

PRML

Evaluation:

[LabWork - In Matlab]

Assignments - 30%

mini Project - 30%

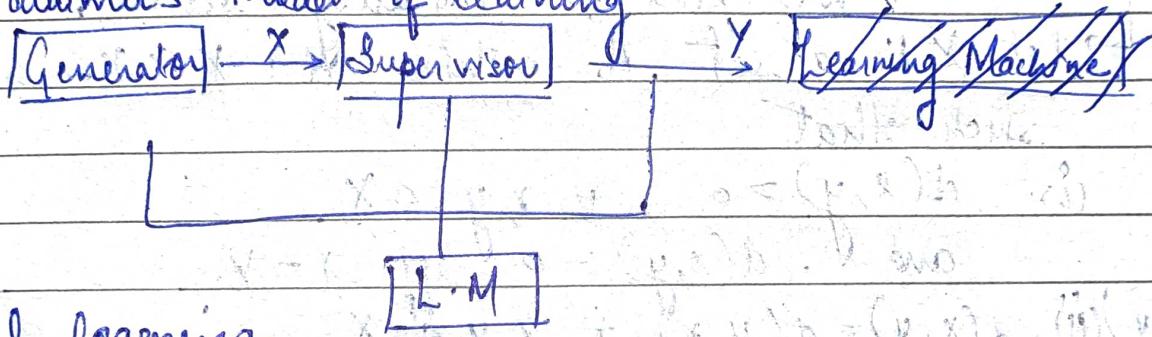
Mid Term - 20%

End Term - 20%

Presentation -

Intro:

Voldimir's Model of learning



Empirical learning
20/01/22 Model

↳ learning machine Zaki - PCA

aleng - SVM

Friedman - Bias

Weisentholt - Mathematics

Classification Problem:

$$T = \{(x_1, y_1), \dots, (x_k, y_k)\}$$

$x_i \in \mathbb{R}^n$, $y_i \in \{1, 2, \dots, K\}$ ↳ K classes.

knowing 'f' that defines relationship b/w X and Y.

We believe mathematics is all about determinism. It is not stochastic.

Number system

— Natural, Whole, Rational, — Real

This will be used.

Operations on Real Numbers

卷之三

$(R, +, *)$ → called field in algebra.

Metric space

e.g. Model studies Signals.

→ We need for similarity scores.

→ Let X be a set. A metric $f_X : X \times X \rightarrow \mathbb{R}$

such that

(i) $d(x, y) > 0$, if $x, y \in X$.

(Symmetry) (ii) $d(x, y) = d(y, x)$ & $x, y \in X$
 (iii) Triangle Inequality.

$A = \{x_1, x_2, \dots, x_m\} \subseteq X_1 \times X_2 \times \dots \times X_n \in \mathbb{R}^n$

Manhattan Distance

$$\sum_{i=1}^n |x_i - y_i| = \sum_{i=1}^n d(x_i, y_i)$$

Euclidean distance

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} = \sqrt{(x-y)^T (x-y)}$$

PAGE NO. 1 DATE: 1/202

Infinity Norm Distance

$$d(x, y) = \max_{i=1, \dots, n} (|x_i - y_i|)$$

(produced by infinity norm)

Norm in Statistics

Mahalanobis Distance $\rightarrow (x - \bar{y})^T \Sigma^{-1} (x - \bar{y})$

(x is far from \bar{y} in terms of std deviation).

Σ is positive semidefinite matrix

$$= [x_1 - \bar{x}_n] \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_{nn} \end{bmatrix} \begin{bmatrix} \bar{y}_1 \\ \bar{y}_2 \\ \vdots \\ \bar{y}_n \end{bmatrix}$$

+ve value is beyond Σ is the semi definite matrix.

$$\sum_{i=1}^n (x_i - \bar{x}_i)^2 = \sum_{i=1}^n \frac{(x_i - \bar{x}_i)^2}{\sigma_i^2}$$

Audio Signal X Audio Signal Y Just to Vector Space.

Addition $\xrightarrow{\quad}$ $\quad + \quad \xrightarrow{\quad}$

Multiplication $\xrightarrow{\quad \mu \quad \rightarrow \quad}$

Scalar $\xrightarrow{\quad x \quad \rightarrow \quad}$

PAGE No
DATE : / 202

PAGE NO
DATE : / / 2002

Def: Vector Space. and prop. of V.S:

- 1. Commutative
- 2. Associative
- 3. Additive identity.
- 4. Multiplicative identity.
- 5. additive inverse.
- 6. scalar multiplication is associative
- 7. scalar multiplication is distributive
- Another eg. V.S - Collection of all $n \times n$ matrix. with real values.

$$F = (\mathbb{R}, +, *)$$

field

Def: Subspace:

- Def: Linear independence
- Prop. \rightarrow A finite set of vectors is linearly dependent if one of them can be expressed as linear combination of others.

- Basis of V.S - Set of vectors which are linearly independent & spans V is called a basis of V.

Def: Span:

Def: Span:

Def: Span:

Def: Span:

25/01/22

L-3

$V = \mathbb{R}^n$ E.g signal

$f+g$ functional data analysis

Classify whether this a disease or not.

E.g. Basis:

- B is basis of V .
- B is minimal gen. set

$$\text{std. Basis for } \mathbb{R}^3 : \mathbb{R}^n \rightarrow \mathbb{R}$$

Linear Mapping \rightarrow Linearity property satisfied.

- \rightarrow (Coordinates) \rightarrow
- \rightarrow Coordinate vector / coordinate representation.

$$\begin{bmatrix} 1 & 2 & 6 \\ 5 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

we are not getting proper inner product (i.e dot product) then we can switch to different inner product.

\rightarrow linear transformation.

\rightarrow Norm

\rightarrow Euclidean Norm (also ℓ_2 norm).

\rightarrow Manhattan (aka ℓ_1 norm).

\rightarrow Infiniti norm.

\rightarrow General inner product.

\rightarrow Bilinear mapping omega.

\rightarrow dot product \equiv inner product.

Distribution & density functions

Regression Problems

mean \rightarrow most commonly occurring value.

$$\hat{y}_i = E[Y_i | \mathbf{x}_i] + \epsilon_i$$

ϵ_i \rightarrow distn of y_i & mean

target function \rightarrow we want to find a function, that approximate this.

distance b/w $x \& y$

$$d(x, y) := \|x - y\|_2 = \sqrt{x - y, x - y}$$

gradient & derivatives

gradient vector \rightarrow vector of partial derivatives w.r.t. component

Convex Sets:

$C \subseteq \mathbb{R}^n$ is convex, if

We can also plot relative freq. []

\rightarrow Group the data
Another plot the histogram.

27/01/22

Assumption behind regression model.

Convex functions:

Unconstrained convex optimization.

Properties:

If f_1, f_2 are convex, then $f_1 + f_2$ is convex, then, max(f_1, f_2) is also convex. It becomes very easy to find convex functions gradients.

gradient is still a vector.



for plotting we divide by bin size as itself.

Area of histogram should represent relative frequency of that class.

R.F. * BinSize.

Not a good method

\rightarrow so, reduce the size of the bin.

Total no.

Gradient descent algorithm.

