

Homework 3

1. Weight Decay

$$\mathcal{J}(\hat{w}) = \frac{1}{2n} \|\hat{w} - t\|_2^2 + \frac{\lambda}{2} \hat{w}^T \hat{w}$$

1.1. Underparametrized Model

if $d \leq n$, $X^T X$ is invertible,

$$\begin{aligned} \frac{\partial \mathcal{J}(\hat{w})}{\partial \hat{w}} &= \frac{1}{2n} (2) X^T (X\hat{w} - t) + \frac{\lambda}{2} (2) \hat{w} \\ &= \frac{1}{n} (X^T X \hat{w} - X^T t) + \lambda \hat{w} = 0 \\ \frac{1}{n} X^T X \hat{w} - \frac{1}{n} X^T t + \lambda \hat{w} &= 0 \end{aligned}$$

$$(X^T X + \lambda I) \hat{w} = X^T t$$

we can show that $X^T X + \lambda I$ is invertible, $X^T X$ is positive semidefinite, λI is definite given that $\lambda \neq 0$, so $X^T X + \lambda I$ is positive definite. All of its eigenvalues are strictly positive, so 0 is not an eigenvalue, now looking at

if this holds for some $v \neq 0$ then 0 is an eigenvalue, contradiction! therefore, $\hat{w} = (X^T X + \lambda I)^{-1} X^T t$.

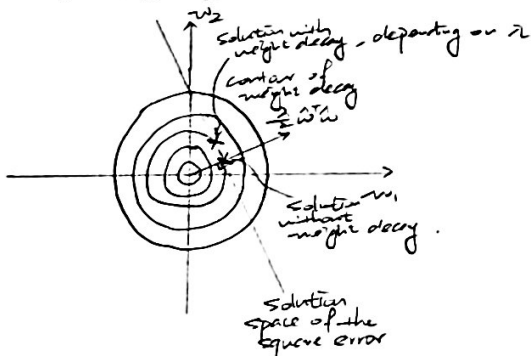
1.2 Overparametrized Model

1.2.1 Warmup: Visualizing Weight Decay

from HW1, we have that empirical risk minimizers satisfy

$$X\hat{w}^* = t$$

$$2w_1^* + w_2^* = 2$$



1.2.2 Gradient Descent and Weight Decay

if $d > n$, $X^T X$ is invertible,

$$\hat{w}(0) = \frac{1}{n} (0 - X^T t) = -\frac{1}{n} X^T t$$

with 0 weight initialization, the gradient direction is $X^T t \in \mathbb{R}^{d \times 1}$.

it cannot be the same as HW1 3.4.1 since the gradient changes along the way.

1.3 Adaptive optimizer and Weight Decay

AdaGrad with weight decay:

$$G_{i,t} = (1 - \delta) G_{i,t-1} + \delta (\nabla_{w_{i,t}} \mathcal{J}(w_{i,t}))^2$$

$$w_{i,t+1} = w_{i,t} - \frac{\eta}{\sqrt{G_{i,t} + \epsilon}} \nabla_{w_{i,t}} \mathcal{J}(w_{i,t})$$

Let's do the 2D toy example first, with $w_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$.

$$\nabla_{w_0} \mathcal{J}(w_0) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} (-2) = \begin{bmatrix} -4 \\ -2 \end{bmatrix}$$

$$G_{1,0} = (1 - \delta) 0 + \delta (-4)^2 = 16\delta$$

$$G_{2,0} = (1 - \delta) 0 + \delta (-2)^2 = 4\delta$$

$$w_{1,1} = 0 - \frac{\eta}{\sqrt{16\delta + \epsilon}} (-4) = \frac{4\eta}{\sqrt{16\delta + \epsilon}} \approx \frac{\eta}{\sqrt{\delta}}$$

$$w_{2,1} = 0 - \frac{\eta}{\sqrt{4\delta + \epsilon}} (-2) = \frac{2\eta}{\sqrt{4\delta + \epsilon}} \approx \frac{\eta}{\sqrt{\delta}}$$

So, the weights still leaves the span of X , since the weight update is not along $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$.

