# Module

### 1 Landscape of Top Eigenvector Problem

**Consider the following optimization problem for computing the top eigenvector of a positive semidefinite matrix $M \in \mathbb{R}^{d \times d}$ :**

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{4} \left\| xx^\top - M \right\|_F^2 .$$

**Let the eigen-decomposition of $M$ be $M = \sum_{i=1}^d \lambda_i v_i v_i^\top$ where $\{\lambda_i\}_{i=1}^d$ are eigenvalues and $\{v_i\}_{i=1}^d$ are the corresponding eigenvectors. Assume $\lambda_1 > \lambda_2 \geq \lambda_3 \geq \cdots \lambda_d \geq 0$. Prove that all saddle points and local maxima are strict, i.e., except the global minima, the Hessian matrices of all other first-order stationary points have a negative eigenvalue.**

First we calculate the gradient of $f$. We will use that to identify stationary points. We rewrite $f$ in a way that makes allows the partial derivative to be computed.

$$f(x) = \frac{1}{4} \sum_{i=1}^d \sum_{j=1}^d (x_i x_j - M_{ij})^2$$

So, using the chain rule (disclosure: classmate, Roger Wong, helped me figure out this derivative)

$$\frac{\partial f}{\partial x_i} = \frac{1}{4} \sum_{j \neq i, j=1}^d 2x_j (x_i x_j - M_{ij}) + \frac{1}{4} \sum_{j \neq i, j=1}^d 2x_j (x_i x_j - M_{ji}) + \frac{1}{4}(2)(2x_i)\left(x_i^2 - M_{ii}\right)$$

$$= \sum_{j=1}^d x_j (x_i x_j - M_{ij})$$

Finding the partial derivative of each $x_i$ leads directly to the the gradient:

$$\nabla f(x) = \begin{bmatrix} \sum_{j=1}^d x_j (x_1 x_j - M_{1j}) \\ \sum_{j=1}^d x_j (x_2 x_j - M_{2j}) \\ . \\ . \\ . \\ \sum_{j=1}^d x_j (x_d x_j - M_{dj}) \end{bmatrix} = \|x\|_2^2 x - Mx$$

The stationary points exist where $\nabla f(x) = 0$, which implies

$$\implies \|x\|_2^2 x = Mx$$

Which is the definition of an eigenvector, where $x$ is an eigenvector with corresponding eigenvalue $\|x\|_2^2$. So, stationary points of $f$ occur when $x$ is a unit eigenvector of $M$, hence we divide by the norm of $x$, yielding $x = \sqrt{\lambda}v$. Using identities on the matrix calculus wikipedia page, the Hessian of $f$ is

$$\nabla^2 f(x) = 2xx^2 + \|x\|_2^2 I - M$$

where $I$ is the identity matrix. Plugging the stationary point in

$$\nabla^2 f(x) = 2(\sqrt{\lambda}v)(\sqrt{\lambda}v)^T + (\sqrt{\lambda}v)^T (\sqrt{\lambda}v)I - M$$
$$= 2\lambda vv^t + \lambda v^T vI - M$$
$$= \lambda(2vv^T + v^T vI) - M$$

Since $\lambda$ will always be less than the largest eigenvalue in $M$, $\lambda_{min}\{\nabla^2 f(x)\} < 0$

**Problem 2.1 Show $L(w_t) \to 0$ as $t \to \infty$**

$L(w)$ can be rewritten as

$$L(w) = \frac{1}{2n} \|X^T w - y\|_2^2$$

We derive the dynamics of the $L(w)$ using the general formula given in class

We can conclude that $\lambda_{min}(H(t)) > 0$

**Problem 2.2 Show that $w$ is always in the span of $(x_1, ..., x_n)$**

We need to show that $w$ is a linear combination of the column vectors of $X$, i.e.

$$w = a_1 x_1 + ... + a_n x_n$$

Observe that each row of $H(t)$ (defined in problem 2.1) represents a linear combination of $x_i$ up to the $n$th dimension (where $n \leq d$).
Let $z = (X^T w - y)$. It is given that $\frac{dw_t}{dt} = -\nabla L(w_t)$, which implies

$$w_{t+1} = w_t + z^T H(0) z$$
$$= \sum_{i,j=1}^{n} [H(0)]_{ij} z_i z_j$$

So, each update of $w$ (start with the 0 vector, which is indeed in the span of $X$) adds a linear combination of $X$ to the previous value of $w$. Since the sum of two linear combinations $X$ is also a linear combination of $X$, $w$ is always in the span of $X$.

**Problem 2.3**

Assume, given the constraints of the problem, that $w_t$ minimizes $L(W)$ but not $\|w\|_2^2$ as $t \to \infty$ (we know $L(W)$ is minimized from the result of problem 2.1). Since $w$ is in the span of $X$, this can't be true: $w$ must also shrink is $L(W)$ shrinks, given $y_i = x_i^T w$.

**Problem 3.1**

We are given

$$k(x, x') = x^\top x' \cdot \mathbb{E}_w \left[ \sigma'(w_1) \sigma' \left( w_1 x^\top x' + w_2 \sqrt{1 - (x^\top x')^2} \right) \right]$$

Letting $x^T x = \cos\theta$ gives

$$k(x, x') = \cos\theta \cdot \mathbb{E}_w \left[ \sigma'(w_1) \sigma'(w_1 \cos\theta + w_2 \sin\theta) \right]$$

Now, let $w_1^* = w_1$ and $w_2^* = w_1 \cos\theta + w_2 \sin\theta$. Solving for $w_1$ yields $w_1 = \frac{w_2^* - w_2 \sin\theta}{\cos\theta}$, which we plug into $k$ as follows:

$$k(x, x') = \cos\theta \cdot \mathbb{E}_w \left[ \sigma'(w_1) \sigma'(w_1 \cos\theta + w_2 \sin\theta) \right]$$
$$= \cos\theta \cdot \mathbb{E}_w \left[ \sigma'(w_1^*) \sigma' \left( \frac{w_2^* - w_2 \sin\theta}{\cos\theta} \cos\theta + w_2 \sin\theta \right) \right]$$
$$= \cos\theta \cdot \mathbb{E}_w \left[ \sigma'(w_1^*) \sigma'(w_2^*) \right]$$

Since $\|x\|, \|x'\| = 1$, both $w^T x^{()}$ and $w_i$ are distributed $N(0, I)$ (and are therefore rotationally invariant), so:

$$k(x, x') = \cos\theta \cdot \mathbb{E}_w \left[ \sigma'(w_1^*) \sigma'(w_2^*) \right]$$
$$= x^T x \cdot \mathbb{E}_w \left[ \sigma'(w_1^*) \sigma'(w_2^*) \right] = x^T x \cdot \mathbb{E}_w \left[ \sigma'(w^T x) \sigma'(w^T x') \right]$$

**Problem 3.2**

Start with

$$k\left(x, x'\right) = \cos\theta \cdot \mathbb{E}_w\left[\sigma'\left(w_1\right)\sigma'\left(w_1\cos\theta + w_2\sin\theta\right)\right]$$

$$= \cos\theta \cdot \int_{-\infty}^{\infty} \sigma'\left(w_1\right)\sigma'\left(w_1\cos\theta + w_2\sin\theta\right)P(w)dw$$

$$= \frac{\cos\theta}{2\pi} \cdot \int_{0}^{\infty} w_1(w_1\cos\theta + w_2\sin\theta)(e^{\frac{-1}{2}(w_1^2+w_2^2)})dw_1dw_2$$

Now, letting $u = w_1$ and $v = w_1\cos\theta + w_2\sin\theta$...

