

Pytorch and Titanic

Now that you've seen the very basics of how to use `pytorch`, we're going to see how to apply it to a machine learning problem. Along the way we'll make connections back to the Titanic dataset and help solidify your understanding of the connection between the math we've been learning and the code we'll be writing using `pytorch`.

To get started, let's load our trusty Titanic dataset.

In [1]:

```
import gdown
import numpy as np
import pandas as pd
```

```
gdown.download('https://drive.google.com/uc?authuser=0&id=1XIFiL3WxxR6M2nWgADi3xWvuRO6A-Ov8&export=download', 'titanic_train.csv', False)
df = pd.read_csv('titanic_train.csv')
df
```

```
Downloading...
From: https://drive.google.com/uc?authuser=0&id=1XIFiL3WxxR6M2nWgADi3xWvuRO6A-Ov8&export=download
To: /Users/pruvolo/Documents/assignments/Module 1/07/titanic_train.csv
100%|██████████| 61.2k/61.2k [00:00<00:00, 2.57MB/s]
```

Out[1]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
5	6	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	NaN	Q
6	7	0	1	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S
7	8	0	3	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	NaN	S
8	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	NaN	S
9	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	NaN	C
10	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 9549	16.7000	G6	S
11	12	1	1	Bonnell, Miss. Elizabeth	female	58.0	0	0	113783	26.5500	C103	S
12	13	0	3	Saunderscock, Mr. William Henry	male	20.0	0	0	A/5. 2151	8.0500	NaN	S
13	14	0	3	Andersson, Mr. Anders Johan	male	39.0	1	5	347082	31.2750	NaN	S
14	15	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	14.0	0	0	350406	7.8542	NaN	S
15	16	1	2	Hewlett, Mrs. (Mary D Kingcome)	female	55.0	0	0	248706	16.0000	NaN	S
16	17	0	3	Rice, Master. Eugene	male	2.0	4	1	382652	29.1250	NaN	Q
17	18	1	2	Williams, Mr. Charles Eugene	male	NaN	0	0	244373	13.0000	NaN	S
18	19	0	3	Vander Planke, Mrs. Julius	female	31.0	1	0	345763	18.0000	NaN	S

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
19	0	3	(Emelia Maria Vande... Masseimani, Mrs. Fatima	female	NaN	0	0	2649	7.2250	NaN	C
20	0	2	Fynney, Mr. Joseph J	male	35.0	0	0	239865	26.0000	NaN	S
21	1	2	Beesley, Mr. Lawrence	male	34.0	0	0	248698	13.0000	D56	S
22	1	3	McGowan, Miss. Anna "Annie"	female	15.0	0	0	330923	8.0292	NaN	Q
23	1	1	Sloper, Mr. William Thompson	male	28.0	0	0	113788	35.5000	A6	S
24	0	3	Palsson, Miss. Torborg Danira	female	8.0	3	1	349909	21.0750	NaN	S
25	1	3	Asplund, Mrs. Carl Oscar (Selma Augusta Emilia...)	female	38.0	1	5	347077	31.3875	NaN	S
26	0	3	Emir, Mr. Farred Chehab	male	NaN	0	0	2631	7.2250	NaN	C
27	0	1	Fortune, Mr. Charles Alexander	male	19.0	3	2	19950	263.0000	C23 C25 C27	S
28	1	3	O'Dwyer, Miss. Ellen "Nellie"	female	NaN	0	0	330959	7.8792	NaN	Q
29	0	3	Todoroff, Mr. Lalio	male	NaN	0	0	349216	7.8958	NaN	S
...
861	0	2	Giles, Mr. Frederick Edward	male	21.0	1	0	28134	11.5000	NaN	S
862	1	1	Swift, Mrs. Frederick Joel (Margaret Welles Ba...)	female	48.0	0	0	17466	25.9292	D17	S
863	0	3	Sage, Miss. Dorothy Edith "Dolly"	female	NaN	8	2	CA. 2343	69.5500	NaN	S
864	0	2	Gill, Mr. John William	male	24.0	0	0	233866	13.0000	NaN	S
865	1	2	Bystrom, Mrs. (Karolina)	female	42.0	0	0	236852	13.0000	NaN	S
866	1	2	Duran y More, Miss. Asuncion	female	27.0	1	0	SC/PARIS 2149	13.8583	NaN	C
867	0	1	Roebing, Mr. Washington Augustus II	male	31.0	0	0	PC 17590	50.4958	A24	S
868	0	3	van Melkebeke, Mr. Philemon	male	NaN	0	0	345777	9.5000	NaN	S
869	1	3	Johnson, Master. Harold Theodor	male	4.0	1	1	347742	11.1333	NaN	S
870	0	3	Balkic, Mr. Cerin	male	26.0	0	0	349248	7.8958	NaN	S
871	1	1	Beckwith, Mrs. Richard Leonard (Sallie Monypeny)	female	47.0	1	1	11751	52.5542	D35	S
872	0	1	Carlsson, Mr. Frans Olof	male	33.0	0	0	695	5.0000	B51 B53 B55	S
873	0	3	Vander Cruyssen, Mr. Victor	male	47.0	0	0	345765	9.0000	NaN	S
874	1	2	Abelson, Mrs. Samuel (Hannah Wizosky)	female	28.0	1	0	P/PP 3381	24.0000	NaN	C
875	1	3	Najib, Miss. Adele Kiamie "Jane"	female	15.0	0	0	2667	7.2250	NaN	C
876	0	3	Gustafsson, Mr. Alfred Ossian	male	20.0	0	0	7534	9.8458	NaN	S
877	0	3	Petroff, Mr. Nedelio	male	19.0	0	0	349212	7.8958	NaN	S
878	0	3	Laleff, Mr. Kristo	male	NaN	0	0	349217	7.8958	NaN	S
879	1	1	Potter, Mrs. Thomas Jr (Lily Alexenia Wilson)	female	56.0	0	1	11767	83.1583	C50	C
880	1	2	Shelley, Mrs. William (Imanita Parrish Hall)	female	25.0	0	1	230433	26.0000	NaN	S
881	0	3	Markun, Mr. Johann	male	33.0	0	0	349257	7.8958	NaN	S
882	0	3	Dahlberg, Miss. Gerda Ulrika	female	22.0	0	0	7552	10.5167	NaN	S
883	0	2	Banfield, Mr. Frederick James	male	28.0	0	0	C.A./SOTON 34068	10.5000	NaN	S
884	0	3	Sutehall, Mr. Henry Jr	male	25.0	0	0	SOTON/OQ 392076	7.0500	NaN	S
885	0	3	Rice, Mrs. William (Margaret Norton)	female	39.0	0	5	382652	29.1250	NaN	Q

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

In the Assignment 6 companion notebook we fit two models:

1. We used the `Age` and `Sex` columns along with a synthetic feature called `is_young_male` as inputs to a logistic regression model in order to predict whether someone survived.
2. We used just the `Age` and `Sex` as inputs to a multilayer perceptron model in order to predict whether someone survived.

In this notebook we'll be implementing both of these models in `pytorch`. To get started, let's perform the following data processing / cleaning steps.

1. Get rid of any passengers where we don't know their age (don't do this in a real machine learning application as it will skew your results).
2. Convert the `Sex` column to a dummy variable called `male` that will take on value 1 if the passenger is male and 0 otherwise.
3. Create the `is_young_male` column that will be 1 for males under the age of 5 and 0 for everyone else.

In [2]:

```
# get rid of null values for age since this is just an illustrative example.
# this would not be a good thing to do if we were trying to evaluate the
# performance of a model.
df_filtered = df[['Age', 'Sex', 'Survived']].dropna()
is_young_male = ((df_filtered['Sex'] == 'male') & (df_filtered['Age'] < 5)).astype(int)
is_young_male.name = 'is_young_male'
experiment_1_data = pd.concat((pd.get_dummies(df_filtered['Sex'], drop_first=True),
df_filtered['Age'], is_young_male), axis=1)
experiment_1_outputs = df_filtered['Survived']
experiment_1_data
```

Out[2]:

	male	Age	is_young_male
0	1	22.0	0
1	0	38.0	0
2	0	26.0	0
3	0	35.0	0
4	1	35.0	0
6	1	54.0	0
7	1	2.0	1
8	0	27.0	0
9	0	14.0	0
10	0	4.0	0
11	0	58.0	0
12	1	20.0	0
13	1	39.0	0
14	0	14.0	0
15	0	55.0	0
16	1	2.0	1
18	0	31.0	0
20	1	35.0	0

21	1	34.0	0
male	Age	is_young_male	
22	0	15.0	0
23	1	28.0	0
24	0	8.0	0
25	0	38.0	0
27	1	19.0	0
30	1	40.0	0
33	1	66.0	0
34	1	28.0	0
35	1	42.0	0
37	1	21.0	0
38	0	18.0	0
...
856	0	45.0	0
857	1	51.0	0
858	0	24.0	0
860	1	41.0	0
861	1	21.0	0
862	0	48.0	0
864	1	24.0	0
865	0	42.0	0
866	0	27.0	0
867	1	31.0	0
869	1	4.0	1
870	1	26.0	0
871	0	47.0	0
872	1	33.0	0
873	1	47.0	0
874	0	28.0	0
875	0	15.0	0
876	1	20.0	0
877	1	19.0	0
879	0	56.0	0
880	0	25.0	0
881	1	33.0	0
882	0	22.0	0
883	1	28.0	0
884	1	25.0	0
885	0	39.0	0
886	1	27.0	0
887	0	19.0	0
889	1	26.0	0
890	1	32.0	0

714 rows × 3 columns

Next, we'll go ahead and build our logistic regression model just as we did in the assignment 6 companion notebook. The only small twist we will introduce is turning off the ridge term of the model so that it will make comparing the results from this analysis to what we do in pytorch easier.

In [3]:

```
from sklearn.linear_model import LogisticRegression

model = LogisticRegression(solver='lbfgs', penalty='none') # setting solve silences annoying
warning
model.fit(experiment_1_data, experiment_1_outputs)
print("coefs", model.coef_)
print("intercept", model.intercept_)

coefs [[-2.67342472  0.00808732  2.38307187]]
intercept [0.9018627]
```

Notebook Exercise 1 (10 minutes)

This is a bit of review. Provide an interpretation for the coefficients learned by the logistic regression model (e.g., what do they mean for the prediction of whether a passenger would survive).

Solution

The negative coefficient for the first feature means that male passengers (\$x_1=1\$) were less likely to survive than female passengers (\$x_1=0\$). The slightly positive coefficient for the second feature (age) means that older passengers were more likely to survive (the effect is not strong though). The large and positive weight for `is young male` means that those passengers had a much higher survival rate in comparison with older males.

Reimplementation in Pytorch

Next, we'll be using `pytorch` to implement our own version of logistic regression! When creating a neural network model (remember, we can think of logistic regression as a perceptron with 2 layers, an input and an output), you create a class that inherits from `nn.Module`. We'll give you an implementation of logistic regression and then give a detailed breakdown of the key lines.

In [4]:

```
from torch import nn
import torch
from torch.autograd import Variable

class LogisticRegressionPytorch(nn.Module):
    def __init__(self):
        super(LogisticRegressionPytorch, self).__init__()
        self.linear = nn.Linear(3,1)

    def forward(self, X):
        """ Propagate data through the network.

        This model first applies the linear layer and then a sigmoid
        """
        X = self.linear(X)
        return torch.sigmoid(X)
```

Here is a breakdown of some of the key lines in this implementation.

Inherit from the super class `nn.Module`:

```
class LogisticRegressionPytorch(nn.Module):
```

Since we're inheriting from `nn.Module`, we need to make sure to call the `__init__` method of `nn.Module` when initializing our class.

```
super(LogisticRegressionPytorch, self).__init__()
```

The linear layer will store the weight vector of our model. The `3` arises from the fact that our model will have 3 inputs (age, male, and is young male). It is very important that you store your layers as attributes of your class. This is how `pytorch` knows about them and can optimize them. If you need to have a list of layers, look into `nn.ModuleList`.

```
self.linear = nn.Linear(3,1)
```

The `forward` function is the heart of the model. It runs input data through the network and returns the output. Writing such

functions usually amounts to passing data between the various layers that were created in the `__init__` method. The syntax for this is a little funny. For instance, in the code below, `self.linear(X)` implicitly calls the `forward` function of the `nn.Linear` class. Yes, we find this kind of weird, but that's how it's done in `pytorch`. Thus is really just a syntactic quirk rather than anything substantive that you need to worry about. The last step of the function involves applying the `sigmoid` and returning the result.