2. Ensembles and Bias-Variance Decomposition

2.1.1 Weight Average or Prediction Average

Suppose there are n models,

$$h_1(x; D_1), h_2(x; D_2), \dots, h_n(x; D_n)$$

using prediction averaging, we have that

$$\begin{aligned} \hat{z}_{\text{avg}} &= \frac{1}{n} (\hat{z}_1 + \hat{z}_2 + \dots + \hat{z}_n) \\ &= \frac{1}{n} (x\hat{w}_1 + x\hat{w}_2 + \dots + x\hat{w}_n + b_1 + \dots + b_n) \end{aligned}$$

where $\hat{w}_1, \dots, \hat{w}_n$ are learned from D_1, \dots, D_n , respectively, using weight averaging, we have that

$$\hat{z}_{\text{avg}} = x \frac{1}{n} (\hat{w}_1 + \hat{w}_2 + \dots + \hat{w}_n) + \frac{1}{n} (b_1 + \dots + b_n)$$

$\hat{z}_{\text{avg}} = \hat{z}_{\text{avg}}'$ - so the ensemble model using prediction averaging and weight averaging with linear activation are the same, so their generalization errors should be the same, too.

2.1.2

No, $\hat{z}_{\text{avg}} \neq \hat{z}_{\text{avg}}'$

2.2 Bagging - Unrelated Models.

$$\bar{h}(x; D) = \frac{1}{K} \sum_{i=1}^K h(x; D_i)$$

2.2.1

we need to show that

$$\mathbb{E}[\bar{h}(x; D) | x] = \mathbb{E}[h(x; D) | x],$$

$$\begin{aligned} \text{LHS} &= \mathbb{E}\left[\frac{1}{K} \sum_{i=1}^K h(x; D_i) | x\right] \\ &= \frac{1}{K} \mathbb{E}\left[\sum_{i=1}^K h(x; D_i) | x\right] \\ &= \frac{1}{K} \sum_{i=1}^K \mathbb{E}[h(x; D_i) | x] \\ &= \frac{1}{K} K \mathbb{E}[h(x; D) | x] \\ &= \mathbb{E}[h(x; D) | x] = \text{RHS}. \end{aligned}$$

[2.2.2]

$$\begin{aligned} \text{Var}(\bar{h}(x; D) | x) &= \mathbb{E}[(\bar{h}(x; D) - \mathbb{E}[\bar{h}(x; D) | x])^2] \\ &= \mathbb{E}\left[\left(\frac{1}{K} \sum_{i=1}^K h(x; D_i) - \mathbb{E}\left[\frac{1}{K} \sum_{i=1}^K h(x; D_i) | x\right]\right)^2\right] \\ &= \frac{1}{K^2} \mathbb{E}\left[\left(\sum_{i=1}^K h(x; D_i) - \sum_{i=1}^K \mathbb{E}[h(x; D_i) | x]\right)^2\right] \\ &= \frac{1}{K^2} \mathbb{E}\left[\sum_{i=1}^K (h(x; D_i) - \mathbb{E}[h(x; D_i) | x])^2\right] \\ &= \frac{1}{K^2} \cdot K \cdot \sigma^2 \\ &= \frac{\sigma^2}{K}. \end{aligned}$$

2.3 Bagging - General case.

$$\rho = \frac{\text{cov}(h(x; D_j), h(x; D_k))}{\sigma_j \sigma_k}, \quad \forall j \neq k.$$

correlation,

[2.3.1] Bias under correlation.

Due to linearity of expectation,

$$\mathbb{E}[aX] = a\mathbb{E}[X]$$

$$\mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

for any $a \in \mathbb{R}$ and random variables X and Y .

(2.2.1)

the same reasoning as before still applies, that $\mathbb{E}\left[\frac{1}{K} \sum_{i=1}^K h(x; D_i) | x\right] = \frac{1}{K} \sum_{i=1}^K \mathbb{E}[h(x; D_i) | x] = \mathbb{E}[h(x; D) | x]$.

So the bias term doesn't change.

2.3.2 Variance under Correlation.

$$\text{variance} = \left(\rho + \frac{1-\rho}{K}\right) \sigma^2$$

[2.3.3] Intuition on Bagging.

Looking at the variance equation, increasing K will lower the variance,

when $\rho = 0$, we recover the variance when ensemble members are not correlated, i.e. we get $\frac{1}{K} \sigma^2$ back.

when $\rho = 1$, we get σ^2 variance, ensemble doesn't decrease variance when ensemble members are fully correlated. This makes sense since in this case we don't get any new information.

