1) Quick recap
2) Overview of GD, SGD, MB GD
3) Polynomial Regression
4) Underfit vs Overfit
5) Bias Variance trade-off

## Today's class

1) Regularization
2) L1 and L2 Regularization
3) Parameter vs Hyper-parameter
4) Cross - Validation
5) K-fold CV

$m = 2^{10}$

$BS = 2^5$

$N.B = \dfrac{2^{10}}{2^5} = 2^5$

$num\_iter = epochs \times N.B$

$\checkmark \qquad = 2 \times 2^5$

Weight
update is
happening
$\Big\{$ grad descent
loop is running

$$x_{ij} \longrightarrow \quad x_1 \quad x_1{}^2 \quad x_1{}^3$$

$$\downarrow \qquad \downarrow \qquad \downarrow$$

$$w_1 \qquad w_2 \qquad w_3$$

$$|w_3| > |w_2| > (w_1)$$

$$SS(x) = x' = \frac{x - \mu}{\sigma}$$

### Ridge Regularization



$$y = 1 + x - 0.1 x^2 + 0.01 x^3$$

$$y = 1 + \boxed{5}(x) - \boxed{5}x^2 + \boxed{5}x^3$$

$$1 + 10x - 10x^2 + 10x^3$$

$$x^2$$

$-x^2$

$x^3$

$0$

$-x^3$

1 to 2 to 4, 8, 16, 32 . . .

regularization proven ⟵ overfitting

| $W_1$ | $W_2$ | $W_3$ | $W_4$ | $\sum W_j$ X |
|-------|-------|-------|-------|--------------|
| 100   | $-50$ | $-25$ | $-25$ | $\to 0$      |

