

#1. ① Linear-dropout - σ 이 σ 이고 ② Linear-dropout - σ 이 σ 이고

$$\textcircled{1} y = \sigma(\tilde{y}') \text{ 이 } \textcircled{2} y_i = \begin{cases} 0 & \text{with prob } p \\ \frac{\sigma(Wx+b)_i}{1-p} & \text{otherwise} \end{cases}$$

$$\left((\tilde{y}')_i = \begin{cases} 0 & \text{with prob } p \\ \frac{y_i}{1-p} = \frac{(Wx+b)_i}{1-p} & \text{otherwise} \end{cases} \right)$$

$$(y)_i \text{ 는 } \frac{1}{1-p} \text{ 이다.}$$

$$\textcircled{1} \text{의 경우, } (y)_i = \sigma((\tilde{y}')_i) = \begin{cases} 0 & \text{with probability } p \\ \sigma\left(\frac{\sum_j W_{ij} X_j + b_i}{1-p}\right) & \text{otherwise} \end{cases}$$

$$\textcircled{2} \text{의 경우, } (y)_i = \begin{cases} 0 & \text{with prob } p \\ \frac{\sigma(\sum_j W_{ij} X_j + b_i)}{1-p} & \text{otherwise} \end{cases}$$

$\frac{1}{1-p}$ ①과 ②가 같으려면, $\sigma(Cx) = C \sigma(x)$ 가 성립하면 된다. ($C > 0$)

$\frac{1}{1-p}$, $\text{ReLU}(a)$ 이 $\text{LeakyReLU}(c)$ 은 equivalent, $\text{sigmoid}(b)$ 은 2가지 있다.

#2. Pytorch = Default weight initialization 이 $k = \frac{1}{\text{in-features}}$ 라고 한다면,

A 이 b 은 $\text{Uniform}(-\sqrt{k}, \sqrt{k})$ 으로 initialize 된다고 한다. ($\frac{1}{3}$ Mean 0, Variance = $\frac{1}{3}k$)

$$x_1, \dots, x_n : \text{i.i.d.}, E[X_i] = 0, \text{Var}(X_i) = 1.$$

$$\text{Then, } y_i = A_i x + b_i, (y)_i = \sum_{j=1}^{n_0} (A_{ij})_j (x_j) + (b_i)_i$$

$$E[(y)_i] = \underbrace{0 + 0 + \dots + 0}_{n_0} + 0 = 0$$

$$\text{Var}[(y)_i] = \text{Var}\left(\sum_{j=1}^{n_0} (A_{ij})_j (x_j) + (b_i)_i\right) = \frac{1}{3} \frac{1}{n_0} + \text{Var}\left(\sum_{j=1}^{n_0} (A_{ij})_j (x_j)\right)$$

$$= \frac{1}{3n_0} + \sum_{j=1}^{n_0} \frac{1}{3} \frac{1}{n_0} \times 1 = \frac{n_0+1}{3n_0} = \frac{1}{3} \left(\frac{1}{n_0} + 1\right)$$

So we can notice that $E[X_L] = 0$

$$\text{And, } \text{Var}((y)_{i-1}) = \frac{1}{3n_{i-1}} + n_{i-1} \times \frac{1}{3n_{i-1}} \times \text{Var}((y)_{i-1})$$

$$= \frac{1}{3} \left(\frac{1}{n_{i-1}} + \text{Var}((y)_{i-1})\right)$$

$$\text{So, } \boxed{\text{Var}(X_L) = \frac{1}{3^L} + \sum_{i=0}^{L-1} \frac{1}{n_i} \frac{1}{3^{L-i}}}$$

#3 <i> $\frac{\partial y_L}{\partial y_{L-1}}$ for $L=2 \dots L$.