- If our function is convex, then we are guaranteed to converge to global minima using gradient descent method.

if any $L \geq 0$,

$$\| \nabla f(x) - \nabla f(y) \|_2 \leq \| x - y \|_2$$

then we have this equality

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

$$\boxed{t = \gamma} \quad \text{step length}$$

$x^k \rightarrow x^{k+1}$ → optimal solution

Iteration.
 $x^0 \rightarrow$ initial.

Stochastic G.D.

Regression Model
→ finding functional dependency of $x \& y$.

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

$$\boxed{\mathbb{E}[Y|X]}$$

This can be obtained when

we have knowledge of all data is given, but in

Mean regression model hence we don't have all data points

Quantized Reg Models

Empirical Risk Minimist
Mean Regression Problem:

$$y = f^*(x) + \epsilon_i \\ \downarrow \\ E[Y|X]$$

Types loss function → 1) $L(f(x_i) - y) = (f(x) - y)^2$

2) $L(f(x) - y) = \| f(x) - y \|$
diff. types of errors

Before modeling we decide the loss function.

$$\hat{L} = L(f(x_i) - y_i) = \int L(f(x_i) - y) dP_{\text{rel}}(x, y)$$

→ sum of all losses should be minimum

Paradigm-1.

Paradigm-2 Parameter distri. Model.

Only for L strain points.

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

$$= \frac{1}{L} \sum (y_i - f(x_i))$$

But for $L \rightarrow \infty$ (By law of large numbers)
 The relative freq. converge to
 ERM defined probability.

$$\text{ERM principle} \quad \int L(Y_i - f(x_i)) dP(x, y)$$

$$= \int (Y_i - f(x_i))^2 dP(x, y).$$

\hookrightarrow SRM

\hookrightarrow Structured Risk minimization.

Original prob.

$f: \mathbb{R}^n \rightarrow \mathbb{R}$

$L \rightarrow$ loss function.

$f_1 =$ linear function

$f_2 =$ quadratic

$f =$ polynomial

$x \in \mathbb{R}^n, y \in \mathbb{R}$

$L \rightarrow$ loss function.
 $f \in \mathcal{F} \equiv$ linear function

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

$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$

$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$b \in \mathbb{R}$

$L = \|w^T x + b\|^2 = w^T w$

where $w = y - f(x)$

\hookrightarrow least squares

\hookrightarrow adv \rightarrow loss function

is convex & also

$$\text{Min } \frac{1}{L} \sum (y_i - (w^T x_i + b))^2$$

$w, b \in \mathbb{R}^n$

least square linear reg. model.

date
01-02-2022

l-4

Normal dist. density functions:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

N. of data $\mu \pm \sigma$
Interval

68% $\rightarrow [\mu - \sigma, \mu + \sigma]$

95% $\rightarrow [\mu - 2\sigma, \mu + 2\sigma]$

99.7% $\rightarrow [\mu - 3\sigma, \mu + 3\sigma]$

(estimate of mean) $\bar{x} = \frac{1}{n} \sum_{i=1}^n (x_i y_i)$

Sample Variance $= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

1.D.

In 2-D, we have to calculate covariance.

\rightarrow Covariance matrix.

$$\Sigma = \begin{bmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{bmatrix}$$

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

$$x = (x_1, x_2)$$

$f(x) = f(x_1) \cdot f(x_2)$.

$$\text{prior} = \frac{1}{\sqrt{2\pi d}} \exp^{-\frac{(x-\mu_1)^2}{2d}} \cdot \frac{1}{\sqrt{2\pi \sigma_2^2}} \exp^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}}$$

$$= \frac{1}{(2\pi)^2} | \Sigma |^{1/2} \exp^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

Variance in particular dir $d \rightarrow$ where d is unit vector.

$$\Rightarrow \begin{bmatrix} d^T \\ d \equiv 1 \end{bmatrix}$$

Want to calc. gradient

$\nabla_d \text{Cov}(x, y)$

Density function of normal dist (Gaussian dist) in n dim.

$$\frac{1}{\sqrt{2\pi}^n | \Sigma |^{1/2}} \exp^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

- follow normal dist.
- jointly

3-D Case

Multivariate normal

$$X = (X_1, X_2, \dots, X_n) \sim N(\mu, \Sigma)$$

$$w^T X \sim N(w^T \mu, w^T \Sigma w)$$

Contours of points when $f(a) = c$

$$\textcircled{1} \quad \begin{cases} x_1 \\ x_2 \end{cases} \quad \begin{array}{l} f(x_1, x_2) \\ f(x_1, x_2) = c \end{array}$$

If $\text{Var}(X) = \text{Var}(Y) \rightarrow$ then there will be co-centric ellipse.

If $X \& Y$ are independent, then - circle should be there.

$$\text{Cov}(X, Y) = 0$$

$$\text{and } \text{Var}(X) = \text{Var}(Y).$$

②

Multivariate Gauss

Gradient descent
Convex Optimization.

$$L(x, y, f) = (y_i - f(x_i))^2 \rightarrow \text{Least square loss function.}$$

Convex

We need to minimize.

$$\text{Min}_x \sum_{i=1}^n (y_i - f(x_i))^2 \rightarrow \text{Empirical loss}$$

$$L(x, y, f) \rightarrow \text{have selected only } L \text{ known points}$$

when we have all data points.

Q2 → Where to select this "f"

Linea function in n dimension can be obtained by

08/02/2022
09/02/2022

Least Square Regression Model

$$\hat{T} = \{(x_1, y_1), (x_2, y_2)\}$$

$x_i \in \mathbb{R}^n$ and $y_i \in \mathbb{R}$

Regression Model

$$y_i = f(x_i) + e_i$$

$\sum (y_i - f(x_i)) \rightarrow$ should be minimum.

$\rightarrow w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b = 0$
by propleme.

$$\rightarrow w_1 x_1 + w_2 x_2 + b = 0.$$

$$\rightarrow w_1 x_1 + w_2 x_2 + w_3 x_3 + b = 0.$$

$$\nabla f(x) = 0$$

Date: 20/02/2022

PAGE NO: 11

PAGE NO: 11

PAGE NO: 11

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_L \end{bmatrix}$$

L → data points.

$$A = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1L} \\ x_{21} & x_{22} & \dots & x_{2L} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nL} \end{bmatrix}$$

$$\text{Added last col.} \rightarrow \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$\Rightarrow L \rightarrow \text{data points.}$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix}$$

$$(Y - Aw)^T (Y - Aw)$$

$$= (Y^T - A^T A) w + A^T A w$$

Equivalent to

$$\text{Min. } \sum_{i=1}^L (y_i - (w^T x_i + b))^2$$

$$y_i - A_i w = y_i - w_1 x_{i1} + w_2 x_{i2} + \dots + w_n x_{in} + b$$

Problem Reduces to.

$$\text{Min. } (Y - Aw)^T (Y - Aw).$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{bmatrix} \in \mathbb{R}^{n+1}$$

where $w \in \mathbb{R}^{n+1}$

$$\nabla_w h(w) = 0.$$

$$\begin{bmatrix} (-A)^T (Y - Aw) \\ + (Y - Aw)^T (-A) \end{bmatrix} = 0.$$

$$= -A^T Y + A^T A w + -Y^T A + w^T A = 0$$

$$\nabla_w h(w) = 0.$$

$$f(w) = w_1 x_1^2 + w_2 x_2^2 + \dots + w_n x_n^2$$

$$= -A^T Y + A^T A w + -Y^T A + w^T A = 0$$

Quadratic function in \mathbb{R}^n .