Next, we'll show how to pass some data into the model.

```
def forward(self, X):
    """ Propagate data through the network.

    This model first applies the linear layer and then a sigmoid
    """
    X = self.linear(X)
    return torch.sigmoid(X)
```

In [5]:

```
# sample_data represents a male passenger who is 10 years old
sample_data = Variable(torch.FloatTensor([1.0, 10.0, 0.0]))
model_pytorch = LogisticRegressionPytorch()
model_pytorch(sample_data)
```

Out[5]:

```
tensor([0.1522], grad_fn=<SigmoidBackward>)
```

The code we computed the probability that the specific passenger would survive. It is very important to realize that right now the model *has not been trained*. This means that we don't expect the output of the model to make any sense (although it might just by chance). If you rerun the code repeatedly, you'll get different results due to the fact that the weights are initialized randomly.

Next, we're going to actually train the network! This is where things get interesting. Run the code and either look through the code on your own or look below for a line-by-line breakdown.

In [6]:

```
model_pytorch = LogisticRegressionPytorch()
model_pytorch.train()

optimizer = torch.optim.SGD(model_pytorch.parameters(), lr=0.01)
criterion = torch.nn.BCELoss()
grad_magnitudes = []

X_data = Variable(torch.Tensor(np.array(experiment_1_data)))
y_data = Variable(torch.Tensor(np.array(experiment_1_outputs)))
for epoch in range(20000):
    optimizer.zero_grad()
    # Forward pass
    y_pred = model_pytorch(X_data)
    # Compute Loss
    loss = criterion(y_pred, y_data)
    # Backward pass
    loss.backward()
    for name, param in model_pytorch.named_parameters():
        if name == 'linear.weight':
            grad_magnitudes.append(np.abs(param.grad.numpy()).mean())

    if epoch % 100 == 0:
        print("epoch", epoch)
        for name, param in model_pytorch.named_parameters():
            print(name, "value", param.data, "gradient", param.grad)
        optimizer.step()

import matplotlib.pyplot as plt
%matplotlib inline
plt.plot(grad_magnitudes)
plt.xlabel('epoch')
plt.ylabel('average magnitude of weight gradient')
plt.yscale('log')
plt.show()
```