$$\left(\frac{\partial y_L}{\partial y_{L-1}}\right)_i = \frac{\partial y_L}{\partial (y_{L-1})_i} = \frac{\partial (A_L y_{L-1} + b_L)}{\partial (y_{L-1})_i} = (A_L)_i \quad \text{Thus, } \boxed{\frac{\partial y_L}{\partial y_{L-1}} = A_L \in \mathbb{R}^{1 \times m}}$$

$$\begin{aligned} \left(\frac{\partial y_L}{\partial y_{L-1}}\right)_{ij} &= \frac{\partial (y_L)_i}{\partial (y_{L-1})_j} = \frac{\partial \left(\nabla (A_L y_{L-1} + b_L) + y_{L-1} \right)_i}{\partial (y_{L-1})_j} = \frac{\partial \left((y_{L-1})_i + \nabla \left(\sum_{j=1}^m (A_L)_{ij} (y_{L-1})_j + (b_L)_i \right) \right)}{\partial (y_{L-1})_j} \\ &\in \mathbb{R}^{m \times m} \\ &= \begin{cases} \nabla' \left(\sum_{j=1}^m (A_L)_{ij} (y_{L-1})_j + (b_L)_i \right) \cdot (A_L)_{ij} & \text{if } i=j \\ 1 + \nabla' \left(\sum_{j=1}^m (A_L)_{ij} (y_{L-1})_j + (b_L)_i \right) \cdot (A_L)_{ij} & \text{else} \end{cases} \end{aligned}$$

$$\boxed{\therefore \left(\frac{\partial y_L}{\partial y_{L-1}}\right) = I + \text{diag} \left(\nabla' (A_L y_{L-1} + b_L) \right) A_L \quad \text{for } L=2, \dots, L-1}$$

$\in \mathbb{R}^{m \times m}$

<ii> $L=1 \dots L$.

$$\frac{\partial y_L}{\partial b_L} = \frac{\partial y_L}{\partial y_{L-1}} \cdot \frac{\partial y_{L-1}}{\partial y_{L-2}} \dots \frac{\partial y_{L-1}}{\partial y_L} \cdot \frac{\partial y_L}{\partial b_L} =$$

$$\frac{\partial y_L}{\partial b_L} = 1 \quad \left(\frac{\partial y_L}{\partial b_L}\right)_{ij} = \frac{\partial (y_L)_i}{\partial (b_L)_j} = \frac{\partial \left(\nabla (A_L y_{L-1} + b_L) + y_{L-1} \right)_i}{\partial (b_L)_j} = \begin{cases} \nabla' \left(\sum_{j=1}^m (A_L)_{ij} (y_{L-1})_j + (b_L)_i \right) & \text{if } i=j \\ 0 & \text{else} \end{cases}$$

$$\Rightarrow \frac{\partial y_L}{\partial b_L} = \text{diag} \left(\nabla' (A_L y_{L-1} + b_L) \right)$$

$$\boxed{\therefore \frac{\partial y_L}{\partial b_L} = 1 \quad \frac{\partial y_L}{\partial b_L} = \frac{\partial y_L}{\partial y_L} \left(\text{diag} \left(\nabla' (A_L y_{L-1} + b_L) \right) \right)}$$

\uparrow
 $\mathbb{R}^{1 \times m}$

* $\frac{\partial y_L}{\partial y_{L-1}}$ is calculated in <i>.

$$\Rightarrow \frac{\partial y_L}{\partial y_L} = \frac{\partial y_L}{\partial y_{L-1}} \dots \frac{\partial y_{L-1}}{\partial y_L}$$

* $\frac{\partial y_L}{\partial b_L} \in \mathbb{R}^m \Rightarrow$ 생각하라고 하면, $\text{diag} \left(\nabla' (A_L y_{L-1} + b_L) \right) \left(\frac{\partial y_L}{\partial y_L} \right)^T \ni$ 생각하라고 하면 $\frac{\partial y_L}{\partial b_L}$.

$$\frac{\partial y_L}{\partial A_L} \quad l=1 \dots L$$

$$\left(\frac{\partial y_L}{\partial A_L} \right) \in \mathbb{R}^{n_L \times n_{L-1}}, \quad \left(\frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}} \quad (\text{clarifying notation}).$$

$$\bullet \quad \frac{\partial y_L}{\partial A_L} \in \mathbb{R}^{1 \times m} \quad \left(\frac{\partial y_L}{\partial A_L} \right)_{1,j} = \frac{\partial y_L}{\partial (A_L)_{1,j}} = (y_{L-1})_j \Rightarrow \boxed{\frac{\partial y_L}{\partial A_L} = y_{L-1}^T}$$

