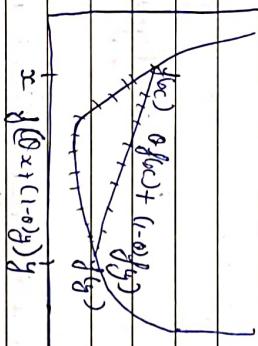


PRM

Date 25-01-2022
Page No.

- Convex sets - $C \subseteq \mathbb{R}^n$ is convex, if $\theta x + (1-\theta)y \in C$ for any $x, y \in C$ and $0 \leq \theta \leq 1$.
i.e. a set is convex if the line connecting any 2 points in the set is entirely inside the set

- Convex functions - $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if $\text{dom}(f)$ (the domain of f) is a convex set and if $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$ for any $x, y \in \text{dom}(f)$ and $0 \leq \theta \leq 1$.



- Finding maxima of convex function is easy
- If $f_1(x)$ is convex and f_2 is convex -
then - (i) $f_1(x) + f_2(x)$ is also convex

(ii) $\lambda f_1(x)$ is also convex

(iii) $\max(f_1(x) + f_2(x))$ is also convex

- If a function is convex function and it is also smooth, then it is differentiable throughout its domain.
- Convexity means when it will decrease it will decrease only and when it will increase it will increase only.

- there is a point whose rate of change in every direction is zero (i.e. gradient), that will give the maximum value

- If function \rightarrow convex & smooth, then necessary and sufficient condition for optimal solution x_0 is $f(x) = 0$

e.g. if function is from $R^2 \rightarrow R$ then the gradient will also be in $R^2 \rightarrow [g]$ \Rightarrow gradient vector direction of gradient is always direction of ascent (moving up) at $x = x_0$.

$$y_i = E(Y/x_i) + \epsilon_i \quad \rightarrow \quad E(\epsilon_i) = 0$$

$\downarrow f(x)$ \rightarrow Mean Regression Model

$$y_i = f(x_i) + \epsilon_i$$

\downarrow if set tells, suppose following 10% of values of y given x , then this type of model is Quantile Regression Model [it gives more information than mean regression model]

- for 2-dimension -
- height
- length will be of pair.
- $[16.0, 6.5], [17.0, 8.0]$
- $f(x,y) =$
- weight
- gradient comes when we are dealing with 2 or more variables otherwise for 1 variable we will subject to the concept of differentiation.

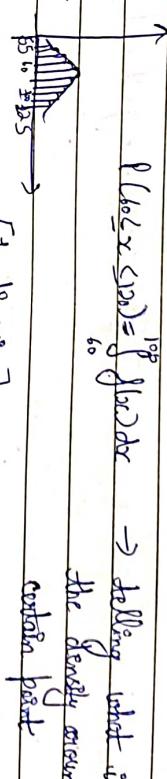
e.g.
 $f: R^2 \rightarrow R$, $f(x,y) = x^2 + y^2$



$$\vec{f} \cdot \vec{v}(a,b)$$

arg Max $(\nabla f(a,b))^\top v$ [where $v = (\nabla f(a,b))^\top$] \downarrow rate of change in v direction

$$\nabla f(x)^\top \nabla f(x) = \|\nabla f(x)\|^2 \quad \text{where } \|v\| \text{ is unit vector}$$



While plotting histogram if we reduce the size of bin then we can join the points and make a curve and then find area of sum of the region we want to evaluate.

- when bin is large, much data is lost because of approximation.

Note. Note of change we can take is the norm of the gradient of $f(x)$

$u = \|\nabla f(x)\| \rightarrow$ this is direction of maximum ascent

$u = -\|\nabla f(x)\| \rightarrow$ direction of minimum descent

* Gradient descent algorithm - [An iterative algorithm]

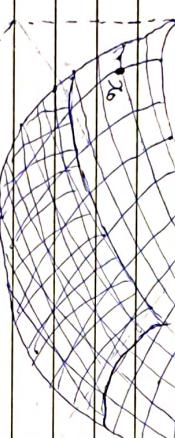
Initialize $x^0 = x_{\text{start}} \in \mathbb{R}^n$

iterate $x^{(k+1)} := x^{(k)} - \gamma_k \nabla f(x^{(k)})$

until $\|\nabla f(x^{(k)})\| \leq \epsilon$

The negative gradient direction is the direction of steepest descent.