epoch 0

linear.weight value tensor([5.0000, 0.1604, 0.2720]) gradient tensor([5.4611e-01

```
linear.weight value tensor([[ -0.2920,  0.1694, -0.2120]]) gradient tensor([[ 4.6414e-01,
1.7580e+01, -9.0821e-03]])
linear.bias value tensor([-0.3314]) gradient tensor([0.5288])
epoch 100
linear.weight value tensor([[ -0.3908,  0.0361, -0.2600]]) gradient tensor([[ 0.2448,  8.4225, -0.0
116]])
linear.bias value tensor([-0.3076]) gradient tensor([0.2097])
epoch 200
linear.weight value tensor([[ -0.4828,  0.0369, -0.2477]]) gradient tensor([[ 0.2393,  8.3909, -0.0
119]])
linear.bias value tensor([-0.2787]) gradient tensor([0.2081])
epoch 300
linear.weight value tensor([[ -0.5711,  0.0376, -0.2351]]) gradient tensor([[ 0.2339,  8.3497, -0.0
122]])
linear.bias value tensor([-0.2494]) gradient tensor([0.2063])
```

```
/anaconda3/lib/python3.6/site-packages/torch/nn/functional.py:2016: UserWarning: Using a target si
ze (torch.Size([714])) that is different to the input size (torch.Size([714, 1])) is deprecated. P
lease ensure they have the same size.
```

```
"Please ensure they have the same size.".format(target.size(), input.size()))
```

```
epoch 400
linear.weight value tensor([[ -0.6559,  0.0382, -0.2222]]) gradient tensor([[ 0.2286,  8.2997, -0.0
124]])
linear.bias value tensor([-0.2198]) gradient tensor([0.2044])
epoch 500
linear.weight value tensor([[ -0.7373,  0.0387, -0.2091]]) gradient tensor([[ 0.2233,  8.2414, -0.0
127]])
linear.bias value tensor([-0.1900]) gradient tensor([0.2023])
epoch 600
linear.weight value tensor([[ -0.8155,  0.0391, -0.1958]]) gradient tensor([[ 0.2182,  8.1754, -0.0
129]])
linear.bias value tensor([-0.1601]) gradient tensor([0.2000])
epoch 700
linear.weight value tensor([[ -0.8906,  0.0395, -0.1823]]) gradient tensor([[ 0.2131,  8.1021, -0.0
131]])
linear.bias value tensor([-0.1302]) gradient tensor([0.1977])
epoch 800
linear.weight value tensor([[ -0.9627,  0.0397, -0.1687]]) gradient tensor([[ 0.2081,  8.0219, -0.0
132]])
linear.bias value tensor([-0.1002]) gradient tensor([0.1952])
epoch 900
linear.weight value tensor([[ -1.0319,  0.0399, -0.1549]]) gradient tensor([[ 0.2031,  7.9352, -0.0
134]])
linear.bias value tensor([-0.0704]) gradient tensor([0.1926])
epoch 1000
linear.weight value tensor([[ -1.0983,  0.0400, -0.1409]]) gradient tensor([[ 0.1982,  7.8422, -0.0
135]])
linear.bias value tensor([-0.0407]) gradient tensor([0.1898])
epoch 1100
linear.weight value tensor([[ -1.1621,  0.0400, -0.1269]]) gradient tensor([[ 0.1934,  7.7432, -0.0
136]])
linear.bias value tensor([-0.0111]) gradient tensor([0.1870])
epoch 1200
linear.weight value tensor([[ -1.2233,  0.0399, -0.1128]]) gradient tensor([[ 0.1885,  7.6384, -0.0
137]])
linear.bias value tensor([0.0182]) gradient tensor([0.1840])
epoch 1300
linear.weight value tensor([[ -1.2819,  0.0398, -0.0985]]) gradient tensor([[ 0.1837,  7.5281, -0.0
138]])
linear.bias value tensor([0.0472]) gradient tensor([0.1809])
epoch 1400
linear.weight value tensor([[ -1.3381,  0.0396, -0.0843]]) gradient tensor([[ 0.1790,  7.4125, -0.0
139]])
linear.bias value tensor([0.0759]) gradient tensor([0.1778])
epoch 1500
linear.weight value tensor([[ -1.3919,  0.0394, -0.0699]]) gradient tensor([[ 0.1742,  7.2917, -0.0
140]])
linear.bias value tensor([0.1044]) gradient tensor([0.1745])
epoch 1600
linear.weight value tensor([[ -1.4434,  0.0391, -0.0555]]) gradient tensor([[ 0.1695,  7.1658, -0.0
140]])
linear.bias value tensor([0.1324]) gradient tensor([0.1711])
epoch 1700
linear.weight value tensor([[ -1.4926,  0.0387, -0.0411]]) gradient tensor([[ 0.1648,  7.0351, -0.0
141]])
```

```
linear.bias value tensor([0.1601]) gradient tensor([0.1676])
epoch 1800
linear.weight value tensor([[ -1.5397,  0.0383, -0.0266]]) gradient tensor([[ 0.1602,  6.8997, -0.0
141]])
linear.bias value tensor([0.1873]) gradient tensor([0.1640])
epoch 1900
linear.weight value tensor([[ -1.5846,  0.0379, -0.0121]]) gradient tensor([[ 0.1555,  6.7598, -0.0
141]])
linear.bias value tensor([0.2142]) gradient tensor([0.1604])
epoch 2000
linear.weight value tensor([[ -1.6274,  0.0374,  0.0024]]) gradient tensor([[ 0.1509,  6.6154, -
0.0141]])
linear.bias value tensor([0.2406]) gradient tensor([0.1566])
epoch 2100
linear.weight value tensor([[ -1.6682,  0.0368,  0.0169]]) gradient tensor([[ 0.1463,  6.4667, -
0.0141]])
linear.bias value tensor([0.2665]) gradient tensor([0.1528])
epoch 2200
linear.weight value tensor([[ -1.7070,  0.0362,  0.0314]]) gradient tensor([[ 0.1417,  6.3138, -
0.0141]])
linear.bias value tensor([0.2919]) gradient tensor([0.1488])
epoch 2300
linear.weight value tensor([[ -1.7439,  0.0356,  0.0459]]) gradient tensor([[ 0.1371,  6.1568, -
0.0141]])
linear.bias value tensor([0.3168]) gradient tensor([0.1448])
epoch 2400
linear.weight value tensor([[ -1.7789,  0.0349,  0.0604]]) gradient tensor([[ 0.1325,  5.9959, -
0.0141]])
linear.bias value tensor([0.3412]) gradient tensor([0.1407])
epoch 2500
linear.weight value tensor([[ -1.8121,  0.0342,  0.0748]]) gradient tensor([[ 0.1279,  5.8310, -
0.0141]])
linear.bias value tensor([0.3651]) gradient tensor([0.1366])
epoch 2600
linear.weight value tensor([[ -1.8435,  0.0334,  0.0892]]) gradient tensor([[ 0.1234,  5.6623, -
0.0141]])
linear.bias value tensor([0.3885]) gradient tensor([0.1323])
epoch 2700
linear.weight value tensor([[ -1.8733,  0.0326,  0.1036]]) gradient tensor([[ 0.1188,  5.4899, -
0.0140]])
linear.bias value tensor([0.4113]) gradient tensor([0.1280])
epoch 2800
linear.weight value tensor([[ -1.9014,  0.0318,  0.1180]]) gradient tensor([[ 0.1143,  5.3138, -
0.0140]])
linear.bias value tensor([0.4335]) gradient tensor([0.1236])
epoch 2900
linear.weight value tensor([[ -1.9279,  0.0309,  0.1323]]) gradient tensor([[ 0.1098,  5.1339, -
0.0140]])
linear.bias value tensor([0.4552]) gradient tensor([0.1191])
epoch 3000
linear.weight value tensor([[ -1.9528,  0.0300,  0.1465]]) gradient tensor([[ 0.1053,  4.9504, -
0.0139]])
linear.bias value tensor([0.4764]) gradient tensor([0.1145])
epoch 3100
linear.weight value tensor([[ -1.9762,  0.0291,  0.1607]]) gradient tensor([[ 0.1008,  4.7632, -
0.0139]])
linear.bias value tensor([0.4969]) gradient tensor([0.1098])
epoch 3200
linear.weight value tensor([[ -1.9983,  0.0282,  0.1748]]) gradient tensor([[ 0.0963,  4.5722, -
0.0138]])
linear.bias value tensor([0.5169]) gradient tensor([0.1051])
epoch 3300
linear.weight value tensor([[ -2.0189,  0.0272,  0.1889]]) gradient tensor([[ 0.0918,  4.3774, -
0.0138]])
linear.bias value tensor([0.5363]) gradient tensor([0.1003])
epoch 3400
linear.weight value tensor([[ -2.0382,  0.0262,  0.2028]]) gradient tensor([[ 0.0872,  4.1785, -
0.0137]])
linear.bias value tensor([0.5551]) gradient tensor([0.0954])
epoch 3500
linear.weight value tensor([[ -2.0563,  0.0251,  0.2167]]) gradient tensor([[ 0.0827,  3.9752, -
0.0136]])
linear.bias value tensor([0.5734]) gradient tensor([0.0903])
epoch 3600
linear.weight value tensor([[ -2.0731,  0.0240,  0.2306]]) gradient tensor([[ 0.0782,  3.7674, -
0.0136]])
linear.bias value tensor([0.5910]) gradient tensor([0.0852])
```



```
epoch 3700
linear.weight value tensor([[ -2.0888,  0.0229,  0.2443]]) gradient tensor([[ 0.0736,  3.5545, -
0.0135]])
linear.bias value tensor([0.6081]) gradient tensor([0.0799])
epoch 3800
linear.weight value tensor([[ -2.1034,  0.0218,  0.2580]]) gradient tensor([[ 0.0690,  3.3359, -
0.0134]])
linear.bias value tensor([0.6247]) gradient tensor([0.0745])
epoch 3900
linear.weight value tensor([[ -2.1169,  0.0206,  0.2715]]) gradient tensor([[ 0.0643,  3.1106, -
0.0133]])
linear.bias value tensor([0.6406]) gradient tensor([0.0689])
epoch 4000
linear.weight value tensor([[ -2.1295,  0.0194,  0.2850]]) gradient tensor([[ 0.0596,  2.8777, -
0.0132]])
linear.bias value tensor([0.6560]) gradient tensor([0.0631])
epoch 4100
linear.weight value tensor([[ -2.1411,  0.0181,  0.2984]]) gradient tensor([[ 0.0547,  2.6355, -
0.0132]])
linear.bias value tensor([0.6708]) gradient tensor([0.0570])
epoch 4200
linear.weight value tensor([[ -2.1518,  0.0167,  0.3117]]) gradient tensor([[ 0.0497,  2.3821, -
0.0131]])
linear.bias value tensor([0.6851]) gradient tensor([0.0507])
epoch 4300
linear.weight value tensor([[ -2.1618,  0.0153,  0.3248]]) gradient tensor([[ 0.0446,  2.1145, -
0.0130]])
linear.bias value tensor([0.6988]) gradient tensor([0.0439])
epoch 4400
linear.weight value tensor([[ -2.1710,  0.0138,  0.3379]]) gradient tensor([[ 0.0391,  1.8291, -
0.0129]])
linear.bias value tensor([0.7120]) gradient tensor([0.0367])
epoch 4500
linear.weight value tensor([[ -2.1795,  0.0122,  0.3509]]) gradient tensor([[ 0.0334,  1.5219, -
0.0128]])
linear.bias value tensor([0.7247]) gradient tensor([0.0288])
epoch 4600
linear.weight value tensor([[ -2.1874,  0.0104,  0.3638]]) gradient tensor([[ 0.0273,  1.1897, -
0.0127]])
linear.bias value tensor([0.7369]) gradient tensor([0.0203])
epoch 4700
linear.weight value tensor([[ -2.1947,  0.0085,  0.3766]]) gradient tensor([[ 0.0211,  0.8380, -
0.0127]])
linear.bias value tensor([0.7486]) gradient tensor([0.0112])
epoch 4800
linear.weight value tensor([[ -2.2017,  0.0067,  0.3892]]) gradient tensor([[ 0.0150,  0.4961, -
0.0126]])
linear.bias value tensor([0.7599]) gradient tensor([0.0024])
epoch 4900
linear.weight value tensor([[ -2.2084,  0.0052,  0.4018]]) gradient tensor([[ 0.0104,  0.2257, -
0.0125]])
linear.bias value tensor([0.7708]) gradient tensor([-0.0046])
epoch 5000
linear.weight value tensor([[ -2.2150,  0.0042,  0.4143]]) gradient tensor([[ 0.0077,  0.0734, -
0.0124]])
linear.bias value tensor([0.7813]) gradient tensor([-0.0084])
epoch 5100
linear.weight value tensor([[ -2.2214,  0.0038,  0.4266]]) gradient tensor([[ 0.0067,  0.0167, -
0.0123]])
linear.bias value tensor([0.7915]) gradient tensor([-0.0096])
epoch 5200
linear.weight value tensor([[ -2.2278,  0.0035,  0.4389]]) gradient tensor([[ 0.0064,  0.0028, -
0.0122]])
linear.bias value tensor([0.8014]) gradient tensor([-0.0097])
epoch 5300
linear.weight value tensor([[ -2.2341,  0.0034,  0.4511]]) gradient tensor([[ 0.0063,  0.0004, -
0.0121]])
linear.bias value tensor([0.8109]) gradient tensor([-0.0094])
epoch 5400
linear.weight value tensor([[ -2.2403,  0.0032,  0.4631]]) gradient tensor([[ 0.0062,  0.0002, -
0.0120]])
linear.bias value tensor([0.8202]) gradient tensor([-0.0091])
epoch 5500
linear.weight value tensor([[ -2.2464,  0.0031,  0.4751]]) gradient tensor([[ 0.0061,  0.0001, -
0.0119]])
linear.bias value tensor([0.8292]) gradient tensor([-0.0088])
epoch 5600
```

```
linear.weight value tensor([[ -2.2525,  0.0030,  0.4870]]) gradient tensor([[ 0.0060,  0.0001, -
0.0118]])
linear.bias value tensor([0.8379]) gradient tensor([-0.0085])
epoch 5700
linear.weight value tensor([[ -2.2584,  0.0028,  0.4988]]) gradient tensor([[ 0.0059,  0.0001, -
0.0117]])
linear.bias value tensor([0.8463]) gradient tensor([-0.0083])
epoch 5800
linear.weight value tensor([[ -2.2643,  0.0027,  0.5105]]) gradient tensor([[ 0.0058,  0.0001, -
0.0117]])
linear.bias value tensor([0.8544]) gradient tensor([-0.0080])
epoch 5900
linear.weight value tensor([[ -2.2701,  0.0026,  0.5221]]) gradient tensor([[ 0.0057,  0.0001, -
0.0116]])
linear.bias value tensor([0.8623]) gradient tensor([-0.0078])
epoch 6000
linear.weight value tensor([[ -2.2758,  0.0025,  0.5336]]) gradient tensor([[ 0.0056,  0.0001, -
0.0115]])
linear.bias value tensor([0.8699]) gradient tensor([-0.0075])
epoch 6100
linear.weight value tensor([[ -2.2814,  0.0024,  0.5451]]) gradient tensor([[ 0.0056,  0.0001, -
0.0114]])
linear.bias value tensor([0.8773]) gradient tensor([-0.0073])
epoch 6200
linear.weight value tensor([[ -2.2869e+00,  2.2549e-03,  5.5641e-01]]) gradient tensor([[ 5.4772e-0
3,  9.9406e-05, -1.1303e-02]])
linear.bias value tensor([0.8845]) gradient tensor([-0.0070])
epoch 6300
linear.weight value tensor([[ -2.2923e+00,  2.1565e-03,  5.6768e-01]]) gradient tensor([[ 5.3926e-0
3,  9.5531e-05, -1.1218e-02]])
linear.bias value tensor([0.8914]) gradient tensor([-0.0068])
epoch 6400
linear.weight value tensor([[ -2.2977e+00,  2.0627e-03,  5.7885e-01]]) gradient tensor([[ 5.3085e-0
3,  9.1121e-05, -1.1134e-02]])
linear.bias value tensor([0.8981]) gradient tensor([-0.0066])
epoch 6500
linear.weight value tensor([[ -2.3029e+00,  1.9733e-03,  5.8994e-01]]) gradient tensor([[ 5.2250e-0
3,  8.6352e-05, -1.1051e-02]])
linear.bias value tensor([0.9046]) gradient tensor([-0.0064])
epoch 6600
linear.weight value tensor([[ -2.3081e+00,  1.8881e-03,  6.0095e-01]]) gradient tensor([[ 5.1421e-0
3,  8.3283e-05, -1.0968e-02]])
linear.bias value tensor([0.9109]) gradient tensor([-0.0062])
epoch 6700
linear.weight value tensor([[ -2.3132e+00,  1.8070e-03,  6.1188e-01]]) gradient tensor([[ 5.0599e-0
3,  7.9408e-05, -1.0887e-02]])
linear.bias value tensor([0.9170]) gradient tensor([-0.0060])
epoch 6800
linear.weight value tensor([[ -2.3182e+00,  1.7299e-03,  6.2273e-01]]) gradient tensor([[ 4.9785e-0
3,  7.6517e-05, -1.0807e-02]])
linear.bias value tensor([0.9229]) gradient tensor([-0.0058])
epoch 6900
linear.weight value tensor([[ -2.3232e+00,  1.6565e-03,  6.3350e-01]]) gradient tensor([[ 4.8978e-0
3,  7.1242e-05, -1.0727e-02]])
linear.bias value tensor([0.9286]) gradient tensor([-0.0056])
epoch 7000
linear.weight value tensor([[ -2.3280e+00,  1.5868e-03,  6.4418e-01]]) gradient tensor([[ 4.8180e-0
3,  6.9574e-05, -1.0649e-02]])
linear.bias value tensor([0.9341]) gradient tensor([-0.0054])
epoch 7100
linear.weight value tensor([[ -2.3328e+00,  1.5206e-03,  6.5479e-01]]) gradient tensor([[ 4.7390e-0
3,  6.5044e-05, -1.0571e-02]])
linear.bias value tensor([0.9394]) gradient tensor([-0.0053])
epoch 7200
linear.weight value tensor([[ -2.3375e+00,  1.4578e-03,  6.6533e-01]]) gradient tensor([[ 4.6608e-0
3,  6.1378e-05, -1.0494e-02]])
linear.bias value tensor([0.9446]) gradient tensor([-0.0051])
epoch 7300
linear.weight value tensor([[ -2.3421e+00,  1.3983e-03,  6.7578e-01]]) gradient tensor([[ 4.5837e-0
3,  5.8785e-05, -1.0418e-02]])
linear.bias value tensor([0.9496]) gradient tensor([-0.0049])
epoch 7400
linear.weight value tensor([[ -2.3467e+00,  1.3419e-03,  6.8616e-01]]) gradient tensor([[ 4.5075e-0
3,  5.3629e-05, -1.0343e-02]])
linear.bias value tensor([0.9544]) gradient tensor([-0.0048])
epoch 7500
linear.weight value tensor([[ -2.3512e+00,  1.2885e-03,  6.9647e-01]]) gradient tensor([[ 4.4322e-0
```

```
3, 5.0962e-05, -1.0269e-02]])
linear.bias value tensor([0.9591]) gradient tensor([-0.0046])
epoch 7600
linear.weight value tensor([[ -2.3555e+00, 1.2380e-03, 7.0670e-01]]) gradient tensor([[ 4.3580e-0
3, 4.9710e-05, -1.0195e-02]])
linear.bias value tensor([0.9636]) gradient tensor([-0.0045])
epoch 7700
linear.weight value tensor([[ -2.3599e+00, 1.1904e-03, 7.1686e-01]]) gradient tensor([[ 4.2847e-0
3, 4.6670e-05, -1.0122e-02]])
linear.bias value tensor([0.9680]) gradient tensor([-0.0043])
epoch 7800
linear.weight value tensor([[ -2.3641e+00, 1.1454e-03, 7.2695e-01]]) gradient tensor([[ 4.2124e-0
3, 4.5165e-05, -1.0050e-02]])
linear.bias value tensor([0.9722]) gradient tensor([-0.0042])
epoch 7900
linear.weight value tensor([[ -2.3683e+00, 1.1031e-03, 7.3696e-01]]) gradient tensor([[ 4.1411e-0
3, 3.9697e-05, -9.9790e-03]])
linear.bias value tensor([0.9763]) gradient tensor([-0.0040])
epoch 8000
linear.weight value tensor([[ -2.3724e+00, 1.0633e-03, 7.4691e-01]]) gradient tensor([[ 4.0709e-0
3, 3.7640e-05, -9.9084e-03]])
linear.bias value tensor([0.9803]) gradient tensor([-0.0039])
epoch 8100
linear.weight value tensor([[ -2.3764e+00, 1.0259e-03, 7.5678e-01]]) gradient tensor([[ 4.0017e-0
3, 3.5569e-05, -9.8386e-03]])
linear.bias value tensor([0.9841]) gradient tensor([-0.0038])
epoch 8200
linear.weight value tensor([[ -2.3804e+00, 9.9089e-04, 7.6659e-01]]) gradient tensor([[ 3.9336e-0
3, 3.2604e-05, -9.7695e-03]])
linear.bias value tensor([0.9878]) gradient tensor([-0.0036])
epoch 8300
linear.weight value tensor([[ -2.3843e+00, 9.5816e-04, 7.7632e-01]]) gradient tensor([[ 3.8665e-0
3, 3.2619e-05, -9.7010e-03]])
linear.bias value tensor([0.9914]) gradient tensor([-0.0035])
epoch 8400
linear.weight value tensor([[ -2.3881e+00, 9.2760e-04, 7.8599e-01]]) gradient tensor([[ 3.8004e-0
3, 2.8446e-05, -9.6333e-03]])
linear.bias value tensor([0.9949]) gradient tensor([-0.0034])
epoch 8500
linear.weight value tensor([[ -2.3919e+00, 8.9915e-04, 7.9559e-01]]) gradient tensor([[ 3.7354e-0
3, 2.7046e-05, -9.5662e-03]])
linear.bias value tensor([0.9982]) gradient tensor([-0.0033])
epoch 8600
linear.weight value tensor([[ -2.3956e+00, 8.7273e-04, 8.0512e-01]]) gradient tensor([[ 3.6714e-0
3, 2.5690e-05, -9.4997e-03]])
linear.bias value tensor([1.0014]) gradient tensor([-0.0032])
epoch 8700
linear.weight value tensor([[ -2.3993e+00, 8.4831e-04, 8.1459e-01]]) gradient tensor([[ 3.6084e-0
3, 2.4661e-05, -9.4339e-03]])
linear.bias value tensor([1.0045]) gradient tensor([-0.0031])
epoch 8800
linear.weight value tensor([[ -2.4028e+00, 8.2576e-04, 8.2399e-01]]) gradient tensor([[ 3.5465e-0
3, 2.1145e-05, -9.3687e-03]])