(Bound by $\pi/2$ since this is the upper range of the domain when $\|x\| = 1$)

$$= \frac{\cos\theta}{2\pi} \cdot \int_{0}^{\pi/2} uv e^{\frac{-1}{2}(u^2+\left(\frac{v-u\cos\theta}{\sin\theta}\right)^2)}dudv$$

$$= \frac{\cos\theta}{2\pi} \cdot \int_{0}^{\pi/2} uv e^{-1(u^2+v^2-2uv\cos\theta)/2\sin^2\theta}dudv$$

Change to polar coordinates where $u = r\cos\phi, v = r\sin\phi$

$$= \frac{\cos\theta}{2\pi} \cdot \int_{0}^{\pi/2} r^2\cos\phi\sin\phi \times e^{-1(u^2+v^2-2uv\cos\theta)/2\sin^2\theta}drd\phi$$
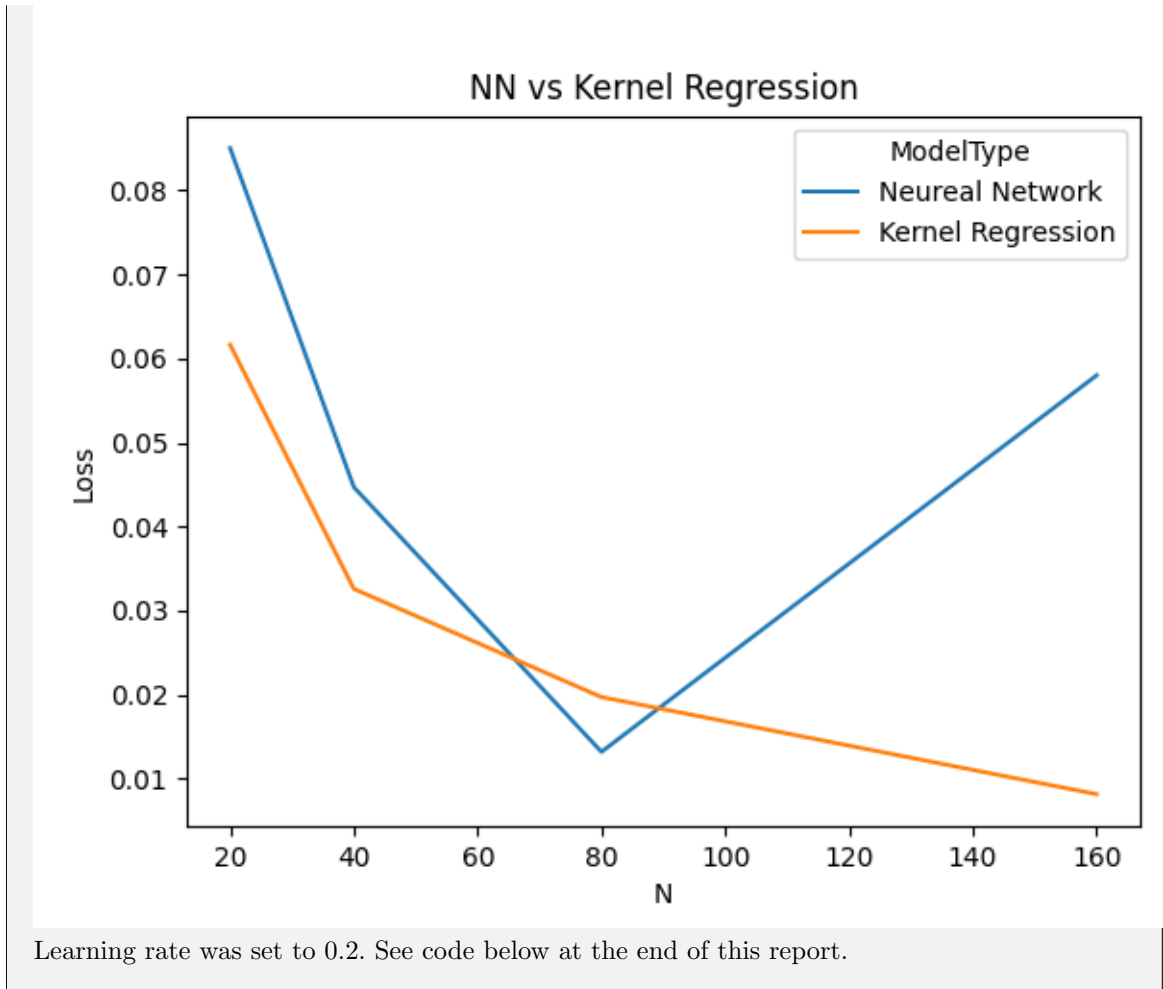
I couldn't figure out how to solve this integral, but I'm sure it evaluates to $\pi - \theta$, so we'd have

$$k\left(x, x'\right) = \frac{\cos\theta}{2\pi} \cdot \int_{0}^{\pi/2} r^2\cos\phi\sin\phi \times e^{-1(u^2+v^2-2uv\cos\theta)/2\sin^2\theta}drd\phi$$

$$= \frac{\cos\theta}{2\pi} \cdot (\pi - \theta)$$

$$= \frac{x^Tx'(\pi - \cos^{-1}(x^Tx'))}{2\pi}.$$

The last step follows from the fact that we initially set $x^Tx' = \cos\theta$

**3.3**

The neural network generalized less way to the test data, especially as sample size increased. It's likely that the NN starts to overfit (high variance) the training data as sample size increases since we have much fewer parameters than data.

NN vs Kernel Regression

Learning rate was set to 0.2. See code below at the end of this report.

**Problem 4.1: Initialization for Leaky ReLU**

Expressing $z^h$ recursively, we have $z^h = w^h x^h$ and $x^h = \sigma(z^{h-1})$. We also know that $E(z_i^h) = 0$ the weights are given to be initialized normally. We want to find a $var(z^h)$ such that

$$\text{var}(z^h) = d_h \text{var}(w^h x^h) = d_{h-1} \text{var}(w^{h-1} x^{h-1}) = \cdots = d_1 \text{var}(w^1 x^1)$$

Which implies

$$\text{var}(z^h) = d_h \text{var}(w_{ij}^h) E[(x_i^h)^2]$$

Now we rewrite $E[(x_i^h)^2]$ in terms of variance

$$
\begin{aligned}
E[(x_i^h)^2] &= E[(\sigma(z_j^{h-1})^2] \\
&= E[(\max(0, z_j^{h-1}) + \alpha \min(0, z_j^{h-1}))^2] \\
&= E[\max(0, z_j^{h-1})^2 + \max(0, z_j^{h-1}) \times \alpha\min(0, z_j^{h-1}) + \alpha^2\min(0, z_j^{h-1})^2] \\
&= E[\max(0, z_j^{h-1})^2] + E[\alpha^2\min(0, z_j^{h-1})^2] \\
&= \int_{-\infty}^{\infty} \max(0, z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} + \alpha^2 \int_{-\infty}^{\infty} \min(0, z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} \\
&= \int_{0}^{\infty} (z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} + \alpha^2 \int_{-\infty}^{0} (z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} \\
&= \frac{1}{2}\int_{-\infty}^{\infty} (z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} + \alpha^2 \frac{1}{2}\int_{-\infty}^{\infty} (z_j^{h-1})^2 P(z_j^{h-1})dz_j^{h-1} \\
&= \frac{1}{2}\mathrm{var}(z_j^{h-1}) + \frac{\alpha^2}{2}\mathrm{var}(z_j^{h-1}) \\
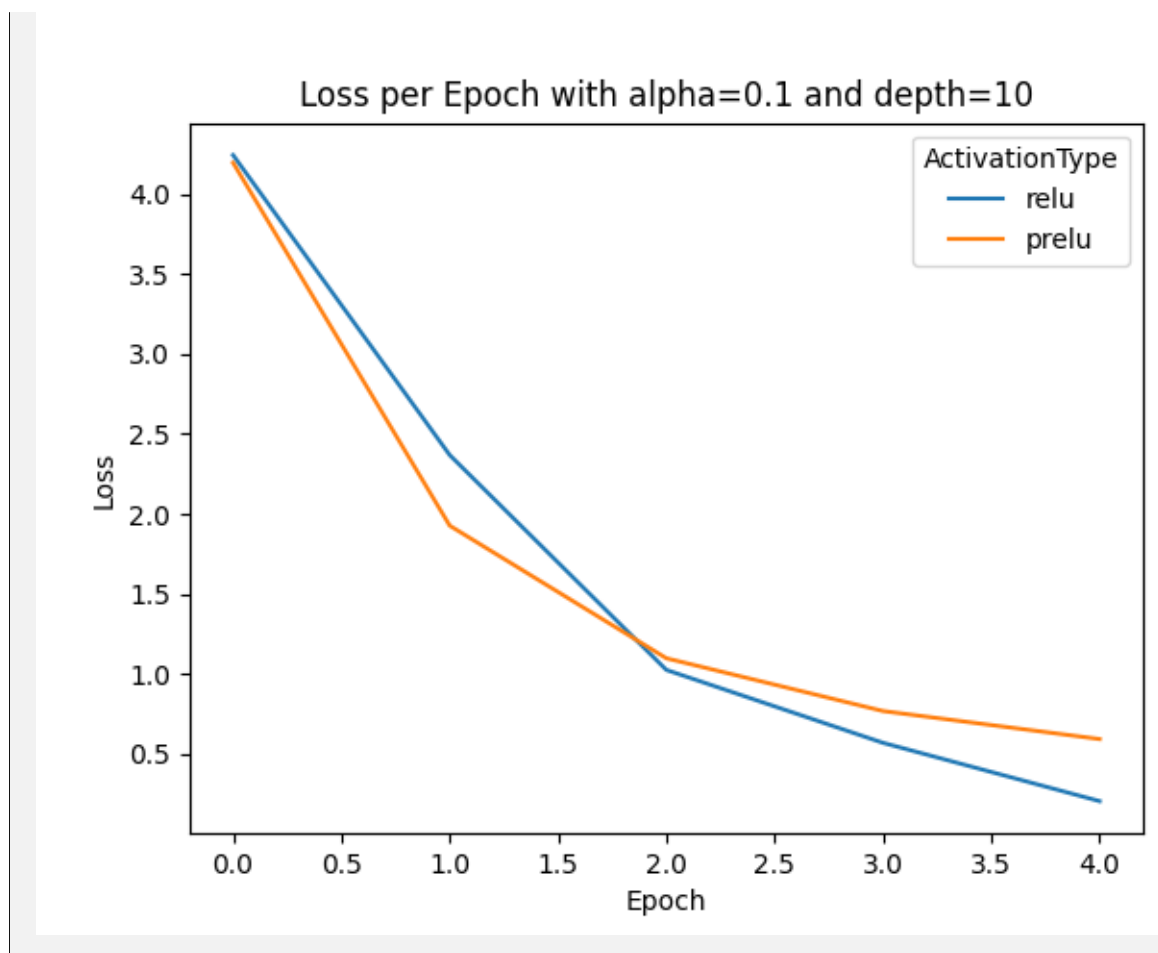&= \frac{1}{2}\mathrm{var}(z_j^{h-1})(1 + \alpha^2)
\end{aligned}
$$