$$f(x_1, x_2)$$



After k th step,

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

where, it is gamma(γ) (step length)

as we increase the value of k , $f(x^{(k)})$ converges to $f(x^*)$

here, $x^k \rightarrow$ at k th step

$x^* \rightarrow$ optimal solution (desired)

$x^0 \rightarrow$ initial position

* loss function - It is perception that what do we mean by error
there are many types of loss functions ($f(x_e) - y_e$)

$$\text{e.g. } (i) (f(x_e) - y_e)^2$$

$$(ii) |f(x_e) - y_e|$$

and gradient is Lipschitz continuous

- If our function is convex, then with this algo we are guaranteed to converge on global local minima

$$f: \mathbb{R}^n \rightarrow \mathbb{R}, \quad \|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|_2$$

L -norm

$$L > 0$$

* F can be anything, lets start with F as a linear function
where $x \in \mathbb{R}^n, y \in \mathbb{R}$

$$= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$$F = \begin{cases} w^T x + b : x, w \in \mathbb{R}^n, b \in \mathbb{R} \end{cases}$$

$$\frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Date _____
Page No. _____

Date _____
Page No. _____

$$L = (Y - f(x))^2 \Rightarrow u^2 \rightarrow \text{Least square loss function, it is smooth and convex, also differentiable at every given point}$$

- Maximization of Least square loss function -

$$\text{Min } J = \sum_{i=1}^n (y_i - (\omega^T x_i + b))^2$$

$$\begin{matrix} \text{weight} \\ \downarrow \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \\ \text{height} \end{matrix} \quad X = (x_1, y_1) \quad \text{with width}$$

$$\rightarrow T \cdot \int_{x_1}^{x_n} f(x, y) dx, \quad (x_1, y_1), \dots, (x_n, y_n)$$

height

$$X = (x_1, y_1)$$

$$X = (x_1, x_2)$$

$$\text{Unbiased Estimate of mean} = \bar{X} = (\bar{x}, \bar{y}) = \frac{1}{n} \sum_{i=1}^n (x_i, y_i) \quad [\text{In 2-Dimensional}]$$

$$\text{Unbiased Estimate of variance} = \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2 \quad [\text{In 1-Dimensional}]$$

If we have large no. of data points, then this will converge to variance

So in 2-D, we have co-variance

$$\Sigma = \begin{bmatrix} \sigma_{xx}^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_{yy}^2 \end{bmatrix} \rightarrow \text{Covariance along x-axis} \quad \begin{bmatrix} \text{constant & positive} \\ \text{semi-definite matrix} \end{bmatrix}$$

$$\text{Cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Variance along a particular direction can be obtained by - $\Sigma \vec{e}_i^T \vec{e}_i$

* On 2-D, if we want to know that our data follows Normal Distribution or not, then we need to project it to a axis. (say $x_i^T \omega$) and then on that axis we will check if 68.3% of data values are falling under $(\mu - \sigma)$ to $(\mu + \sigma)$ and 95.4% are falling under $(\mu - 2\sigma)$ to $(\mu + 2\sigma)$ and 99.7% are falling under $(\mu - 3\sigma)$ to $(\mu + 3\sigma)$. Then it is normal distribution.

$$Y = (x, y) \sim N(\mu, \Sigma) \quad [\text{Density function of } Y \text{ is independent}]$$

Then we can say x & y are jointly normal

$$f(x) = f(x_1, x_2) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp^{-\frac{1}{2} \frac{(x_1 - \mu_1)^2}{\sigma_1^2}} * \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp^{-\frac{1}{2} \frac{(x_2 - \mu_2)^2}{\sigma_2^2}}$$

$$\downarrow$$

$$(\sqrt{2\pi})^2 |\Sigma|^{-\frac{1}{2}}$$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp^{-\frac{1}{2} \frac{(x_1 - \mu_1)^2}{\sigma_1^2}}, \quad X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

density function of Gaussian \downarrow
normal distribution in n-dimensional space are independent, covariance of them is Σ