linear.bias value tensor([1.0075]) gradient tensor([-0.0029])
epoch 8900
linear.weight value tensor([[ -2.4063e+00, 8.0505e-04, 8.3333e-01]]) gradient tensor([[ 3.4856e-0
3, 2.0593e-05, -9.3042e-03]])
linear.bias value tensor([1.0104]) gradient tensor([-0.0028])
epoch 9000
linear.weight value tensor([[ -2.4098e+00, 7.8612e-04, 8.4260e-01]]) gradient tensor([[ 3.4257e-0
3, 1.7241e-05, -9.2402e-03]])
linear.bias value tensor([1.0132]) gradient tensor([-0.0027])
epoch 9100
linear.weight value tensor([[ -2.4132e+00, 7.6890e-04, 8.5181e-01]]) gradient tensor([[ 3.3669e-0
3, 1.5795e-05, -9.1769e-03]])
linear.bias value tensor([1.0159]) gradient tensor([-0.0026])
epoch 9200
linear.weight value tensor([[ -2.4165e+00, 7.5334e-04, 8.6095e-01]]) gradient tensor([[ 3.3091e-0
3, 1.5214e-05, -9.1141e-03]])
linear.bias value tensor([1.0185]) gradient tensor([-0.0026])
epoch 9300
linear.weight value tensor([[ -2.4198e+00, 7.3939e-04, 8.7004e-01]]) gradient tensor([[ 3.2522e-0
3, 1.4067e-05, -9.0519e-03]])
linear.bias value tensor([1.0210]) gradient tensor([-0.0025])
epoch 9400
linear.weight value tensor([[ -2.4230e+00, 7.2696e-04, 8.7906e-01]]) gradient tensor([[ 3.1964e-0
3, 1.2666e-05, -8.9903e-03]])
```

```
linear.bias value tensor([1.0234]) gradient tensor([-0.0024])
epoch 9500
linear.weight value tensor([[-2.4262e+00,  7.1604e-04,  8.8802e-01]]) gradient tensor([[ 3.1415e-0
3,  1.0729e-05, -8.9293e-03]])
linear.bias value tensor([1.0258]) gradient tensor([-0.0023])
epoch 9600
linear.weight value tensor([[-2.4293e+00,  7.0657e-04,  8.9692e-01]]) gradient tensor([[ 3.0875e-0
3,  8.9407e-06, -8.8688e-03]])
linear.bias value tensor([1.0280]) gradient tensor([-0.0022])
epoch 9700
linear.weight value tensor([[-2.4324e+00,  6.9845e-04,  9.0576e-01]]) gradient tensor([[ 3.0346e-0
3,  5.4091e-06, -8.8088e-03]])
linear.bias value tensor([1.0302]) gradient tensor([-0.0021])
epoch 9800
linear.weight value tensor([[-2.4354e+00,  6.9172e-04,  9.1454e-01]]) gradient tensor([[ 2.9826e-0
3,  7.3016e-06, -8.7494e-03]])
linear.bias value tensor([1.0322]) gradient tensor([-0.0020])
epoch 9900
linear.weight value tensor([[-2.4384e+00,  6.8627e-04,  9.2326e-01]]) gradient tensor([[ 2.9315e-0
3,  4.5300e-06, -8.6905e-03]])
linear.bias value tensor([1.0342]) gradient tensor([-0.0020])
epoch 10000
linear.weight value tensor([[-2.4413e+00,  6.8209e-04,  9.3192e-01]]) gradient tensor([[ 2.8814e-0
3,  3.3677e-06, -8.6322e-03]])
linear.bias value tensor([1.0361]) gradient tensor([-0.0019])
epoch 10100
linear.weight value tensor([[-2.4441e+00,  6.7912e-04,  9.4052e-01]]) gradient tensor([[ 2.8321e-0
3,  3.0696e-06, -8.5743e-03]])
linear.bias value tensor([1.0380]) gradient tensor([-0.0018])
epoch 10200
linear.weight value tensor([[-2.4469e+00,  6.7732e-04,  9.4907e-01]]) gradient tensor([[ 2.7838e-0
3,  9.6858e-07, -8.5170e-03]])
linear.bias value tensor([1.0398]) gradient tensor([-0.0017])
epoch 10300
linear.weight value tensor([[-2.4497e+00,  6.7667e-04,  9.5756e-01]]) gradient tensor([[ 2.7363e-0
3, -1.9372e-07, -8.4602e-03]])
linear.bias value tensor([1.0415]) gradient tensor([-0.0017])
epoch 10400
linear.weight value tensor([[-2.4524e+00,  6.7712e-04,  9.6599e-01]]) gradient tensor([[ 2.6897e-0
3, -1.4603e-06, -8.4038e-03]])
linear.bias value tensor([1.0431]) gradient tensor([-0.0016])
epoch 10500
linear.weight value tensor([[-2.4551e+00,  6.7862e-04,  9.7437e-01]]) gradient tensor([[ 2.6440e-0
3, -2.7418e-06, -8.3479e-03]])
linear.bias value tensor([1.0446]) gradient tensor([-0.0015])
epoch 10600
linear.weight value tensor([[-2.4577e+00,  6.8116e-04,  9.8269e-01]]) gradient tensor([[ 2.5992e-0
3, -2.8014e-06, -8.2926e-03]])
linear.bias value tensor([1.0461]) gradient tensor([-0.0015])
epoch 10700
linear.weight value tensor([[-2.4603e+00,  6.8468e-04,  9.9095e-01]]) gradient tensor([[ 2.5551e-0
3, -4.0829e-06, -8.2376e-03]])
linear.bias value tensor([1.0476]) gradient tensor([-0.0014])
epoch 10800
linear.weight value tensor([[-2.4628e+00,  6.8917e-04,  9.9916e-01]]) gradient tensor([[ 2.5119e-0
3, -4.1574e-06, -8.1832e-03]])
linear.bias value tensor([1.0489]) gradient tensor([-0.0013])
epoch 10900
linear.weight value tensor([[-2.4653e+00,  6.9458e-04,  1.0073e+00]]) gradient tensor([[ 2.4694e-0
3, -6.5267e-06, -8.1292e-03]])
linear.bias value tensor([1.0503]) gradient tensor([-0.0013])
epoch 11000
linear.weight value tensor([[-2.4677e+00,  7.0087e-04,  1.0154e+00]]) gradient tensor([[ 2.4279e-0
3, -6.5565e-06, -8.0757e-03]])
linear.bias value tensor([1.0515]) gradient tensor([-0.0012])
epoch 11100
linear.weight value tensor([[-2.4701e+00,  7.0803e-04,  1.0235e+00]]) gradient tensor([[ 2.3870e-0
3, -9.2387e-06, -8.0226e-03]])
linear.bias value tensor([1.0527]) gradient tensor([-0.0012])
epoch 11200
linear.weight value tensor([[-2.4725e+00,  7.1605e-04,  1.0315e+00]]) gradient tensor([[ 2.3469e-0
3, -9.0450e-06, -7.9700e-03]])
linear.bias value tensor([1.0539]) gradient tensor([-0.0011])
epoch 11300
linear.weight value tensor([[-2.4748e+00,  7.2485e-04,  1.0394e+00]]) gradient tensor([[ 2.3077e-0
3, -9.6560e-06, -7.9177e-03]])
linear.bias value tensor([1.0549]) gradient tensor([-0.0011])
```

```
epoch 11400
linear.weight value tensor([[ -2.4771e+00,  7.3443e-04,  1.0473e+00]]) gradient tensor([[ 2.2691e-03, -1.0207e-05, -7.8660e-03]])
linear.bias value tensor([1.0560]) gradient tensor([-0.0010])
epoch 11500
linear.weight value tensor([[ -2.4794e+00,  7.4478e-04,  1.0551e+00]]) gradient tensor([[ 2.2312e-03, -1.1116e-05, -7.8146e-03]])
linear.bias value tensor([1.0570]) gradient tensor([-0.0010])
epoch 11600
linear.weight value tensor([[ -2.4816e+00,  7.5582e-04,  1.0629e+00]]) gradient tensor([[ 2.1942e-03, -1.2428e-05, -7.7637e-03]])
linear.bias value tensor([1.0579]) gradient tensor([-0.0009])
epoch 11700
linear.weight value tensor([[ -2.4838e+00,  7.6762e-04,  1.0707e+00]]) gradient tensor([[ 2.1577e-03, -1.2413e-05, -7.7132e-03]])
linear.bias value tensor([1.0588]) gradient tensor([-0.0009])
epoch 11800
linear.weight value tensor([[ -2.4859e+00,  7.8003e-04,  1.0784e+00]]) gradient tensor([[ 2.1221e-03, -1.3798e-05, -7.6631e-03]])
linear.bias value tensor([1.0596]) gradient tensor([-0.0008])
epoch 11900
linear.weight value tensor([[ -2.4880e+00,  7.9317e-04,  1.0860e+00]]) gradient tensor([[ 2.0870e-03, -1.2353e-05, -7.6134e-03]])
linear.bias value tensor([1.0604]) gradient tensor([-0.0008])
epoch 12000
linear.weight value tensor([[ -2.4901e+00,  8.0689e-04,  1.0936e+00]]) gradient tensor([[ 2.0527e-03, -1.3217e-05, -7.5641e-03]])
linear.bias value tensor([1.0612]) gradient tensor([-0.0007])
epoch 12100
linear.weight value tensor([[ -2.4921e+00,  8.2122e-04,  1.1011e+00]]) gradient tensor([[ 2.0190e-03, -1.5214e-05, -7.5151e-03]])
linear.bias value tensor([1.0619]) gradient tensor([-0.0007])
epoch 12200
linear.weight value tensor([[ -2.4941e+00,  8.3616e-04,  1.1086e+00]]) gradient tensor([[ 1.9859e-03, -1.5035e-05, -7.4666e-03]])
linear.bias value tensor([1.0625]) gradient tensor([-0.0006])
epoch 12300
linear.weight value tensor([[ -2.4961e+00,  8.5164e-04,  1.1161e+00]]) gradient tensor([[ 1.9536e-03, -1.6063e-05, -7.4185e-03]])
linear.bias value tensor([1.0632]) gradient tensor([-0.0006])
epoch 12400
linear.weight value tensor([[ -2.4980e+00,  8.6771e-04,  1.1235e+00]]) gradient tensor([[ 1.9217e-03, -1.6347e-05, -7.3708e-03]])
linear.bias value tensor([1.0637]) gradient tensor([-0.0006])
epoch 12500
linear.weight value tensor([[ -2.4999e+00,  8.8428e-04,  1.1308e+00]]) gradient tensor([[ 1.8905e-03, -1.6108e-05, -7.3234e-03]])
linear.bias value tensor([1.0643]) gradient tensor([-0.0005])
epoch 12600
linear.weight value tensor([[ -2.5018e+00,  9.0135e-04,  1.1381e+00]]) gradient tensor([[ 1.8600e-03, -1.7583e-05, -7.2764e-03]])
linear.bias value tensor([1.0648]) gradient tensor([-0.0005])
epoch 12700
linear.weight value tensor([[ -2.5037e+00,  9.1893e-04,  1.1454e+00]]) gradient tensor([[ 1.8300e-03, -1.7986e-05, -7.2298e-03]])
linear.bias value tensor([1.0653]) gradient tensor([-0.0005])
epoch 12800
linear.weight value tensor([[ -2.5055e+00,  9.3701e-04,  1.1526e+00]]) gradient tensor([[ 1.8005e-03, -1.7449e-05, -7.1835e-03]])
linear.bias value tensor([1.0657]) gradient tensor([-0.0004])
epoch 12900
linear.weight value tensor([[ -2.5073e+00,  9.5549e-04,  1.1597e+00]]) gradient tensor([[ 1.7717e-03, -1.7747e-05, -7.1376e-03]])
linear.bias value tensor([1.0661]) gradient tensor([-0.0004])
epoch 13000
linear.weight value tensor([[ -2.5090e+00,  9.7440e-04,  1.1668e+00]]) gradient tensor([[ 1.7434e-03, -1.9163e-05, -7.0920e-03]])
linear.bias value tensor([1.0665]) gradient tensor([-0.0004])
epoch 13100
linear.weight value tensor([[ -2.5107e+00,  9.9376e-04,  1.1739e+00]]) gradient tensor([[ 1.7157e-03, -1.8731e-05, -7.0469e-03]])
linear.bias value tensor([1.0668]) gradient tensor([-0.0003])
epoch 13200
linear.weight value tensor([[ -2.5124e+00,  1.0135e-03,  1.1809e+00]]) gradient tensor([[ 1.6884e-03, -1.8850e-05, -7.0020e-03]])
linear.bias value tensor([1.0671]) gradient tensor([-0.0003])
epoch 13300
```

```
linear.weight value tensor([[ -2.5141e+00,  1.0337e-03,  1.1879e+00]]) gradient tensor([[ 1.6617e-03, -2.0579e-05, -6.9575e-03]])
linear.bias value tensor([1.0674]) gradient tensor([-0.0003])
epoch 13400
linear.weight value tensor([[ -2.5158e+00,  1.0543e-03,  1.1948e+00]]) gradient tensor([[ 1.6355e-03, -2.0877e-05, -6.9134e-03]])
linear.bias value tensor([1.0676]) gradient tensor([-0.0002])
epoch 13500
linear.weight value tensor([[ -2.5174e+00,  1.0752e-03,  1.2017e+00]]) gradient tensor([[ 1.6098e-03, -2.0877e-05, -6.8695e-03]])
linear.bias value tensor([1.0678]) gradient tensor([-0.0002])
epoch 13600
linear.weight value tensor([[ -2.5190e+00,  1.0964e-03,  1.2086e+00]]) gradient tensor([[ 1.5846e-03, -2.0891e-05, -6.8261e-03]])
linear.bias value tensor([1.0680]) gradient tensor([-0.0002])
epoch 13700
linear.weight value tensor([[ -2.5206e+00,  1.1180e-03,  1.2154e+00]]) gradient tensor([[ 1.5599e-03, -2.1428e-05, -6.7829e-03]])
linear.bias value tensor([1.0682]) gradient tensor([-0.0001])
epoch 13800
linear.weight value tensor([[ -2.5221e+00,  1.1399e-03,  1.2222e+00]]) gradient tensor([[ 1.5357e-03, -2.2277e-05, -6.7401e-03]])
linear.bias value tensor([1.0683]) gradient tensor([-0.0001])
epoch 13900
linear.weight value tensor([[ -2.5236e+00,  1.1621e-03,  1.2289e+00]]) gradient tensor([[ 1.5119e-03, -2.2352e-05, -6.6976e-03]])
linear.bias value tensor([1.0684]) gradient tensor([-9.0903e-05])
epoch 14000
linear.weight value tensor([[ -2.5251e+00,  1.1846e-03,  1.2356e+00]]) gradient tensor([[ 1.4886e-03, -2.2709e-05, -6.6554e-03]])
linear.bias value tensor([1.0685]) gradient tensor([-6.6032e-05])
epoch 14100
linear.weight value tensor([[ -2.5266e+00,  1.2074e-03,  1.2422e+00]]) gradient tensor([[ 1.4657e-03, -2.2054e-05, -6.6135e-03]])
linear.bias value tensor([1.0685]) gradient tensor([-4.1848e-05])
epoch 14200
linear.weight value tensor([[ -2.5281e+00,  1.2304e-03,  1.2488e+00]]) gradient tensor([[ 1.4433e-03, -2.3231e-05, -6.5720e-03]])
linear.bias value tensor([1.0686]) gradient tensor([-1.8389e-05])
epoch 14300
linear.weight value tensor([[ -2.5295e+00,  1.2538e-03,  1.2553e+00]]) gradient tensor([[ 1.4213e-03, -2.3708e-05, -6.5307e-03]])
linear.bias value tensor([1.0686]) gradient tensor([4.4077e-06])
epoch 14400
linear.weight value tensor([[ -2.5309e+00,  1.2773e-03,  1.2618e+00]]) gradient tensor([[ 1.3998e-03, -2.2814e-05, -6.4898e-03]])
linear.bias value tensor([1.0686]) gradient tensor([2.6582e-05])
epoch 14500
linear.weight value tensor([[ -2.5323e+00,  1.3011e-03,  1.2683e+00]]) gradient tensor([[ 1.3786e-03, -2.3529e-05, -6.4491e-03]])
linear.bias value tensor([1.0685]) gradient tensor([4.8071e-05])
epoch 14600
linear.weight value tensor([[ -2.5337e+00,  1.3251e-03,  1.2747e+00]]) gradient tensor([[ 1.3579e-03, -2.5123e-05, -6.4088e-03]])
linear.bias value tensor([1.0685]) gradient tensor([6.8910e-05])
epoch 14700
linear.weight value tensor([[ -2.5350e+00,  1.3494e-03,  1.2811e+00]]) gradient tensor([[ 1.3375e-03, -2.3857e-05, -6.3688e-03]])
linear.bias value tensor([1.0684]) gradient tensor([8.9214e-05])
epoch 14800
linear.weight value tensor([[ -2.5363e+00,  1.3739e-03,  1.2875e+00]]) gradient tensor([[ 1.3175e-03, -2.3723e-05, -6.3291e-03]])
linear.bias value tensor([1.0683]) gradient tensor([0.0001])
epoch 14900
linear.weight value tensor([[ -2.5376e+00,  1.3985e-03,  1.2938e+00]]) gradient tensor([[ 1.2978e-03, -2.5496e-05, -6.2896e-03]])
linear.bias value tensor([1.0682]) gradient tensor([0.0001])
epoch 15000
linear.weight value tensor([[ -2.5389e+00,  1.4234e-03,  1.3001e+00]]) gradient tensor([[ 1.2786e-03, -2.4974e-05, -6.2505e-03]])
linear.bias value tensor([1.0680]) gradient tensor([0.0001])
epoch 15100
linear.weight value tensor([[ -2.5402e+00,  1.4484e-03,  1.3063e+00]]) gradient tensor([[ 1.2598e-03, -2.4959e-05, -6.2116e-03]])
linear.bias value tensor([1.0679]) gradient tensor([0.0002])
epoch 15200
linear.weight value tensor([[ -2.5414e+00,  1.4736e-03,  1.3125e+00]]) gradient tensor([[ 1.2414e-03, -2.4959e-05, -6.2116e-03]])
```

```
3, -2.4855e-05, -6.1730e-03]])
linear.bias value tensor([1.0677]) gradient tensor([0.0002])
epoch 15300
linear.weight value tensor([[ -2.5427e+00,  1.4989e-03,  1.3186e+00]]) gradient tensor([[ 1.2231e-0
3, -2.5943e-05, -6.1347e-03]])
linear.bias value tensor([1.0675]) gradient tensor([0.0002])
epoch 15400
linear.weight value tensor([[ -2.5439e+00,  1.5245e-03,  1.3247e+00]]) gradient tensor([[ 1.2053e-0
3, -2.5690e-05, -6.0967e-03]])
linear.bias value tensor([1.0673]) gradient tensor([0.0002])
epoch 15500
linear.weight value tensor([[ -2.5451e+00,  1.5501e-03,  1.3308e+00]]) gradient tensor([[ 1.1879e-0
3, -2.6405e-05, -6.0589e-03]])
linear.bias value tensor([1.0671]) gradient tensor([0.0002])
epoch 15600
linear.weight value tensor([[ -2.5463e+00,  1.5758e-03,  1.3369e+00]]) gradient tensor([[ 1.1708e-0
3, -2.5555e-05, -6.0214e-03]])
linear.bias value tensor([1.0668]) gradient tensor([0.0002])
epoch 15700
linear.weight value tensor([[ -2.5474e+00,  1.6018e-03,  1.3429e+00]]) gradient tensor([[ 1.1539e-0
3, -2.6450e-05, -5.9842e-03]])
linear.bias value tensor([1.0666]) gradient tensor([0.0003])
epoch 15800
linear.weight value tensor([[ -2.5486e+00,  1.6278e-03,  1.3488e+00]]) gradient tensor([[ 1.1374e-0
3, -2.6554e-05, -5.9473e-03]])
linear.bias value tensor([1.0663]) gradient tensor([0.0003])
epoch 15900
linear.weight value tensor([[ -2.5497e+00,  1.6540e-03,  1.3548e+00]]) gradient tensor([[ 1.1213e-0
3, -2.5719e-05, -5.9106e-03]])
linear.bias value tensor([1.0660]) gradient tensor([0.0003])
epoch 16000
linear.weight value tensor([[ -2.5508e+00,  1.6802e-03,  1.3607e+00]]) gradient tensor([[ 1.1053e-0
3, -2.7016e-05, -5.8742e-03]])
linear.bias value tensor([1.0657]) gradient tensor([0.0003])
epoch 16100
linear.weight value tensor([[ -2.5519e+00,  1.7066e-03,  1.3665e+00]]) gradient tensor([[ 1.0897e-0
3, -2.6673e-05, -5.8380e-03]])
linear.bias value tensor([1.0654]) gradient tensor([0.0003])
epoch 16200
linear.weight value tensor([[ -2.5530e+00,  1.7330e-03,  1.3723e+00]]) gradient tensor([[ 1.0745e-0
3, -2.7150e-05, -5.8021e-03]])
linear.bias value tensor([1.0651]) gradient tensor([0.0003])
epoch 16300
linear.weight value tensor([[ -2.5541e+00,  1.7596e-03,  1.3781e+00]]) gradient tensor([[ 1.0593e-0
3, -2.6211e-05, -5.7665e-03]])
linear.bias value tensor([1.0648]) gradient tensor([0.0003])
epoch 16400
linear.weight value tensor([[ -2.5551e+00,  1.7862e-03,  1.3839e+00]]) gradient tensor([[ 1.0448e-0
3, -2.6628e-05, -5.7311e-03]])
linear.bias value tensor([1.0644]) gradient tensor([0.0004])
epoch 16500
linear.weight value tensor([[ -2.5562e+00,  1.8129e-03,  1.3896e+00]]) gradient tensor([[ 1.0303e-0
3, -2.7776e-05, -5.6960e-03]])
linear.bias value tensor([1.0641]) gradient tensor([0.0004])
epoch 16600
linear.weight value tensor([[ -2.5572e+00,  1.8397e-03,  1.3953e+00]]) gradient tensor([[ 1.0160e-0
3, -2.7955e-05, -5.6611e-03]])
linear.bias value tensor([1.0637]) gradient tensor([0.0004])
epoch 16700
linear.weight value tensor([[ -2.5582e+00,  1.8665e-03,  1.4009e+00]]) gradient tensor([[ 1.0022e-0
3, -2.6688e-05, -5.6264e-03]])
linear.bias value tensor([1.0633]) gradient tensor([0.0004])
epoch 16800
linear.weight value tensor([[ -2.5592e+00,  1.8935e-03,  1.4065e+00]]) gradient tensor([[ 9.8833e-0
4, -2.5928e-05, -5.5921e-03]])
linear.bias value tensor([1.0629]) gradient tensor([0.0004])
epoch 16900
linear.weight value tensor([[ -2.5602e+00,  1.9203e-03,  1.4121e+00]]) gradient tensor([[ 9.7510e-0
4, -2.6777e-05, -5.5579e-03]])
linear.bias value tensor([1.0625]) gradient tensor([0.0004])
epoch 17000
linear.weight value tensor([[ -2.5611e+00,  1.9474e-03,  1.4176e+00]]) gradient tensor([[ 9.6183e-0
4, -2.5794e-05, -5.5240e-03]])
linear.bias value tensor([1.0621]) gradient tensor([0.0004])
epoch 17100
linear.weight value tensor([[ -2.5621e+00,  1.9744e-03,  1.4231e+00]]) gradient tensor([[ 9.4897e-0
4, -2.7373e-05, -5.4903e-03]])
```

