

1.2.1  $x_1 = [2, 1]$ ,  $t_1 = 2$ .  $J_{\text{reg}} = J + \lambda R$   $J = \frac{1}{2n} \|XW - t\|_2^2 + \frac{\lambda}{2} W^T W$ .

$$W^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{\partial J}{\partial W} = \frac{2}{n} X^T (X\hat{W} - t) + \lambda W$$

$$W^1 \leftarrow (1 - \alpha \lambda) \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \alpha \begin{bmatrix} 2 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 2 \right) + \lambda \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$= 4\alpha \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

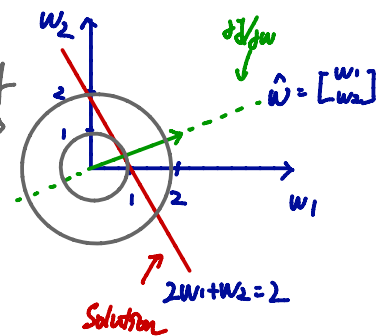
$$W^2 \leftarrow (1 - \alpha \lambda) 4\alpha \begin{bmatrix} 2 \\ 1 \end{bmatrix} - \alpha \left[ 2 \begin{bmatrix} 2 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 2 \\ 1 \end{bmatrix} 4\alpha \begin{bmatrix} 2 \\ 1 \end{bmatrix} - 2 \right) + 4\alpha \lambda \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right]$$

$$\hat{W} = \text{ro}(\begin{bmatrix} 2 \\ 1 \end{bmatrix}). \text{ direction of grad stays in the row space of } \begin{bmatrix} 2 \\ 1 \end{bmatrix}.$$

$$= a \begin{bmatrix} 2 \\ 1 \end{bmatrix}, a \in \mathbb{R}$$

$$\frac{\lambda}{2} W^T W = \frac{\lambda}{2} a [2, 1] a \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \frac{\lambda}{2} a^2 = \frac{\lambda}{2} [w_1, w_2] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \frac{\lambda}{2} (w_1^2 + w_2^2)$$

$$W^T X_1 = t_1 \Rightarrow [2, 1] \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = 2 \quad \left\{ \begin{array}{l} 2w_1 + w_2 = 2 \Rightarrow \text{solution dir.} \\ w_1^2 + w_2^2 = 2a^2 \Rightarrow \text{circle} \end{array} \right.$$



The weight decay will not change the direction of gradient descent, but the solution will not be on the solution space:  $2w_1 + w_2 = 2$ . Since the weight decay term perturbs the solution off this line.

1.3. Yes. as we showed before, the weight decay term guides/adds to update the weight in direction of row space of weight, while other terms shifts the offset. Without weight decay forms, solution resides on any lines parallel to solution, while weight decay term penalize the weight to smaller value such that it generalize to the solution in the end.

2.1.1.  $h_{\text{weight}}(x; D) = \left( \frac{1}{m} \sum_{i=1}^m w_i^T \right) x = \frac{1}{m} \left( \sum_{i=1}^m w_i^T x \right) = h_{\text{pred}}(x; D)$ .

2.2.2.  $\text{Var}[\bar{h}(x; D)] = \text{Var}\left[\frac{1}{k} \sum_{i=1}^k h(x; D_i)\right] = \frac{1}{k^2} \text{Var}\left[\sum_{i=1}^k h(x; D_i)\right] = \frac{1}{k^2} \sum_{i=1}^k \text{Var}(h(x; D_i)) = \frac{1}{k^2} \cdot k \sigma^2 = \frac{1}{k} \sigma^2$

2.3.1. No.  $E[\bar{h}(x; D) | x] = E\left[\frac{1}{k} \sum_{i=1}^k h(x; D_i) | x\right] = \frac{1}{k} \sum_{i=1}^k E[h(x; D_i) | x] = E[h(x; D) | x]$   
 $\Rightarrow E[|E[\bar{h}(x; D) | x] - y_x(x)|^2] = E[|E[h(x; D) | x] - y_x(x)|^2]$ ,  $\Rightarrow$  not changed.

