

8/01/21

Assignment → 30%

Mini-project → 30%

Mid-sem. → 20%

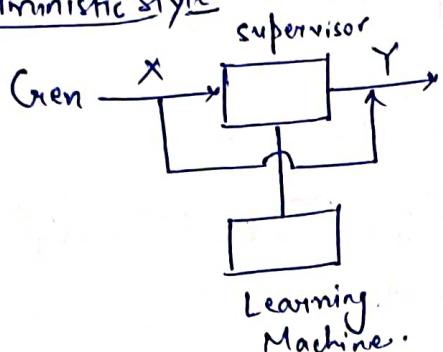
End-sem. → 20%

Presentation →

Pattern-Recognition & Machine Learning

10/01/21

Deterministic style →



features & labels →

$f_1, f_2, f_3, \dots, f_m$ → features
 $x_1 \rightarrow \begin{matrix} x \\ x \\ x \\ \vdots \end{matrix}$ → labels
 \dots
 $x_n \rightarrow \begin{matrix} x \\ x \\ x \\ \vdots \end{matrix}$ → coal → ①.
 \dots
 \dots → Diamond → ②.
 $y_i \in \{1, 2\}$

→ Here labels are represented in
euclidean space of features in \mathbb{R}^m .

so → $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$.
 where $x_i \in \mathbb{R}^n$; and $y_i \in \{1, 2\} \Rightarrow$ Example of
a stone supervised learning,

Ex. ①. fish Recognition problem in daada Hart. } \Rightarrow classification problem.
 ②. Classification of Handwritten Numbers. }

③. Problem to find Expected income. } \Rightarrow Regression problem.

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}$$

$x_i \in \mathbb{R}^n, y_i \in \mathbb{R}; i=1, 2, 3, \dots, L$

so → find a f_n for y_i for $x \in \mathbb{R}^n$

→ if output data is not labelled, then these problems are generally known as unsupervised learning. { if input data is not mapped to output labels. } \Rightarrow { clustering problems. }

④. Semi-supervised learning.

→ Regression prob. are more general & class. problem is a particular prob. of regression.

20/01/2022

Linear Algebra →

①. Metric Space →

let X be a set. A metric $d: X \times X \rightarrow \mathbb{R}$ such that →

- ①. $d(x, y) \geq 0$, $\forall x, y \in X$ and $d(x, y) = 0$ iff $x = y$.
- ②. $d(x, y) = d(y, x)$, $\forall x, y \in X$
- ③. $d(x, y) + d(y, z) \geq d(x, z)$, $\forall x, y, z \in X$

$$\text{So} \rightarrow A = \{(x_1, \dots, x_n) : x_1, x_2, \dots, x_n \in \mathbb{R}\}$$

Distance definition in A →

$$①. d(x, y) = \sum_{i=1}^n |x_i - y_i| ; \{ \text{Manhattan Distance} \}$$

②.

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

③.

$$d(x, y) = \max_{i=1, \dots, n} (|x_i - y_i|) ; \{ \text{Distance induced by infinity Norm} \}$$

④.

$$d(x, y) = X^T \Sigma^{-1} Y ; \{ \text{Mahalanobis Distance} \}$$

$$\Sigma = [x_1, \dots, x_m] \cdot \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1m} \\ \vdots & \ddots & \vdots \\ \sigma_{m1} & \dots & \sigma_{mm} \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

→ we can also write this as →

$$\Rightarrow \sum_{i=1}^n \frac{(x_i - y_i)^2}{\sigma_i^2}$$

②. Vector Space → all prop. of v.s. should satisfy.

- * closed under addition
- * " " scalar multiplication
- * additive & multiplicative inverse should exist in set.

Ex → ①. set of all n -ordered tuples.

$$R^n = \{(x_1, x_2, \dots, x_n); x_1, \dots, x_n \in R\}.$$

$$\rightarrow (P_1, P_2, \dots, P_n)$$

$$\rightarrow (P'_1, P'_2, \dots, P'_n)$$

Ex → ②. set of all 2×2 matrix over the field of R .

③. → Subspaces →

④. → linearly dependence & independence →

A finite set of vectors is linearly dependent if at least one of them can be expressed as L.C. of others.