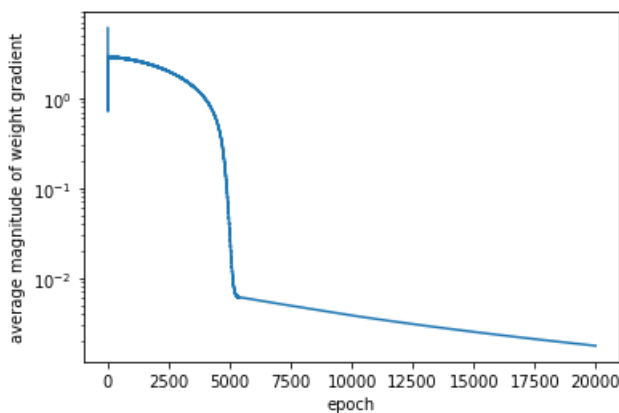
```
, 1.000000e+00, 1.000000e+00]],
linear.bias value tensor([1.0617]) gradient tensor([0.0004])
epoch 17200
linear.weight value tensor([[-2.5630e+00, 2.0014e-03, 1.4286e+00]]) gradient tensor([[ 9.3625e-04, -2.8342e-05, -5.4569e-03]])
linear.bias value tensor([1.0613]) gradient tensor([0.0004])
epoch 17300
linear.weight value tensor([[-2.5640e+00, 2.0285e-03, 1.4341e+00]]) gradient tensor([[ 9.2377e-04, -2.8476e-05, -5.4237e-03]])
linear.bias value tensor([1.0608]) gradient tensor([0.0004])
epoch 17400
linear.weight value tensor([[-2.5649e+00, 2.0557e-03, 1.4395e+00]]) gradient tensor([[ 9.1158e-04, -2.6211e-05, -5.3907e-03]])
linear.bias value tensor([1.0604]) gradient tensor([0.0005])
epoch 17500
linear.weight value tensor([[-2.5658e+00, 2.0828e-03, 1.4448e+00]]) gradient tensor([[ 8.9957e-04, -2.7061e-05, -5.3579e-03]])
linear.bias value tensor([1.0599]) gradient tensor([0.0005])
epoch 17600
linear.weight value tensor([[-2.5667e+00, 2.1100e-03, 1.4502e+00]]) gradient tensor([[ 8.8777e-04, -2.7522e-05, -5.3254e-03]])
linear.bias value tensor([1.0594]) gradient tensor([0.0005])
epoch 17700
linear.weight value tensor([[-2.5676e+00, 2.1372e-03, 1.4555e+00]]) gradient tensor([[ 8.7617e-04, -2.7761e-05, -5.2931e-03]])
linear.bias value tensor([1.0590]) gradient tensor([0.0005])
epoch 17800
linear.weight value tensor([[-2.5684e+00, 2.1643e-03, 1.4608e+00]]) gradient tensor([[ 8.6485e-04, -2.6777e-05, -5.2610e-03]])
linear.bias value tensor([1.0585]) gradient tensor([0.0005])
epoch 17900
linear.weight value tensor([[-2.5693e+00, 2.1916e-03, 1.4660e+00]]) gradient tensor([[ 8.5368e-04, -2.7731e-05, -5.2292e-03]])
linear.bias value tensor([1.0580]) gradient tensor([0.0005])
epoch 18000
linear.weight value tensor([[-2.5701e+00, 2.2187e-03, 1.4712e+00]]) gradient tensor([[ 8.4263e-04, -2.8536e-05, -5.1975e-03]])
linear.bias value tensor([1.0575]) gradient tensor([0.0005])
epoch 18100
linear.weight value tensor([[-2.5710e+00, 2.2459e-03, 1.4764e+00]]) gradient tensor([[ 8.3201e-04, -2.6971e-05, -5.1661e-03]])
linear.bias value tensor([1.0570]) gradient tensor([0.0005])
epoch 18200
linear.weight value tensor([[-2.5718e+00, 2.2732e-03, 1.4816e+00]]) gradient tensor([[ 8.2121e-04, -2.6822e-05, -5.1349e-03]])
linear.bias value tensor([1.0565]) gradient tensor([0.0005])
epoch 18300
linear.weight value tensor([[-2.5726e+00, 2.3003e-03, 1.4867e+00]]) gradient tensor([[ 8.1108e-04, -2.7061e-05, -5.1039e-03]])
linear.bias value tensor([1.0560]) gradient tensor([0.0005])
epoch 18400
linear.weight value tensor([[-2.5734e+00, 2.3275e-03, 1.4918e+00]]) gradient tensor([[ 8.0066e-04, -2.7135e-05, -5.0731e-03]])
linear.bias value tensor([1.0555]) gradient tensor([0.0005])
epoch 18500
linear.weight value tensor([[-2.5742e+00, 2.3546e-03, 1.4968e+00]]) gradient tensor([[ 7.9083e-04, -2.7269e-05, -5.0426e-03]])
linear.bias value tensor([1.0550]) gradient tensor([0.0005])
epoch 18600
linear.weight value tensor([[-2.5750e+00, 2.3817e-03, 1.5018e+00]]) gradient tensor([[ 7.8100e-04, -2.7150e-05, -5.0122e-03]])
linear.bias value tensor([1.0544]) gradient tensor([0.0005])
epoch 18700
linear.weight value tensor([[-2.5758e+00, 2.4089e-03, 1.5068e+00]]) gradient tensor([[ 7.7121e-04, -2.7582e-05, -4.9821e-03]])
linear.bias value tensor([1.0539]) gradient tensor([0.0005])
epoch 18800
linear.weight value tensor([[-2.5765e+00, 2.4359e-03, 1.5118e+00]]) gradient tensor([[ 7.6194e-04, -2.6911e-05, -4.9521e-03]])
linear.bias value tensor([1.0534]) gradient tensor([0.0005])
epoch 18900
linear.weight value tensor([[-2.5773e+00, 2.4631e-03, 1.5167e+00]]) gradient tensor([[ 7.5226e-04, -2.7850e-05, -4.9224e-03]])
linear.bias value tensor([1.0528]) gradient tensor([0.0005])
epoch 19000
linear.weight value tensor([[-2.5781e+00, 2.4900e-03, 1.5217e+00]]) gradient tensor([[ 7.4339e-04, -2.7224e-05, -4.8928e-03]])
linear.bias value tensor([1.0523]) gradient tensor([0.0005])
```



```

linear.bias value tensor([1.0000]), gradient tensor([0.0000]),
epoch 19100
linear.weight value tensor([[ -2.5788e+00,  2.5169e-03,  1.5265e+00]]) gradient tensor([[ 7.3445e-0
4, -2.6673e-05, -4.8634e-03]])
linear.bias value tensor([1.0517]) gradient tensor([0.0006])
epoch 19200
linear.weight value tensor([[ -2.5795e+00,  2.5441e-03,  1.5314e+00]]) gradient tensor([[ 7.2535e-0
4, -2.5779e-05, -4.8343e-03]])
linear.bias value tensor([1.0512]) gradient tensor([0.0006])
epoch 19300
linear.weight value tensor([[ -2.5802e+00,  2.5709e-03,  1.5362e+00]]) gradient tensor([[ 7.1693e-0
4, -2.6569e-05, -4.8053e-03]])
linear.bias value tensor([1.0506]) gradient tensor([0.0006])
epoch 19400
linear.weight value tensor([[ -2.5810,  0.0026,  1.5410]]) gradient tensor([[ 7.0832e-04, -2.5690e-
05, -4.7765e-03]])
linear.bias value tensor([1.0500]) gradient tensor([0.0006])
epoch 19500
linear.weight value tensor([[ -2.5817,  0.0026,  1.5458]]) gradient tensor([[ 6.9982e-04, -2.7284e-
05, -4.7480e-03]])
linear.bias value tensor([1.0495]) gradient tensor([0.0006])
epoch 19600
linear.weight value tensor([[ -2.5824,  0.0027,  1.5505]]) gradient tensor([[ 6.9176e-04, -2.6807e-
05, -4.7196e-03]])
linear.bias value tensor([1.0489]) gradient tensor([0.0006])
epoch 19700
linear.weight value tensor([[ -2.5830,  0.0027,  1.5552]]) gradient tensor([[ 6.8351e-04, -2.5615e-
05, -4.6914e-03]])
linear.bias value tensor([1.0483]) gradient tensor([0.0006])
epoch 19800
linear.weight value tensor([[ -2.5837,  0.0027,  1.5599]]) gradient tensor([[ 6.7562e-04, -2.6718e-
05, -4.6634e-03]])
linear.bias value tensor([1.0477]) gradient tensor([0.0006])
epoch 19900
linear.weight value tensor([[ -2.5844,  0.0027,  1.5645]]) gradient tensor([[ 6.6790e-04, -2.7508e-
05, -4.6356e-03]])
linear.bias value tensor([1.0472]) gradient tensor([0.0006])