Therefore

$$
\mathrm{var}(z^h) = d_h \mathrm{var}(w_{ij}^h)\frac{1}{2}\mathrm{var}(z_j^{h-1})(1 + \alpha^2)
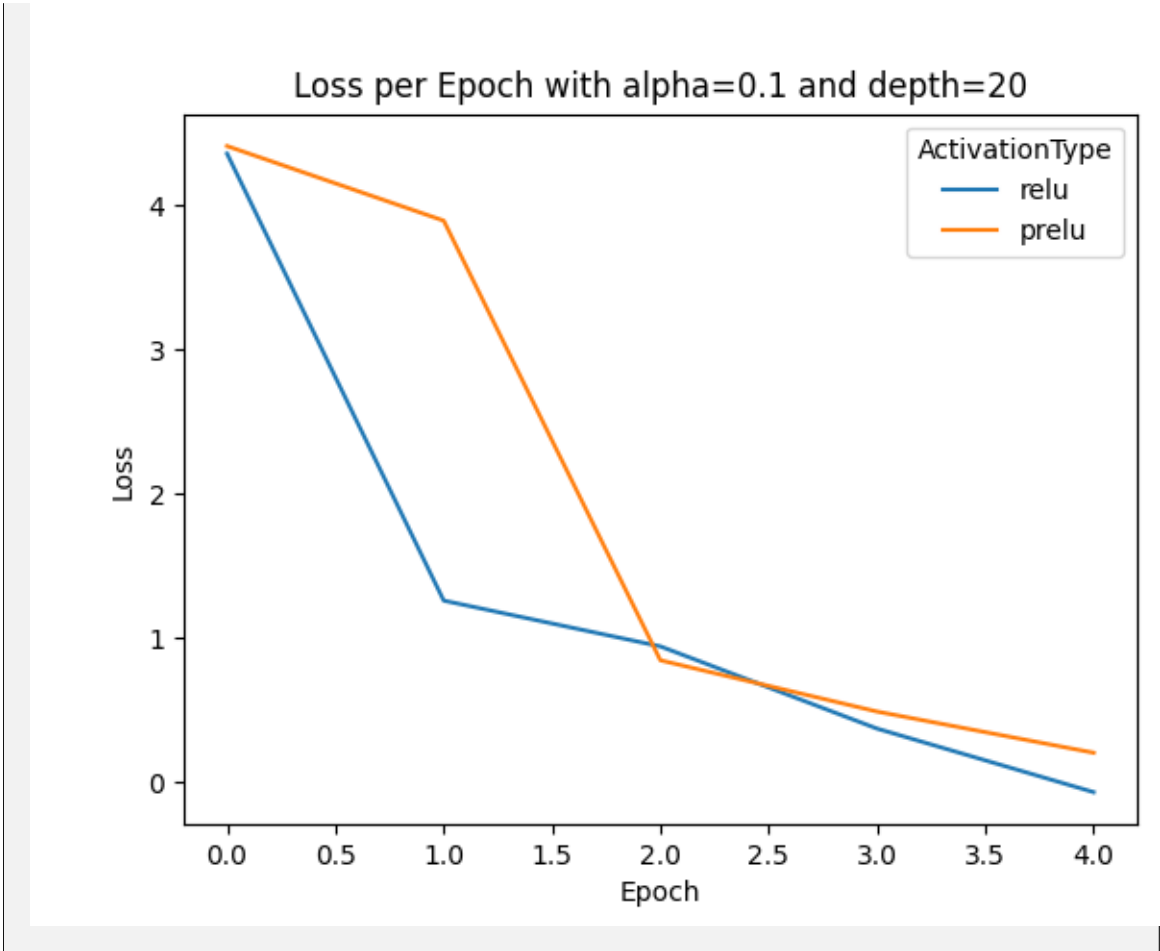$$

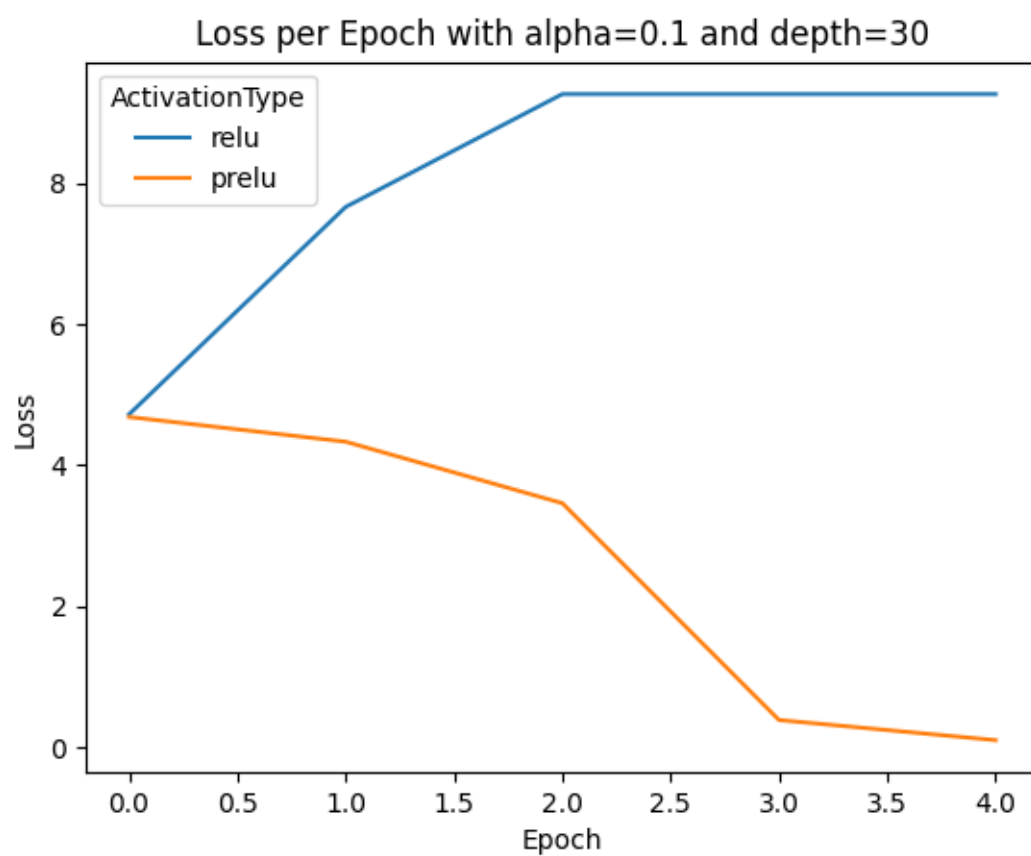And now solving for $\mathrm{var}(w_{ij}^h)$, where all $z$ has the same variance at all layers