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

Biased estimate.

$$Y = \begin{bmatrix} 1 & 2 & 1 & 2 & 2 & 2 & 1 & 2 & 1 & 2 & 1 & 2 \end{bmatrix}$$

$$\begin{aligned} A &= \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 4 & 9 & 16 & 25 & 36 & 49 & 64 & 81 & 100 & 121 & 144 \end{bmatrix} \\ &\quad \text{and } X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \end{bmatrix} \\ &\quad \text{and } Y = \begin{bmatrix} 78 \\ 112 \\ 176 \\ 240 \\ 304 \\ 368 \\ 432 \\ 500 \\ 564 \\ 632 \\ 700 \\ 768 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} &\text{AT } (Y - Aw) = 0 \\ &\Rightarrow A^T(Y - Aw) = 0 \\ &\Rightarrow A^T A w = A^T Y \\ &\Rightarrow w = (A^T A)^{-1} A^T Y \\ &\Rightarrow \text{Least square solution.} \end{aligned}$$

These are same.

Quadratic terms
Linear terms

PAGE NO: 1 / 202

$$f(u) = \underbrace{x^T H x}_{\text{Quadratic}} + \cancel{\underbrace{C^T X + b}_{\text{Linear}}}$$

$$H = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nn} \end{bmatrix}$$

$$\text{Matrix } \phi \rightarrow \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_{n-1}(x) \end{bmatrix}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

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

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

Least square quadratic regression model.

$$\text{Min}(u) = \text{Min}(u) - y$$

$$\phi_1(x) = 1$$

$$\phi_2(x) = x$$

$$\phi_1(x_2) = x^3 x_2$$

M-degree polynomial fit in N-dimension.

Want to obtain polynomial function of degree m
and $x_i \in \mathbb{R}^n$

$$x_i \in \mathbb{R}^1$$

N
o
r
e

$$y_i = \sin(2\pi x_i) + \epsilon_i$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n-1} (y_i - \hat{y}_i)^2}$$

calculation

$$x_i x_{i+1}$$

$$x_1 x_2 \dots x_{m-2}$$

$$x_1 \\ x_2 \\ \vdots \\ x_m$$

PAGE NO: 2 / 202

Bishop → 2nd Chapter

10/62

Linear Estimate

$$f(w) = w^T x + b$$

find $n+1$ parameterized $w \in \mathbb{R}^n$

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

Non-linear estimate

(M) $f(x) = w^T \phi(x) + b$,
 $\phi(x) \rightarrow$ matrix consisting of basis
 $\phi(x) = [\phi_1(x) \quad \phi_2(x) \quad \dots \quad \phi_{n-1}(x) \quad \phi_n(x)]$

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

$$\text{N. } f(x) = w^T \phi(x)$$

when $\phi(x) = [b, \phi_1(x), \dots, \phi_{n-1}(x), \phi_n(x)]^T$

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

What happens when we move from

linearity to non-linearity.

$\phi(x) \rightarrow$ functions of x .

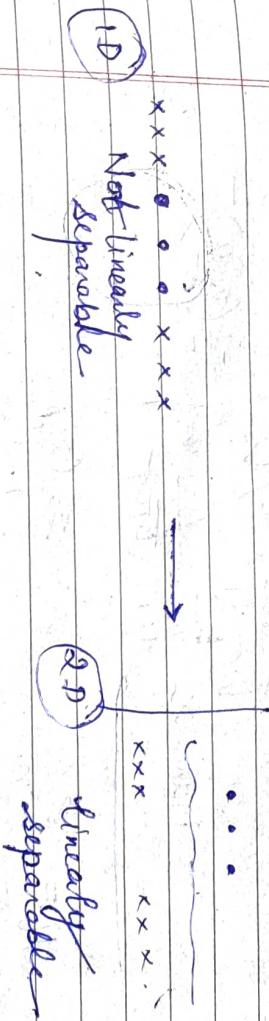
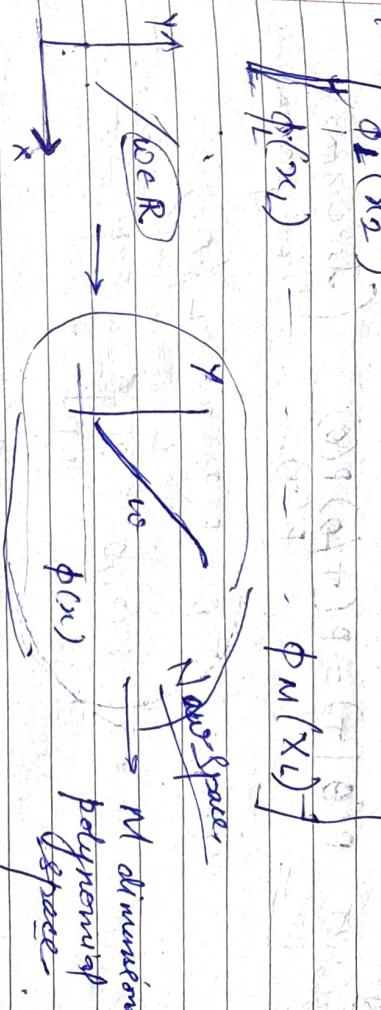
All x 's are mapped to a new space
where $f(x)$ is linearly separable.

$\phi(x) \rightarrow$ one function.

Infinite dimensional Hilbert spaces.
and learn how we will study about
kernels.

$$\phi = [\phi_1(x_1) \quad \phi_2(x_1) \quad \dots \quad \phi_M(x_1)]$$

$$\phi = [\phi_1(x_1) \quad \phi_2(x_1) \quad \dots \quad \phi_M(x_1)]$$



$$\phi(x) = [1 \quad x \quad x^2]^T$$

$$f(x) = w_0 x_0^2 + w_1 x_1 + w_2$$

Gaussian Basis function

$$\phi_j(x) = e^{-\frac{1}{2s^2} (x - x_j)^2}$$

We can obtain any curve by combining basis with appropriate parameters.

x_1, x_2, \dots, x_L

$$f(x_i) = E[Y/x_i] \quad \text{Some}$$

$$P(\theta | T) = P(T|\theta) P(\theta) \quad (\text{Likelihood})$$

prior - probability
what here?

$$P(T|\theta) = P(x_1, x_2, \dots, x_L | \theta)$$

Note:
Take derivative w.r.t. θ

$$= \frac{1}{\prod_{i=1}^L} P(x_i | \theta) = \prod_{i=1}^L \frac{1}{2\pi\sigma^2} \exp^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

Maximization of $\log P(\theta | T)$

$$= \text{Max. } \log \prod_{i=1}^L \frac{1}{2\pi\sigma^2} \exp^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}}$$

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

Maximum likelihood estimate
of θ w.r.t. T

$\hat{P}_S \rightarrow E[Y|X] \rightarrow$ same as $f(x)$?