```



Code breakdown

Put the model into training mode (this only affects certain models that behave differently during training and testing). Even though it doesn't affect our logistic regression model, calling the `.train` function is a good habit to get into.

```
model_pytorch.train()
```

Create an optimizer that will tune the parameters of the network. Here we are using an algorithm called *stochastic gradient descent*. It's a very popular choice for an optimizer. Similarly to gradient descent it optimizes a function by stepping down the gradient (step size is given by learning rate or `lr`). Where it differs from normal gradient descent is that it doesn't necessarily use all of the data to compute the gradient. If you have a big dataset, it might only use a small number of data points, compute gradient over those, and then step down that estimate of the gradient. The collection of datapoints used for estimating the gradient is called a *mini batch*. In this example we are using the entire dataset each time, so we are really doing normal gradient descent.

```
optimizer = torch.optim.SGD(model_pytorch.parameters(), lr=0.01)
```

The criterion defines the loss function we are minimizing. When the training outputs are binary, the `BCELoss` function is equivalent to the log loss that we have been using in this course.

```
criterion = torch.nn.BCELoss()
```

In order to operate on data in pytorch, you have to convert any matrix or vector data into a pytorch variable. This should be familiar based on the tutorial you went through earlier.

```
X_data = Variable(torch.Tensor(np.array(experiment_1_data)))
y_data = Variable(torch.Tensor(np.array(experiment_1_outputs)))
```

An *epoch* in neural network training is a single pass through the data. In this case we are taking a single gradient step on the whole dataset, so the number of epochs is the same as the number of gradient steps.

```
for epoch in range(200):
```

This tells the optimizer to throw away any gradients it has accumulated from previous data (do not forget to call this!!!).

```
optimizer.zero_grad()
```

Apply the forward model to get predictions.

```
y_pred = model_pytorch(X_data)
```

Calculate the loss of the model by comparing its predictions with the actual outputs.

```
loss = criterion(y_pred, y_data)
```

Use backpropagation to compute the gradient of all of the model parameters with respect to the loss.

```
loss.backward()
```

Calculate the gradient magnitudes so we can make a plot after training.

```
for name, param in model_pytorch.named_parameters():
    if name == 'linear.weight':
        grad_magnitudes.append(np.abs(param.grad.numpy()).mean())
```

Print out the values of the parameters and the gradient of the parameters with respect to the loss every 50 epochs.

```
if epoch % 50 == 0:
    print("epoch", epoch)
    for name, param in model_pytorch.named_parameters():
        print(name, "value", param.data, "gradient", param.grad)
```

Perform the gradient step.

```
optimizer.step()
```

Notebook Exercise 2 (30 minutes)

(a) Explain the output you see when you run the previous code cell. How are the weights changing over time? How is the gradient changing over time? Is the algorithm close to converging (i.e., computing the optimal solution)? How would you know if it has converged?