$$\frac{\partial y_L}{\partial A_L} \in \mathbb{R}^{m \times m} \quad \left(\frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}} = \left(\frac{\partial y_L}{\partial y_{L-1}} \dots \frac{\partial y_{L+1}}{\partial y_L} \right) \frac{\partial y_L}{\partial (A_L)_{ij}}$$

$$\left(\frac{\partial y_L}{\partial A_L} \right)_{ij} = \frac{\partial y_L}{\partial (A_L)_{ij}} = \left(\frac{\partial y_L}{\partial y_{L-1}} \dots \frac{\partial y_{L+1}}{\partial y_L} \right) \frac{\partial y_L}{\partial (A_L)_{ij}}$$

$$\frac{\partial y_L}{\partial (A_L)_{ij}} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ r' \left(\sum_{k=1}^m (A_L)_{ik} (y_{L-1})_k + (b_L)_i \right) (y_{L-1})_j \\ \vdots \\ 0 \end{pmatrix} \leftarrow i\text{-번째 행}$$

$$\therefore \left(\frac{\partial y_L}{\partial A_L} \right)_{ij} = \left(\frac{\partial y_L}{\partial y_L} \right) \begin{pmatrix} 0 \\ 0 \\ \vdots \\ r' \left(\sum_{k=1}^m (A_L)_{ik} (y_{L-1})_k + (b_L)_i \right) (y_{L-1})_j \\ \vdots \\ 0 \end{pmatrix} \leftarrow i\text{-번째 행}$$

$$\Rightarrow \left(\frac{\partial y_L}{\partial A_L} \right) = \underset{\substack{\cap \\ \mathbb{R}^{m \times m}}}{\text{diag} \left(r' (A_L y_{L-1} + b_L) \right)} \underset{\substack{\cap \\ \mathbb{R}^{m \times 1}}}{\left(\frac{\partial y_L}{\partial y_L} \right)^T} \underset{\substack{\cap \\ \mathbb{R}^{1 \times m}}}{(y_{L-1})^T}$$

(iii) $\frac{\partial y_L}{\partial b_i}$ 와 $\frac{\partial y_L}{\partial A_i}$ 는 $i \neq L$ 인, $i \leq L$ $\left(\frac{\partial y_L}{\partial y_L} \right)^T$ term 이 존재함. 즉, residual equation 이 있는 경우,

$$\left(\frac{\partial y_L}{\partial y_L} \right) = (A_L) \left(I + \text{diag} \left(r' (A_{L-1} y_{L-2} + b_{L-1}) \right) A_{L-1} \right) \dots \left(I + \text{diag} \left(r' (A_{L+1} y_L + b_{L+1}) \right) A_{L+1} \right)$$

라 같이 재귀적으로 표현, A_i 는 $r' (A_i y_{i-1} + b_i) = 0$ 이 되는 상황

4대, $\text{diag} \left(r' (A_{L-1} y_{L-2} + b_{L-1}) \right) A_{L-1}$ 가 0이 되더라도 Identity Matrix가 되기 된다.

2번째, 4번째, 6번째 MLP 라는 것은, $\frac{\partial y_L}{\partial b_i}, \frac{\partial y_L}{\partial A_i}$ 가 0으로 vanish 되기 때문이다.

#4

(a) Parameter Count

$$\text{<i>Naive structure: } 128 \times (256 \times 1 \times 1 + 1) + 128 \times (128 \times 3 \times 3 + 1) + 256 \times (128 \times 1 \times 1 + 1) \\ = 213504 \text{ \>}$$

<i>Split-transform-Merge

$$: 32 \times \left(4 \times (256 \times 1 \times 1 + 1) + 4 \times (4 \times 3 \times 3 + 1) + 256 \times (4 \times 1 \times 1 + 1) \right) \\ = 78592 \text{ \>}$$

(b) Implement

* (init \>

self.convlayers1 = nn.ModuleList([nn.Conv2d(256, 4, 1) for i in range(32)])

self.convlayers2 = nn.ModuleList([nn.Conv2d(4, 4, 3, padding=1) for i in range(32)])

self.convlayers3 = nn.ModuleList([nn.Conv2d(4, 256, 1) for i in range(32)])