$\mathcal{D}(Y|X) \rightarrow$ Distributions

Probability ~ In reality most of

$y = f(x) + \epsilon$ the things follow normal estimate

target function which loss function we will use for which distribution

(Gaussian, normal distribution)

$$y_i \sim N(f(x), \beta) \quad \begin{matrix} \text{mean} \\ \text{variance} \end{matrix}$$

$$P(Y|X, \omega, \beta) = N(Y | \omega^T \phi(x), \beta) \quad \begin{matrix} f(x) = \omega^T \phi(x) \text{ linear} \\ \text{mean} \\ \text{variance} \end{matrix}$$

Suppose
L - training points

$$(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)$$

$$P(Y_1 | X_1, \omega, \beta) = N(Y_1 | X_1, \omega, \beta) \quad \text{Maximize}$$

We need to maximize this (prob. of getting y_1, y_2, \dots, y_L respectively given θ)

$$\prod_{i=1}^L N(Y_i | \omega^T \phi(x_i), \beta) \quad \begin{matrix} \text{Independent samples} \\ \text{mean} \end{matrix}$$

$$= \prod_{i=1}^L \frac{1}{\sqrt{2\pi\beta}} \exp^{-\frac{1}{2\beta} (x_i - \omega^T \phi(x))^2}$$

$$\text{Max. } -x \approx \text{Max. } \log x$$

$$\Rightarrow \text{Max. } \log \prod_{i=1}^L \frac{1}{\sqrt{2\pi\beta}} \exp^{-\frac{1}{2\beta} (x_i - \omega^T \phi(x))^2}$$

$$= \text{Max. } \sum_{i=1}^L \log \beta - \frac{1}{2} \log(2\pi) - \beta \sum_{i=1}^L (y_i - \omega^T \phi(x_i))^2$$

$$\text{Max} \rightarrow \sum_{i=1}^L (\gamma_i - w^\top \phi(x_i))^2$$

$$\text{Min}_{w, b} \rightarrow \sum_{i=1}^L (\gamma_i - (w^\top \phi(x_i) + b))^2$$

Normal dis. \rightarrow min least square loss function.

Drawback

Least square loss function - sensitive to outliers.

\rightarrow error is generalized form of square sum.

2022/03/15

Least square reg. with normalization

$$\text{Loss}(L) = \min_w \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^L (\gamma_i - (w^\top \phi(x_i) + b))^2$$

$$\frac{\partial L}{\partial t} = 0$$

obtain

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

Now we will consider only this optimization problem.

$L(\lambda) = -18$ \rightarrow Better regularization

$\lambda = 0 \rightarrow$ straight line.

$(\ln \lambda = -\infty \Rightarrow \lambda = 0)$ (No regularization)

If we increase the value of lambda, then regularization there is too much regularization.

Tuning - finding optimal value of λ .

so that there is no overfit.

As the value of λ increase, the typical magnitude of the coefficient get smaller.

$$\ln \lambda = -\infty \quad \ln \lambda = -18$$

$$\lambda = e^{-\infty} \quad \lambda = e^{-18}$$

$$\lambda = \frac{1}{e^\infty} = \frac{1}{\infty} \Rightarrow \lambda = \frac{1}{e^{18}}$$

$$\Rightarrow \lambda = 0$$

$$\Rightarrow [1.52 \times 10^{-8}]$$

too much regularization

No regularization



Testing error

This value of λ is req'd.

$$\min_{w, b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^L (\gamma_i - (w^\top \phi(x_i) + b))^2$$

Lab 2 Week

PAGE NO. / DATE: / 2022

NON-UNI

PAGE NO. / DATE: / 2022

In linear estimate, the variance is much less than variance in non-linear estimate.

- 1) plots of testing & training without regularization
- 2) plots of \hat{y} after regularization
plots of \hat{y} after regularization
prior $[u]$ \rightarrow prior coeff. for this

- 3) Find optimal value of λ .

~~Notes~~
~~Date~~
14/02/2022

$$f(x) = w^T \phi(x) + b$$

We obtain this from

training set:

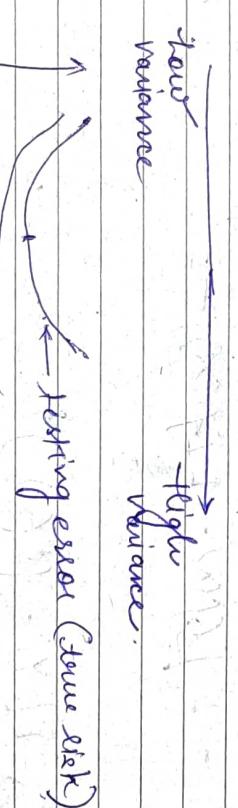
$$\begin{bmatrix} w \\ b \end{bmatrix} = [A^T A + \lambda I]^{-1} A^T Y$$

- A contains basis function
- evaluation

For training set $T_1 = \{$

$$\text{we get } \begin{bmatrix} w_1 \\ b_1 \end{bmatrix}$$

These' w_1 values with b_1 differ from w_2 and b_2 drawing set

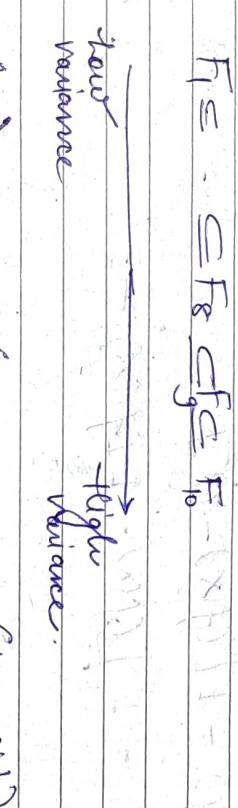


low variance \rightarrow fitting error (true risk)
high variance \rightarrow training error.

$$F_1 = \dots = F_8 \subset F_9 \subset F_{10}$$

make more complex functions to outliers as we

model sensitive



underfitting \rightarrow underfitting
overfitting \rightarrow overfitting (empirical risk)

① Modeling with fixed basis functions.

high deg Poly.
high deg.

high bias (fits poorly)
high variance
underfitting
underfitting

$$\text{Avg. Sys. Error} = (\text{Bias Error})^2 + \text{Variance}$$

Actual ~~Estimated~~ Risk → Should be estimated on entire dataset.

$$\int (Y - f(x))^2 dP(x, y)$$

True Risk → on infinite data points

$$R(f) = E[(f(x) - y)^2]$$

$$= E \left[(f(x) - E[f(x)])^2 \right] = E$$

$$\text{Expected error} = E[(e_i - f^*(x))^2]$$

Irreducible error } will always exist

$$= \underline{0} - \left(\underline{\underline{E[f(x)]}}^2 + \left(E[f(x)] \right)^2 \right) = f$$

