

به نام خدا

## تمرین ششم یادگیری عمیق

غزل زمانی نژاد

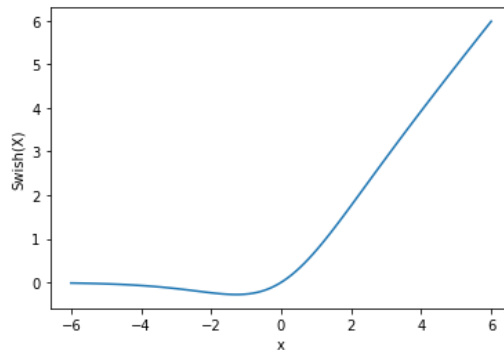
۹۷۵۲۲۱۶۶

### 1. الف) Swish:

$$f(x) = x * \text{sigmoid}(x)$$

تابع sigmoid و swish را تعریف میکنیم و بعد به کمک matplotlib آن را رسم میکنیم.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import math
4
5 def sigmoid(x):
6     return 1 / (1 + np.exp(-x))
7
8 def swish(x):
9     return x * sigmoid(x)
10
11 x = np.linspace(-6, 6, 100)
12 y = swish(x)
13
14 plt.plot(x, y)
15 plt.xlabel("x")
16 plt.ylabel("Swish(X)")
17
18 plt.show()
```



### Mish:

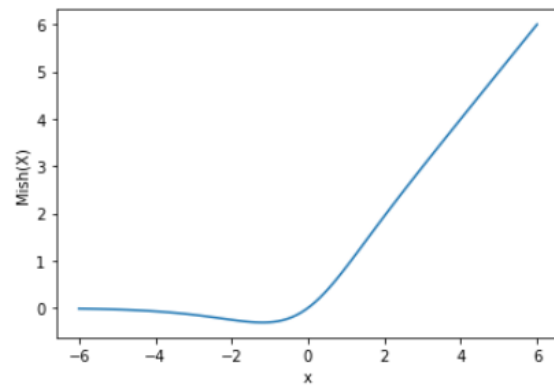
$$f(x) = x * \tanh(\text{softplus}(x)) = x * \tanh(\ln(1 + e^x))$$

تابع softplus و mish را تعریف میکنیم و بعد به کمک matplotlib آن را رسم میکنیم.

```

1 def softplus(x):
2     return np.log(1 + np.exp(x))
3
4 def mish(x):
5     return x * np.tanh(softplus(x))
6
7 # x2 = np.linspace(-3, 3, 100)
8 y2 = mish(x)
9
10 plt.plot(x, y2)
11 plt.xlabel("x")
12 plt.ylabel("Mish(X)")
13
14 plt.show()

```



ب) ابتدا مشتق swish را محاسبه میکنیم:

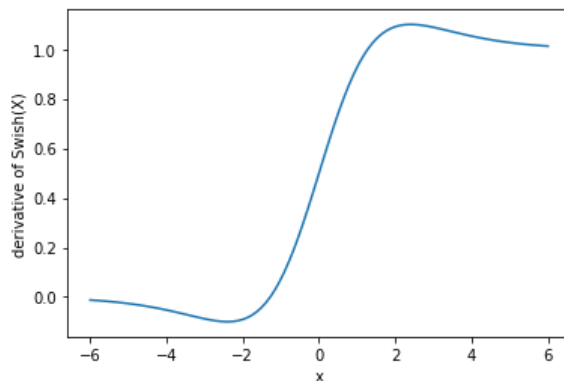
$$\begin{aligned}
 \text{Swish}(x) &= x \times \sigma(x) \\
 \xrightarrow{\text{derivative}} \quad 1 \cdot \sigma(x) + x (\sigma(x) (1 - \sigma(x))) &= \sigma(x) + x \sigma(x) + x \sigma^2(x) \\
 &= \underbrace{x \sigma(x)}_{\text{swish}(x)} + \sigma(x) \underbrace{[1 - x \sigma(x)]}_{\text{swish}(x)} \\
 &= \text{Swish}(x) + \sigma(x) (1 - \text{swish}(x))
 \end{aligned}$$

سپس آن را رسم میکنیم:

```

1 y2 = swish(x) + sigmoid(x) * (1 - swish(x))
2
3 plt.plot(x, y2)
4 plt.xlabel("x")
5 plt.ylabel("derivative of Swish(X)")
6
7 plt.show()

```

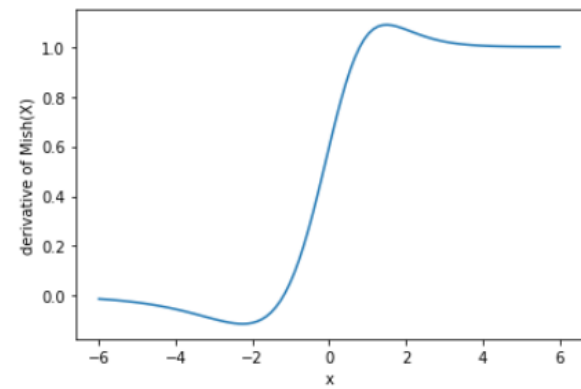


ابتدا مشتق mish را محاسبه میکنیم:

$$\begin{aligned} \text{mish}(x) &= x \tanh(\text{softplus}(x)) \\ \xrightarrow{\text{derivative}} \quad & \tanh(\text{softplus}(x)) + x \underbrace{(1 - \tanh^2(\text{softplus}(x)))}_{\text{sech}^2(\text{softplus}(x))} \frac{e^x}{1+e^x} = \\ &= \frac{\text{mish}(x)}{x} + x \cdot \text{sech}^2(\text{softplus}(x)) \text{sigmoid}(x) \end{aligned}$$

سپس آن را رسم میکنیم:

```
1 y4 = mish(x) / x + x * (1 - np.tanh(softplus(x))**2) * sigmoid(x)
2
3 plt.plot(x, y4)
4 plt.xlabel("x")
5 plt.ylabel("derivative of Mish(X)")
6
7 plt.show()
```



ت) مزیت ReLU نسبت به sigmoid و tanh: در شبکه های عمیق، ReLU نسبت به sigmoid و tanh عملکرد بهتری دارد زیرا با مشکل vanishing gradients یعنی کوچک شدن بیش از حد مقادیر گرادیان ها رو به رو نمیشویم و برای مقادیر مثبت، گرادیان در شبکه جریان پیدا میکند. به دلیل اینکه ReLU از بالا بیکران است، از saturation برای مقادیر مثبت ورودی جلوگیری میکند. دلیل دیگر برتری ReLU آن است که تعمیم بهتری دارد و سرعت همگرایی آن بیشتر است. همچنین بسیاری از متخصصین سادگی و قابل اطمینان بودن ReLU را به سایر توابع فعالسازی ترجیح میدهند زیرا در بسیاری از موارد بهبود عملکرد سایر توابع به مدل و دیتاست وابسته است.

Swish کدام مشکلات ReLU را رفع کرده و چه شباهت هایی با آن دارد؟: این تابع همچون ReLU از بالا بیکران و از پایین کران دار است. اما برخلاف ReLU غیریکنواخت است. تمامی ویژگی های swish یعنی کران دار از پایین، بی کران از بالا، غیریکنواختی و هموار بودن آن مزیت محسوب میشوند. به بررسی آنها میپردازیم: بیکران بودن از بالا یک ویژگی مهم محسوب میشود چون مشکل saturation آموزش کند برای گرادین های نزدیک 0 را برطرف کرده.

کران دار بودن از پایین اثر قوی regularization دارد. توابعی که حد آنها 0 است تاثیر به سزای regularization دارند زیرا مقادیر منفی بزرگ در آنها فراموش میشود. swish برخلاف ReLU برای ورودی های منفی کوچک، مقداری منفی است و باعث غیریکنواختی میشود. همین امر باعث تقویت جریان یافتن گرادین ها میشود. در تابع ReLU، مقدار خروجی به ازای تمامی مقادیر منفی 0 است. اگر شبکه به جایی برسد که همیشه مقادیر منفی داشته باشد یعنی ReLU به طور مداوم مقدار 0 را برگرداند، باعث بروز پدیده dying ReLU میشود و شبکه قدرت طبقه بندی بین داده ها را از دست میدهد. با استفاده از swish از این پدیده جلوگیری میشود. همچنین باعث بالا بردن robustness در برابر مقداردهی های اولیه پارامترها و نرخ یادگیری های متفاوت میشود.

هموار بودن (smoothness) نقش مهمی در بهینه سازی و تعمیم دارد. همواری خروجی بر همواری میزان ضرر تاثیر مستقیم میگذارد. اگر ضرر هموار باشد در بهینه سازی راحت تر است و حساسیت مقداردهی اولیه پارامترها و نرخ یادگیری کاهش می یابد. مقایسه در batch size های مختلف: همانطور که انتظار میرود با افزایش اندازه batch، میزان عملکرد ReLU و swish کاهش می یابد. با این وجود برای batch size های متفاوت (چه کوچک و چه بزرگ) swish عملکرد بهتری دارد.

Mish کدام مشکلات ReLU را رفع کرده و چه شباهت هایی با آن دارد؟: تابع mish هموار، پیوسته، self-regularized و غیریکنواخت است. همچون swish از بالا بی کران و از پایین کران دار است (تمامی مزیت های swish را دارد). Mish با داشتن یک پیش شرط یعنی همان  $\Delta(x)$  از dying ReLU جلوگیری می کند. و باعث جریان یافتن بهتر گرادین در شبکه می شود.

بی کران بودن از بالا از saturation جلوگیری می کند (وجود saturation موجب آهسته شدن فرآیند یادگیری در گرادین های نزدیک 0 می شود).

با وجود این موارد، در صورت استفاده از تابع فعال سازی mish یک trade-off میان بالا رفتن دقت و بالا رفتن هزینه های محاسباتی وجود دارد.

ث) اگر پارامترهای swish را به شکل زیر تغییر دهیم به ReLU شباهت بیشتری پیدا میکند:

$$f(x, \beta) = 2x * \text{sigmoid}(\beta * x)$$

در این تابع اگر  $\beta = 0$  باشد،  $\text{swish} = x$  و خطی میشود.

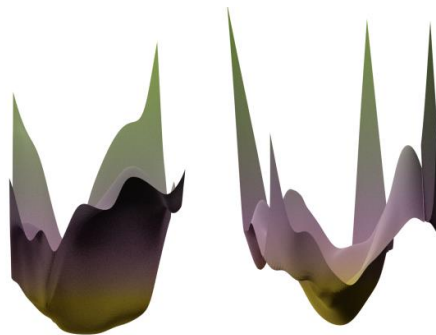
اگر  $\beta \rightarrow \infty$ ، سیگموئید به تابع 1-0 تبدیل میشود و swish شبیه ReLU میشود.

اگر  $\beta$  به عنوان یک پارامتر قابل آموزش تعیین شود، میتوانیم درجه برون یابی مدل را با آن کنترل کنیم. به این گونه از تابع swish بدون ضریب 2، swish-beta میگوییم. مزیت تابع swish-beta این است که میتواند میان تابع خطی و ReLU تغییر کند.

ج) تابع mish بسیار به تابع swish شباهت دارد. در مشتق اول تابع mish، ضریبی از تابع swish وجود دارد:

$$f'(x) = \text{sech}^2(\text{softplus}(x))x\text{sigmoid}(x) + \frac{f(x)}{x}$$
$$= \Delta(x)\text{swish}(x) + \frac{f(x)}{x}$$

$\Delta(x)$  موجود در این رابطه باعث هموارتر شدن گرادیان ها نسبت به گرادیان های تابع swish می شود. این ضریب اثر regularization دارد و باعث میشود بهینه سازی کانتور بهتر انجام شود. و در نهایت موجب می شود تابع mish در شبکه های عمیق و پیچیده نسبت به swish برتری داشته باشد. شکل زیر به خوبی نشان دهنده این برتری است. در سمت راست تابع ضرر با تابع فعال سازی swish و در سمت چپ تابع ضرر با تابع فعال سازی mish دیده می شود. تصویر سمت چپ (mish) هموارتر است.



left: mish, right: swish

2. الف) فرمول تابع ضرر MSE:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

فرمول تابع ضرر Binary Cross Entropy:

$$Loss = -\frac{1}{\text{output size}} \sum_{i=1}^{\text{output size}} y_i \cdot \log \hat{y}_i + (1 - y_i) \cdot \log (1 - \hat{y}_i)$$

همانطور که در فرمول های فوق مشاهده می شود، مقدار تابع MSE حداکثر می تواند 1 باشد و به دلیل کوچک بودن مقادیر گرادیان دیرتر بهینه می شود. اما Binary Cross Entropy نسبت به خطا بسیار حساس است و حتی مقدار loss می تواند به بی نهایت میل کند. به همین دلیل در epoch اول میان این دو مقدار اختلاف دیده می شود. در مسائل دسته بندی هدف بالا بردن مقدار دقت است و نه کم کردن میزان ضرر. یعنی با مقایسه نمودارهای loss نمی توانیم بگوییم کدام یک عملکرد بهتری داشته است.

در ابتدای کار که هنوز شبکه پارامترها را یاد نگرفته، خروجی 0.5 است (مستقل از داده شبکه احتمال 50 درصد را پیش بینی می کند). خطای MSE را محاسبه می کنیم (اختلاف بین پیش بینی شبکه و مقدار واقعی در هر صورت 0.5 است):

$$MSE = (1 / n) (0.5)^2 = 0.25$$

خطای Binary Cross Entropy را محاسبه می کنیم:

$$Y = 1: \text{Binary Cross Entropy} = -1 * \ln(0.5) - 0 = \ln(2) = 0.69$$

$$Y = 0: \text{Binary Cross Entropy} = 0 - (1 - 0) * \ln(1 - 0.5) = \ln(2) = 0.69$$

مشاهده می شود این مقادیر با مقادیر نمودارها در epoch نخست یکی است.

ب) همانطور که در نمودار دیده می شود، تابع MSE پس از 100 epoch همچنان آموزش می بیند و به همگرایی نرسیده به همین دلیل اختلاف میان نمودار اعتبارسنجی و آموزش کمتر است. اما در نمودار Binary Cross Entropy بعد از آموزش دیدن در تعدادی epoch، شبکه به سمت overfit شدن می رود. یعنی الگوها را بیشتر حفظ می کند و آنها را به درستی یاد نمی گیرد تا تعمیم دهد. به همین دلیل در نمودار سمت راست، بعد از مدت خطای آموزش کاهش می باید اما خطای اعتبارسنجی تقریباً ثابت می ماند و میان این دو نمودار فاصله ایجاد می شود.

ج) Binary Cross Entropy: در این نمودار تقریباً بعد از 60 epoch مشاهده می شود که مقدار خطای اعتبارسنجی رو به افزایش می رود. یعنی در آن نقطه شبکه به سمت overfit شدن می رود. پس بهتر است آموزش را متوقف کنیم تا شبکه الگوها را حفظ نکند و بتواند تعمیم دهد.

MSE: در این نمودار هم خطای آموزش و هم خطای اعتبارسنجی همچنان در حال کاهش هستند. یعنی شبکه هنوز underfit است و مشکل بایاس داریم. پس بهتر است همچنان به آموزش ادامه دهیم تا زمانی که خطای اعتبارسنجی به کمترین مقدار خود برسد (سرعت یادگیری در MSE نسبتاً کم است و به همین دلیل بعد از epoch100 شبکه هنوز در حال آموزش دیدن است).

3. سلول اول: دیتاست MNIST که شامل رقم های 0 تا 9 است از طریق کراس load می کنیم و شکل داده های آزمون و تست را چاپ می کنیم.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import random
4 from keras.datasets import mnist
5
6 (x_train, y_train), (x_test, y_test) = mnist.load_data()
7
8 print("train input shape", x_train.shape)
9 print("train label shape", y_train.shape)
10 print("test input shape", x_test.shape)
11 print("test label shape", y_test.shape)
```

