	ALZIMIE(AITO) classmate
	Date
4	Page
	Foundations of Machine Learning
	Assignment - 4 8 stuttedus del
	1)35191/merie
10	of the section
0.1	Non- Unitorm Weights in Linear Regression.
9	Non-Uniform weights in Linear Regression. EDIW) = 1 \(\frac{2}{2} \) \(g_n \) \(\tan \)^2 2 \(n=1 \)
	2 n=1
	i) Taking gradient of above function and
	equation it to zero.
	$\frac{1}{2} \sum_{n=1}^{N} \sum_{n=1}^{N} \frac{1}{2} \frac{1}$
	i) Taking gradient of above function and equating it to zero. a FD(W) = - > rn {tn - w d(xn) { d(xn) = 0} ow n=1
	Solving for w
	N
	$\sum_{n=1}^{\infty} r_n t_n \phi(x_n) = \left(\sum_{n=1}^{\infty} r_n \phi(x_n) \phi(x_n)^{\top}\right) \omega$
	V=1
	ii) Matrix form: Rewriting the error function in terms of matrix products ED(W)= 1 > rn ftn-w p(xn) f -1 (diw-t) R (dw-t).
	in terms of matrix products
	[(1) 1 F K C I WITH CO C 2 Products
	to (w)= 1 > in stn-w quan)s
	1 (1) 1) 7 2 (1 , 1)
	$=\frac{1}{2}(qw-t)R(qw-t)$
	$= \frac{1}{2} \left(w^{T} \phi^{T} R \phi w - w^{T} \phi^{T} R t - t^{T} R \phi w \right)$ $= \frac{1}{2} \left(w^{T} \phi^{T} R \phi w - w^{T} \phi^{T} R t - t^{T} R \phi w \right)$ $= \frac{1}{2} \left(w^{T} \phi^{T} R \phi w - w^{T} \phi^{T} R t - t^{T} R \phi w \right)$
	++TR+)
	$= \frac{1}{2} \left(\omega^{T} \phi^{T} R \phi \omega - 2 t^{T} R \phi \omega + t^{T} R t \right)$
	2
	Padiga (r. ru)
	R=) diag. (r, rN) : Gradient > OTROW- +TRO
	0 0000
	of enor function .: w=(dRd) tTRD
	·· W-(4 K4) TRP
	$\cdots W = (\phi^T R \phi)^T \phi^T R t$
The state of the s	

Scanned with CamScanner

Define of as NxN matrix, $\phi(i,j) = \sqrt{9}i \phi_j(x_i)$ The interpretation in 2 folds

Define of as NxN matrix, $\phi(i,j) = \sqrt{9}i \phi_j(x_i)$ The interpretation in 2 folds

Define of as NxN matrix, $\phi(i,j) = \sqrt{9}i \phi_j(x_i)$ The interpretation in 2 folds

In P(F/w,B) = $\frac{N}{N}$ In N (tn (w^T ϕ (xn),B^T))

 $= \frac{N \ln \beta - N \ln(2\pi)}{2} - \beta E_{D}(\omega)$

Ist we substitute Bt by InBthen expression would be,

 $E_{D}(\omega) = \frac{1}{2} \sum_{n=1}^{N} \frac{1}{n} \sum_{n=1}^{N} \frac{1}{n} \left(\frac{1}{n} \right)^{\frac{N}{2}}$

ii) Replicated Datapoints,

In is effective number of observations of (xn, tn) as repeatedly occurring Xntimes.

6e2 = (1 - 8xy) 6y2

Sxy) correlation coefficient for population.

0.2) Bayes Optimal Classifier Here, ∑ P(F/hi). P(hi/D)= 0.4 ZP(Uhi). P(hilD)= 0.2+0.1+0.2=0.5 E P (R/hi). P(hilD)= 0.1 ... Must probable hypothesis, (L). P(L/hi)=0.5. lie. Bayes Optimal Estimat Map Hypothesis:

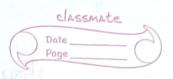
HMAP 3 Argmax P(DIh). P(h). For Forward, 3 argmax (0.4 XI) = 0.4 For Left = argmax (0.2x1, 0.1x1), 0.2x1)

=) argmax (0.2, 0.1, 0.2) and found commend form For Right = orgmax (0.1x1) = 01

. For MAP Hypothesis, Forward has maximum Value.

Hence, both are not same

0.3	H= SP, 95. classified as I iff p <x<9.< th=""></x<9.<>
	Light to a later than 129 V AZ
	The VC dimension of a set of hypothesis H is the size of the largest set C = X such that C is shattered by H.
	is the gize of the largest set C = X such
	that C is shattered by H.
20	a) Let the training points are in sphere
	At Radius D'
	Let G(x) = Sign(f(x)) = Sign (Bx+Bo) b) So, class of functions,
	b) so, class of functions
	215 24 CO 170 310 DY CAX CO 750 XI
1. Edhmate	G(x), IIBII & A & it has VC dimension H satisfying & R2A2
	H satisfying \$\leq R^2A^2
	Map Hypothesis
	: VC dimension of I-D data of hypothesis
	chare His
XG ELV	Dio = (IX Lie) XAHLER2A2 THUYAN NOT
1	for Left 3 argmax (22x1 0.1x1) 0.2x1)
4.9	y(x,w)= Wo + \(\frac{1}{2}\) Wk Xk \(\to\) ID dota
	c = 1
	sum of squares error function for N data samples
	[[1]]= 1 \(\frac{1}{2}\) \(\frac{1}{2}\) \(\frac{1}{2}\) \(\frac{1}{2}\)
	$\frac{1}{2} = \frac{1}{2} = \frac{1}$
- aron	$F_{0}(w) = 1 \sum_{i=1}^{N} (y(x_{i}, w) - t_{i})^{2} - II$ Substituting I in II with noise E
	No CO HOUSE E
	[(w)= 1 Σ [wo + Σ wi (xi+εi)]-ti ξ²
	2 1=1) [1=1
	= 1 \frac{1}{2} \left\{ \text{Wo} + \frac{1}{2} \text{Wi} \text{Xi} \right\} - \frac{1}{1} + \frac{1}{2} \text{Wi} \text{Ei} \left\{ \text{2}}
	2 = 1



$$= \frac{1}{2} \sum_{i=1}^{N} \left\{ y(x_i, \omega) - t_i \right\} + \sum_{i=1}^{N} \omega_i \varepsilon_i \right\}^2$$

$$= \frac{1}{2} \sum_{i=1}^{N} \left\{ y(x_i, \omega) - t_i \right\}^2 + \left(\sum_{i=1}^{N} \omega_i \varepsilon_i \right)^2 + \left(\sum_{i=1}^{N} \omega_i \omega_i \varepsilon_i$$

of Machine Legranda

$$\Rightarrow 2(y(x_i, \omega) - t_i) \stackrel{>}{\geq} E_{\varepsilon}[\omega_i \varepsilon_i] = 0$$

Therefore if we calculate the expectation of FD(W), with respect to Gaussian noise E, we can obtain N = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 = 1 > (y (xi, w)-ti)^2 + 62 > wi^2 > (y (xi, w)-ti)^2 + 62 > (y (xi, w)-ti)^2 + (y (xi, w)-ti)^2