• Consider first one the points where value of $f(x_1, x_2)$ is same, i.e. $f(x_1, x_2) = c$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}, A = \begin{bmatrix} x_{11}, x_{12}, \dots, x_{1m}, 1 \end{bmatrix}, w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix}$$

$$\text{Min}_{w,b} (Y - Aw)^T (Y - Aw)$$

$$\text{Min}_{w,b} h(w)$$

$$h(w) = 0$$

* Least squares regression problem

$$w = (A^T A)^{-1} A^T Y$$

$$b$$

* Least square quadratic regression model - Refer previous L assignment

$$h_{\text{quad}}(w) = \text{Min}_w (C(Y - Aw)^T (Y - Aw))$$

$$1. \quad \nabla_w (2(Y - Aw)^T (Y - Aw))$$

$$2. \quad A^T (Y - Aw) = 0$$

$$A^T Y - A^T Aw = 0$$

$$w = (A^T A)^{-1} A^T Y$$

where, $\phi =$

$$\begin{pmatrix} \phi_1(x_1) & \dots & \phi_{n+1}(x_1) \\ \vdots & \ddots & \vdots \\ \phi_1(x_m) & \dots & \phi_{n+1}(x_m) \end{pmatrix}$$

$$\phi(x) = w^T \phi(x)$$

$$\phi(x) = w^T \phi(x)$$

* Polynomial basis function = [Non-linear estimate]

$$\left[1, x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10}, x_{11}, x_{12}, x_{13}, x_{14}, x_{15}, x_{16} \right]$$

$$\text{Base function } \Phi_m(x)_{m=1}^6$$

$$\Phi_1(x) = 1, \Phi_2(x) = x_1, \Phi_3(x) = x_2, \dots, \Phi_6(x) = x_1 x_2 \dots x_5$$

Independent of choice of basis functions, the regression parameters were calculated using the all-zero equations for linear regression.

$$f(x) = \sum_{i=1}^{n+1} w_i \phi_i(x) + b \Rightarrow w_0 + \sum_{j=1}^{n+1} w_j \phi_j(x)$$

(M) is the total no. of parameters in the model

$$\text{Min}_{w, w_0, \dots, w_m} \sum_{i=1}^m (y_i - (w_0 + \sum_{j=1}^{n+1} w_j \phi_j(x_i)))^2$$

$$\text{Min}_w \sum_{i=1}^m (y_i - w^T \phi(x_i))^2$$

$$\text{Min}_w (Y - w^T \phi)^T (Y - w^T \phi)$$

$$w = (\phi^T \phi)^{-1} \phi^T Y$$

- * Maximum likelihood estimate - It estimates the parameter θ of a distribution, suppose we have normal distribution then our parameter (θ) to be estimated will be (μ, σ)
- sample - $\{x_1, x_2, \dots, x_n\}$

$$P(\theta | T) = P(T|\theta) \frac{P(\theta)}{P(T)}$$

$\text{Max } P(\theta | T) = \text{Max } P(T|\theta) \quad [\text{because we don't have info about } P(\theta) \text{ & } P(T)]$

$$= \prod_{i=1}^n P(x_i | \theta)$$

let us suppose θ is coming from normal distribution

$$\prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2}$$

$$\text{Max Log } P(\theta | T) = \text{Max Log} \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2}$$

Now, we can take derivative with respect to μ & σ & minimize it

$$\mu = \frac{\sum x_i}{n}, \quad \sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n-1}}$$

$$\text{Max Log} \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2}$$

$$\text{Max} -\frac{1}{2} \log \sigma - \frac{1}{2} \log (2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Max} - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2$$

If sample follows normal distribution then maximum likelihood estimation will be equal to the minimization of least square loss function.

- Disadvantage of least square loss function is that it is sensitive to presence of outliers

- Overfitting occurs mostly when our training set is small.

Introducing Regularization -

$$\text{Min } \frac{1}{2} \|\omega\|_2^2 + \sum_{i=1}^n (y_i - (\omega^\top \phi(x_i) + b))^2$$

$$\frac{\partial L}{\partial \omega} = 0 \quad \frac{\partial L}{\partial b} = 0$$

$$\begin{bmatrix} \omega \\ b \end{bmatrix} = (\phi^\top \phi + \lambda I)^{-1} \phi^\top y$$

- Real feature size - 9 (17-62-2000) — [PCA, variance, etc.]
- Variances of complex data sets decreases if working with large data sets