3. Generalization and Dropout.

3.1.1 Regress from X_1 .

$$\hat{y} = w_1 x_1$$

$$\mathcal{J} = \mathbb{E}[(Y - w_1 x_1)^2]$$

$$= \mathbb{E}[Y^2] - 2w_1 \mathbb{E}[x_1 Y] + w_1^2 \mathbb{E}[x_1^2].$$

$$\frac{\partial \mathcal{J}}{\partial w_1} = -2 \mathbb{E}[x_1 Y] + 2w_1 \mathbb{E}[x_1^2] = 0$$

$$\mathbb{E}[x_1 Y] = w_1 \mathbb{E}[x_1^2].$$

$$\mathbb{E}[x_1] \mathbb{E}[Y] + \text{cov}(x_1, Y) = w_1 \mathbb{E}[x_1^2]$$

$$0 \quad w_1 \mathbb{E}[x_1^2] = \text{cov}(x_1, Y)$$

$$w_1 \mathbb{E}[(x_1 - 0)^2] = \text{cov}(x_1, Y)$$

$$w_1 \text{cov}(x_1, x_1) = \text{cov}(x_1, Y)$$

$$w_1 = \frac{\text{cov}(x_1, Y)}{\text{cov}(x_1, x_1)}$$

$$= \frac{\text{cov}(x_1, x_1 + N(0, \sigma^2))}{\sigma^2} = \frac{\text{cov}(x_1, x_1) + \text{cov}(x_1, N(0, \sigma^2))}{\sigma^2} = \frac{\sigma^2 + 0}{\sigma^2} = 1.$$

[3.2.2] Regress from x_2

$$\hat{y} = w_2 x_2$$

$$\mathcal{J} = \mathbb{E}[(Y - w_2 x_2)^2]$$

$$= \mathbb{E}[Y^2] - 2w_2 \mathbb{E}[x_2 Y] + w_2^2 \mathbb{E}[x_2^2]$$

$$\frac{\partial \mathcal{J}}{\partial w_2} = -2 \mathbb{E}[x_2 Y] + 2w_2 \mathbb{E}[x_2^2] = 0$$

$$w_2 = \frac{\mathbb{E}[x_2 Y]}{\mathbb{E}[x_2^2]} = \frac{\mathbb{E}[Y] \mathbb{E}[x_2] + \text{cov}(x_2, Y)}{\sigma^2}$$

$$= \frac{0 + \text{cov}(x_2, Y)}{\sigma^2}$$

$$= \text{cov}(x_2, x_1 + N(0, \sigma^2)) = \text{cov}(x_2, x_1) + \text{cov}(x_2, N(0, \sigma^2))$$

6 terms.

$$\begin{aligned} &= \text{cov}(x_2, x_1) + \text{cov}(x_2, N(0, \sigma^2)) + \text{cov}(N(0, \sigma^2), x_1) \\ &\quad + \text{cov}(N(0, \sigma^2), N(0, \sigma^2)) + \text{cov}(N(0, \sigma^2), x_1) \\ &\quad + \text{cov}(N(0, \sigma^2), N(0, \sigma^2)) + \text{cov}(N(0, \sigma^2), x_1) \end{aligned}$$

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

$$= 2\sigma^2$$

3.1.3 Regress from (x_1, x_2) .

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

$$J = E[(y - w_1 x_1 - w_2 x_2)^2]$$

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

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

$$= \text{cov}(y, y) + w_1^2 \text{cov}(x_1, x_1) + w_2^2 \text{cov}(x_2, x_2) - 2w_1 \text{cov}(x_1, y) - 2w_2 \text{cov}(x_2, y) + 2w_1 w_2 \text{cov}(x_1, x_2)$$

$$\frac{\partial J}{\partial w_1} = 2w_1 \text{cov}(x_1, x_1) - 2 \text{cov}(x_1, y) + 2w_2 \text{cov}(x_1, x_2)$$

$$= 2w_1 \sigma^2 - 2\sigma^2 + 2w_2 \sigma^2$$

$$= (2w_1 + 2w_2 - 2)\sigma^2 = 0, \quad \begin{matrix} 2w_1 = 2 - 2w_2 \\ w_1 = 1 - w_2 \end{matrix}$$

$$\frac{\partial J}{\partial w_2} = 2w_2 \text{cov}(x_2, x_2) - 2 \text{cov}(x_2, y) + 2w_1 \text{cov}(x_1, x_2)$$