Downloading data from <https://storage.googleapis.com/tensorflow/tf-11493376/11490434> [=====] - 0s 0us/step  
11501568/11490434 [=====] - 0s 0us/step  
train input shape (60000, 28, 28)  
train label shape (60000,)  
test input shape (10000, 28, 28)  
test label shape (10000,)

سلول دوم: برای اینکه داده های آزمون بیشتری در اختیار داشته باشیم، از data augmentation استفاده می کنیم. در اینجا رنگ هر پیکسل را invert می کنیم. یعنی اگر سیاه باشد آن را سفید، و اگر سفید باشد آن را سیاه می کنیم. با استفاده از این کار تعداد داده ها دو برابر می شود. پس از تولید تصاویر جدید، آنها را به دیتاست اصلی اضافه می کنیم (concatenate).

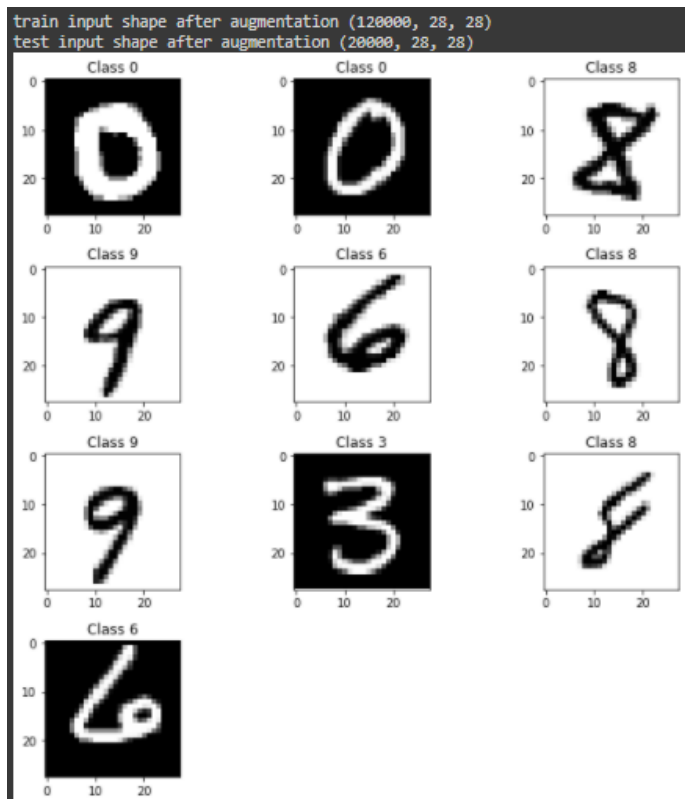
```
1 # data augmentation
2 x_train_aug = np.subtract(255, x_train)
3 x_test_aug = np.subtract(255, x_test)
4
5 new_x_train = np.concatenate((x_train, x_train_aug), axis=0)
6 new_y_train = np.concatenate((y_train, y_train), axis=0)
7 new_x_test = np.concatenate((x_test, x_test_aug), axis=0)
8 new_y_test = np.concatenate((y_test, y_test), axis=0)
```