$$
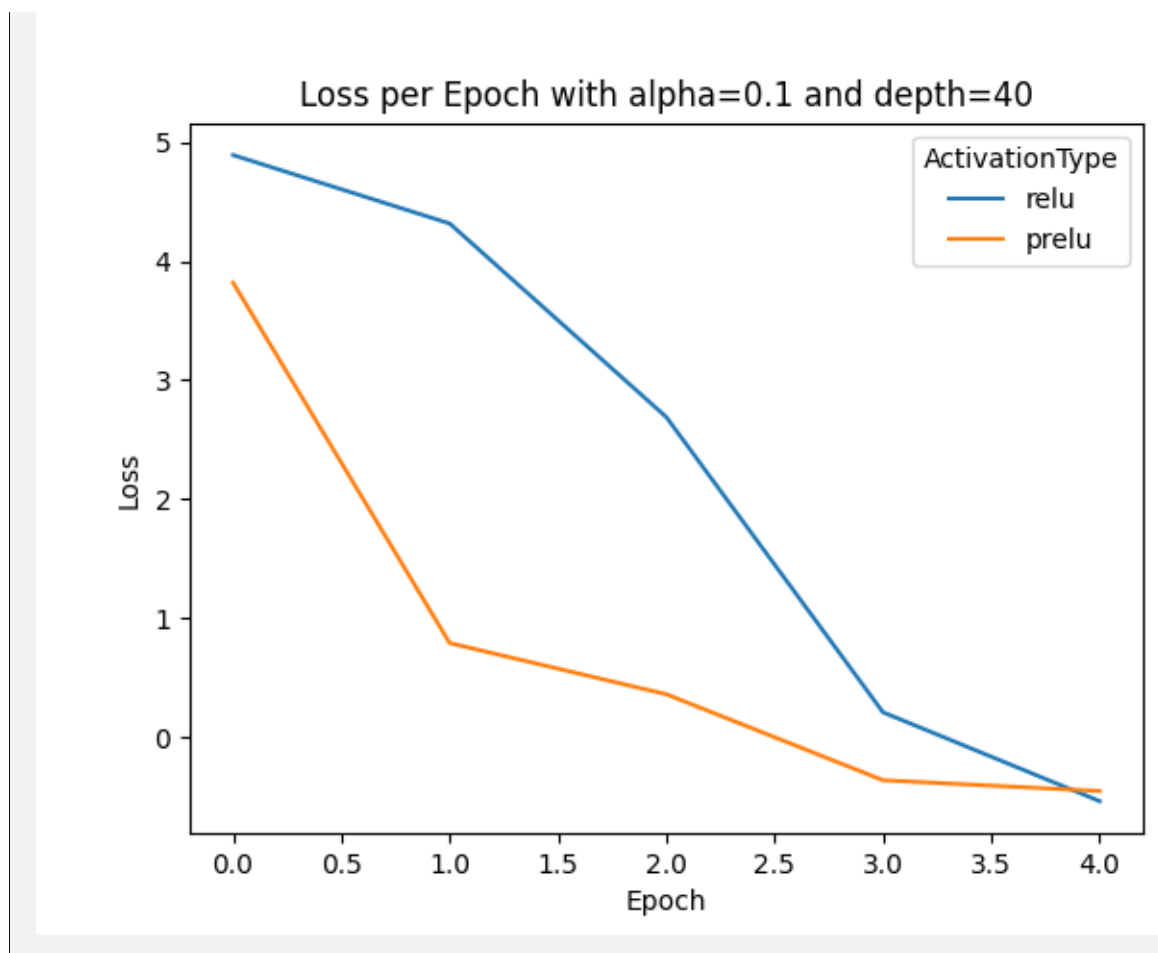\mathrm{var}(w_{ij}^h) = \beta_h = \frac{2}{d_h(1 + \alpha^2)}
$$