→ i.e., that vector will lie in the plane of other Hyper

vectors:

⑤. Span of a set $\rightarrow \text{span}(\text{set}) \Rightarrow$ Total no. of vectors that can be expressed using set.

⑥. Basis of a set \rightarrow Min. no. of elements which can span over whole space.

⑦. Inner Product space \rightarrow Ref 3 properties related to I.P.S. \rightarrow

25 | 01 | 22
all Rules & properties of basis \rightarrow

⑧. Linear mapping \rightarrow linear transformation \rightarrow

$$\forall x, y \in V, \forall \lambda, \psi \in \mathbb{R}; \boxed{\phi(\lambda x + \psi y) = \lambda \cdot \phi(x) + \psi \cdot \phi(y)}$$

⑨. Norm of a vector \rightarrow

(first Norm)
①. Manhattan Norm \rightarrow

②. Euclidean Norm \rightarrow

L_1 -Norm
 L_2 -Norm

⑩. Bilinear Mapping \rightarrow

(General inner product) \rightarrow

⑪. standard inner product (Dot product) \rightarrow

symmetric

Note \rightarrow Every Cr.I.P. can be written in form of a (two) semi-definite matrix.

$$\langle \mathbf{x}, \mathbf{y} \rangle = \hat{\mathbf{x}}^T \cdot A \cdot \hat{\mathbf{y}}$$

symmetric (+ve)- semi-definite matrix.

distance \rightarrow

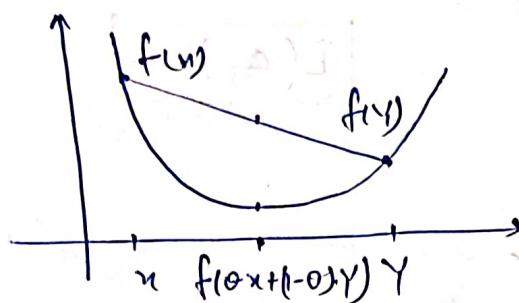
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle} = \|\mathbf{x} - \mathbf{y}\|$$

⑫. Gradients & derivatives →

⑬. Convex set →

⑭. Convex function → {concept of convexity?}.

for all values of 'x' & 'Y', this curve should have a lesser value than the line.



$$f(\theta x + (1-\theta)Y) \leq \theta \cdot f(x) + (1-\theta) \cdot f(Y)$$

* Unconstrained Convex Optimization → easier to find the minima.

(Ex) → find a f_n to check how weight & height are related?

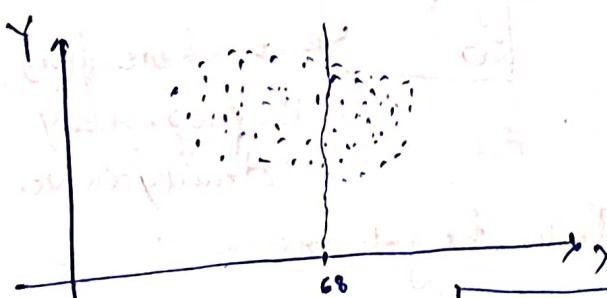
weight	\rightarrow	height
68	\rightarrow	170
72	\rightarrow	168
110	\rightarrow	180
69	\rightarrow	170
65	\rightarrow	180
68	\rightarrow	159

Ans → step ① → Plot them first.

Density & Distribution of data

check for a point

checks in a interval.



So, we can calc. $F(Y/x)$

To calc. $Y @ x = 68$, we'll take mean of all the values.

→ B/w 2 mean is a pts. where sum of square of deviations from ~~mean~~ that pt. is zero. $\Rightarrow \text{MIN. } \sum_{i=1}^n (x_i - c)^2$