(b) Increase the number of epochs until you get convergence (you may want to make it so the gradient prints out less by changing the line `if epoch % 100 == 0` to `if epoch % 1000 == 0`). Roughly how many did it take?

(c) Tune the learning rate to some other values. How does this change the algorithm's behavior?

(d) Change the optimizer to the ADAM optimizer by swapping out the previous optimizer with this new line of code.

```
optimizer = torch.optim.Adam(model_pytorch.parameters())
```

Roughly many epochs does it take to reach convergence now?

Solution

(a) (this might vary based on the random initialization of the weights in the network). In this run the weights slowly adjusted to the values computed by the sklearn model. The size of the gradient remained largely constant over time (as seen both in the output and in the plot). This indicates that the algorithm has not converged yet. If the gradient magnitude was close to 0, that would indicate convergence.

- (b) It took about 150,000 epochs to reach convergence!
- (c) Setting it too high causes divergence of the algorithm. Setting it too low causes slow convergence.
- (d) It seems to take about 10,000 epochs to converge (that's a huge speedup!)

Multilayer Perceptron

Now that we've shown you how to implement the logistic regression model, we want you to implement the MLP model from the previous companion notebook. Remember, the MLP had 2 input features (we didn't use `is young male` as an input) and 2 hidden units. We'll provide you with the skeleton of the code as well as some code to generate the visualization from the previous notebook.

In [7]:

```
# start from this and modify it
class LogisticRegressionPytorch(nn.Module):
    def __init__(self):
        super(LogisticRegressionPytorch, self).__init__()
        self.linear = nn.Linear(3,1)

    def forward(self, X):
        """ Propagate data through the network.

        This model first applies the linear layer and then a sigmoid
        """
        X = self.linear(X)
        return torch.sigmoid(X)
```

Notebook Exercise 3 (20 minutes + 20 minutes of optional work)

Non-optional: modify the starter code listed above (where it says `class LogisticRegressionPytorch`) to create a class called `TitanicMLP` that has 2 input units and 2 hidden units. Train your network on the Titanic dataset (we recommend using the Adam optimizer you learned about in the last problem). We have defined a function called `visualize_model_probs` for visualizing the probability plot that we saw in the last companion notebook (this code should be run after the model is done training).

Optional: visualize the hidden unit representations in the network (similar to what we did in the companion notebook last time). Hint: you can define additional functions (rather than just `forward` to compute the hidden units in the network).

In [8]:

```
def visualize_model_probs(model):
    xx, yy = np.mgrid[-.1:1.1:.01, 0:85:.1]
    grid = np.c_[xx.ravel(), yy.ravel()]
    probs = model(torch.Tensor(grid)).detach().numpy().reshape(xx.shape)

    f, ax = plt.subplots(figsize=(8, 6))
    contour = ax.contourf(xx, yy, probs, 25, cmap="RdBu",
                          vmin=0, vmax=1)
    ax_c = f.colorbar(contour)
    ax_c.set_label("$P(\text{survived})$")
    ax_c.set_ticks([0, .25, .5, .75, 1])

    ax.scatter(experiment_1_data['male'], experiment_1_data['Age'], c=experiment_1_outputs, s=50,
               cmap="RdBu", vmin=-.2, vmax=1.2,
               edgecolor="white", linewidth=1)

    ax.set(xlim=(-.1, 1.1),
           ylim=(0, 85),
           xlabel="is male", ylabel="age (years)")
    plt.show()
```

In [9]:

```
class TitanicMLP(nn.Module):
    def __init__(self):
        super(TitanicMLP, self).__init__()
        self.linear_1 = nn.Linear(2,2)
        self.linear_2 = nn.Linear(2,1)
```

```

def forward(self, X):
    """ Propagate data through the network.

    This model first applies the linear layer, a sigmoid, a linear
    layer, and finally a sigmoid
    """
    X = self.linear_1(X)
    X = torch.sigmoid(X)
    X = self.linear_2(X)
    return torch.sigmoid(X)
# Note: the hidden function is only needed if you want to visualize hidden
# layers. If you didn't do that part, you wouldn't have this in your
# solution.
def hidden(self, X):
    """ Propagate data to the hidden layer """
    X = self.linear_1(X)
    return torch.sigmoid(X)

```

In [10]:

```

mlp_pytorch = TitanicMLP()
mlp_pytorch.train()
optimizer = torch.optim.Adam(mlp_pytorch.parameters())
criterion = torch.nn.BCELoss()
grad_magnitudes = []

X_data = Variable(torch.Tensor(np.array(experiment_1_data.drop('is_young_male',axis=1))))
y_data = Variable(torch.Tensor(np.array(experiment_1_outputs)))
for epoch in range(20000):
    optimizer.zero_grad()
    # Forward pass
    y_pred = mlp_pytorch(X_data)
    # Compute Loss
    loss = criterion(y_pred, y_data)
    # Backward pass
    loss.backward()
    for name, param in mlp_pytorch.named_parameters():
        if name == 'linear_1.weight':
            grad_magnitudes.append(np.abs(param.grad.numpy()).mean())

    if epoch % 1000 == 0:
        print("epoch", epoch)
        for name, param in mlp_pytorch.named_parameters():
            print(name, "value", param.data, "gradient", param.grad)
    optimizer.step()

plt.plot(grad_magnitudes)
plt.show()
visualize_model_probs(mlp_pytorch)

```

/anaconda3/lib/python3.6/site-packages/torch/nn/functional.py:2016: UserWarning: Using a target size (torch.Size([714])) that is different to the input size (torch.Size([714, 1])) is deprecated. Please ensure they have the same size.
 "Please ensure they have the same size.".format(target.size(), input.size()))

```

epoch 0
linear_1.weight value tensor([[ 0.0783, -0.3369],
                             [-0.6310, -0.5063]]) gradient tensor([[1.6124e-04, 3.6440e-04],
                             [1.7096e-05, 5.6956e-05]])
linear_1.bias value tensor([-0.4244, -0.3620]) gradient tensor([3.4531e-04, 3.6274e-05])
linear_2.weight value tensor([[ -0.1838, -0.0255]]) gradient tensor([[ -0.0027, -0.0018]])
linear_2.bias value tensor([0.0716]) gradient tensor([0.1109])
epoch 1000
linear_1.weight value tensor([[ 1.0951, -0.4031],
                             [-1.7972, -0.0024]]) gradient tensor([[ -3.5508e-04, 3.5121e-05],
                             [ 2.8090e-02, -2.2486e-04]])
linear_1.bias value tensor([0.6286, 0.4005]) gradient tensor([-0.0009, -0.0002])
linear_2.weight value tensor([[0.9174, 1.5132]]) gradient tensor([[ -0.0032, -0.0218]])
linear_2.bias value tensor([-0.7331]) gradient tensor([0.0457])
epoch 2000
linear_1.weight value tensor([[ 2.1469, -0.5015],
                             [-2.8846, 0.0098]]) gradient tensor([[ -3.7360e-04, -8.1380e-07],
                             [ 1.7378e-02, -2.3938e-05]])
linear_1.bias value tensor([0.9809, 0.7351]) gradient tensor([ 2.5607e-04, -2.0111e-05])

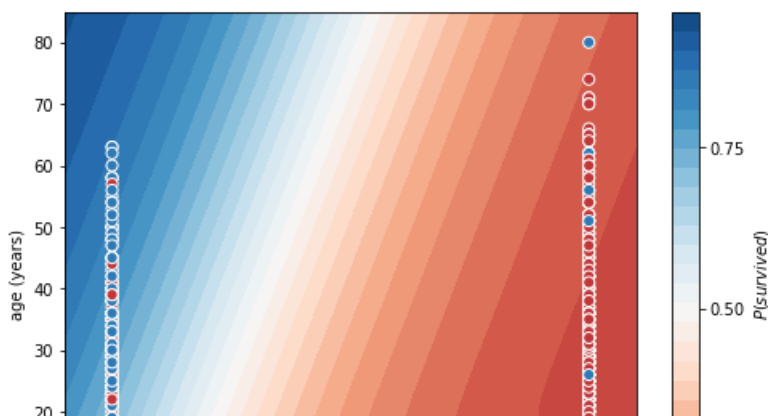
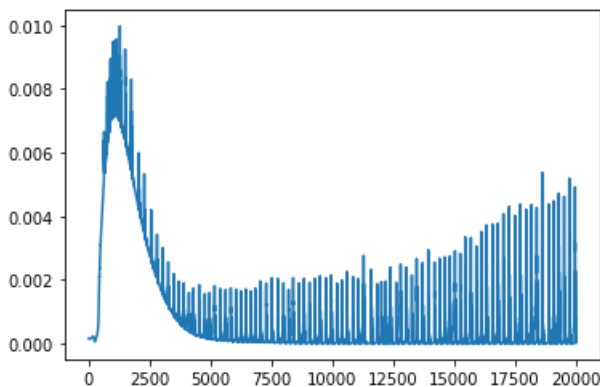
```