**Problem 4.2**

Learning rate for all computations was 0.1. All other hyper-parameters are assigned as specified in the problem statement.
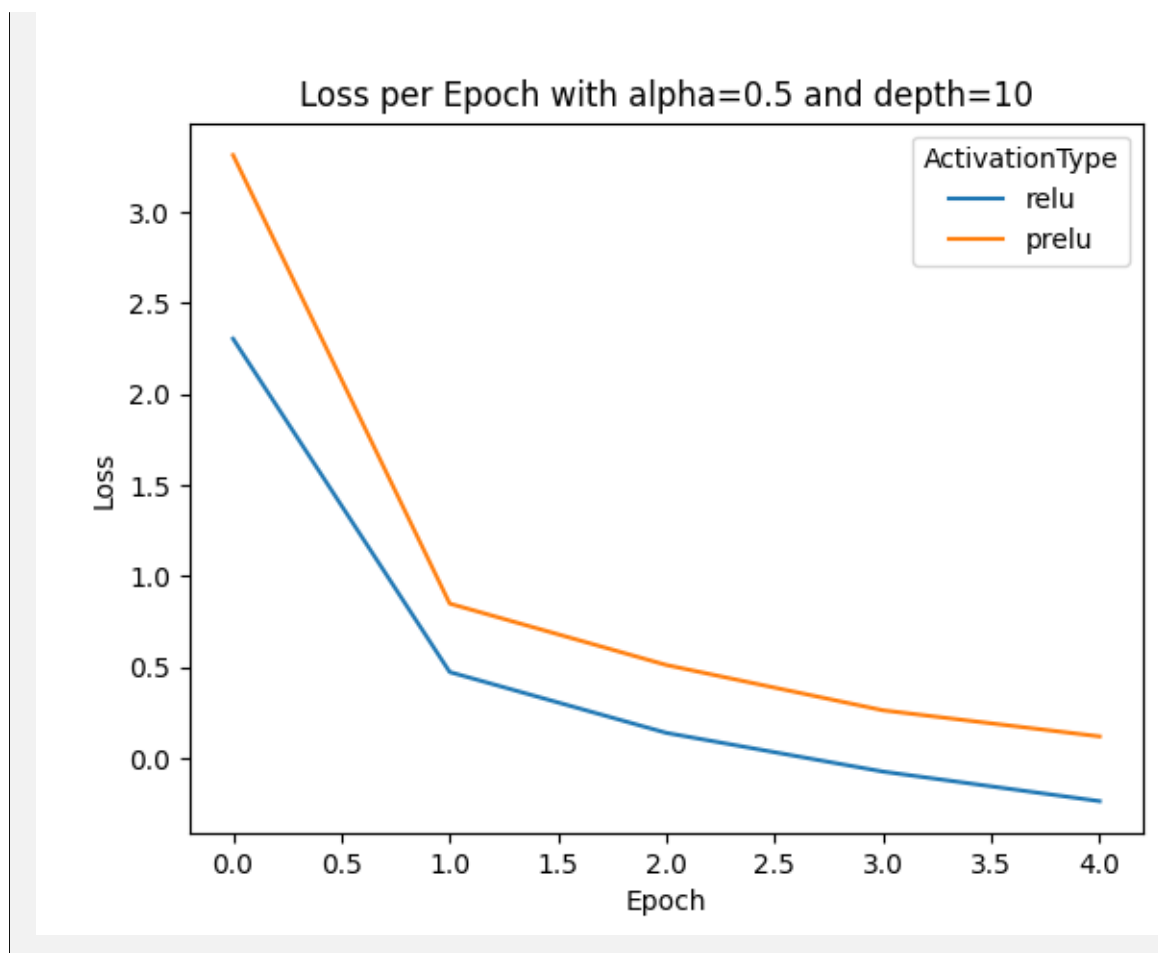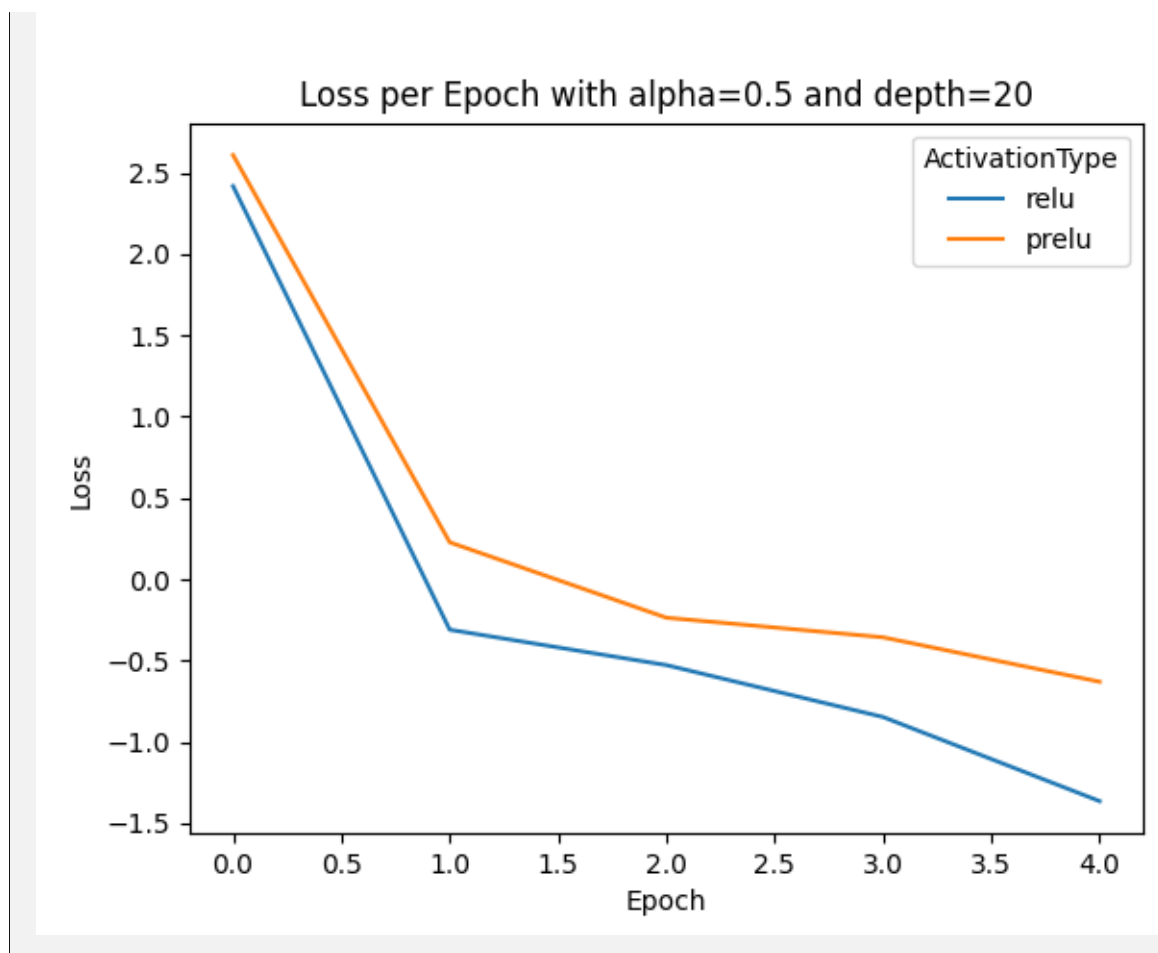
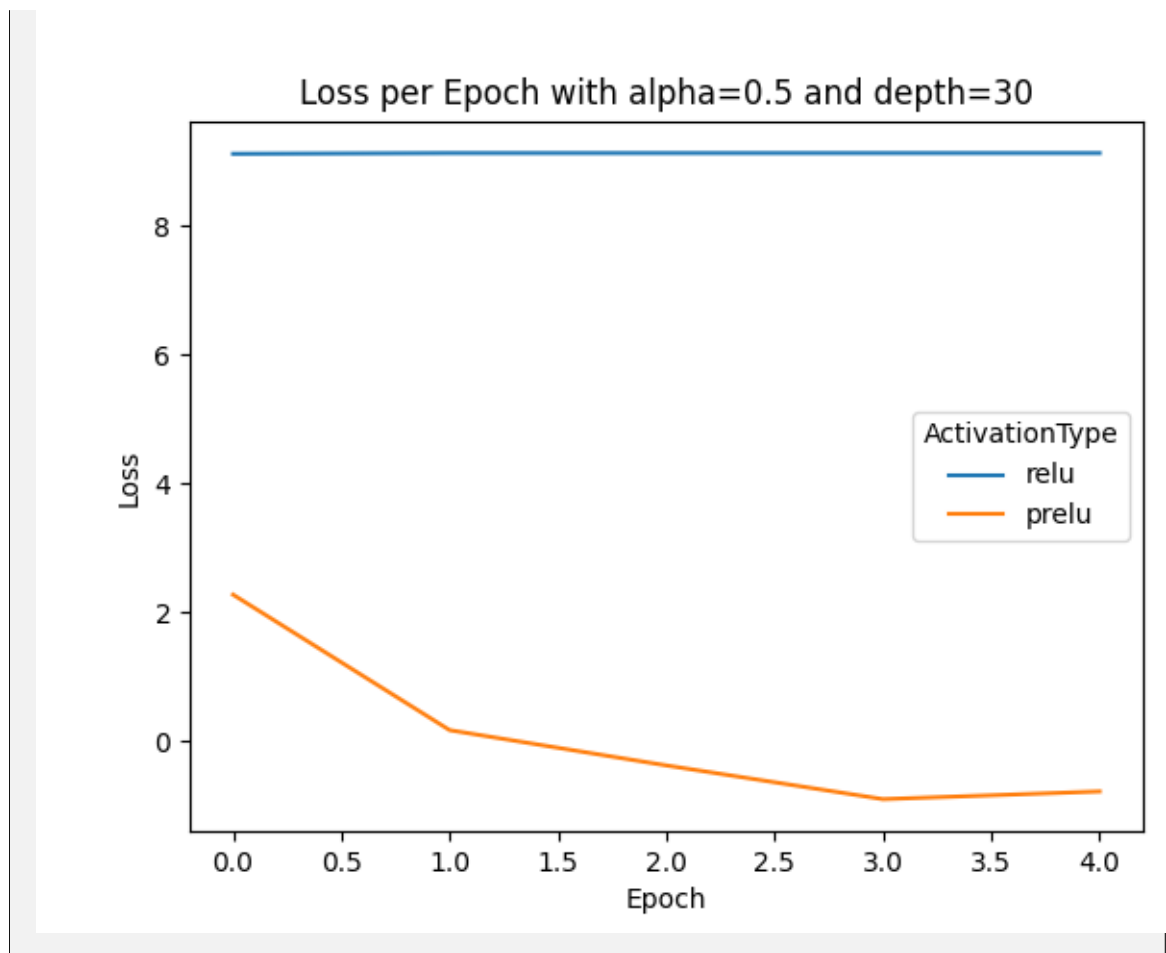Loss per Epoch with alpha=0.1 and depth=10