سلول سوم: داده های آموزشی و آزمون را shuffle می کنیم تا در هر مجموعه داده ها ترتیب خاصی نداشته باشند. در اینجا آرگومان random\_state را عدد خاصی تنظیم می کنیم تا برچسب متناظر با هر داده اشتباه نشود و لیست برچسب ها هم به همان صورت ترکیب شوند.

```
1 from sklearn.utils import shuffle
2 new_x_train, new_y_train = shuffle(new_x_train, new_y_train, random_state=10)
3 new_x_test, new_y_test = shuffle(new_x_test, new_y_test, random_state=8)
```

سلول چهارم: پس از مخلوط کردن داده ها شکل آنها و همچنین 10 تا از داده های ورودی با برچسب اصلی شان را چاپ می کنیم.

```
1 print("train input shape after augmentation", new_x_train.shape)
2 print("test input shape after augmentation", new_x_test.shape)
3
4 plt.rcParams['figure.figsize'] = (9,9) # Make the figures a bit bigger
5
6 for i in range(10):
7     plt.subplot(4,3,i+1)
8     plt.imshow(new_x_train[i], cmap='gray', interpolation='none')
9     plt.title("Class {}".format(new_y_train[i]))
10
11 plt.tight_layout()
```





سلول پنجم: برای اینکه داده های ورودی را از لایه کانولوشنی عبور دهیم، باید شکل آنها را تغییر دهیم. برای این کار از reshape استفاده می کنیم.

```
1 shape1 = new_x_train.shape
2 shape2 = new_x_test.shape
3 new_x_train = np.reshape(new_x_train, (shape1[0], shape1[1], shape1[2], 1))
4 new_x_test = np.reshape(new_x_test, (shape2[0], shape2[1], shape2[2], 1))
5
6 print("train input shape after reshape", new_x_train.shape)
7 print("test input shape after reshape", new_x_test.shape)
```

train input shape after reshape (120000, 28, 28, 1)  
test input shape after reshape (20000, 28, 28, 1)

سلول ششم: برای اینکه پراکندگی داده های ورودی زیاد نباشد، آنها را normalize می کنیم. یعنی مقدار تمامی پیکسل های تمام تصاویر را بر 255 تقسیم می کنیم. سپس برچسب هر تصویر را به شکل one hot با استفاده از to\_categorical انکود می کنیم. ورودی اول این تابع برداری است که می خواهیم آن را انکود کنیم و ورودی دوم تعداد کل کلاس هاست که در اینجا 10 رقم داریم پس تعداد کلاس ها نیز 10 است.

```
1 # normalized = (x - min) / (max - min)
2 # in this example = x / 255
3 import tensorflow as tf
4
5 new_x_train = new_x_train / 255
6 new_x_test = new_x_test / 255
7
8 new_y_train = tf.keras.utils.to_categorical(new_y_train, num_classes=10)
9 new_y_test = tf.keras.utils.to_categorical(new_y_test, num_classes=10)
```

سلول هفتم: تابع sequential\_network را تعریف می کنیم. این تابع آلفا را به عنوان ورودی می گیرد و در خروجی مدل از نوع sequential را برمی گرداند. این مدل دارای لایه های زیر است:

لایه کانولوشن اول که دارای تابع فعال سازی Leaky ReLU با پارامتر آلفا است.

لایه کانولوشن دوم که مشابه لایه کانولوشن اول است و تنها kernel size آن با لایه اول متفاوت است.

لایه Flatten که خروجی های لایه قبل را flat می کند.

لایه Dense که خروجی نهایی شبکه (10 نرون) را تشکیل می دهد. دارای تابع فعال سازی softmax است.

```

1 from keras.models import Sequential
2 from keras.layers import Dense, Conv2D, LeakyReLU, Flatten, PReLU
3
4 def sequential_network(alpha):
5     leaky_relu = LeakyReLU(alpha=alpha)
6     model = Sequential()
7     model.add(Conv2D(8, 7, activation=leaky_relu))
8     model.add(Conv2D(8, 5, activation=leaky_relu))
9     model.add(Flatten())
10    model.add(Dense(10, activation='softmax'))
11
12    return model
13

```

برای اینکه از درست بودن لایه های شبکه مطمئن شویم یک نمونه از آن را می سازیم و آن را build می کنیم.  
سپس خلاصه شبکه را چاپ می کنیم.

```

14 test_model = sequential_network(0.3)
15 test_model.build(new_x_train.shape)
16 print(test_model.summary())

```

Model: "sequential\_26"

Layer (type)	Output Shape	Param #
conv2d_52 (Conv2D)	(120000, 22, 22, 8)	400
conv2d_53 (Conv2D)	(120000, 18, 18, 8)	1608
flatten_26 (Flatten)	(120000, 2592)	0
dense_26 (Dense)	(120000, 10)	25930