$$= 2w_2 (2\sigma^2 + 1) - 4\sigma^2 + 2w_1 \sigma^2$$

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

$$4w_2 \sigma^2 + 2w_2 - 4\sigma^2 + 2(1 - w_2)\sigma^2 = 0$$

$$4w_2 \sigma^2 + 2w_2 - 4\sigma^2 + 2\sigma^2 - 2w_2 \sigma^2 = 0$$

$$4w_2 \sigma^2 + 2w_2 - 2\sigma^2 - 2w_2 \sigma^2 = 0$$

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

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

$$w_2 = \frac{\sigma^2}{\sigma^2 + 1}, \quad w_1 = 1 - \frac{\sigma^2}{\sigma^2 + 1} = \frac{1}{\sigma^2 + 1}$$

At test time, $\hat{y} = w_1 x_1 + w_2 x_2$

$$= \frac{1}{\sigma^2 + 1} x_1 + \frac{\sigma^2}{\sigma^2 + 1} x_2 \text{ is a weighted sum of } x_1 \text{ and } x_2$$

Suppose σ^2 gets ^{bigger} smaller during the test time,
^(decrease) w_1 will increase and w_2 will decrease.

So x_1 has ^(more) more effect on the output.
 and x_2 has ^(less) less effect on the output

In the extreme case, if $\sigma^2 = 0$ for training, then $w_1 = 1, w_2 = 0$.

$$\hat{y} = x_1$$

if $\sigma^2 \neq 0$ for test, then prediction is off. more generalizable.

3.3 Effect on Dropout

$$E[J] = E[(\hat{y} - t)^2] + \frac{1}{2} \sum_j \text{Var}[x_j] w_j^2$$

$$= E[(w_1 x_1 + w_2 x_2 - y)^2] + (w_1^2 \sigma^2 + w_2^2 (\sigma^2 + 1))$$

$$= E[w_1^2 x_1^2 + w_2^2 x_2^2 + 2w_1 w_2 x_1 x_2 - 2w_1 x_1 y - 2w_2 x_2 y] + w_1^2 \sigma^2 + w_2^2 (\sigma^2 + 1)$$

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

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

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

$$\frac{\partial E[J]}{\partial w_2} = 8w_2 \sigma^2 + 4w_2 + 2w_1 + 2 - 4\sigma^2 = 0$$

$$w_1 = \frac{2\sigma^2 + 2}{\sigma^2 + 1}$$

$$w_2 = \frac{\sigma^2}{\sigma^2 + 1}$$

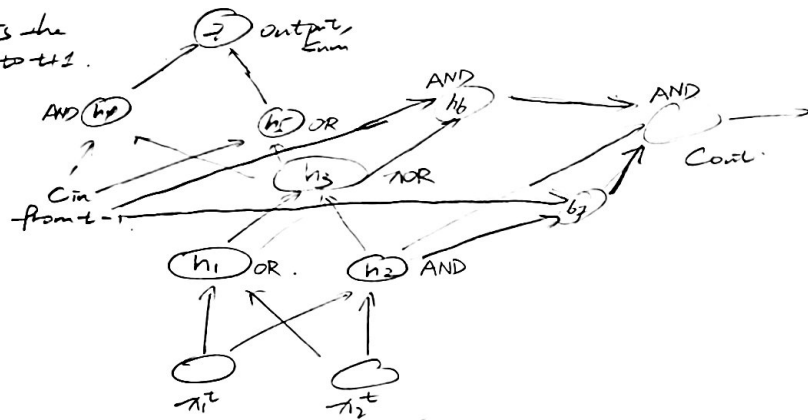
this is better because the weights don't change a lot with σ^2 . Since both the numerator and denominator are second orders, this will result in more generalizability.

4. Hard-coding Recurrent Neural Networks.

Truth table - this is the memory from $t-1$.

x_1^t	x_2^t	C_{in}^t	Sum^t	$Count^t$
0	0	0	0	0
0	0	1	1	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

this is the output.



$$Sum^t = \neg XOR(x_1^t, x_2^t, C_{in}^t)$$

$$= NOT[XOR(XOR(x_1^t, x_2^t), C_{in}^t)]$$

$$Count^t = OR[AND(x_1^t, x_2^t), AND(x_1^t, C_{in}^t), AND(x_2^t, C_{in}^t)]$$

$$= C_{in}^{t+1}$$