Loss per Epoch with alpha=0.1 and depth=20

Loss per Epoch with alpha=0.1 and depth=30

Loss per Epoch with alpha=0.1 and depth=40

Loss per Epoch with alpha=0.5 and depth=10

Loss per Epoch with alpha=0.5 and depth=20

Loss per Epoch with alpha=0.5 and depth=30

Loss per Epoch with alpha=0.5 and depth=40

Loss per Epoch with alpha=1 and depth=10

Loss per Epoch with alpha=1 and depth=20

Loss per Epoch with alpha=1 and depth=30

Loss per Epoch with alpha=1 and depth=40

Loss per Epoch with alpha=2 and depth=10

Loss per Epoch with alpha=2 and depth=20

Loss per Epoch with alpha=2 and depth=30

Loss per Epoch with alpha=2 and depth=40

```
 1  import torch
 2  import torch.nn as nn
 3  import torch.nn.functional as F
 4  import torch.optim as optim
 5
 6  from sklearn.kernel_ridge import KernelRidge
 7
 8  import numpy as np
 9  import pandas as pd
10  import seaborn as sns
11
12  DEBUG_MODE = True
13  def pp(s):
14      if(DEBUG_MODE):
15          print(s)
16
17  class Net(nn.Module):
18
19      def __init__(self, depth, width, initVariance):
20          super(Net, self).__init__()
21
22          self.preLayers = nn.ModuleList()
23          for _ in range(depth):
24              newLayer = nn.Linear(width, width)
25              nn.init.normal_(newLayer.weight, 0, np.sqrt(initVariance))
26              self.preLayers.append(newLayer)
27
28          self.outputLayer = nn.Linear(width, 1)
29
30      def forward(self, x):
31          pp(f"x dims: {x.shape}")
32          for l in self.preLayers:
33              x = F.relu(l(x))
34
35          y = F.relu(self.outputLayer(x))
36          return y
37
38
39      def fit(self, xs, ys, learningRate, nEpics=5): # TODO: change back to 10
40
41          optimizer = optim.SGD(self.parameters(), lr=learningRate)
42          lossFn = nn.MSELoss()
43
44          lossPerEpoch = []
45
46          for _ in range(nEpics):
47              self.train()
48
49              optimizer.zero_grad()
50              y_hat = self(xs)
```
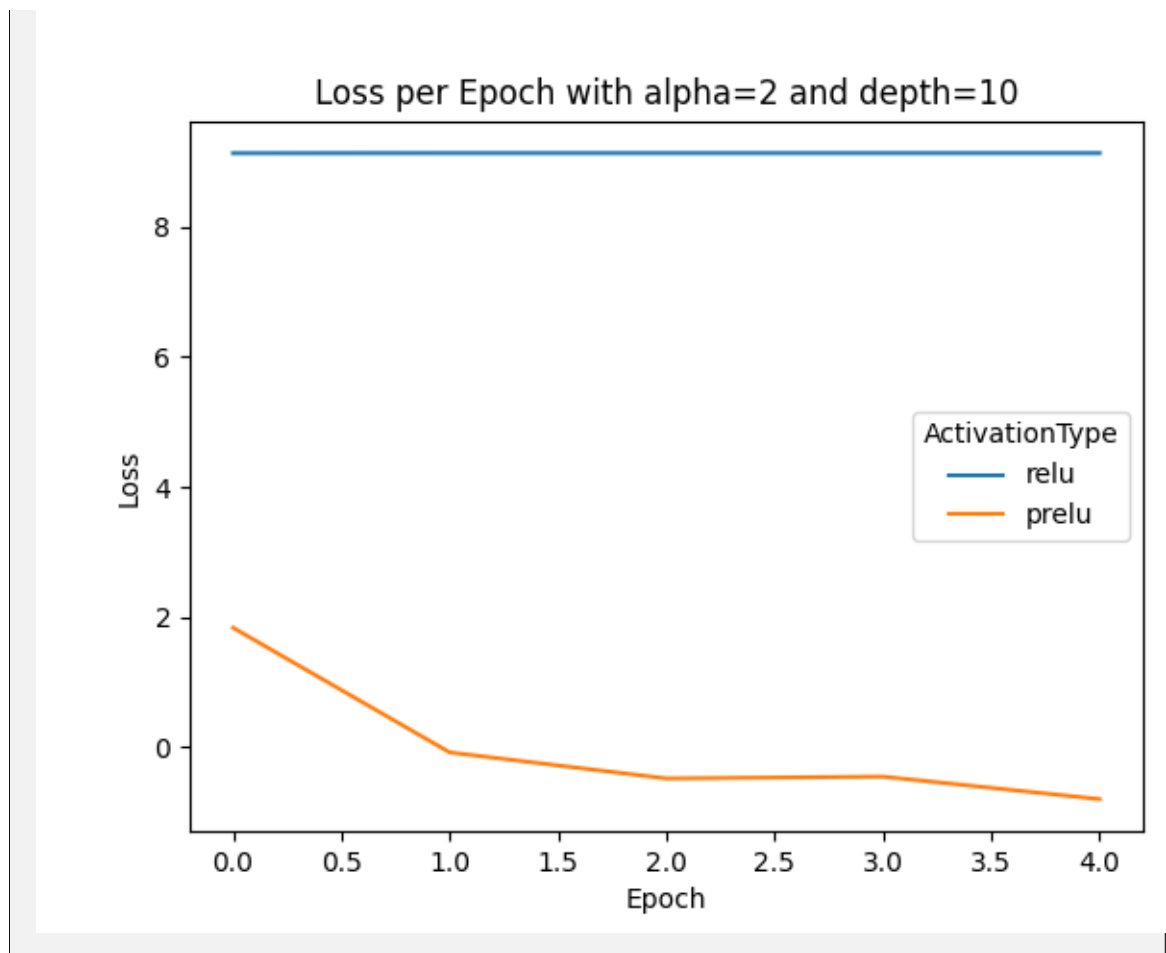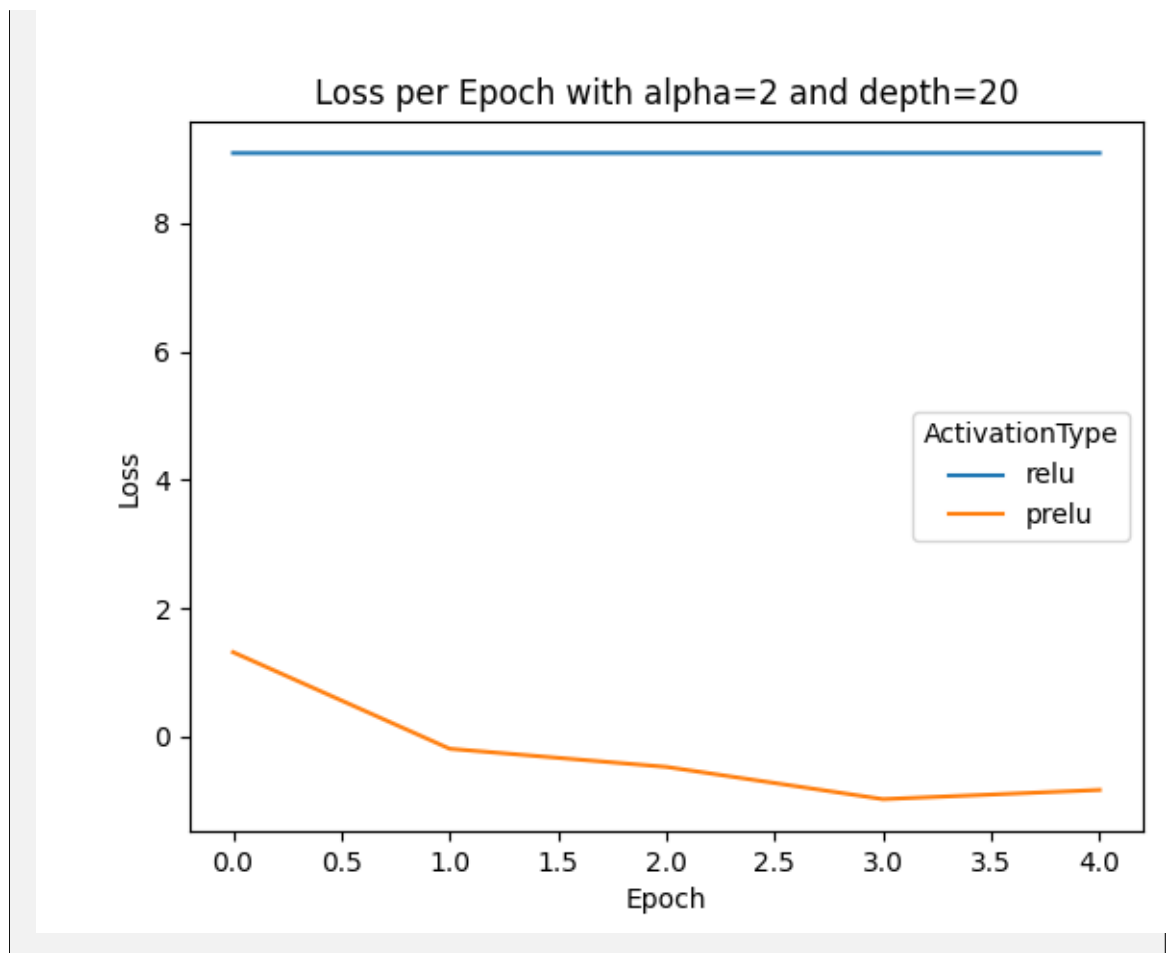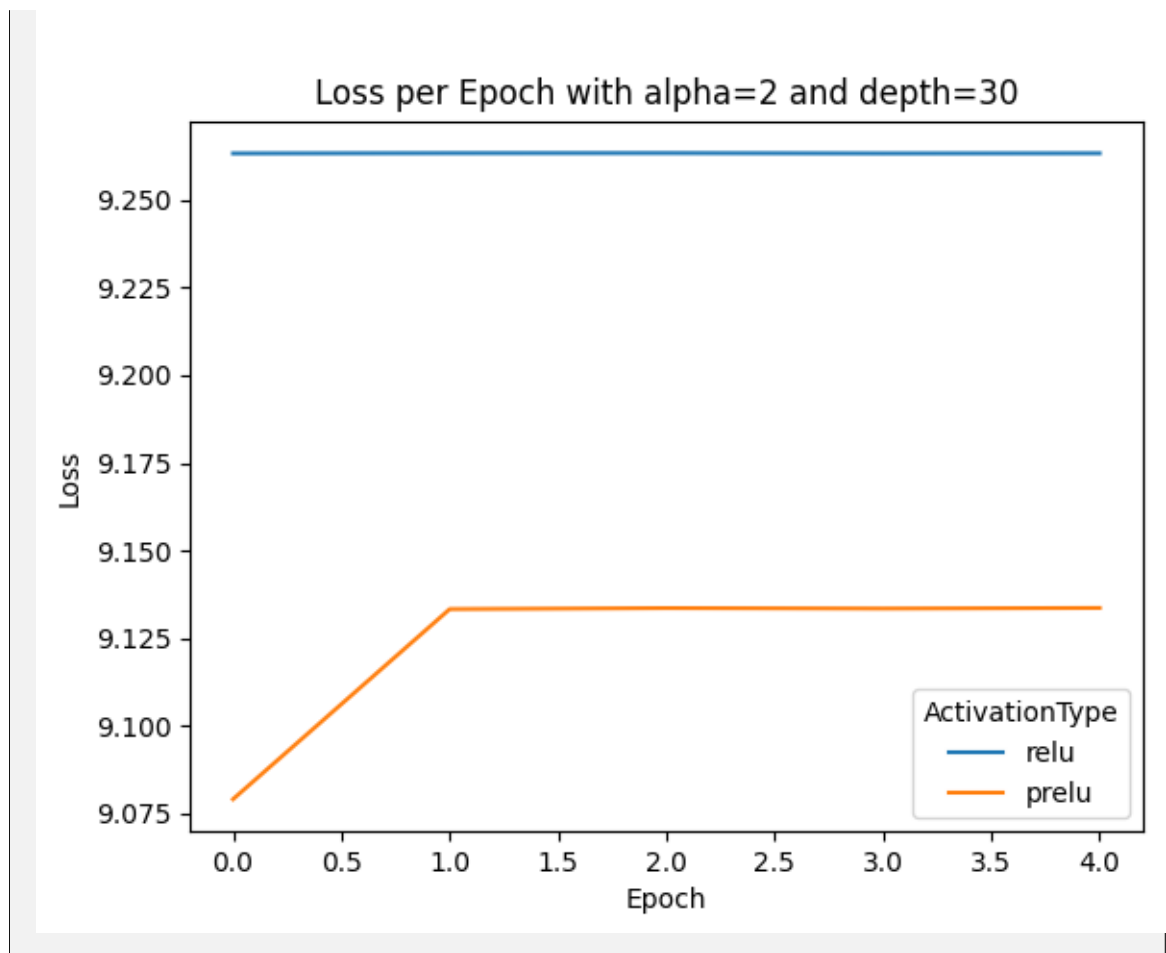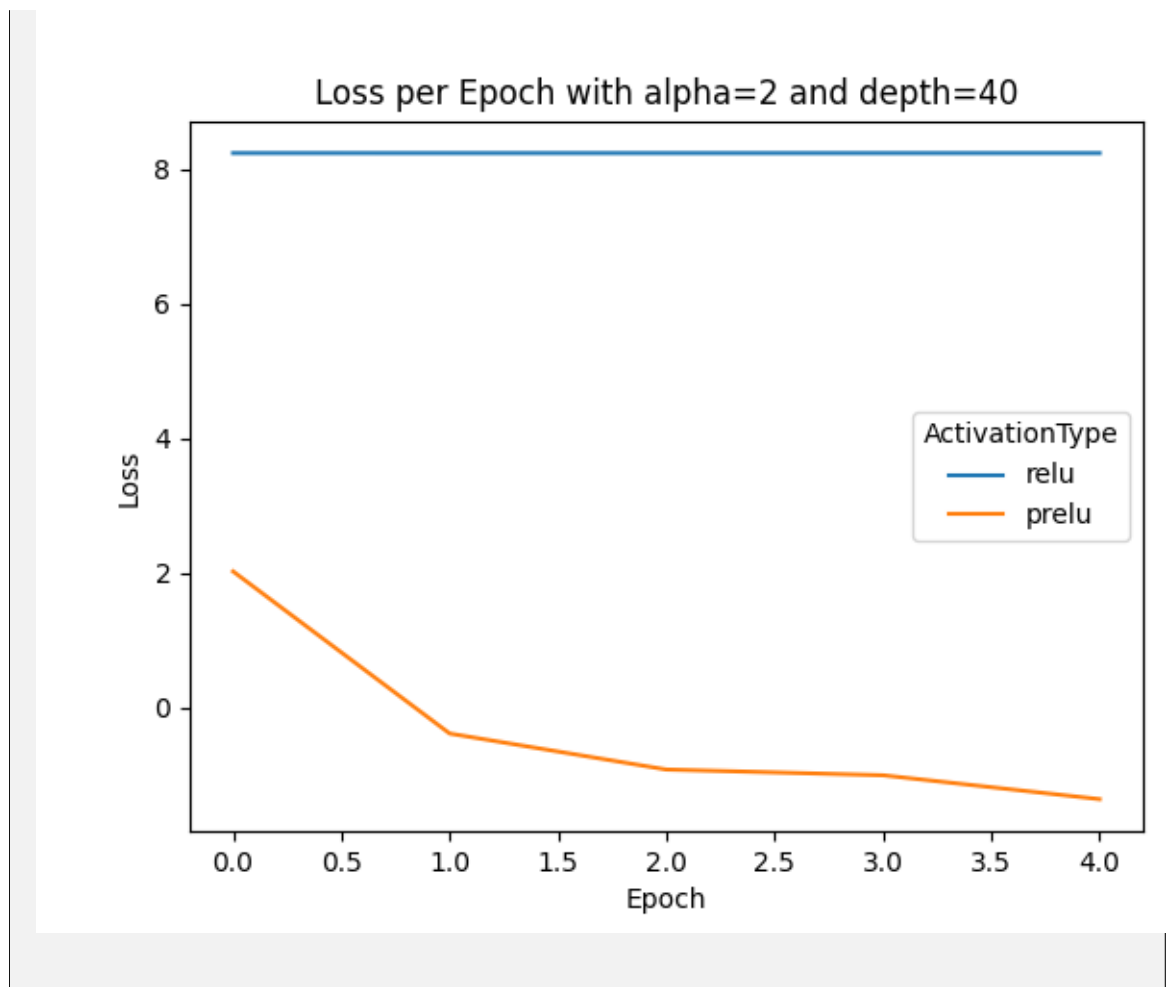
```python
51                    pp(f"y_hats : {y_hat}")
52                    pp(f"ys : {ys}")
53                    loss_output = lossFn(y_hat, ys)
54                    loss_output.backward()
55                    optimizer.step()
56                    lossPerEpoch.append(loss_output.item()/ys.shape[0])
57
58              return lossPerEpoch
59
60  def genData(n, m, d):
61      def relu(num):
62          return 0 if num < 0 else num
63      def genNSphere(r, d):
64          v = np.random.normal(0, r, d)
65          d = np.sum(v**2) **(0.5)
66          return v/d
67
68      xss = [genNSphere(1, d) for _ in range(n)]
69      ys = []
70      for xs in xss:
71          ys.append(sum(map(relu, xs))/m)
72
73      return (np.array(xss), np.array(ys))
74
75  def calcError(preds, targets):
76      pp(f"predictions: {preds}")
77      pp(f"targets: {targets}")
78      return ((preds - targets)**2).mean()
79
80  def arcKernel(x1: np.array, x2: np.array):