2.3.3. Variance =  $\left(1 + \frac{1}{k}\right) \sigma^2$  if  $k \uparrow \Rightarrow \text{var} \downarrow$ . Then the variance decreases

$$\rho = 0 \Rightarrow \text{Cov}(h(x; D_i), h(x; D_j)) = 0 \quad \forall i \neq j$$

$\Leftrightarrow$  all ensemble members are not related to each other i.e. mutually independent.

Then.  $\text{Var} = \sigma^2/k$

$$\rho = 1 \Rightarrow \text{Cov}(h(x; D_i), h(x; D_j)) = \sigma_i \sigma_j \quad \forall i \neq j$$

$\Leftrightarrow$  The correlations between each ensemble members are very strong. i.e. perfect positive relationship.

Then  $\text{Var} = \sigma^2$  same as without bagging.

We see that the more correlations between data, the less effect bagging has.

3.1.2. Predicting  $Y$  only with  $X_2$ .

$$\hat{y}^{(0)} = w_2 x_2^{(0)}$$

$$J = E_{(X_1, X_2, Y) \sim (X_1, X_2, Y)} [(y^{(0)} - w_2 x_2^{(0)})^2]$$

$$E(G(0, 2\sigma^2) (G(0, 2\sigma^2 + 1)))$$

$$= E_{(X_1, X_2, Y) \sim (X_1, X_2, Y)} (y^{(0)^2} - 2y^{(0)} w_2 x_2^{(0)} + w_2^2 x_2^{(0)^2})$$

$$= E_{(X_1, X_2, Y) \sim (X_1, X_2, Y)} y^{(0)^2} - 2w_2 E_{(X_1, X_2, Y) \sim (X_1, X_2, Y)} y^{(0)} x_2^{(0)} + w_2^2 E x_2^{(0)^2}$$

$$= E Y^2 - 2w_2 E Y X_2 + w_2^2 E X_2^2$$

$$= 2\sigma^2 - 2w_2 \cdot 2\sigma^2 + w_2^2 (2\sigma^2 + 1)$$

$$= 2\sigma^2 (1 - 2w_2 + w_2^2) + w_2^2$$

$$= 2\sigma^2 (1 - w_2)^2 + w_2^2$$

$$\begin{aligned} E Y X_2 &= E(Y(Y + \text{Gaussian}(0, 1))) \text{ in short } \rightarrow G(0, 1) \\ &= E Y^2 + E Y G(0, 1). \quad Y \perp G(0, 1) \\ &= E Y^2 + E Y E G(0, 1) \\ &= 2\sigma^2 \end{aligned}$$

$$\text{Var } Y = E Y^2 - E^2 Y$$

$$\begin{aligned} E Y^2 &= E^2 Y + \text{Var } Y \\ &= 0 + 2\sigma^2 = 2\sigma^2 \end{aligned}$$

$$E X_2^2 = E^2 X_2 + \text{Var } X_2$$

$$= 0 + 2\sigma^2 + 1$$

$$= 2\sigma^2 + 1$$

$$\frac{\partial J}{\partial w_2} = -4\sigma^2 + 4w_2\sigma^2 + 2w_2 = 0 \Rightarrow w_2 = \frac{2\sigma^2}{2\sigma^2 + 1}$$