So, mathematically \rightarrow

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

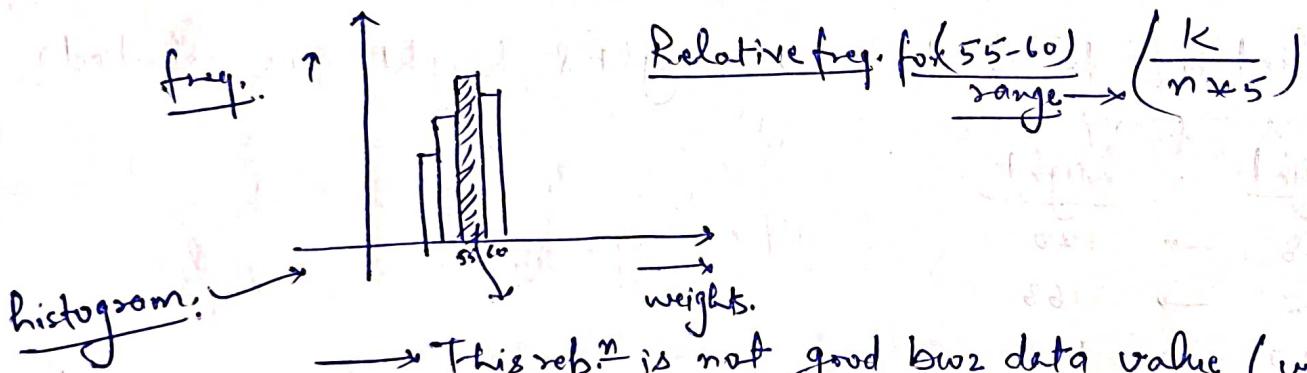
distance from mean.

\rightarrow So, we want to find a fm that can generalise that.

$$E(\epsilon_i) = 0$$

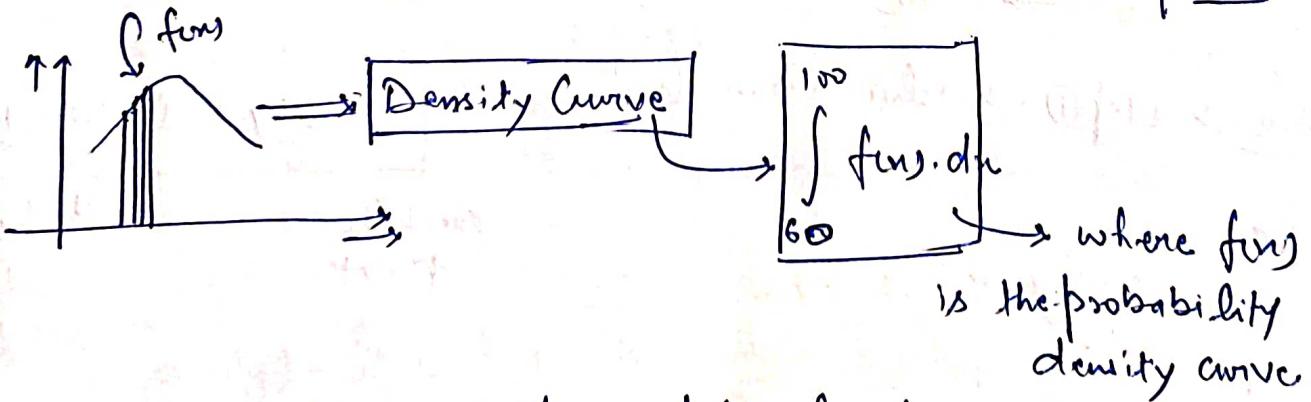
Assumptions of Regression Model →

27/01/2022
Ex(1) Weight of students \rightarrow 100 weights.



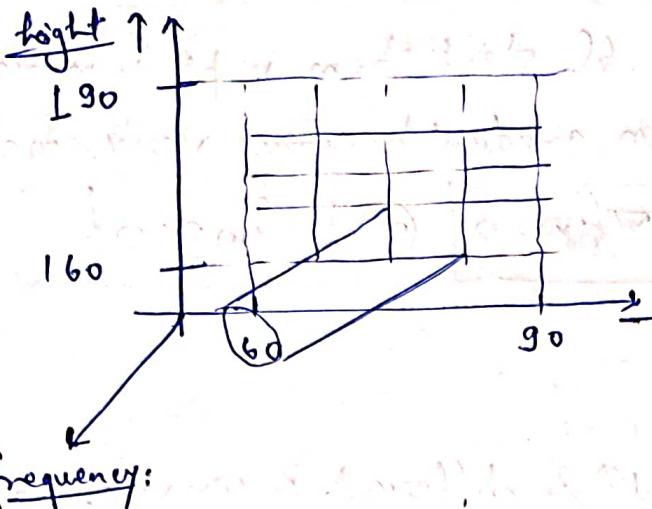
\rightarrow This rep. is not good b/w 2 data value (what it is exactly) will be lost in the range.