81      x1Tx2 = np.dot(x1, x2)
82      x1Tx2 = 1.0 if x1Tx2 > 1.0 else x1Tx2
83      x1Tx2 = -1.0 if x1Tx2 < -1.0 else x1Tx2
84      return x1Tx2*(np.pi - np.arccos(x1Tx2))/(2*np.pi)
85
86  def plotResults(df):
87      sns.color_palette("Set2")
88
89      #df['Loss'] = np.log(df['Loss'].astype(float))
90
91      title = f"NN vs Kernel Regression"
92      plot = sns.lineplot(
93          data=df, x="N", y="Loss", hue="ModelType", ci=None
94      ).set_title(title)
95      plot.figure.savefig(f"{title}.png")
96      plot.figure.clf()
97
98  def addHistoryRow(df, n, loss, modelType):
99      print(df.head())
100     new_row = pd.DataFrame(columns=df.columns)
```

```python
101        new_row.loc[0] = [n, loss, modelType]
102        df = pd.concat([df, new_row], ignore_index=True)
103        return df
104
105 ############ BEGIN APPLICATION ############
106
107 ns = [20, 40, 80, 160]
108 d = 10
109 m = 5
110 lr = 0.1
111
112 historyDf = pd.DataFrame(columns=[
113        "N", "Loss", "ModelType"
114 ])
115
116 for n in ns:
117        trainXs, trainYs = genData(n, m, d)
118        testXs, testYs = genData(100, m, d)
119
120        #pp(trainXs)
121        #pp(trainYs)
122        net = Net(depth=m, width=d, initVariance=1/m)
123        lossPerEpoch = net.fit(torch.Tensor(trainXs), torch.Tensor(trainYs), lr)
124        pp(f"Loss per eopch: {lossPerEpoch}")
125        nnError = calcError(net(torch.Tensor(testXs)).detach().numpy().flatten(), testYs)
126
127        historyDf = addHistoryRow(historyDf, n, nnError, "Neureal Network")
128
129        krr = KernelRidge(kernel=arcKernel)
130        print(trainXs)
131        krr.fit(trainXs, trainYs)
132        krrError = calcError(krr.predict(testXs), testYs)
133        historyDf = addHistoryRow(historyDf, n, krrError, "Kernel Regression")
134
135
136 plotResults(historyDf)
```

```python
1  import torch
2  import torch.nn as nn
3  import torch.nn.functional as F
4  import torch.optim as optim
5  from torch.utils.data import DataLoader, Subset
6
7  import torchvision
8  from torchvision.datasets import MNIST
9  from torchvision import transforms
10
11 import numpy as np
12 import pandas as pd
13 import seaborn as sns
14
15 DEBUG_MODE = True
16 def pp(s):
17     if(DEBUG_MODE):
18         print(s)
19
20 class Net(nn.Module):
21
22     def __init__(self, depth, alpha, width, initVariance):
23         super(Net, self).__init__()
24         self.alpha = alpha
25
26         #self.conv1 = nn.Conv2d(1, 6, 5)
27         #self.conv2 = nn.Conv2d(6, 16, 5)
28
29         self.preLayers = nn.ModuleList()
30         for _ in range(depth):
31             newLayer = nn.Linear(width, width)
32             #pp(f"initVariance: {initVariance}")
33             nn.init.normal_(newLayer.weight, 0, np.sqrt(initVariance))
34             self.preLayers.append(newLayer)
35
36         self.outputLayer = nn.Linear(width, 1)
37         self.outSig = nn.Sigmoid()
38
39     def forward(self, x):
40
41         x = torch.flatten(x, 2, 3)
42
43         for l in self.preLayers:
44             x = F.leaky_relu(l(x), self.alpha)
45