3.1.3.  $\hat{y} = w_1 x_1 + w_2 x_2$

$$\hat{y} = w_1 x_1 + w_2 x_2$$

$$J = E[(Y - \hat{y})^2] = E[(Y - w_1 x_1 - w_2 x_2)^2]$$

$$= E[Y^2 - 2w_1 x_1 Y - 2w_2 x_2 Y + w_1^2 x_1^2 + w_2^2 x_2^2 + 2w_1 w_2 x_1 x_2]$$

$$= E Y^2 - 2w_1 E(x_1 Y) - 2w_2 E(x_2 Y) + w_1^2 E(x_1^2) + w_2^2 E(x_2^2) + 2w_1 w_2 E(x_1 x_2)$$

$$= 2\sigma^2 - 2w_1 \sigma^2 - 2w_2 2\sigma^2 + w_1^2 \sigma^2 + w_2^2 (2\sigma^2 + 1) + 2w_1 w_2 \sigma^2$$

$$= \sigma^2 (2 - 2w_1 - 4w_2 + w_1^2 + 2w_2^2 + 2w_1 w_2) + w_2^2$$

$$\frac{\partial J}{\partial w_1} = -2\sigma^2 + 2w_1 \sigma^2 + 2w_2 \sigma^2 = 0$$

$$\frac{\partial J}{\partial w_2} = -4\sigma^2 + 4w_2 \sigma^2 + 2w_1 \sigma^2 + 2w_2 = 0$$

$$\Rightarrow \begin{cases} w_1 = \frac{1}{\sigma^2 + 1} \\ w_2 = \frac{\sigma^2}{\sigma^2 + 1} \end{cases}$$

$$E X_1^2 = E^2 X_1 + \text{Var}(X_1)$$

$$= 0 + \sigma^2 = \sigma^2$$

$$E(X_1 X_2) = E(X_1 (X_1 + G(0, \sigma^2) + G(0, 1)))$$

$$= E(X_1^2) + E(X_1)E(G(0, \sigma^2)) + E(X_1)E(G(0, 1))$$

$$= \sigma^2 \quad X_1 \perp G(0, \sigma^2), X_1 \perp G(0, 1)$$

$$E(X_1 Y) = E(X_1 (X_1 + G(0, \sigma^2)))$$

$$= E X_1^2 + E X_1 E(G(0, \sigma^2))$$

$$= E X_1^2 = \sigma^2$$

$$\Rightarrow \text{as } \sigma^2 \downarrow, \Rightarrow \begin{cases} w_1 \uparrow \\ w_2 \downarrow \end{cases}$$

we have causal effect that  $Y$  is more related to  $X_1$  than  $X_2$ .

Therefore the model will generalize more to  $x_1$  in test, which could potentially destroy the generalisation in all.

3.3.  $\hat{y} = 2(m_1 w_1 x_1 + m_2 w_2 x_2)$

$$E[J] = E[(Y - \hat{y})^2] = E[(Y - 2m_1 w_1 x_1 - 2m_2 w_2 x_2)^2]$$

$$E(m_1) = E(m_2) = 1/2$$

$$E(m_1^2) = \text{Var}(m_1) + E(m_1)^2 = (\frac{1}{4}) + (\frac{1}{2})^2 = \frac{1}{2}$$

$$= E[Y^2 + 4m_1^2 w_1^2 x_1^2 + 4m_2^2 w_2^2 x_2^2 - 4m_1 w_1 x_1 Y - 4m_2 w_2 x_2 Y + 8m_1 m_2 w_1 w_2 x_1 x_2]$$

$$= E Y^2 + 4E m_1^2 w_1^2 E x_1^2 + 4E m_2^2 w_2^2 E x_2^2 - 4E m_1 w_1 E x_1 Y - 4E m_2 w_2 E x_2 Y + 8E m_1 m_2 w_1 w_2 E x_1 x_2$$

