```
Gearly posterior expected loss is
             R\left(\hat{y}=0|X\right) = \lambda_0 |P(y=1|X) = \lambda_0 |P| P_{,=}P(y=1|X)
(9)
           and R ( 5=0/x1= 10 + (y=0/x) = 10%
                So We will predict y=0 = 110 (1-P1)
            4 RI9=01X) < R(9=11X)
               \lambda_{0} | P_{1} < \lambda_{10} (1 - P_{1})
                      P1 < 110 = 0
                                \lambda v + \lambda 10
                16 210 = 0.1 = 10 1+9

10 1+210

Hen 10=1 and 201=9 (Not. unique)
```

	clearly loss matrix will bo
	rodicted Tone y 0 0 9
	(Note any multiple of I and 9 will also
	give some threshold or
<u>- 2</u>	posterior expected loss/Risk
7.7	postegos expected 1055/Risk
9	Cost of nejecting is ly
•	COST of picking most probable class is
	J = arg max. P (y=ie/x) is
	\(\frac{\gamma}{\pi}\) \[\frac{\cost \phi}{\pi} \text{is picking right} \\ \(\pi\) \(\p
	(+j) Class is 07
	So pick 1 if
	$\lambda_{Y} \geq \lambda_{S} P(y=i X)$
	1- P(y=) X) (Sum to
	$\sigma_{s} P(y=1 x) > 1 - \frac{\lambda y}{\lambda s}$
	15

otherwise choose réject.

Mote if a we decide to choose 9 class we has to choose J = arginax P (4=i/X) It we choose other class K =)
we will incur more wet. ie cost of choosing k will be 5 /s P(y=i (x) = /s (1-P(y=K/x) > 1s (1-P(4=J1X) because j = arg max p(y=i|x)

16 dr = 0 there i (no cost of nejecting. $\frac{2}{\sqrt{3}} \rightarrow 1$ ost of négerting in (reases-Above inequality of for most probable dows is satisfied more and more, We always a rept the most Probable class. Total points

From section 5.7.2

Pick $\hat{y}=1$ Follow upto equation 5.114

Pick $\hat{y}=1$ Follow upto equation 5.114

P(y=1|X) > LFP P(y=1|X)=Pthen $\frac{P_1}{1-P_1}$ > $\frac{1}{C}$ $\frac{1$

=> P,> 1+1.

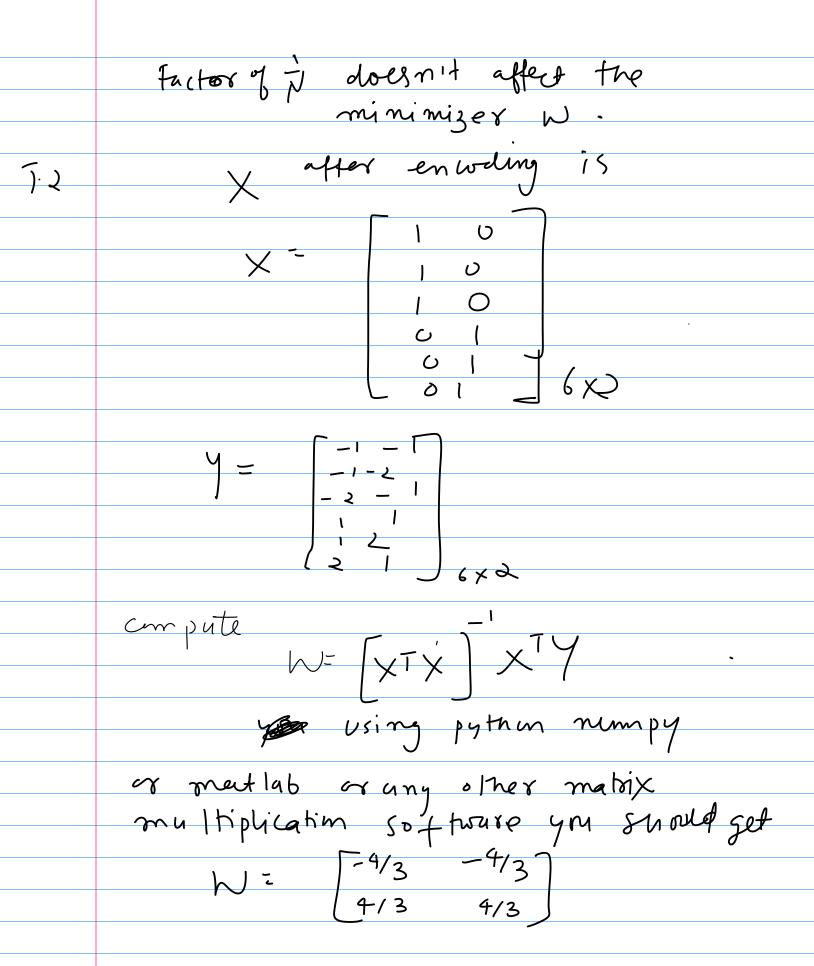
```
ridge regression derivation
                         arg max \( \sigma \) \( \sigma 
                                     expression inside ary max will be
                          2 lg 1 exp (-1 (y;-wx)2) + 2 lg 1 exp (-1 (wi)2/2)

ει (2π)/2 6
= 2 W (317/26 - 2 191-WIX)2 + 5 [ W(21/26 - 2 1/2)]
                                             we cam i gnore constant form

max mizahan

[4:-wix]2 - \frac{1}{2} \frac{1}{2}
                                                                         Hence maximizing above objective is some as minimizing

\frac{(-1)^{2}}{(-1)^{2}} = \frac{(-1)^{2}}{(-1)^{2}} + \frac{(-1)^{2}}{(-1)^{2}} = \frac{(-1)^{2}}{(-1)^{2}
                                                                                                                                                                                                      com ignorp+ve constant infrant
                                                             of whole objetive as some w will mini mize
                                                                                                                                                          > (yi-WTx)2+ > 11W112
```



using equalion 7.8 We have 7.4 l(0)-l(w,6) = -1 (4i - w[xi) 2 - N by (21162) for my MIE estimate of 62 les tuke derivative of CIW) with nespect to 6 $d(w,b) = -(2) \frac{2}{2} (y_i - w^T x_i)^T - \frac{N}{2} \frac{1}{2} \frac{(y_i - w^T x_i)^T}{2} \frac{1}{2} \frac{$ 9 - EN (9:- WTX,)L We wheaty know how to get MLE estimate of W. 8.3,41 use 1D Calculus to show J 6(a) - ((a) (1-6(a))

We know
$$N = \sum_{i=1}^{N} \left[y_{i} \log M_{i} + (I-y_{i}) \log (I-M_{i}) \right]$$

Hence $NL(w) = -\sum_{i=1}^{N} \left[y_{i} M_{i}(I-M_{i}) \times i + (I-y_{i}) \log (I-M_{i}) \times i \right]$
 $= -\sum_{i=1}^{N} \left[y_{i} - y_{i} M_{i} - M_{i} + y_{i} M_{i} \right] \times i$
 $= -\sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$
 $= \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i + \sum_{i=1}^{N} \left[y_{i} - M_{i} \right] \times i$

 $\sum_{y} = \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(y - \overline{y} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x - \overline{x} \right) / \left(x - \overline{x} \right)$ $= \frac{1}{m} \left(x$

$$E[Y|X] = Y - W^{T}X + W^{T}X = W_{0} + W^{T}X$$

$$W = [X_{0}^{T}X_{0}] \times [X_{0}^{T}X_{0}]$$

$$W = [X_{0}^{T}X_{0}] \times [X_{0}^{T}X_{0}]$$