46         x = F.leaky_relu(self.outputLayer(x), self.alpha)
47
48         return torch.squeeze(self.outSig(x))
49
50
```

```python
51      def fit(self, batches, learningRate, epochs_n=5): # TODO: change back to 10
52          loss_per_epoch = []
53
54          optimizer = optim.SGD(self.parameters(), lr=learningRate)
55          lossFn = nn.BCELoss()
56
57          for i in range(epochs_n):
58              batch_loss = 0
59
60              self.train()
61              for image, label in batches:
62                  optimizer.zero_grad()
63                  y_hat = self(image)
64                  #pp(f"y_hat: {y_hat}" )
65                  #pp(f"label: {label}" )
66                  loss = lossFn(y_hat, label.float())
67                  loss.backward()
68                  optimizer.step()
69                  #pp(f"loss: {loss}")
70                  batch_loss += loss.item()
71
72              loss_per_epoch.append(batch_loss)
73
74          return loss_per_epoch
75
76 # setup data
77 trainDataRaw = MNIST(
78     root='data',
79     train=True,
80     transform=transforms.Compose([
81         transforms.Resize(16),
82         transforms.ToTensor(),
83     ]),
84     download=True
85 )
86
87
88 def addToDf(df: pd.DataFrame, history, actType):
89     for i, loss in enumerate(history):
90         new_row = pd.DataFrame(columns=df.columns)
91         new_row.loc[0] = [i, loss, actType]
92         df = pd.concat([df, new_row], ignore_index=True)
93
94     return df
95
96 def plotResults(df, d, a, lr):
97     sns.color_palette("Set2")
98
99     df['Loss'] = np.log(df['Loss'].astype(float))
100
```

```
101        title = f"Loss per Epoch with alpha={a} and depth={d}"
102        plot = sns.lineplot(
103            data=df, x="Epoch", y="Loss", hue="ActivationType", ci=None
104        ).set_title(title)
105        plot.figure.savefig(f"{title}_{d}_{a}.png")
106        plot.figure.clf()
107
108 ########### DO THE STUFF ###########
109 batchSize = 64
110 width = 256
111 depths = [10, 20, 30, 40]
112 alphas = [2, 1, 0.5, 0.1]
113
114 keepIdxs = (trainDataRaw.targets==0) | (trainDataRaw.targets==1)
115 trainDataRaw.targets = trainDataRaw.targets[keepIdxs]
116 trainDataRaw.data = trainDataRaw.data[keepIdxs]
117
118 trainBatches = DataLoader(trainDataRaw, batch_size=batchSize, shuffle=True)
119
120 initFuncs = {
121     "relu": (lambda alpha: 2/width),
122     "prelu": (lambda alpha: 2/(width*(1 + alpha**2))),
123 }
124
125 history = []
126 learningRate = 0.01
127
128 for depth in depths:
129     for alpha in alphas:
130
131         historyDf = pd.DataFrame(columns=[
132             "Epoch", "Loss", "ActivationType"
133         ])
134
135         for activationType, initFunc in initFuncs.items():
136             net = Net(depth, alpha, width, initFunc(alpha))
137             losses = net.fit(trainBatches, learningRate, epochs_n=5)
138             pp(f"Losses: {losses}")
139             historyDf = addToDf(historyDf, losses, activationType) # TODO
140
141         print(historyDf.head())
142         plotResults(historyDf, d=depth, a=alpha, lr=learningRate)
```