\hookrightarrow To solve this, we can reduce the sizes of 'bin'.



Ex(2) \rightarrow if in addition to prev. data, heights are also given.

→ We can map in 3-D.



→ We'll get a landscape structure.

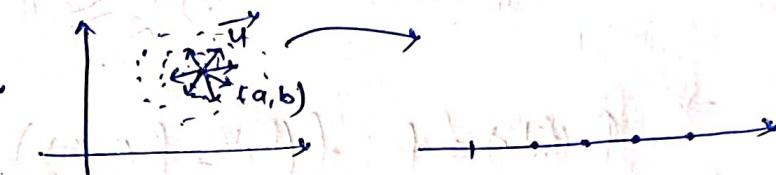
, that $f(x,y)$ will be the density f_n .

$$f(x,y) = \frac{1}{\pi}$$

→ Convex Optimization →

(Ex) $f: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(x,y) = x^2 + y^2 \Rightarrow$$



→ we want to move in the dir. u such that change is max^m.

So $\vec{u} \cdot \nabla f(a,b)$

$$\arg \max_u \nabla f(x)^T \cdot u$$

$$u = \frac{\nabla f(x)}{\|\nabla f(x)\|}$$

{ Dir^m of Max^m ascent }.

→ x
find y , s.t. $x^T y \rightarrow$
 $y = \frac{x}{\|x\|}$

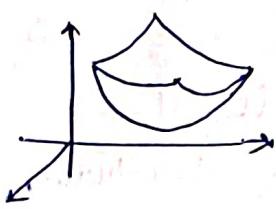
So $\frac{\text{Max.}^m \text{ rate of change}}{\rightarrow \nabla f(x)^T \cdot \nabla f(x)} \frac{\|\nabla f(x)\|}{\|\nabla f(x)\|}$

$$\Rightarrow \|\nabla f(x)\|$$

$$u = -\frac{\nabla f(x)}{\|\nabla f(x)\|}$$

{ Dir^m of Max^m Descent }

Gradient Descent Algorithm



→ we'll start from 1 pt. & move until change in gradient becomes very small.

$$\boxed{\nabla f(x) = 0} \quad \text{or} \quad \boxed{\nabla f(x) \approx 0}$$

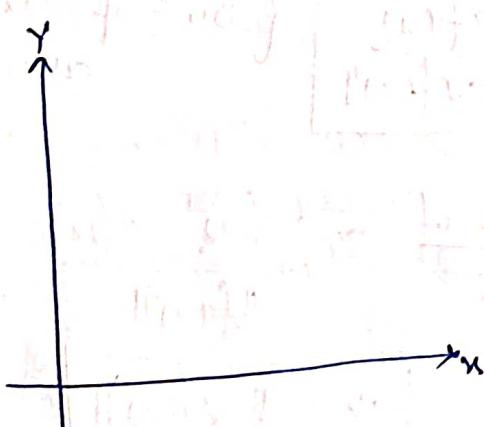
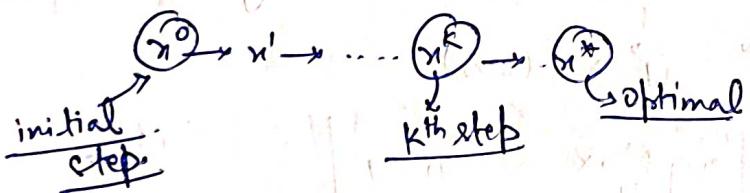
(But $f(x)$ should be convex).

Definition → if f' is contⁿ & different. & convex, such that

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

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

$$\Rightarrow f(x^{(k)}) - f(x^{(k)} - \gamma x^*) \leq \frac{\|x^{(0)} - x^*\|_2^2}{2k}; \quad \begin{array}{l} L \geq 0 \\ t \rightarrow \gamma \end{array}$$



$$f: x \rightarrow y$$

$$\text{Target: } E(Y/x) \rightarrow$$

$$T = \{x_1, y_1, \dots, x_n, y_n\}.$$

$$Y_i = f(x_i) + \epsilon_i; \quad E(\epsilon) = 0$$

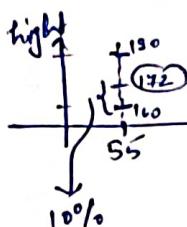
$$E(Y/x_i)$$

→ it is called Mean Regression Model.