Total params: 27,938  
 Trainable params: 27,938  
 Non-trainable params: 0

None

سلول هشتم: در آن 3 متد تعریف می کنیم:

Train: به عنوان ورودی مدل را دریافت می کند. آن را با شرایط خواسته شده یعنی تابع ضرر Categorical Cross Entropy و بهینه ساز Adam کامپایل می کنیم. سپس آن را در پنج epoch و با اندازه batch 128 با استفاده از تابع fit آموزش می دهیم.

```

1 epoch_size = 5
2 batch_size = 128
3
4
5 def train(model):
6     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
7     history = model.fit(new_x_train, new_y_train, epochs=epoch_size, batch_size=batch_size)
8     return history

```

Test: مقدار دقت و خطای روی داده های آزمون را محاسبه کرده و باز می گرداند.

```

9
10 def test(model):
11     loss, acc = model.evaluate(new_x_test, new_y_test, verbose = 0)
12     return (loss, acc)

```

Plot: به عنوان ورودی تاریخچه ای از مدل را دریافت می کند تا بتواند نمودار میزان دقت و خطای مدل را در 2 subplot رسم کند.

```

14 def plot(history):
15     figure, (ax1, ax2) = plt.subplots(1, 2)
16
17     ax1.plot(history.history['accuracy'])
18     ax1.set_title('model accuracy')
19
20     ax2.plot(history.history['loss'])
21     ax2.set_title('model loss')
22     figure.tight_layout()
23     plt.show()

```

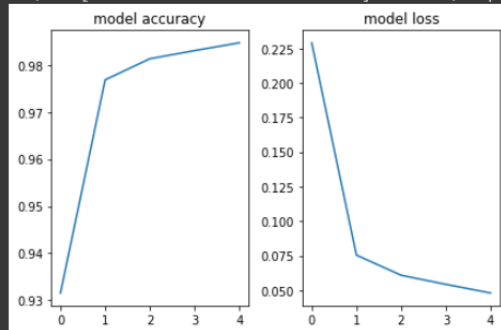
سلول نهم: مقادیر آلفا داده شده را در یک لیست میریزیم تا به کمک یک حلقه مدل ها با آلفای متفاوت را آموزش دهیم و نتیجه را چاپ کنیم. در این حلقه: ابتدا با استفاده از تابع sequential\_network ابتدا یک مدل با آلفای مشخص می سازیم. سپس با استفاده از تابع train آن را آموزش می دهیم. بعد از تابع test استفاده کرده و مقدار دقت و خطای مدل روی داده آزمون را می یابیم. بعد با استفاده از plot نمودارها را رسم می کنیم. نتایج در تصاویر زیر دیده می شود:

مدل اول با  $\alpha = -0.5$

```

model with alpha = -1
Epoch 1/5
938/938 [=====] - 7s 7ms/step - loss: 0.2291 - accuracy: 0.9316
Epoch 2/5
938/938 [=====] - 7s 7ms/step - loss: 0.0755 - accuracy: 0.9769
Epoch 3/5
938/938 [=====] - 7s 7ms/step - loss: 0.0609 - accuracy: 0.9814
Epoch 4/5
938/938 [=====] - 7s 7ms/step - loss: 0.0543 - accuracy: 0.9831
Epoch 5/5
938/938 [=====] - 7s 7ms/step - loss: 0.0480 - accuracy: 0.9847

```



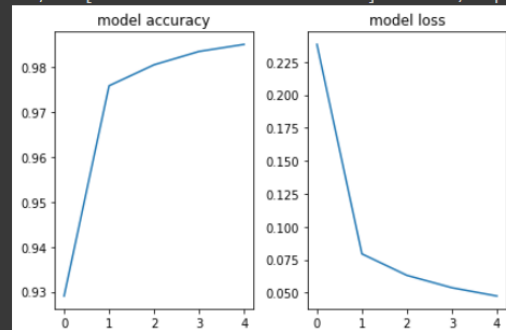
Accuracy on test data: 97.96  
Loss on test data: 0.07

مدل دوم با  $\alpha = -1$

```

model with alpha = -0.5
Epoch 1/5
938/938 [=====] - 7s 7ms/step - loss: 0.2384 - accuracy: 0.9291
Epoch 2/5
938/938 [=====] - 7s 7ms/step - loss: 0.0795 - accuracy: 0.9758
Epoch 3/5
938/938 [=====] - 7s 7ms/step - loss: 0.0632 - accuracy: 0.9805
Epoch 4/5
938/938 [=====] - 7s 7ms/step - loss: 0.0538 - accuracy: 0.9835
Epoch 5/5
938/938 [=====] - 7s 7ms/step - loss: 0.0475 - accuracy: 0.9850

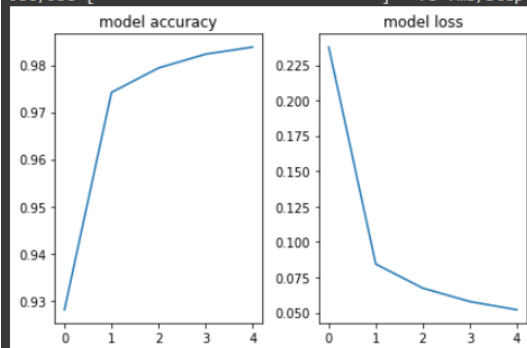
```