* forward \>

answer = torch.zeros(x.size())

for i in range(32):

out = torch.nn.functional.relu(self.convlayers1[i](x))

out = torch.nn.functional.relu(self.convlayers2[i](x))

out = torch.nn.functional.relu(self.convlayers3[i](x))

answer += out

out = answer

#5

다음과 같이 짤다

```
ddescent.py x
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5 """
6 Step 1 : Generate Toy data
7 """
8
9 d = 35
10 n_train, n_val, n_test = 300, 60, 30
11 np.random.seed(0)
12 beta = np.random.randn(d)
13 beta_true = beta / np.linalg.norm(beta)
14 # Generate and fix training data
15 X_train = np.array([np.random.multivariate_normal(np.zeros(d), np.identity(d)) for _ in range(n_train)])
16 Y_train = X_train @ beta_true + np.random.normal(loc = 0.0, scale = 0.5, size = n_train)
17 # Generate and fix validation data (for tuning lambda).
18 X_val = np.array([np.random.multivariate_normal(np.zeros(d), np.identity(d)) for _ in range(n_val)])
19 Y_val = X_val @ beta_true
20 # Generate and fix test data
21 X_test = np.array([np.random.multivariate_normal(np.zeros(d), np.identity(d)) for _ in range(n_test)])
22 Y_test = X_test @ beta_true
23
24
25 """
26 Step 2 : Solve the problem
27 """
28 def f(X,W):
29     L=[]
30     for i in range(len(X)):
31         L.append((W @ X[i]) * (W @ X[i]>0))
32     return np.array(L)
33
34 fixed_lambda = 0.01
35 lambda_list = [2 ** i for i in range(-6, 6)]
36 num_params = np.arange(1,1501,10)
37
38 errors_opt_lambda = []
39 errors_fixed_lambda = []
40 for p in num_params :
41     # fix W, calculate Xtilda based on fixed W
42     W=np.random.normal(0, 1/np.sqrt(p), size=(p,d))
43     X_tilda = f(X_train, W)
44     X_val_tilda = f(X_val, W)
45     X_test_tilda = f(X_test, W)
46
47     # theta based on fixed lambda
48     theta_fixed = np.linalg.inv(X_tilda.T @ X_tilda + fixed_lambda * np.identity(p)) @ X_tilda.T @ Y_train
49
50     # find optimal lambda using validation data
51     minloss = 1000000000000
52     theta_optimal = np.zeros(p)
53     for l in lambda_list:
54         theta = np.linalg.inv(X_tilda.T @ X_tilda + l * np.identity(p)) @ X_tilda.T @ Y_train
55         loss = np.linalg.norm(X_val_tilda@theta - Y_val)
56         if(loss < minloss):
57             minloss = loss
58             theta_optimal = theta
59
60
61     #use test set to calculate accuracy
62     fixed_loss = np.linalg.norm(X_test_tilda@theta_fixed - Y_test)
63     errors_fixed_lambda.append(fixed_loss)
64
65     optimal_loss = np.linalg.norm(X_test_tilda@theta_optimal - Y_test)
66     errors_opt_lambda.append(optimal_loss)
67
68     #just for debugging
69     print(p, " done")
70
```

```

70
71 """
72 Step 3 : Plot the results
73 """
74
75 plt.figure(figsize = (24, 8))
76 plt.rc('text', usetex = False)
77 plt.rc('font', family = 'serif')
78 plt.rc('font', size = 24)
79
80
81 plt.scatter(num_params, errors_fixed_lambda, color = 'black',
82             label = r"Test error with fixed  $\lambda = 0.01$ ",
83             )
84 plt.legend()
85
86 plt.plot(num_params, errors_opt_lambda, 'k', label = r"Test error with tuned  $\lambda$ ")
87 plt.legend()
88 plt.xlabel(r' $\lambda$  parameters')
89 plt.ylabel('Test error')
90 plt.title(r'Test error vs.  $\lambda$  params : SeoyongLee')
91
92 plt.savefig('double_descent.png')
93 plt.show()

```

그러면 똑같은 결과를 얻게 된다.