* Kernel

$$\begin{array}{c} y \\ \vdots \\ \phi(x_1) \\ \vdots \\ \phi(x_n) \end{array} \xrightarrow{\omega^\top \phi(x) + b} \begin{array}{c} y \\ \vdots \\ \omega^\top \phi(x_1) + b \\ \vdots \\ \omega^\top \phi(x_n) + b \end{array}$$

mapping to higher dimensional space

$$x_1 \\ x_2 \\ \vdots \\ x_n \\ \phi_1(x) \\ \phi_2(x)$$

$$K(x_1, x_2) = \phi(x_1)^\top \phi(x_2)$$

positive semi-definite

$$\left[\omega = \sum_{i=1}^n w_i \phi(x_i) \right] \quad f(x) = \omega^\top \phi(x) + b$$

$$= \left(\sum_{i=1}^n w_i \phi(x_i)^\top \phi(x) \right)^\top \phi(x) + b$$

$$= \sum_{i=1}^n w_i K(x_i, x) w_i + b$$

• Least square regularized kernel regression -

$$J(\omega, b) = \sum_{i=1}^n \|y_i - (\sum_{j=1}^k K(x_j, x) w_j + b)\|^2$$

$$\frac{\partial J}{\partial \omega} = 0$$

$$\begin{bmatrix} \omega \\ b \end{bmatrix} = (\sum_{i=1}^n K(x_i, x)^\top K(x_i, x))^{-1} \sum_{i=1}^n y_i K(x_i, x)$$

• We can also call this as Least square support vector regression because it is all about working with kernel generated surfaces

$$\|\omega\|^2 = \omega^\top \omega$$

$$\nabla_{\omega} J(\omega) = \left[\frac{\partial J}{\partial \omega_1}, \frac{\partial J}{\partial \omega_2}, \dots, \frac{\partial J}{\partial \omega_n} \right] = 2 \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = 2w$$

$$\nabla_{\omega} J(\omega) = \frac{1}{2} \times 2 \omega + \frac{1}{2} \sum_{i=1}^n (y_i - (\omega^\top \phi(x_i) + b))^2$$

$$= \Delta \omega + -\sum_{i=1}^n (y_i - (\omega^\top \phi(x_i) + b)) \phi(x_i)$$

$$\nabla_b J(\omega) = 0 - \sum_{i=1}^n (y_i - (\omega^\top \phi(x_i) + b))$$

Gradient-descent \rightarrow do while

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

$$\text{until } \|\nabla f(x^k)\|_2 \leq \epsilon$$

$$\begin{bmatrix} \omega^{k+1} \\ b^{k+1} \end{bmatrix} = \begin{bmatrix} \omega^k \\ b^k \end{bmatrix} - \eta (\nabla J(\omega^k, b^k))$$

* Stochastic gradient — we take K-random points instead of whole data set, that's why it is also known as K-min batch stochastic gradient.

* bringing back robustness in least square estimate by introducing weights, dense points will have more weightage than the sparse ones

$$\text{Min}_{w,b} \frac{1}{2} \sum_{i=1}^n \frac{1}{\alpha_i} |y_i - (w^T \phi(x_i) + b)|^2$$

$$J(w) = \alpha w^T b$$

K-NN can also be used for assigning weights, as dense points will have lot of neighbors than the sparse ones

* Hausser, we made a new loss function — (using ℓ_1 -norm)

~~$$\text{Min}_{w,b} \frac{1}{2} \sum_{i=1}^n \| (y_i - (w^T \phi(x_i) + b)) \|_1$$~~

In ℓ_1 -norm, we don't take squares, so it is less sensitive in presence of outliers

* Disadvantage of ℓ_1 -norm — function is not smooth, we can not differentiate at every point

$$\rightarrow \text{Min}_{w,b} \frac{1}{2} \sum_{i=1}^n (y_i - (w^T \phi(x_i) + b))^2$$

↓
dense regression

In this we are considering ℓ_1 -norm for regularization and ℓ_2 -norm for loss function to overcome the problem of singularity of solutions

$$\text{Min}_{w,b} \frac{1}{2} \sum_{i=1}^n (y_i - (w^T \phi(x_i) + b))^2 + \lambda (w^T \phi(x_i) + b)^2$$

* $w = w^T \phi + b$
if we have sparse vector (w), then we are dropping some features
 $(w_1, 0, w_2, 0, 0, \dots, 0)$

thus, a simpler model because of usage of less features, so linear regression model will perform better when we will have large data points. [Simpson's paradox generalizes well] as variance of unwanted features does not effect our estimate

$$\text{Min}_{w,b} \frac{1}{2} \| w \|_2^2 + \frac{\lambda}{2} \sum_{i=1}^n |y_i - (w^T \phi(x_i) + b)|$$

↓
sparse

* most common way to handle in linear regression by taking gradient and equating it to zero

* Solution of ℓ_1 -norm Regression model — (Optimization problem)

$$\text{Min}_{w,b} \frac{1}{2} \| w \|_2^2 + \frac{\lambda}{2} \sum_{i=1}^n |y_i - (w^T \phi(x_i) + b)|$$

We will solve optimization problem by sub-gradient

we have 2 terms, so we cannot use normal gradient

$$\nabla_w f(x) = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \end{cases}$$

$$\nabla_w f(x_j) = (w^T \phi(x_j) + b) = 1(\phi(x_j))$$

$$= \phi(x_j), \text{ if } (y_j - (w^T \phi(x_j) + b)) > 0$$

$$= \phi(x_j), \text{ if } (y_j - (w^T \phi(x_j) + b)) < 0$$

for kernel regression -

$$\left[\gamma_i - \left(\sum_{r=1}^R K(x_r, x_i) w_r + b \right) \right] = - \sum_{r=1}^R K(x_r, x_i) \cdot \delta(\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b))$$

$$= \sum_{r=1}^R K(x_r, x_i) \cdot \delta(\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b))$$

Support Vector Regression model -

$$\text{Min}_{w, b} \frac{1}{2} w^T w + \frac{1}{2} \sum_{i=1}^n \delta(\gamma_i - (w^T x_i + b))$$

L1 norm loss is also known as softplus loss

* C-sensisitive loss function -

L1 norm may be robust but may not be sparse, so to overcome this we have C-insensitive loss function

$$l(w) = \begin{cases} 0 & \text{if } |w| < c \\ |w| - c & \text{otherwise} \end{cases}$$

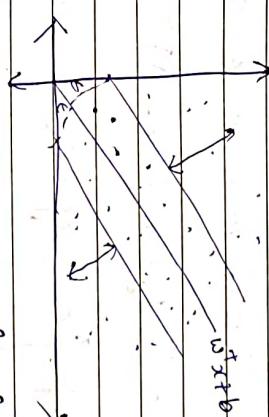
if absolute error is less than c, then it is insensitive, if it crosses c, then it gets converted to L1 norm loss

$$\text{Min}_{w, b} \frac{1}{2} w^T w + \frac{1}{2} \sum_{i=1}^n \delta(\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b))$$

$$\delta(\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b)) = \begin{cases} 0 & \text{if } (\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b)) < c \\ |(\gamma_i - (\sum_{j=1}^R K(x_j, x_i) w_j + b))| - c & \text{otherwise} \end{cases}$$

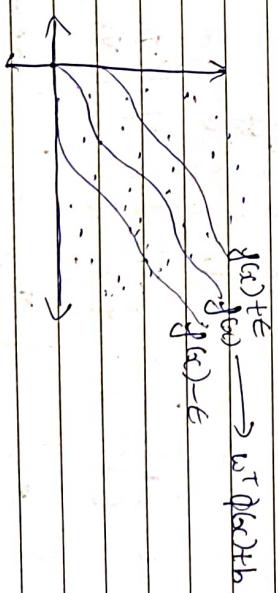
* Linear case -
 $\text{Min}_{w, b} \frac{1}{2} w^T w + \frac{1}{2} \sum_{i=1}^n \delta(\gamma_i - (w^T x_i + b))$

$$\delta(\gamma_i - (w^T x_i + b)) = \begin{cases} 0 & \text{if } |\gamma_i - (w^T x_i + b)| < c \\ |\gamma_i - (w^T x_i + b)| - c & \text{otherwise} \end{cases}$$



minimization of C-insensitive loss function is equal to drawing a C-insensitive tube around the data point and ignoring all the data points which are falling inside the C-insensitive tube. It is not sensitive to outliers and solution is going to be sparse

* Non-linear case -



for non-linear
Non
 $\frac{1}{2} w^T w + c \sum_{i=1}^n |y_i - (w^T x_i + b)|$

$$w, b \quad \text{for linear}$$

$$\frac{1}{2} w^T w + c \sum_{i=1}^n |y_i - (w^T x_i + b)|$$

$L(w) = \text{rule.}$

for non-linear
Non

$$f(w, b) = \frac{1}{2} w^T w + c \sum_{i=1}^n |y_i - (y_i - \sum_{j=1}^n K(x_j, x_i) w_j + b)|$$

$$f(w, b) = w + c \sum_{i=1}^n \begin{cases} 0 & \text{if } -c \leq y_i - (w^T x_i + b) \leq c \\ -\frac{2(y_i - (w^T x_i + b))}{c} & \text{if } y_i - (w^T x_i + b) > c \\ 1 & \text{if } y_i - (w^T x_i + b) < -c \end{cases}$$

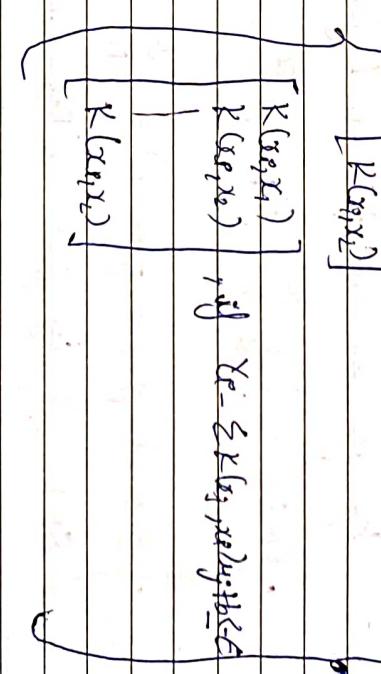
Support Vector regression model is also known as C-SVM.

$$f(w, b) = \frac{1}{2} w^T w + c \sum_{i=1}^n |y_i - (y_i - (w^T x_i + b))|$$

$$\text{The } f(w, b) = w + c \sum_{i=1}^n \begin{cases} 0 & \text{if } -c \leq y_i - (w^T x_i + b) \leq c \\ -\frac{2(y_i - (w^T x_i + b))}{c} & \text{if } y_i - (w^T x_i + b) > c \\ 1 & \text{if } y_i - (w^T x_i + b) < -c \end{cases}$$

$$f(w, b) = \frac{1}{2} w^T w + c \sum_{i=1}^n \begin{cases} 0 & \text{if } -c \leq y_i - (w^T x_i + b) \leq c \\ -\frac{2(y_i - (w^T x_i + b))}{c} & \text{if } y_i - (w^T x_i + b) > c \\ 1 & \text{if } y_i - (w^T x_i + b) < -c \end{cases}$$

* Resampling method —



$$\begin{bmatrix} w^{k+1} \\ b^{k+1} \end{bmatrix} = \begin{bmatrix} w^k \\ b^k \end{bmatrix} - \eta \begin{bmatrix} \nabla_w f(w, b) \\ \nabla_b f(w, b) \end{bmatrix}$$

until $\left\| \begin{bmatrix} \nabla_w f(w, b) \\ \nabla_b f(w, b) \end{bmatrix} \right\| \leq \text{some value}$

K-fold cross validation — will give K RMSE values and we will take mean of that along with \pm standard deviation. It is good for large dataset

$$RMSE = \pm S.D$$

A good regression model will obtain lesser mean RMSE value and lesser variance of RMSE values.

- leave-one-out validation - Here we leave one data point for testing and rest are taken for training.

e.g. go data points \rightarrow skip for training $k-1$ for testing

we will have k MSE's & k variances.

In k -fold of $k=2$, then all becomes leave-one-out
Particular we'll be less likely have leave-one-out, but will be
high after less value of k in k -fold.

$$RMSE = \sqrt{\frac{1}{k} \sum_{i=1}^k (f(x_i) - y_i)^2}$$

$$SSE \text{ (Sum of square of error)} = \sum_{i=1}^n (f(x_i) - y_i)^2$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 \quad [\text{Variance in } Y]$$

$$NMSE = \frac{SSE}{SST} = \frac{\sum_{i=1}^n (f(x_i) - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$MAE = \frac{1}{k} \sum_{i=1}^k |(f(x_i) - y_i)|$$

$$R^2 = \frac{\sum_{i=1}^n (f(x_i) - \bar{f}(x))^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \text{ if } R^2 \leq 1, \text{ that means model is good}$$

$$R^2 = \frac{\sum_{i=1}^n (f(x_i) - \bar{f}(x))^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \text{ if } R^2 \geq 1, \text{ that means model is good}$$

For model to be good — $R^2 \leq 1$
 one '0' of solution
 \Rightarrow $NMSE \leq 0$
 $MAE \leq 0$

\Rightarrow Total No. of us