Accuracy on test data: 97.94  
Loss on test data: 0.06

مدل سوم با  $\alpha = 0$

```
model with alpha = 0
Epoch 1/5
938/938 [=====] - 7s 7ms/step - loss: 0.2377 - accuracy: 0.9281
Epoch 2/5
938/938 [=====] - 7s 7ms/step - loss: 0.0843 - accuracy: 0.9742
Epoch 3/5
938/938 [=====] - 7s 7ms/step - loss: 0.0673 - accuracy: 0.9794
Epoch 4/5
938/938 [=====] - 7s 7ms/step - loss: 0.0579 - accuracy: 0.9823
Epoch 5/5
938/938 [=====] - 7s 7ms/step - loss: 0.0520 - accuracy: 0.9838
```

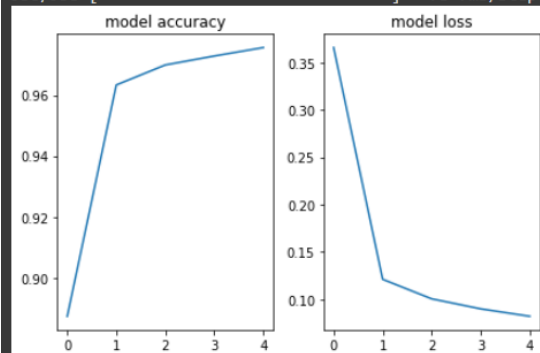


Accuracy on test data: 98.22  
Loss on test data: 0.06

\*\*\*\*\*

مدل چهارم با  $\alpha = 0.5$

```
model with alpha = 0.5
Epoch 1/5
938/938 [=====] - 8s 7ms/step - loss: 0.3658 - accuracy: 0.8875
Epoch 2/5
938/938 [=====] - 7s 7ms/step - loss: 0.1212 - accuracy: 0.9633
Epoch 3/5
938/938 [=====] - 7s 8ms/step - loss: 0.1005 - accuracy: 0.9699
Epoch 4/5
938/938 [=====] - 7s 8ms/step - loss: 0.0899 - accuracy: 0.9728
Epoch 5/5
938/938 [=====] - 7s 7ms/step - loss: 0.0821 - accuracy: 0.9756
```

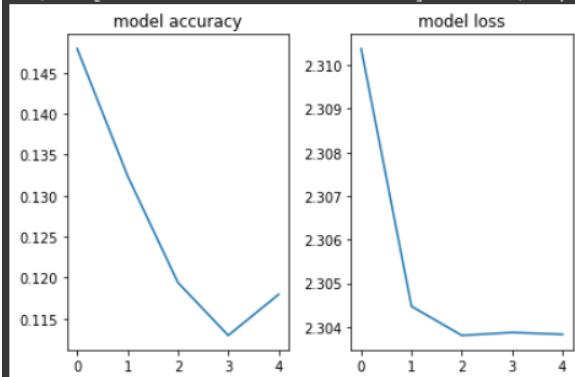


Accuracy on test data: 97.29  
Loss on test data: 0.08

\*\*\*\*\*

مدل پنجم با  $\alpha = 1$

```
model with alpha = 1
Epoch 1/5
938/938 [=====] - 7s 7ms/step - loss: 2.3104 - accuracy: 0.1479
Epoch 2/5
938/938 [=====] - 7s 7ms/step - loss: 2.3045 - accuracy: 0.1324
Epoch 3/5
938/938 [=====] - 7s 7ms/step - loss: 2.3038 - accuracy: 0.1194
Epoch 4/5
938/938 [=====] - 7s 7ms/step - loss: 2.3039 - accuracy: 0.1129
Epoch 5/5
938/938 [=====] - 7s 7ms/step - loss: 2.3038 - accuracy: 0.1179
```



Accuracy on test data: 11.69

Loss on test data: 2.30

\*\*\*\*\*

برای پیدا کردن بهترین مقدار آلفا در میان مقادیری داده شده، دقت مدل ها را روی داده آزمون مقایسه می کنیم. مدل سوم با  $\alpha = 0$  بهترین دقت آزمون را دارد.