$\exists y \forall x$

$$f_{\text{true}}(n) = E(Y - E(f(x)) + E_{\alpha}(f'(x)) - f(x))^2.$$

$$Y = E_x(f(x))^2 + E((E_x(f(x)) - f(x))^2)$$

$$E \left[f(a) \right]$$

$$(f_1 + w_2 + w_n) \oplus (a) + (b) = b_2$$

① we will run to 'O' constitutive

$$E(Y) = E(f(x)) + f(\bar{x}) \cdot E(Y - f(x))$$

$$E\{Y E(f(x))\} = E\{Eg(x) E(f(x))\}$$

- E { yfowf

$$= 2 \left[E[f(x)] - \left(E[f(x)] \right)^2 \right] =$$

$$E[Y/X] = f(x)$$

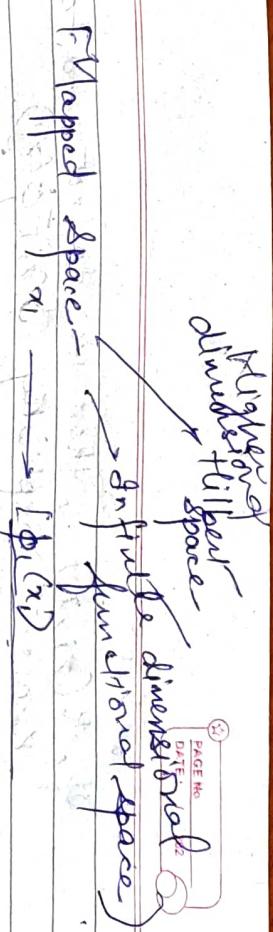
$$\text{LSE} \rightarrow E[Y - f(x)]^2 = E[(Y - E_0[f(x)])]^2$$

$$+ E \{ E_D [f^{(n)}] - f^{(n)} \}$$

$$E_0 \left[(y - f^*(x))^2 + E_0 f^*(x) - E_0 (f^*)^2 \right] + \text{variance}$$

Highly linear models
Neural networks → low Bias
But still perform better. Why?
→ Because we have large dataset, then variance is comparatively low

High variance
Highly non-linear models
Linear → high variance



Kernel Trick - Finding linear function in higher dimensional functional space is similar to finding non-linear func. in original space.

22/02/2022

$$\text{Linear } f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \quad \mathbf{w} \in \mathbb{R}^n \quad \mathbf{x} = (x_1, y_1, x_2, y_2)$$

Non-linear $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$.

$$\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_M(\mathbf{x})]$$

base function,

$$\rightarrow \sum_{i=1}^M w_i \phi_i(\mathbf{x}) + b.$$

Kernel
representing
space

analogous prop. (constraints) that \mathbf{k} and ϕ should follow.

$$\mathbf{\Phi}(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \end{bmatrix}$$

Suppose for every training point, we have a basis function.

$$f(\mathbf{x}) = \sum_{i=1}^L \phi_i(\mathbf{x}_i) w_i + b.$$

Kernel
trick
tricks

$$f(\mathbf{x}) = \sum_{i=1}^L w_i \phi_i(\mathbf{x}_i) + b.$$

Mapped
Space

$$\rightarrow \sum_{i=1}^L w_i \phi_i(\mathbf{x}) + b$$

(Input) Original
space

(1) Functional
space

$$= \sum_{i=1}^L k_{\mathbf{x}}(\mathbf{x}, \mathbf{x}_i) w_i + b \rightarrow \text{Kernel generated surface}$$

we know
if $K(x_i, x_j)$, then we do not need to
know explicitly $\phi(x_i)^\top \phi(x_j)$.

Some popular kernel → should be the semidefinite
matrix → symmetric Matrix.

Moser condition

\hookrightarrow

$x_1, x_2 \in \mathbb{R}^n$.

$$K(x_1, x_2) = k(x_1, x_2)$$

$$k(x_1, x_1)$$

$$k(x_2, x_2)$$

$$k(x_1, x_2)$$

$$k(x_2, x_1)$$

$$k(x_1, x_1)$$

$$k(x_2, x_2)$$

$$k(x_1, x_2)$$

$$k(x_2, x_1)$$

$$k(x_1, x_1)$$

$$k(x_2, x_2)$$

$$V_b(f(x), b) = 0 - \sum_{i=1}^L (y_i(\mathbf{w}^\top \phi(x_i) + b))$$

Gradient descent

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

→ data

step length

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

Norm of gradient should be very very close to 0.

Select α_1, α_2

Reduce value of η

\mathbf{x}^{k+1} Norm of gradient should be very very close to 0.

every

iteration

evaluate.

$$\eta_{(i)} = \frac{\eta}{i}$$

Better convergence

Mid Sem

Exam

Once done we move

$$\mathbf{g}(\mathbf{x}, \mathbf{u}, \mathbf{A}, \mathbf{y})$$

$$\begin{cases} x^2 \\ y^2 \end{cases}$$

$$\mathbf{A} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \\ b \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\mathbf{u} = (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{y}$$

$$\text{gradient } \mathbf{u} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \\ b \end{bmatrix}$$

$$\text{Optim} \quad \min \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \sum_{i=0}^L (y_i - (\mathbf{w}^\top \phi(x_i) + b))^2$$

$$\rightarrow \min \frac{1}{2} \frac{(\mathbf{w}^\top \mathbf{w})^2}{\mathbf{u}^\top \mathbf{u}} + (\mathbf{y} - \mathbf{A}\mathbf{u})^\top (\mathbf{y} - \mathbf{A}\mathbf{u})$$

format.

$$\mathbf{u}^\top \mathbf{Q}_0 \mathbf{u}$$

When we add b^2 → set stick to some frame of reference

do this until $J_\theta(\mathbf{u}) < \text{tolerance}$

$$\mathcal{J}_\theta(\mathbf{u}) = \mathbf{u}^\top \mathbf{Q}_0 \mathbf{u} - \mathbf{u}^\top \mathbf{A}^\top (\mathbf{y} - \mathbf{A}\mathbf{u})$$

$$\mathbf{u} = \mathbf{u} - \eta \left(\mathbf{A} \mathbf{u} - \mathbf{A}^\top (\mathbf{y} - \mathbf{A} \mathbf{u}) \right)$$

do this until $J_\theta(\mathbf{u}) < \text{tolerance}$

O + O O

DATE: / 202

24/02/2022

line 162

$\boxed{2 \times 20}$

$$M = 1$$

$$\eta = 1$$

$\begin{bmatrix} x & y \\ 1 & 1 \end{bmatrix}$

$$M = 2$$

$$\begin{bmatrix} x^2 & y^2 & xy & x & y \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\begin{bmatrix} x^2 & y^2 & xy & x & y \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$$x^2 y^0 \cdot x^1 y^1 \cdot y^2 x^0 \cdot M = M+1$$

$$\frac{M(M+1)}{2}$$

(M) $\Rightarrow 5 \times 5^2 = 25$

(M) Distance from the predicted hyperplane

$$\left(N - \frac{N-1}{2} \right)$$

$$\left(M - \frac{M-1}{2} \right)$$

K-NN based

Also we need some loss functions

$$\text{Min}_{w,b} \| w^T x_i + b - y_i \|_2^2$$

L-1 norm \rightarrow we are not squaring there

Suppose -

$$\begin{array}{ccc} \vdots & \vdots & \vdots \\ f(x) = ax + b & & \text{good estimate} \end{array}$$

But of some noise \rightarrow deviated from outliers.

When we have

least square loss estimate

It is very much prone to outliers.

How to overcome this?

Certain local measure \rightarrow to calc. density of a point.

If this is too less, then we give less weight to that data point.

Dense points are given more weightage compared to outliers.

$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|.$$

For kernel method

$$\min_{(\omega, b)} \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \left[y_i - \left(\sum_{j=1}^n K(x_j, x_i) \omega_j + b \right) \right]^+$$

Check

L1 norm - (Not smooth
C₀af)

C₁ Not differentiable

L2 norm

Smooth &

differentiable

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

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

L1 norm of reg. (3)
regular

- for D use converge use gradient descent

- for B - use convex

Use concept of sub-gradient.