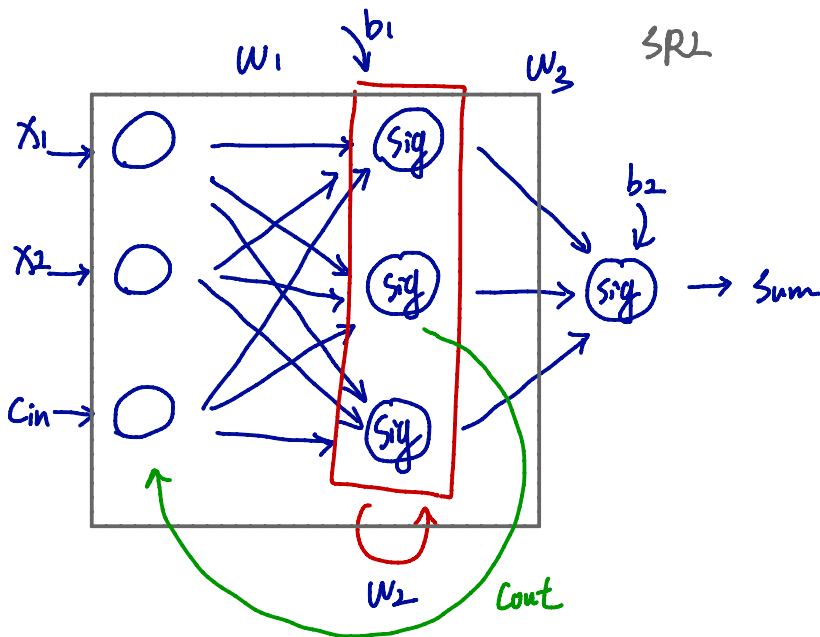
$$= 2\sigma^2 + 2w_1^2 \sigma^2 + 2(2\sigma^2 + 1)w_2^2 - 2w_1 \sigma^2 - 2w_2 2\sigma^2 + 2w_1 w_2 \sigma^2$$

$$\frac{\partial E[J]}{\partial w_1} = 4w_1 \sigma^2 - 2\sigma^2 + 2w_2 \sigma^2 = 0 \Rightarrow \begin{cases} w_1 = \frac{2\sigma^2 + 2}{7\sigma^2 + 4} \\ w_2 = \frac{3\sigma^2}{7\sigma^2 + 4} \end{cases}$$

$$\frac{\partial E[J]}{\partial w_2} = 4(2\sigma^2 + 1)w_2 - 4\sigma^2 + 2w_1 \sigma^2 = 0$$

Applying dropout changes would definitely help on the generalisation, since  $Y$  is more evenly influenced by  $w_1$  and  $w_2$ , and less affected by  $\sigma^2$ .

# A "Simple Recurrent Layer" Component (SRL) :



$$w_1 = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b_1 = \begin{bmatrix} -0.5 \\ -1.5 \\ -2.5 \end{bmatrix}$$

$$w_2 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad w_2 h = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \end{bmatrix} = \begin{bmatrix} h_{12} \\ h_{12} \\ h_{12} \end{bmatrix}$$

⇒ each time,  $h_{12}$  output is working as a carry-on bit.

$$w_3 = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} \quad b_2 = -0.1$$

$$\begin{aligned} & \text{Sig}(x_1 + x_2 + \text{cin} - 0.5) \\ & \downarrow \\ & \text{Sig}(x_1 + x_2 + \text{cin} - 1.5) \\ & \downarrow \\ & \text{Sig}(x_1 + x_2 + \text{cin} - 2.5) \end{aligned}$$

| $x_1$ | $x_2$ | $\text{cin}$ | Sum | P-Sum | Count | $h_{11}$ | $h_{12}$ | $h_{13}$ |
|-------|-------|--------------|-----|-------|-------|----------|----------|----------|
| 0     | 0     | 0            | 0   | 0     | 0     | 0        | 0        | 0        |
| 0     | 1     | 0            | 1   | 1     | 0     | 0        | 0        | 1        |
| 1     | 0     | 0            | 1   | 1     | 0     | 0        | 0        | 1        |
| 0     | 0     | 1            | 1   | 1     | 0     | 0        | 0        | 1        |
| 0     | 1     | 1            | 2   | 0     | 1     | 0        | 1        | 1        |
| 1     | 1     | 0            | 2   | 0     | 1     | 0        | 1        | 1        |
| 1     | 0     | 1            | 2   | 0     | 1     | 0        | 1        | 1        |
| 1     | 1     | 1            | 3   | 1     | 1     | 1        | 1        | 1        |