Quantile Regression Model

we are interested in finding 'f', such that for a given value of (m) \rightarrow $P(Y \leq f(m)) = T$

\curvearrowleft given weight, what are lowest 10% values \curvearrowleft top 10% values.



$$P(172 \leq f(55)) = 0.1$$

\rightarrow This ~~is not~~ model will give more information.

(i) in Mean Regression Model \rightarrow mean height of all people having weight 55.

(ii) in Quantile Regression Model \rightarrow it can give percentage boundary as well based on value of T .

Empirical \Rightarrow Empirical Loss minimization \rightarrow

$$f: x \rightarrow Y \quad ; \text{ such that: } Y = f^*(x_i) + \epsilon_i$$

$E(Y/x_i)$

To find the 'f', we can say that \rightarrow

$$L(f(x_i) - Y_i) \text{ should be minimum } \forall i.$$

$$\stackrel{\text{So} \rightarrow}{L(f(x_i) - Y)} = (f(x_i) - Y)^2 ; \quad ; \quad x_i \in \mathbb{R}^n, Y \in \mathbb{R}$$

$$L(f(x_i) - Y) = |f(x_i) - Y| ;$$

→ Considering L is an arbitrary loss fn., then →

$$\sum_i L(f(x_i) - y_i) = \int L(f(x) - y) \cdot dP(x, y)$$

This should be
minimized.

- ↓
①. Parametric M.L. modelling.
②. Non-Parametric M.L. "

⇒ if we have L p.t.s. →

Empirical Risk minimization $\frac{1}{L} \cdot \sum_{i=1}^L L(y_i - f(x_i))$ $= \frac{1}{L} \cdot \sum (y_i - f(x_i))^2$

structured Risk Minimization $\int L(y_i - f(x_i)) \cdot dP(x, y) = \int (y_i - f(x_i))^2 \cdot dP(x, y)$

This is to be minimized.

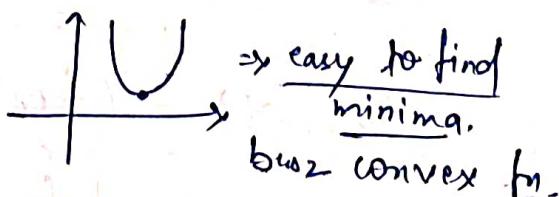
$$f \in F_1 = \{ \text{set of linear fn} \} = w^T x + b$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$F = \{ w^T x + b; w \in \mathbb{R}^n, b \in \mathbb{R} \};$$

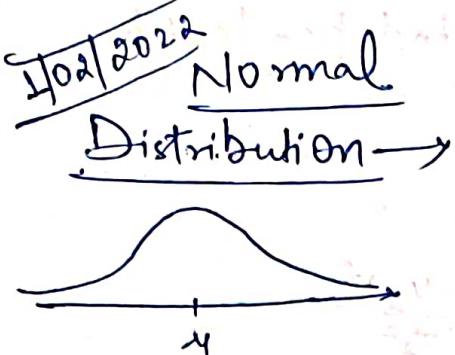
least square loss $L = (y - f(x))^2 = u^2$

where $u = y - f(x)$



Now, we want to minimise →

$$\underset{\substack{w, \\ \mathbb{R}^n}}{\text{Min}} \quad \frac{1}{L} \cdot \sum_{i=1}^L (y_i - (w^T x_i + b))^2$$



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

μ, σ

$$\left. \begin{array}{l} 68\% \text{ Data points.} \rightarrow (\mu - \sigma, \mu + \sigma) \\ 95\% \text{ " } \rightarrow (\mu - 2\sigma, \mu + 2\sigma) \\ 99.7\% \text{ " } \rightarrow (\mu - 3\sigma, \mu + 3\sigma) \end{array} \right\}$$



$$x = (x, y)$$

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

Covariance of $(x, y) \Rightarrow \Sigma = \begin{bmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{bmatrix}$

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

Ques → How can we say that given set of data points will follow Normal Distribution?