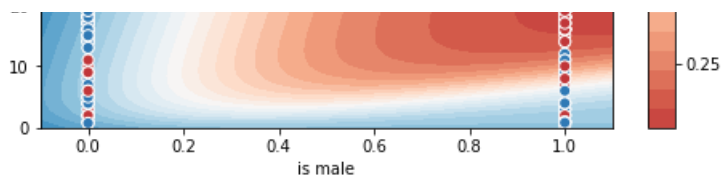
```
linear_2.weight value tensor([[1.6476, 2.7293]]) gradient tensor([[ -0.0014, -0.0179]])
linear_2.bias value tensor([-1.3760]) gradient tensor([0.0188])
epoch 3000
linear_1.weight value tensor([[ 3.2709, -0.4568],
                               [-3.4559,  0.0198]]) gradient tensor([[ -3.4939e-04, -2.4040e-06],
                               [ 6.0881e-03, -7.9927e-05]])
linear_1.bias value tensor([0.0602, 0.6662]) gradient tensor([4.2575e-04, 2.2223e-05])
linear_2.weight value tensor([[2.2241, 3.4301]]) gradient tensor([[ -0.0008, -0.0076]])
linear_2.bias value tensor([-1.7176]) gradient tensor([0.0052])
epoch 4000
linear_1.weight value tensor([[ 3.9437, -0.4280],
                               [-3.7260,  0.0250]]) gradient tensor([[ -1.1966e-04, -4.4906e-07],
                               [ 1.6162e-03,  2.8851e-05]])
linear_1.bias value tensor([-0.6584,  0.4778]) gradient tensor([1.7836e-04, 1.8110e-05])
linear_2.weight value tensor([[2.5620, 3.8046]]) gradient tensor([[ -0.0002, -0.0025]])
linear_2.bias value tensor([-1.8323]) gradient tensor([0.0004])
epoch 5000
linear_1.weight value tensor([[ 4.2070, -0.4196],
                               [-3.8430,  0.0255]]) gradient tensor([[ -2.4529e-05, -6.5224e-07],
                               [ 4.9488e-04, -2.5686e-04]])
linear_1.bias value tensor([-0.9889,  0.3030]) gradient tensor([3.9215e-05, 1.8510e-06])
linear_2.weight value tensor([[2.6310, 3.9967]]) gradient tensor([[ -9.6581e-07, -8.6432e-04]])
linear_2.bias value tensor([-1.8130]) gradient tensor([-0.0004])
epoch 6000
linear_1.weight value tensor([[ 4.2803, -0.4215],
                               [-3.9182,  0.0248]]) gradient tensor([[ -2.1080e-06,  3.7682e-08],
                               [ 2.4280e-04, -1.0977e-05]])
linear_1.bias value tensor([-1.0569,  0.1765]) gradient tensor([-2.0427e-06,  4.6357e-06])
linear_2.weight value tensor([[2.6129, 4.1229]]) gradient tensor([[ 1.0606e-05, -3.9720e-04]])
linear_2.bias value tensor([-1.7739]) gradient tensor([-0.0002])
epoch 7000
linear_1.weight value tensor([[ 4.2658, -0.4250],
                               [-3.9823,  0.0243]]) gradient tensor([[ 2.7147e-06,  3.7136e-08],
                               [ 1.2849e-04, -1.8030e-06]])
linear_1.bias value tensor([-1.0126,  0.0853]) gradient tensor([-4.4638e-06,  3.3051e-06])
linear_2.weight value tensor([[2.5884, 4.2237]]) gradient tensor([[ 4.2278e-06, -2.0000e-04]])
linear_2.bias value tensor([-1.7457]) gradient tensor([-7.6671e-05])
epoch 8000
linear_1.weight value tensor([[ 4.2150, -0.4255],
                               [-4.0372,  0.0237]]) gradient tensor([[2.1690e-06, 1.5692e-06],
                               [6.5628e-05, 4.2954e-04]])
linear_1.bias value tensor([-0.9729,  0.0188]) gradient tensor([-1.0133e-06,  1.6265e-05])
linear_2.weight value tensor([[2.5781, 4.3074]]) gradient tensor([[ 7.1179e-07, -9.0927e-05]])
linear_2.bias value tensor([-1.7255]) gradient tensor([-1.4759e-05])
epoch 9000
linear_1.weight value tensor([[ 4.1725, -0.4245],
                               [-4.0799,  0.0233]]) gradient tensor([[7.3228e-07, 2.1322e-08],
                               [2.7750e-05, 4.5947e-06]])
linear_1.bias value tensor([-0.9589, -0.0281]) gradient tensor([-2.0795e-07,  1.9030e-06])
linear_2.weight value tensor([[2.5762, 4.3755]]) gradient tensor([[ 1.1726e-07, -4.8460e-05]])
linear_2.bias value tensor([-1.7114]) gradient tensor([-1.3222e-05])
epoch 10000
linear_1.weight value tensor([[ 4.1528, -0.4239],
                               [-4.1057,  0.0228]]) gradient tensor([[ 3.7247e-08, -1.5200e-06],
                               [ 6.7470e-06, -4.2325e-04]])
linear_1.bias value tensor([-0.9557, -0.0596]) gradient tensor([-1.8007e-07, -1.2192e-05])
linear_2.weight value tensor([[2.5747, 4.4282]]) gradient tensor([[ -5.5341e-08, -3.0997e-05]])
linear_2.bias value tensor([-1.7027]) gradient tensor([-2.1252e-05])
epoch 11000
linear_1.weight value tensor([[ 4.1479, -0.4238],
                               [-4.1099,  0.0225]]) gradient tensor([[ -3.7065e-09, -2.7059e-08],
                               [-2.3399e-06, -8.7875e-06]])
linear_1.bias value tensor([-0.9570, -0.0794]) gradient tensor([2.9888e-08, 4.6006e-07])
linear_2.weight value tensor([[2.5735, 4.4680]]) gradient tensor([[ 2.3152e-08, -1.0761e-05]])
linear_2.bias value tensor([-1.6990]) gradient tensor([-9.2937e-07])
epoch 12000
linear_1.weight value tensor([[ 4.1521, -0.4239],
                               [-4.0898,  0.0221]]) gradient tensor([[ -1.7622e-06, -2.1967e-05],
                               [-1.4893e-05, -6.3051e-03]])
linear_1.bias value tensor([-0.9614, -0.0921]) gradient tensor([-2.5697e-06, -1.9980e-04])
linear_2.weight value tensor([[2.5732, 4.5027]]) gradient tensor([[ -2.1108e-06, -1.3055e-04]])
linear_2.bias value tensor([-1.6997]) gradient tensor([-0.0002])
epoch 13000
linear_1.weight value tensor([[ 4.1638, -0.4240],
                               [-4.0440,  0.0219]]) gradient tensor([[ -4.6930e-08,  3.7338e-08],
                               [-6.4580e-06,  1.0225e-05]])
linear_1.bias value tensor([-0.9694, -0.1043]) gradient tensor([4.8225e-08, 1.0827e-06])
```

```

linear_2.weight value tensor([[2.5736, 4.5471]]) gradient tensor([[ -1.2660e-09, -6.1914e-06]])
linear_2.bias value tensor([-1.7040]) gradient tensor([1.3856e-06])
epoch 14000
linear_1.weight value tensor([[ 4.1836, -0.4243],
                               [-3.9671,  0.0213]]) gradient tensor([[ -1.0517e-07, -6.4353e-07],
                               [-6.5725e-06, -1.8544e-04]])
linear_1.bias value tensor([-0.9825, -0.1228]) gradient tensor([-3.7626e-08, -4.4789e-06])
linear_2.weight value tensor([[2.5744, 4.6186]]) gradient tensor([[ -8.3979e-08, -1.0333e-05]])
linear_2.bias value tensor([-1.7119]) gradient tensor([-5.6618e-06])
epoch 15000
linear_1.weight value tensor([[ 4.2132, -0.4246],
                               [-3.8503,  0.0205]]) gradient tensor([[ -1.8978e-08,  3.3623e-07],
                               [-5.3742e-06,  1.0696e-04]])
linear_1.bias value tensor([-1.0018, -0.1521]) gradient tensor([7.9998e-08, 5.5730e-06])
linear_2.weight value tensor([[2.5760, 4.7343]]) gradient tensor([[ 2.7678e-08, -6.0736e-06]])
linear_2.bias value tensor([-1.7249]) gradient tensor([5.2253e-06])
epoch 16000
linear_1.weight value tensor([[ 4.2509, -0.4251],
                               [-3.6980,  0.0194]]) gradient tensor([[ -2.6985e-08,  1.3053e-07],
                               [-3.9564e-06,  4.4197e-05]])
linear_1.bias value tensor([-1.0265, -0.1933]) gradient tensor([5.1067e-08, 4.6999e-06])
linear_2.weight value tensor([[2.5783, 4.9019]]) gradient tensor([[ 8.8476e-09, -8.3302e-06]])
linear_2.bias value tensor([-1.7437]) gradient tensor([3.2760e-06])
epoch 17000
linear_1.weight value tensor([[ 4.2892, -0.4255],
                               [-3.5372,  0.0182]]) gradient tensor([[ -1.0974e-06, -1.2230e-05],
                               [-1.0328e-05, -4.2157e-03]])
linear_1.bias value tensor([-1.0518, -0.2429]) gradient tensor([-1.5030e-06, -1.2597e-04])
linear_2.weight value tensor([[2.5810, 5.1075]]) gradient tensor([[ -1.5310e-06, -6.9310e-05]])
linear_2.bias value tensor([-1.7656]) gradient tensor([-0.0001])
epoch 18000
linear_1.weight value tensor([[ 4.3218, -0.4258],
                               [-3.3945,  0.0173]]) gradient tensor([[ -8.6642e-08, -6.9989e-07],
                               [-1.9062e-06, -2.4824e-04]])
linear_1.bias value tensor([-1.0736, -0.2932]) gradient tensor([-6.3623e-08, -3.1157e-06])
linear_2.weight value tensor([[2.5836, 5.3223]]) gradient tensor([[ -9.5140e-08, -1.2591e-05]])
linear_2.bias value tensor([-1.7869]) gradient tensor([-5.5695e-06])
epoch 19000
linear_1.weight value tensor([[ 4.3477, -0.4260],
                               [-3.2758,  0.0165]]) gradient tensor([[ -2.0641e-08, -9.0367e-09],
                               [-9.3695e-07,  7.0594e-07]])
linear_1.bias value tensor([-1.0912, -0.3406]) gradient tensor([1.7099e-08, 4.5355e-06])
linear_2.weight value tensor([[2.5859, 5.5302]]) gradient tensor([[ -6.1555e-09, -8.8378e-06]])
linear_2.bias value tensor([-1.8060]) gradient tensor([1.7219e-06])

```





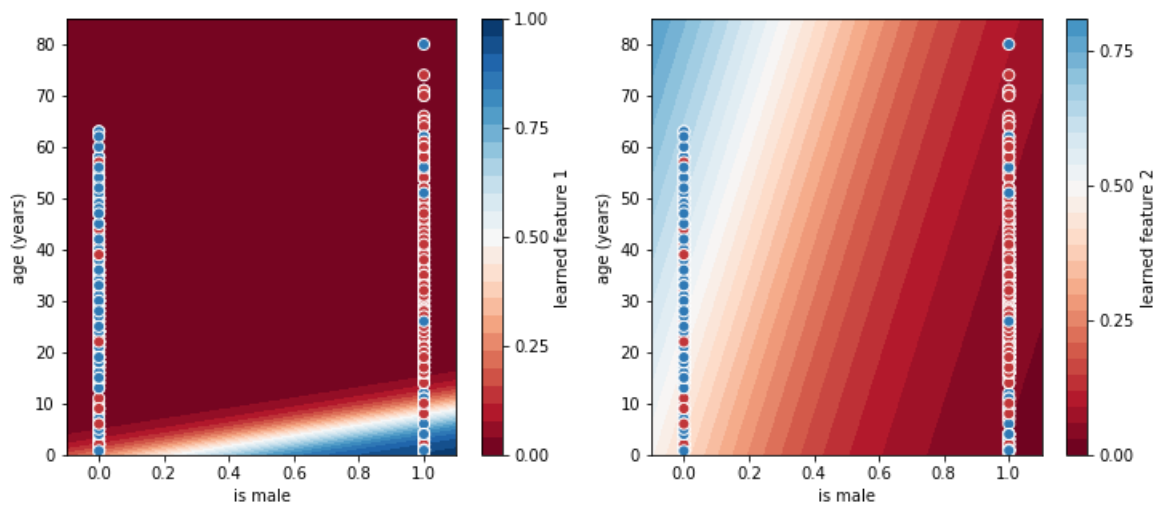
The next code block visualizes the hidden layer activations.

In [11]:

```
f = plt.figure(figsize=(12, 5))
xx, yy = np.mgrid[-.1:1.1:.01, 0:85:.1]
grid = np.c_[xx.ravel(), yy.ravel()]
hidden_units = mlp_pytorch.hidden(Variable(torch.Tensor(grid))).detach().numpy()
for i in range(2):
    ax = f.add_subplot(1,2,i+1)
    contour = ax.contourf(xx, yy, hidden_units[:,i].reshape(xx.shape), 25, cmap="RdBu",
                          vmin=0, vmax=1)

    ax.scatter(experiment_1_data['male'], experiment_1_data['Age'], c=experiment_1_outputs, s=50,
               cmap="RdBu", vmin=-.2, vmax=1.2,
               edgecolor="white", linewidth=1)
    ax_c = f.colorbar(contour)
    ax_c.set_label("learned feature %d" % (i+1))
    ax_c.set_ticks([0, .25, .5, .75, 1])

    ax.set(xlim=(-.1, 1.1),
           ylim=(0, 85),
           xlabel="is male", ylabel="age (years)")
plt.show()
```



In [0]: