)}}$

$$B = [\vec{v}_1 \ \vec{v}_2]$$

we then calculate projection matrix P ,

$$P = BB^T$$

where $\dim(B) = 4 \times M$

$\dim(P) = 4 \times 4$.

$$X_{\text{reconstruct}} = \underbrace{BB^T}_{\substack{\text{projection} \\ \text{matrix } P \\ (4 \times 4)}} \cdot \underbrace{X_{\text{norm}}^T}_{\substack{(4 \times N) \\ \uparrow \\ \text{total data points}}}$$

$$\underbrace{X_{\text{reconstruct}}}_{(N \times 4)} = \underbrace{X_{\text{reconstruct}}^T}_{(4 \times N)}$$

↑
This matrix is reconstructed feature matrix whose dimension in our case = $(435742, 4)$ (same as X_{norm}) but with lower intrinsic dimensionality.

Further, we calculate y (A&I values) from this reconstructed feature matrix & compare the new A&I values with the actual A&I values.

◦ Support vector Machines

(handwritten Mathematical Analysis)

→ The equation of hyperplane in 'M' dimension is —

$$\begin{aligned} y &= w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots \\ &= w_0 + \sum_{i=1}^m w_i x_i = w_0 + w^T X \\ &= \underset{\substack{\uparrow \\ \text{biased} \\ \text{term}}}{b} + w^T X \end{aligned}$$

→ Hard margin SVM :

we conclude, for any point X_i ,

if $Y_i (w^T X_i + b) \geq 1$ then $X_i \rightarrow$ correctly classified

else $X_i \rightarrow$ incorrectly classified.

If points are linearly separable then only our hyperplane is able to distinguish between them and if any outlier is introduced then it is not able to separate them. So these type of SVM is called hard margin SVM.

→ Soft margin SVM :

If the points are not linearly separable then new slack variable is introduced (ξ) ~~which is called~~.

new equation:

$$Y_i (w^T X_i + b) \geq 1 - \xi_i$$

if $\xi_i = 0$, points can be correctly classified

else if $\xi_i > 0$, points are incorrectly classified.

mean ξ_i = error term & the average error is,

$$\text{avg error} = \frac{1}{n} \sum_{i=1}^n \xi_i$$

hence, our mathematical objective is,

$$\text{minimize}_{w, b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \xi_i$$

such that $y_i(w^T x_i + b) \geq 1 - \xi_i$, for all $i=1, 2, \dots, n$

→ Dual form SVM:

with SVM, we can separate each data point by projecting it onto a higher dimension.

alternate method is dual form SVM which uses Lagrange's multiplier to solve the constraints optimization problem.

$$\text{maximise}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i^T x_j)$$

subject to $\alpha_i \geq 0$ for all $i=1, 2, \dots, n$ &
 $\sum_{i=1}^n \alpha_i y_i = 0$.

If $\alpha_i > 0$ then $x_i \rightarrow$ support vector &
 when $\alpha_i = 0$ then $x_i \rightarrow$ not a support vector.

→ Kernel trick

new equation with kernel function

$$\text{maximise}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{K(x_i, x_j)}_{\text{kernel function}}$$

linear kernel
 $K(x_1, x_2) = x_1^T x_2$

poly kernel
 $K(x_1, x_2) = (a + x_1^T x_2)^b$

gaussian kernel
 $K(x_1, x_2) = \exp(-\gamma \|x_1 - x_2\|^2)$

Gaussian mixture model

↳ clustering algorithm

$$p(x|\mu, \Sigma, \pi) = \sum_{j=1}^n \pi_j \cdot \mathcal{N}(x|\mu_j, \Sigma_j)$$

subject to $\sum_{j=1}^n \pi_j = 1$. ↑
weighted sum of
different gaussian
distributions

pdf of gaussian distribution,
$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x-\mu)^2}{2\sigma^2} \right)$$

for multivariate gaussian distribution,

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^K |\Sigma|}} \exp \left(-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2} \right)$$

Expectation-maximization algorithm is applied as the parameters cannot be estimated in the closed form.

⑦ Step 2:
we know, $p(x|\theta) = \prod_{n=1}^N p(x_n|\theta)$

= $\sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$ and

$$p(x|\theta) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$$

$$\log(p(x|\theta)) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k) \right)$$

log likelihood \mathcal{L}

To find optimal parameters μ_k, Σ_k, π_k ,

$$\frac{dL}{d\mu_k} = 0^T, \quad \frac{dL}{d\Sigma_k} = 0^{\oplus}, \quad \frac{dL}{d\pi_k} = 0 \quad \text{--- (1)}$$

Responsibility \rightarrow

$$r_{nk} = \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}$$

Responsibility of

k th mixture component for the
 n th data point

on applying eqn (1),
we get

$$\mu_k^{\text{new}} = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}} \quad \text{--- (A)}$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N r_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

$$\text{where } N_k = \sum_{n=1}^N r_{nk}$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}, \quad k=1, 2, \dots, K.$$

$N \rightarrow$ no of data points