Obtain gradient of

$$\text{Reg. L1} \quad \min_{(\omega, b)} \sum_{i=1}^n (y_i - (\omega^\top \phi(x_i) + b))^+$$

* More simple exp. is guaranteed to generalize well

Response Sol.

Thus, variance of unselected feature doesn't affect

$$\begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Check for every training point.

We have to calc. gradient for every sample point individually

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum_{i=1}^N (y_i - \sum_{j=1}^L k(x_j, x_i) w_j + b)^2$$

Laplace Reg. Model (we just solved)

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{1}{N} \sum_{i=1}^N |y_i - \sum_{j=1}^L k(x_j, x_i) w_j + b|$$

Laplace Support Vector Reg. Model.

L1-norm loss also called as Laplace loss.

Similarly also for

$$y_j - \left(\sum_{i=1}^L k(x_i, x_j) w_i + b \right)$$

$$u = \begin{bmatrix} w_0 \\ w \end{bmatrix}$$

vector form

$$\boxed{y - \mathbf{w}^T \mathbf{u}}$$

+ in sensitive loss function

$$L(u) = \begin{cases} 0 & \text{if } |u| \leq e \\ |u| - e & \text{otherwise} \end{cases}$$

beyond e it becomes L1 norm loss

$$\boxed{\sum_{i=1}^L k(x_i, x_j)} \text{ if } y_j - \left(\sum_{i=1}^L k(x_i, x_j) w_i + b \right) > 0$$

$$\boxed{0} \text{ if } y_j - \left(\sum_{i=1}^L k(x_i, x_j) w_i + b \right) \leq 0$$

$$A = \int \int \int \text{col} \times \text{col}$$

$K \rightarrow$ size
L \times L
data points

K - Kernel Matrix
L data points

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^l y_i (w^T \phi(x_i) + b)$$

$$U = (K + \lambda I)^{-1} Y$$

Without Regularisation,

Since K does not have 1's in last col for bias term, we make $R = [k | 1]$

$$K(x, y) = \exp\left(-\frac{\|x-y\|^2}{\sigma^2}\right)$$

Suppose we have

$$(x_1, y_1), \dots, (x_{400}, y_{400})$$

and each $x_i^\circ = (x_{i1}^\circ, x_{i2}^\circ)$.

$$K(x_i^\circ, x_j^\circ) = \exp\left(-\frac{\|x_i^\circ - x_j^\circ\|^2}{\sigma^2}\right)$$

$$U = (A^T \cdot A) + \lambda I \cdot A^T \cdot Y$$

$$y_{\text{pred}} = A \cdot U$$

$$y_{\text{pred}} =$$

$$(L+1) \times 1$$

$$R^T \times R$$

$$(L+1) \times (L+1)$$

$$I \rightarrow (L+1) \times (L+1)$$

$$\text{Identity matrix}$$

With Regularisation

$$U = (R^T \cdot R + \lambda I)^{-1} R^T \cdot Y$$

Part 1

Part 2

Part 3

Part 4

Part 5

Part 6

Part 7

Part 8

Part 9

Part 10

Part 11

Part 12

Part 13

Part 14

Part 15

Part 16

Part 17

Part 18

Part 19

Part 20

Part 21

Part 22

Part 23

Part 24

Part 25

Part 26

Part 27

Part 28

Part 29

Part 30

Part 31

Part 32

Part 33

Part 34

Part 35

Part 36

Part 37

Part 38

Part 39

Part 40

Part 41

Part 42

Part 43

Part 44

Part 45

Part 46

Part 47

Part 48

Part 49

Part 50

Part 51

Part 52

Part 53

Part 54

Part 55

Part 56

Part 57

Part 58

Part 59

Part 60

Part 61

Part 62

Part 63

Part 64

Part 65

Part 66

Part 67

Part 68

Part 69

Part 70

Part 71

Part 72

Part 73

Part 74

Part 75

Part 76

Part 77

Part 78

Part 79

Part 80

Part 81

Part 82

Part 83

Part 84

Part 85

Part 86

Part 87

Part 88

Part 89

Part 90

Part 91

Part 92

Part 93

Part 94

Part 95

Part 96

Part 97

Part 98

Part 99

Part 100

Part 101

Part 102

Part 103

Part 104

Part 105

Part 106

Part 107

Part 108

Part 109

Part 110

Part 111

Part 112

Part 113

Part 114

Part 115

Part 116

Part 117

Part 118

Part 119

Part 120

Part 121

Part 122

Part 123

Part 124

Part 125

Part 126

Part 127

Part 128

Part 129

Part 130

Part 131

Part 132

Part 133

Part 134

Part 135

Part 136

Part 137

Part 138

Part 139

Part 140

Part 141

Part 142

Part 143

Part 144

Part 145

Part 146

Part 147

Part 148

Part 149

Part 150

Part 151

Part 152

Part 153

Part 154

Part 155

Part 156

Part 157

Part 158

Part 159

Part 160

Part 161

Part 162

Part 163

Part 164

Part 165

Part 166

Part 167

Part 168

Part 169

Part 170

Part 171

Part 172

Part 173

Part 174

Part 175

Part 176

Part 177

Part 178

Part 179

Part 180

Part 181

Part 182

Part 183

Part 184

Part 185

Part 186

Part 187

Part 188

Part 189

Part 190

Part 191

Part 192

Part 193

Part 194

Part 195

Part 196

Part 197

Part 198

Part 199

Part 200

Part 201

Part 202

Part 203

Part 204

Part 205

Part 206

Part 207

Part 208

Part 209

Part 210

Part 211

Part 212

Part 213

Part 214

Part 215

Part 216

Part 217

Part 218

Part 219

Part 220

Part 221

Part 222

Part 223

Part 224

Part 225

Part 226

Part 227

Part 228

Part 229

Part 230

Part 231

Part 232

Part 233

Part 234

Part 235

Part 236

Part 237

Part 238

Part 239

Part 240

Part 241

Part 242

Part 243

Part 244

Part 245

Part 246

Part 247

Part 248

Part 249

Part 250

Part 251

Part 252

Part 253

Part 254

Part 255

Part 256

Part 257

Observation

RMSD

$$\sigma = 0.05 \quad (0.2\%)$$

$$g = 0.01 \quad \underline{0.00}$$

High value of σ , higher value of RMSE.
Poor predictions.

~~Peter~~

12

~~401 x 1~~ → coefficient

$$K\text{-test} = \begin{cases} K & K \geq \text{x-test}^*, \\ x - \text{brunng} \end{cases}$$

$$Y_{\text{test}} - \hat{y}_{\text{pred}} = R_{\text{test}} - R_{\text{pred}}$$

Yipped R.D. 5

1

gradient

20 marks

$$\sqrt{3^2 + 4^2}$$

→ Mid term exam

$$Q_2 = \lambda U - A^T(\gamma - A)U$$

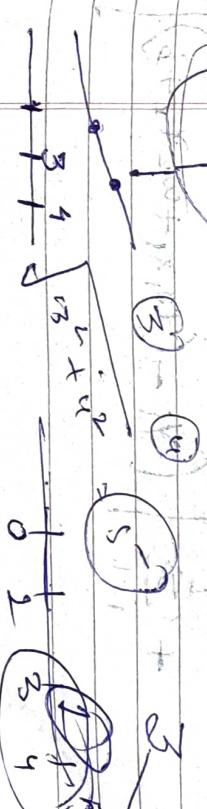


$$A(p, q), B(x, y)$$

$$(p-x)^2 + (q-y)^2$$



$$\frac{1}{x} \cdot \frac{1}{\theta}$$



WTX + B

$$\frac{\partial L}{\partial w_1} = 1w_1 + \sum_{i=1}^L \begin{cases} 1 & \text{if } Y_i - (w_1x_1 + w_2x_2) > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial L}{\partial w_2} = 1w_2 + \sum_{i=1}^L \begin{cases} -1 & \text{if } Y_i - (w_1x_1 + w_2x_2) < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\frac{\partial L}{\partial b} = 1b + \sum_{i=1}^L \begin{cases} 1 & \text{if } Y_i - (w_1x_1 + w_2x_2 + b) > 0 \\ -1 & \text{if } Y_i - (w_1x_1 + w_2x_2 + b) < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$(5) \quad \mathbf{v}^{k+1} = \mathbf{v}^k - \eta \mathbf{f}(\mathbf{A}^T(\mathbf{y} - \mathbf{H}\mathbf{v}) - \mathbf{H}^T(\mathbf{v} - \mathbf{H}\mathbf{v}))$$

L1 norm loss

$$f(w, b) = \min_{w, b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^L \left| \left(Y_i - (w^T x_i + b) \right) \right|$$

$$f(w, b) = \lambda w + \sum_{i=1}^L \begin{cases} x_i^* & \text{if } (Y_i - (w^T x_i + b)) > 0 \\ 0 & \text{if } (Y_i - (w^T x_i + b)) \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\nabla_b f(w, b) = 1b + \sum_{i=1}^L \begin{cases} 1 & \text{if } (Y_i - (w^T x_i + b)) > 0 \\ -1 & \text{if } (Y_i - (w^T x_i + b)) < 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\nabla_w f(w, b) = \lambda w + \sum_{i=1}^L \begin{cases} x_i^* & \text{if } (Y_i - (w^T x_i + b)) > 0 \\ 0 & \text{if } (Y_i - (w^T x_i + b)) \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

$$u = \begin{cases} u^* & \text{if } Y - Hu > 0 \\ H^T u & \text{if } Y - Hu \leq 0 \end{cases}$$

$$f(w_1, w_2, b) = \min_{w_1, w_2, b} (w_1^2 + w_2^2 + b^2)$$

$$\nabla_b f(w_1, w_2, b) = \sum_{i=1}^L \begin{cases} 1 & \text{if } (Y_i - (w_1x_1 + w_2x_2 + b)) > 0 \\ -1 & \text{if } (Y_i - (w_1x_1 + w_2x_2 + b)) < 0 \\ 0 & \text{otherwise} \end{cases}$$

Mid Sem - 2nd

$$H = \begin{cases} x_{11} & y_{11} \\ x_{11} & y_{12} \\ \dots & \dots \end{cases}$$

\hat{y}_i

\hat{y}_i^2

\hat{y}_i

$$\hat{y}_i^2 f(\hat{y}_i) \geq 0 \quad \forall i \quad \hat{y}_i^2 > 0.$$

$$H(u) = \begin{cases} w^T x_i + b & \hat{y}_i^2 u \\ w^T x_i + b & \hat{y}_i u \end{cases}$$

$$w^T x_i + b$$

$$H(u)^{(i)} = w^T x_i + b$$

Minimize least square

$$\min_{w_1, w_2, w} \sum_{j=1}^k (y_j - \sum_{i=1}^n (w_1^T x_i + b) \beta_i)^2$$

Just minimize this $\beta_1, \beta_2 = -\beta_L$, to

get the optimal solution

Only this also suffice

(β_1, β_2)

Activation

Extreme ML

$e(w^T x + b)$

$$y^o = \sum_{j=1}^L q(w_j^T x + b) \beta_j$$

$$\text{estimate} = q(w_1^T x + b) \beta_1 + q(w_2^T x + b) \beta_2$$

$$\left. \begin{array}{l} \gamma - (w^T x_i + b) > 0 \\ \gamma - H(u)^{(i)} > 0 \end{array} \right\} \quad \left. \begin{array}{l} \gamma - (w^T x_i + b) > 0 \\ \gamma - H(u)^{(i)} > 0 \end{array} \right\}$$

$$\left. \begin{array}{l} \gamma - (w^T x_i + b) > 0 \\ \gamma - H(u)^{(i)} > 0 \end{array} \right\}$$

$$\left. \begin{array}{l} \gamma - (w^T x_i + b) > 0 \\ \gamma - H(u)^{(i)} > 0 \end{array} \right\}$$

Gradient

$\frac{\partial L}{\partial w} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) x_i$

$\frac{\partial L}{\partial b} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)$

$$\text{Min}_{(\beta_0, \beta_1, \dots, \beta_L)} \frac{1}{2} \sum_{k=1}^K \left(Y_k - \sum_{j=1}^L \gamma_j (\omega_j^T x_k + b_j) \right)^2$$

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_K \end{bmatrix} \quad H = \begin{bmatrix} \gamma_1(\omega_1^T x_1 + b_1) & \gamma_2(\omega_2^T x_1 + b_2) \\ \vdots & \vdots \\ \gamma_1(\omega_1^T x_K + b_1) & \gamma_2(\omega_2^T x_K + b_2) \end{bmatrix}$$

Least square

When noise follows $\sim N(\mu, \sigma^2)$

\rightarrow when noise follows uniform L1 norm.

$L = \# \text{ of points}$
~~hidden nodes~~
~~weights~~

$$L \cdot \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_L \end{bmatrix}$$

$\propto X^T \beta$

$$\text{Min}_{\beta} \left\{ \frac{1}{2} (Y - H\beta)^T (Y - H\beta) \right\}$$

$$\therefore \hat{\beta} = (H^T H)^{-1} H^T Y$$

(+) for regularization



Thus, we need to make loss function adaptive.
 \rightarrow Which is not easy.

adaptive.

\hookrightarrow Only this is done is extreme machine learning.

Arriving at a powerful SVM with piecewise linear loss function.

Recall $A = \begin{bmatrix} \phi(x_1) \\ \vdots \\ \phi(x_L) \end{bmatrix}$

$$= \begin{bmatrix} \phi_1(x_1) & \dots & \phi_L(x_1) \\ \vdots & \ddots & \vdots \\ \phi_1(x_L) & \dots & \phi_L(x_L) \end{bmatrix}$$

Tuesday

Assignments - 4

PAGE NO: / 2022
DATE: / 2022

- (1) Univariate - $50 X \rightarrow \text{train}$.
(Uniform dist)

$$(2) Y_{\text{train}} \rightarrow \sin(2\pi X_i) + \epsilon_i$$

$$Y_{\text{test}} \rightarrow \sin(2\pi X_{\text{test}}) + \epsilon_i$$

- (3) 500 - testing points - $(X_{\text{test}}, Y_{\text{test}})$

- (4) Regularised least square kernel regression

$$\text{Kernel RBF} \quad K_{\text{train}} = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}_{50 \times 50}$$

$$R_{\text{train}} = \begin{bmatrix} K_{\text{train}} & 1 \\ 1^T & 50 \times 50 \end{bmatrix}_{50 \times 51}$$

$$R_{\text{test}} = \begin{bmatrix} K_{\text{test}} & 1 \end{bmatrix}_{500 \times 51}$$

and

$$K_{\text{test}}(i,j) = \exp\left(-\frac{\|X_i - X_j\|}{\sigma}\right)^2$$

where $X_i \in X_{\text{test}}$.

and $X_j \in X_{\text{train}}$

$$U = (R_{\text{train}}^T R_{\text{train}} + \lambda I)^{-1} R_{\text{train}}^T Y$$

$$(\underbrace{[51 \times 50] \times [50 \times 51]}_{[51 \times 50]} + \underbrace{[51 \times 51]}_{[51 \times 50]} \cdot \underbrace{[51 \times 50]}_{[50 \times 51]} \cdot \underbrace{[51 \times 51]}_{[51 \times 51]} = [51 \times 51])$$

Q5

$X_{\text{train-modified}}$.