$$\begin{array}{c} x = (x, y) \\ x \sim N(\mu, \Sigma) \\ w^T x \sim N \end{array}$$

$$x = (x_1, x_2, \dots, x_n) \sim N(\mu, \Sigma)$$

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

$$\Rightarrow \begin{cases} X = (x, y) \\ x = (x_1, x_2) \end{cases} \quad \left\{ \begin{array}{l} \text{if } X \text{ & } Y \text{ are independent} \\ \text{if } X \text{ & } Y \text{ are not independent} \end{array} \right.$$

$\xrightarrow{\text{So}} f_{(x,y)} = f_{(x_1)} \cdot f_{(x_2)}$

$$= \frac{1}{\sigma_1 \sqrt{2\pi}} \cdot e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \cdot \frac{1}{\sigma_2 \sqrt{2\pi}} \cdot e^{-\frac{(y-\mu_2)^2}{2\sigma_2^2}}$$

$\Rightarrow \frac{1}{(\sqrt{2\pi})^d \cdot |\Sigma|^{\frac{1}{2}}} \cdot \exp \left\{ -\frac{1}{2} (x-\mu)^T \cdot \Sigma^{-1} \cdot (x-\mu) \right\};$

$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$

$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$

$\boxed{\text{cov}(x, y) = 0}$

Density

Probability Distribution fn of Gaussian Distribution in 'd'-dim

$$f_{(x)} = \frac{1}{(\sqrt{2\pi})^d \cdot |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \cdot \Sigma^{-1} \cdot (x-\mu) \right\};$$

$$x = (x_1, x_2, x_3, \dots, x_m);$$

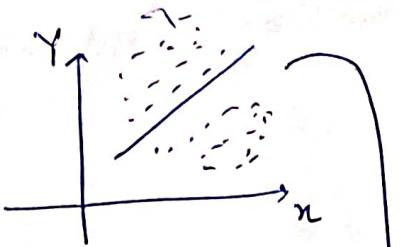
Contours \rightarrow Contour pts. are where f_x takes the constant value.

06/02/2022

$$\Rightarrow f_{\text{new}} = \omega^T n + b$$

$$\omega \in \mathbb{R}^m$$

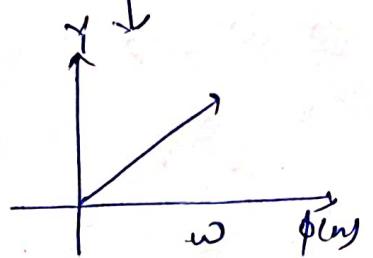
$$b \in \mathbb{R}$$



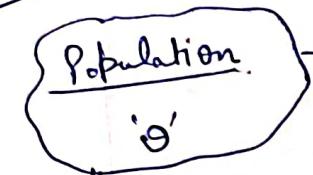
$$\Rightarrow f_{\text{new}} = \omega^T Q_{\text{new}}$$

\Rightarrow

$$Q_{\text{new}} = \begin{bmatrix} Q_1(n) \\ Q_2(n) \\ \vdots \\ Q_m(n) \end{bmatrix} \rightarrow \perp \quad \perp$$



15/02/2022



$$\{n_1, n_2, n_3, \dots, n_L\}.$$

$$N(\mu, \sigma)$$

\Rightarrow

$$P(\theta/T) = \frac{P(T/\theta) \cdot P(\theta)}{P(T)}$$

$$\approx P(\theta/T) = P(T/\theta) = P(n_1, n_2, \dots, n_L/\theta)$$

$$\Rightarrow \prod_{i=1}^L P(n_i/\theta)$$

$$\Rightarrow \prod_{i=1}^L \frac{1}{\sigma \sqrt{2\pi}} \cdot \exp \left\{ -\frac{1}{2} \frac{(n_i - \mu)^2}{\sigma^2} \right\}.$$

$$\text{Max log } P(\theta/T) = \text{Max log } \left\{ \frac{1}{\sigma \sqrt{2\pi}} \cdot \exp \left\{ -\frac{1}{2} \frac{(n_i - \mu)^2}{\sigma^2} \right\} \right\}$$