در مدل پنجم که  $\alpha = 1$  است، مشاهده می شود که دقت مدل از epoch اول بسیار کم و تقریباً نزدیک به حالت پیش بینی تصادفی است. و در طی epoch های بعدی در حال کاهش است. بررسی دلیل:

می دانیم فرمول تابع فعال سازی Leaky ReLU به صورت زیر است:

$$f(x) = \alpha * x \text{ if } x < 0$$
$$f(x) = x \text{ if } x \geq 0$$

اگر در این فرمول مقدار آلفا را برابر 1 قرار دهیم به شکل زیر در می آید:

$$f(x) = x \text{ if } x < 0$$

$$f(x) = x \text{ if } x \geq 0$$

$$\Rightarrow f(x) = x$$

که این فرمول عملاً یک تابع فعال سازی خطی است و هیچ تغییری در خروجی نوروں ها حاصل نمی شود. یعنی نوروں های خارج شده از لایه کانولوشن به همان صورت به لایه بعدی منتقل می شوند. به همین دلیل شبکه عملکرد بسیار ضعیفی داشته است.

سلول دهم: در این سلول تابعی تعریف می کنیم که به عنوان خروجی یک مدل برمی گرداند. مشابه مدل های قبل است با این تفاوت که برای لایه های کانولوشنی آن از تابع فعال سازی PReLU استفاده می کنیم. برای اینکه هر فیلتر فقط یک مجموعه از پارامترها را داشته باشد مقدار shared\_axes را در آن به [1,2] ست می کنیم.

```
1 def sequential_network_Prelu():
2     prelu = PReLU(shared_axes=[1,2])
3     model = Sequential()
4     model.add(Conv2D(8, 7, activation=prelu))
5     model.add(Conv2D(8, 5, activation=prelu))
6     model.add(Flatten())
7     model.add(Dense(10, activation='softmax'))
8
9     return model
```

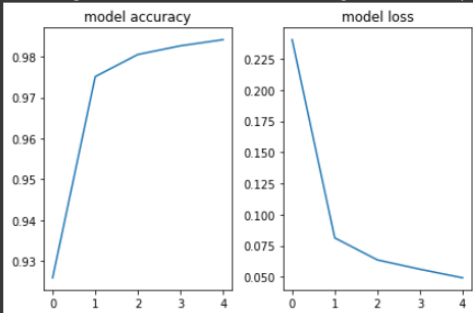
سلول یازدهم: مشابه مراحل، ابتدا یک مدل با تابع فعال سازی PReLU می سازیم. سپس آن را آموزش می دهیم. بعد مقدار دقت و خطا را روی داده آزمون محاسبه می کنیم. و در نهایت نمودارها را رسم می کنیم.

```
1 print("PReLU model")
2 model_prelu = sequential_network_Prelu()
3 hist = train(model_prelu)
4 test_loss, test_acc = test(model_prelu)
5 plot(hist)
6 print('Accuracy on test data: %.2f' % (test_acc*100))
7 print('Loss on test data: %.2f' % (test_loss))
```

```

PReLU model
Epoch 1/5
938/938 [=====] - 8s 8ms/step - loss: 0.2403 - accuracy: 0.9259
Epoch 2/5
938/938 [=====] - 7s 8ms/step - loss: 0.0813 - accuracy: 0.9751
Epoch 3/5
938/938 [=====] - 7s 8ms/step - loss: 0.0637 - accuracy: 0.9806
Epoch 4/5
938/938 [=====] - 7s 8ms/step - loss: 0.0562 - accuracy: 0.9827
Epoch 5/5
938/938 [=====] - 7s 8ms/step - loss: 0.0493 - accuracy: 0.9842

```



Accuracy on test data: 98.39  
Loss on test data: 0.05

سلول دوازهم: می دانیم تابع PReLU تقریباً مشابه Leaky ReLU است. با این تفاوت که در آن پارامتر آلفا از مقادیر trainable شبکه است و شبکه خودش آن را آموزش می بیند. اما در Leaky ReLU مقدار آلفا به صورت دستی tune می شود.

پارامترهای آلفایی که شبکه پس از آموزش به آنها رسیده را چاپ می کنیم.

```

1 layers = model_prelu.layers
2
3 W1 = layers[0].get_weights()
4 print(W1[2])
5
6 W2 = layers[1].get_weights()
7 print(W2[2])

[[[-0.09735187 -0.2429723 -0.36114842 -0.15082116 -0.23535672
  -0.02661474 -0.2351685  0.391844  ]]]
[[[-0.09735187 -0.2429723 -0.36114842 -0.15082116 -0.23535672
  -0.02661474 -0.2351685  0.391844  ]]]

```

در حالت قبل، بهترین آلفا 0 بود. در اینجا نیز اکثر مقادیر آلفا (به جز یکی) که برای هر فیلتر یاد گرفته، کمترین فاصله را با 0 دارند (نسبت به فاصله با آلفاهای دیگر). پس بهترین مقدار آلفا نتیجه این آزمایش را تایید می کند. اگر بخواهیم مقدار آلفا در تابع فعال سازی را Leaky ReLU را به صورت دستی tune کنیم، بهتر است از مقدار منفی کوچک تا 0 را امتحان کنیم.