→ Select 5 training points

↑ dimension

modifying their

Training $\hat{Y}_{\text{train}} = \hat{U} \cdot R_{\text{train}} + U$
estimated $\hat{Y}_{\text{train-modified}} = \hat{U} \cdot R_{\text{train-modified}} + U$
train-modified

$$[50 \times 51] \quad [51 \times 1]$$

↓

$$[50 \times 1]$$

For test RMSE.

$$\hat{Y}_{\text{test-modified}} = R_{\text{test}} * U$$

$$[500 \times 51] \quad [51 \times 1]$$

Performance comp. Q4 & Q5.

L1-norm Loss Kernel Regression Model

PAGE NO: / DATE: / 202

$$\Delta Y = Y_{\text{train}}(i^*) - R_u(i^*) \cdot R_{\text{train}}$$

$$R_{\text{train}} \rightarrow [50 \times 51]$$

gradient descent

$$U = \begin{bmatrix} 1 \\ \vdots \\ 51x1 \end{bmatrix}$$

$$U^{k+1} = U^k - \eta (\text{grad}_k)$$

ini grad-vec = gradient

gradient $\nabla_{Y_{\text{train}}, R_{\text{train}}} U$

[As per was doubt session] $f_{R_u} = R_{\text{train}} \cdot U$

grad-vec op? \rightarrow $i - ?$ how to set?

$$\Rightarrow U + \underbrace{0}_{R_{\text{train}}^T(i)} - R_{\text{train}}^T(i^*) - R_u(i^*) > 0$$

$$+ R_{\text{train}}^T(i) \quad \text{otherwise}$$

Since we have

To consider every \forall every single
drawing point, how

we will i select?

\rightarrow Can this missing area be replaced ($\Sigma_{i=1}^{50}$)

$i \rightarrow 1 \quad 2 \quad 3 \quad \dots \quad 50$

$$\begin{cases} \Delta Y \rightarrow 0 \quad \text{else } \alpha_i = 1 \\ \alpha_i = \begin{bmatrix} 1 \\ \vdots \\ 50 \end{bmatrix} \end{cases}$$

$$R_{\text{train}}^T(i) \rightarrow (51 \times 50)$$

$$L \begin{bmatrix} 1 \\ \vdots \\ 51 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ \vdots \\ 50 \end{bmatrix}$$

$$\text{grad} = U + \underbrace{R_{\text{train}}^T}_{R_{\text{train}}} \alpha$$

$$AU + R_{\text{train}}^T \alpha$$

23/03/2022

PAGE No: / 202
DATE: / 2022

Epsilon insensitive loss function.

$$\text{Le}(u) = \begin{cases} 0 & \text{if } |u| \leq \epsilon \\ |u| - \epsilon & \text{otherwise} \end{cases}$$

$u = y_i^0 - f(x_i^0)$

$$\text{Le}(y - f(x_i)) = \begin{cases} 0 & \text{if } |y - f(x_i)| \leq \epsilon \\ |y - f(x_i)| - \epsilon & \text{otherwise} \end{cases}$$

What will happen if we want to minimize epsilon sensitive loss

$$\begin{cases} f(x) = w^T x + b \end{cases}$$

Book

Duda Hartz - Chap 1 & Chap 2 - Classification.

Fish classification problem.

Class - (w_1, w_2) .

$x = (x_1, x_2)$

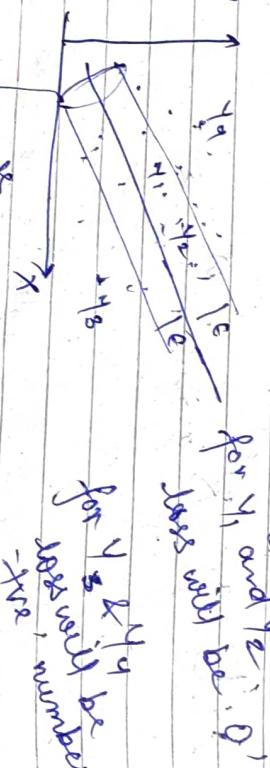
Given w_1 , class & what is the distribution of

(x/w_1)

for $y \in \mathbb{N}$

for $y \in \mathbb{N}$ be number.

$-x/w_1$



Data points lying inside epsilon insensitive region will contribute to loss function.

Parametric Bayesian methods

→ maximum likelihood methods.

Fisher discriminant analysis.

Neural networks bcoz assumption on distribution were not working.

Distribution of X in
w₁ class called likelihood
evidence

$$P(w_1 | x) = P(x | w_1) \propto P(w_1)$$

Given feature X what is the prob. that
belongs to w_1 class.
 $P(w_1 | x)$.

w_2 class

$$P(w_2 | x).$$

$$\text{Prob. of } x \text{ being } w_1 \rightarrow P(w_1).$$

$$P(w_2 | x) = P(x | w_2).$$

Density function of distribution of w_1 class
 \rightarrow how much dense our points in particular region



That means lot of data points are accumulated.

e.g. length \rightarrow

$$P(w_j | x) = \phi(x | w_j) P(w_j)$$

Minimize error. (The avg. prob. of error).

$$P(\text{error}) = \int_0^\infty P(\text{error}(x)) P(x) dx.$$

\Rightarrow This integral will result in smaller error.

$$P(x) = \sum_{j=1}^2 P(x | w_j) P(w_j)$$

posterior = likelihood \times prior evidence.

$$P(x) = \int_0^\infty f(x | k) f(k) dk$$

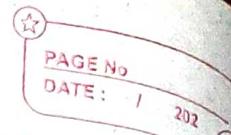
dependent variable k dependent on x .

FRMC

5:30

1 hr exam -

Mid Sem -



also

$$P(\text{error}|x) = \min(P(w_1|x), P(w_2|x))$$

$$P(w_1|x) > P(w_2|x) \Rightarrow \begin{array}{l} \text{Bayesian decision} \\ \text{Rule} \end{array}$$

w_1 class, otherwise
 w_2 class

$$\frac{P(x|w_1) P(w_1)}{P(x)} \geq \frac{P(x|w_2) P(w_2)}{P(x)}$$

Decide: w_1 : if $\frac{P(x|w_1) P(w_1)}{P(x)} \geq \frac{P(x|w_2) P(w_2)}{P(x)}$

Statistically
most powerful
But difficult
should be
known

Thus,
likelihood
prior prob factors
are important

New decision
criteria.

Here 'evidence' can be
omitted.

Multiclass classification

↳ Concept of loss function

But usually entire universe of data is
not known.