$\sum W_j^4 \to \sum W_j^2$   $\sum W_j^3$ X

Occam razor

$W = 1.5, \quad 2.0 \quad 3.5 \quad 4 \quad \Big] \; L2$

$\downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow$

$0.1 \qquad 0.2 \qquad 0.3 \qquad 0.4$

$W = 1.5, \quad 2.0, \quad 3.5 \quad 4 \quad \Big]$ 1.5+2.7

$\downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow$ 3.5+4

$\qquad \qquad \qquad = 11$

$\bigcirc \qquad 2.0 \qquad 0 \qquad 4.0$

$\sum\limits_{j=1}^{d} |w_j|$ $\qquad \qquad L> 2+4=6$



$\sqrt{0.5^2 + 0.5^2}$

$= \sqrt{0.25 + 0.25} =$

$= \sqrt{0.5} =$

$=$

$\sqrt{1^2 + 0^2}$

$= \sqrt{1} = 1$

$1 > \sqrt{0.5} \qquad \times$

$$y = w_j^r$$



$$y = |w_j|$$



**L1**

$w_j = 0.5$

$$\boxed{\eta = 0.1}$$

$w_j:$

$$0.5 - 0.1 \times 1$$
$$= 0.4$$

$$\frac{\partial |w_j|}{\partial w_j}$$

$$0.4 - 0.1 \times 1$$
$$= 0.3$$

$$0.3 - 0.1 \times 1 = 0.2$$
$$0.2 - 0.1 \times 1 = 0.1$$

$$0.1 - 0.1 \times 1 = 0$$

$$\frac{\partial w_j^2}{\partial w_j} = 2 w_j$$

**L2:**

$$w_j = 0.5 - 0.1 \times (2 \times 0.5) \downarrow$$

$$= 0.5 - 0.1$$
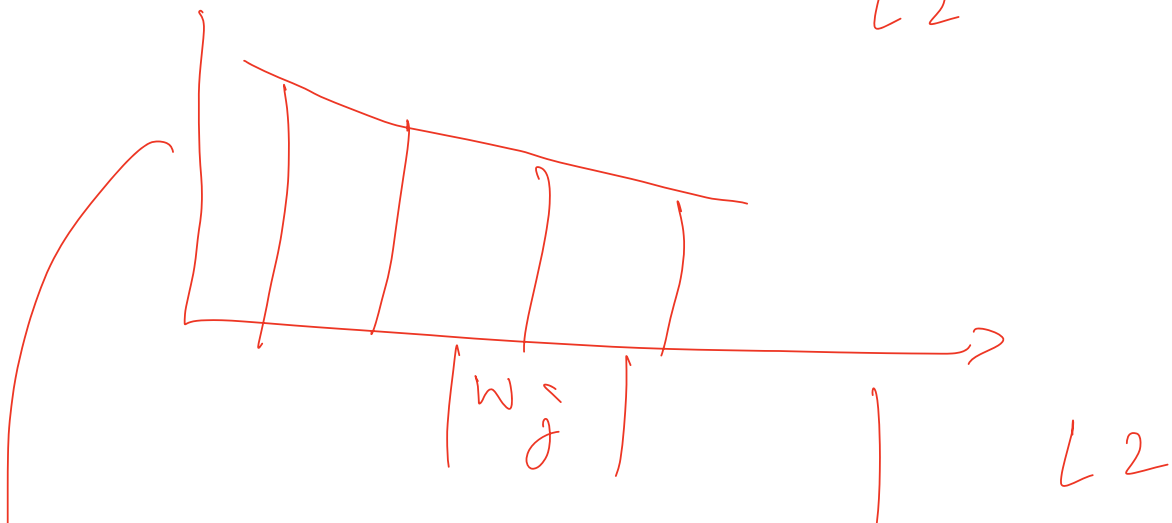
$$= 0.4$$

$$0.4 - 0.1 \times (2 \times 0.4) \downarrow$$

$$0.4 - 0.08 = 0.32$$

$$0.32 - 0.1 \times (2 \times 0.32) \downarrow$$

$$= 0.32 - 0.064$$

$$\approx 0.25$$

Elastic Net: combination of L1 & L2

L2



$|w_j|$

L2

$L_1$

$( |w_j| )$

$( |w_j| )$

$Loss =$ MSE

① ②

$$\sum (\hat{y} - y)^2 + \lambda_2 \left( \sum_{j=1}^{d} w_j{}^2 \right)$$

$$\sum (\hat{y} - y)^2 + \lambda_1 \sum_{j=1}^{d} |w_j|$$

$< 1$

$= 0.1$

$$\boxed{\lambda_2 = 10^3}$$   $$\boxed{\lambda_2 = 10^6}$$

$\downarrow \lambda$ low

$\lambda$ high $\rightarrow$ Under fitting   Over fitting

very low   too much
learning    learning

$\lambda_1 / \lambda_2 = \{ 0.01, 0.05, 0.1,$
$\{ 0.5, 1.0 \}$

$\lambda_1 \rightarrow$ min RMSE error
on test data

parameters $\longrightarrow$   $\underset{b,}{w_0}, \underset{w}{w_1, w_2, \cdots w_d}$

hyper-params $\longrightarrow$ $n$, no. of epochs for gradient descent, $\lambda_1, \lambda_2$, degree $d$ of model

$\searrow \quad \downarrow$
regularization

## Big data

$\longmapsto$ 1 M data-points

60% $\longrightarrow$ train
20% $\longrightarrow$ val     randomly
20% $\longrightarrow$ test

1000 data-points

600 $\longrightarrow$ train
200 $\longrightarrow$ val
200 $\longrightarrow$ test

$\chi$

1000 data points    0 to 999

↓ 5 folds

5-fold CV

200    399    600    799

① | ② | ③ | ④ | ⑤

0    199    400    599    800    999

↓ val

train

train

1 | 2 | 3 | 4 | 5

↓ val

$$[\lambda_1', \lambda_2', \eta', d', epoch']$$

↓ hyper-param set

for every

avg of $5/K$ val fold errors

$\swarrow$

1) Train on entire 1000 train data with optimal hyper-param ✓

2) Train on 4 folds individualy [4 out of 5 folds] with optimal hyper-param, & then finally avg out the lin reg coefficients. ✗