

① Supervised Learning

So given $X \Rightarrow$ predict $y \Rightarrow$ input $\xrightarrow{\text{model}}$ output
(label)

↙ ↘

Regression
↳ output will be in a continuous range

Classification
↳ output will be in classes (discrete)

② Unsupervised learning

So, given $X \rightarrow$ do sth with it (No label)

So, technically we only detect patterns!

↳ kmeans → k++ means

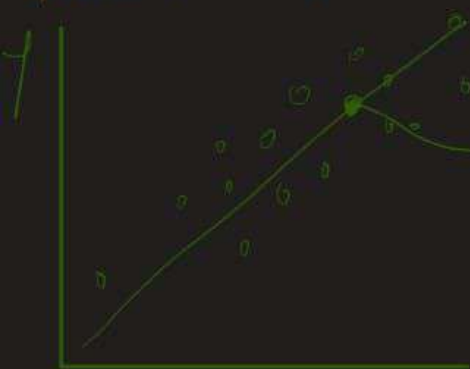
Regression and classification

14 December 2025 13:40

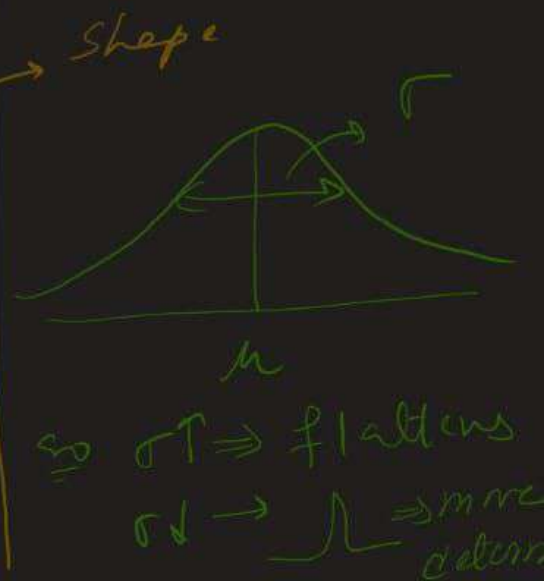
MSE Loss :- ① Normal Distribution

$$N(\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

for linear Regression \Rightarrow Assumptⁿ \rightarrow



our samples are normally distributed @
 mean = $w^T x$!



\Rightarrow to calculate the optimal value we do MLE :-
 so, if we have a distribution y we want to calculate its parameter, as we only know distributⁿ but don't know param

\Rightarrow histogram if we check the quantity



$$M(\theta, y) = \prod P(y|\theta)$$

$$\downarrow$$

$$M(2) > M(1), M(2) > M(3)$$

because probab is more where points are more so \prod will have more value!! so, that's why we define

$$MLE \Rightarrow \boxed{\hat{\theta} = \arg\max_{\theta} M(\theta, y)}$$

Regression and classification

14 December 2025

13:40

MLE for Linear regression \Rightarrow

$$\hat{w} = \underset{w}{\operatorname{argmax}} \prod_i P(y_i^o / w^T x_i)$$

$$\Rightarrow \underset{w}{\operatorname{argmax}} \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\left(\frac{y_i^o - w^T x_i}{\sigma_i}\right)^2}; \text{ assume } \sigma_i = \text{const}$$

$$\underset{w}{\operatorname{argmax}} L = \left(\frac{1}{\sigma \sqrt{2\pi}}\right)^n e^{-\sum_i \left(\frac{y_i - w^T x_i}{\sigma}\right)^2}$$

$$\underset{w}{\operatorname{argmax}} \log L = -\frac{1}{\sigma^2} \sum_i \left(\frac{y_i - w^T x_i}{\sigma}\right)^2 = -\frac{N(\text{MSE})}{\text{const}} + \text{const}$$

$$\underline{\underline{\underset{w}{\operatorname{argmax}} \log L = \underset{w}{\operatorname{argmin}} \text{MSE}}}}$$

MLE for Classification

\Rightarrow binary model $\Rightarrow p_i \rightarrow \text{probability of } y \in \{0\} (1-p_i) \Rightarrow \underline{\underline{SE}}$

$$\underline{\underline{L = p_i^m (1-p_i)^n}}; \begin{array}{l} m \rightarrow \# \text{ samples } \in 0 \\ n \rightarrow \# \text{ samples } \in 1 \end{array}$$

$$\underline{\underline{\underset{w}{\operatorname{argmax}} \log L = m \log p_i^o + n \log (1-p_i^o)}}$$

Regression and classification

14 December 2025

13:40

quick calculation

$$\nabla_{\text{loss}} = \nabla \left(\frac{1}{N} \sum (y_i - \underbrace{w^T x_i}_{{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}})^2 \right) = \nabla \ell$$

$w^T = [w_1, w_2, \dots, w_n]$ $\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

$$\frac{\partial \ell}{\partial w_i} = -\frac{2}{N} \left(\sum (y_i - w^T x_i) x_i \right)$$

if is a

$$X = \begin{bmatrix} x_1^1 & x_1^2 & \dots & x_1^m \\ x_2^1 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & \ddots & \vdots \\ x_n^1 & x_n^2 & \dots & x_n^m \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} \quad y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix}$$

then

$$\textcircled{1} \text{ MSE} = \frac{1}{N} (y - Xw)^T (y - Xw)$$
$$\textcircled{2} \nabla \ell = \frac{2}{N} (X^T) [y - Xw]$$

Is linear regression really linear?

14 December 2025

13:40

$$\hat{y} = w^T x$$

→ let's look at a polynomial f^n

$$y = ax^3 + bx^2 + cx + d \rightarrow$$

but if we map features

$$\hookrightarrow \vec{x} = [1, x, x^2, x^3]$$



then,

$$y = [d, c, b, a] \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \Rightarrow \underline{\text{so becomes linear!!}}$$

→ so $w^T x \rightarrow$ hyperplane in feature dim \Rightarrow but if features are correlated \Rightarrow can become non-linear

Q:- $\hat{y} = w^T x$ is a polynomial of degree n if training dataset contains ' m ' points where

$m < n \Rightarrow$ what is the min loss possible

$$\{y = w^T x\}??$$

Is linear regression really linear?

14 December 2025

13:40

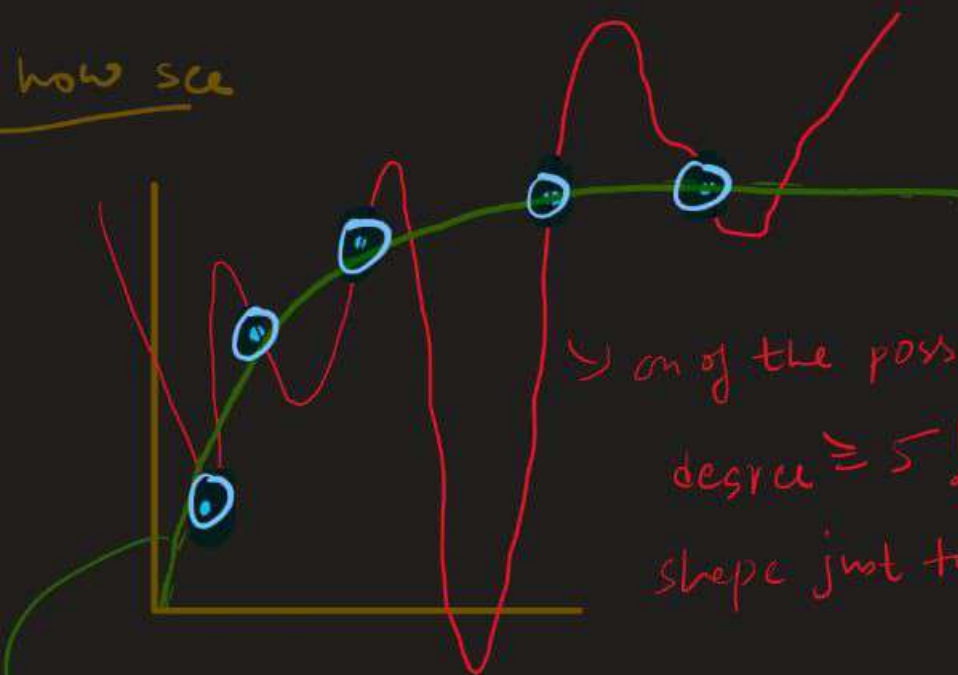
for the problem since $\exists m < n$ points $\Rightarrow (y - w^T x)$ have
'm' predefined roots now choose $n - m$ roots = 0 i.e

$$p(n) = (n - \alpha_1)(n - \alpha_2) \dots (n - \alpha_{n-m}) x^{n-m} \rightarrow n \text{ degree poly}$$

$$G = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 \Rightarrow \boxed{MSE = 0}$$

==

how see



one of the possible polynomials of
degree ≥ 5 ! \Rightarrow giving random
shape just to set $MSE = 0$

while for this $MSE \neq 0$ but it represent the
correct trend!

in the real case its observed that w_i become much
large $\approx 10^5 \rightarrow$ so we make $\boxed{\text{new loss} = MSE + \lambda \|w\|^2}$

\rightarrow this $\lambda \|w\|^2$ additⁿ is called Regularizatiⁿ

problem with normal Kmeans \Rightarrow

its possible that the randomly initialized means are too close to each other



instead of



Kmeans++ \Rightarrow for the randomly chosen points

① choose 1 random pt. (distributed probability uniformly)

② make probabilities st. $P(n_i) \propto \|x_i - \mu_0\|^2$

so, new points will be further away!

no closer possible

③ do it like that!

not used in the project