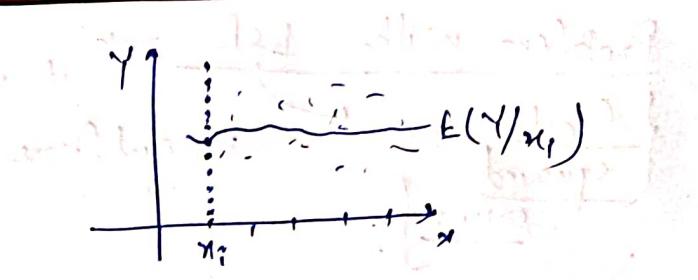
$$\mu = \frac{\sum_{i=1}^m n_i}{m}$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^m (n_i - \mu)^2}{m-1}}$$

$$Y_i = f(x_i) + \epsilon$$

\downarrow
few

$E(\epsilon) = 0$



→ Which loss fn to use in which case →

① So, we can say that (Y/x_i) is following Normal distribution with mean $f(x_i)$. $(Y/x_i) \sim N(f(x_i), \beta)$

$$P(Y/x_i, w) = N(Y_i / w^T \phi(x_i), \beta)$$

Suppose → datapoints → $\{(x_1, Y_1), (x_2, Y_2), \dots, (x_L, Y_L)\}$.

$$\text{Max } P(Y_i/x_i, w, \beta) = N(Y_i/x_i, w, \beta)$$

$$\Rightarrow \text{Max } \prod_{i=1}^L N(Y_i / w^T \phi(x_i), \beta)$$

$$\Rightarrow \text{Max } \prod_{i=1}^L \frac{1}{\beta \sqrt{2\pi}} \cdot \exp \left\{ -\frac{(Y_i - w^T \phi(x_i))^2}{2\beta^2} \right\}$$

$$\Rightarrow \text{Max } \log \prod_{i=1}^L \frac{1}{\beta \sqrt{2\pi}}$$

$$\Rightarrow \text{Max } \frac{L}{2} \log \beta - \frac{L}{2} \cdot \log(2\pi) - \beta \cdot \sum_{i=1}^L (Y_i - w^T \phi(x_i))^2$$

$$\Rightarrow \text{Max } - \sum_{i=1}^L (Y_i - w^T \phi(x_i))^2$$

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

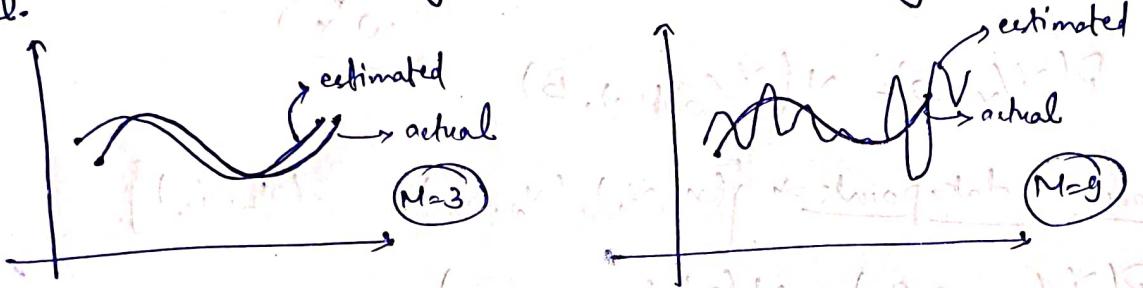
Problems with LSE → it is sensitive to the existence of outliers.

{ least squared error }.

15th feb/2022

①. for higher degree, overfitting exist. bcoz we have made our curve (2) model more flexible, so, it is fitting on the noise values as well.

②. Overfitting occurs when our training sample size is small.



so, to correct this, we minimize as following →

$$\Rightarrow \text{Min}_{w,b} \frac{\lambda}{2} \|w\|_2^2 + \sum_{i=1}^L (\gamma_i - (w^T \phi(u_i) + b))^2$$

we basically try to regularize the model.

so from now onwards, we'll only this eq.

22/02/22 → Training set $\rightarrow T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}$.

and $\rightarrow f(x) = (\omega^T x + b)$; $\omega \in \mathbb{R}^n$, $b \in \mathbb{R}$

$$x \in \mathbb{R}^L = \sum_{i=1}^L$$

→ Gradient Descent →

→ Stochastic Gradient Descent →