

شبکه‌های عصبی  
دانشگاه فردوسی مشهد  
تمرین دو

نیمسال دوم تحصیلی ۱۴۰۳-۱۴۰۲

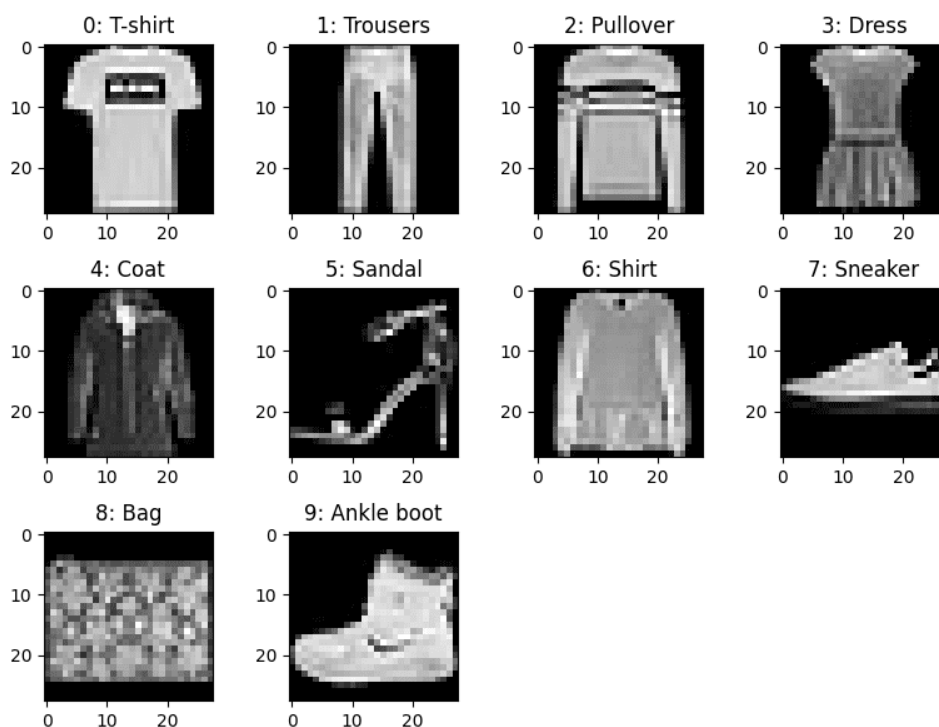
مهلت ارسال: ۲۳:۵۹ ۱۴۰۳/۱/۲۰

گروه مهندسی کامپیوتر

در این تمرین قصد داریم از شبکه MLP که در تمرین قبل کد آن را آماده کرده اید برای طبقه بندی دیتاست Fashion MNIST استفاده کنیم و به بررسی تاثیرات مفاهیم مختلف شبکه های عصبی از جمله توابع متفاوت فعالیت، وزن دهی اولیه وزن ها، mini-batch gradient descent و بهینه سازهای متفاوت اسن. برای پیاده سازی از توابع استاندارد Python و Numpy استفاده کنید و از Toolbox های مختلف مربوط به دیپ لرنینگ استفاده نکنید. در ابتدا به توضیح دیتاست می پردازیم و سپس مواردی که باید پیاده شوند توضیح داده می شود.

### دیتاست Fashion MNIST

دیتاست Fashion MNIST این دیتاست همانند دیتاست MNIST شامل ۶۰۰۰۰ عکس به عنوان داده آموزش و ۱۰۰۰۰ عکس به عنوان داده تست است و داده ها مربوط به ۱۰ نوع لباس مختلف همانند کفش، شلوار و ... می باشد. عکس ها به صورت grayscale هستند و ابعاد آنها ۲۸ در ۲۸ می باشد و مقادیر پیکسل ها از ۰ تا ۲۵۵ است.



شکل ۱: چند نمونه از داده های Fashion MNIST.

برای این تمرین، از آنجایی که این دیتاست شامل ۱۰ کلاس است، دو کلاس ابتدایی آن که پیراهن و شلوار است

را جدا شده اند و با فرمت npz در اختیار شما قرار داده شده است. برای لود کردن داده ها از دستورات زیر استفاده کنید.

```
x_train = np.load('x_train.npz')['arr_0']
x_test = np.load('x_test.npz')['arr_0']
y_train = np.load('y_train.npz')['arr_0']
y_test = np.load('y_test.npz')['arr_0']
```

## توابع فعالیت

هدف در این قسمت بررسی دقیق تر توابع فعالیت و ویژگی های آنها می باشد.

### ۱.۱ تابع فعالیت همانی

تابع فعالیت همانی، ورودی نورون را به همان مقدار برای خروجی نورون نسبت می دهد.

$$f(x) = x$$

استفاده از این تابع فعالیت برای داده های غیر جداپذیر خطی ممکن نیست. علت این امر این است که شبکه ای که نورون های لایه های مخفی آن از این تابع استفاده کنند، صرف نظر از تعداد لایه ها و تعداد نورون های هر لایه همواره تابعی خطی را یاد می گیرد.

### ۲.۱ Sigmoid (Logistic Activation Function)

$$f(x) = \frac{1}{1 + e^{-x}}$$

این تابع به ازای ورودی های منفی تر و مثبت از مقداری، شیب کمی دارد. این امر باعث می شود که در فرآیندهای اپدیت گرادیانی که نیاز به مشتق گرفتن از توابع فعالیت است، مشتق ها کوچک باشند و فرآیند همگرایی به کندی برای رنج زیادی از مقادیر انجام شود. استفاده معمول تر از تابع در لایه آخر برای طبقه بندی دو کلاسه است.

### ۳.۱ Tanh or hyperbolic tangent Activation Function

$$f(x) = 2 \times \text{sigmoid}(2x) - 1$$

این تابع فعالیت نیز همانند سیگموید  $s$  شکل است. از آنجایی که خروجی این تابع بین ۰ و ۱ قرار دارد و یک تابع فعالیت zero-centered است، مقادیر خروجی نورون ها را اطراف ۰ نگاه می‌دارند و به عبارتی از بزرگ شدن یا کوچک شدن گرادیان‌ها جلوگیری می‌کند. به همین علت از این تابع بیشتر در لایه‌های مخفی و از سیگموید در لایه آخر استفاده می‌شود.

این تابع همچنان مشکل تابع سیگموید در کند کردن فرآیند آموزش و vanishing gradient به خصوص در شبکه‌های عمیق را دارد.

## ۴.۱ ReLU

$$f(x) = \max(0, x)$$

این تابع سرعت خوبی برای رسیدن به همگرایی دارد و از مشکل اشباع گرادیان برای ورودی‌های مثبت جلوگیری می‌کند و به مشکل vanishing gradients را نیز کمک می‌کند. این تابع مشکلی دارد که آن را dying ReLU problem می‌نامند. زمانی که اکثر داده ورودی در رنج منفی قرار بگیرد چون این تابع خروجی ۰ می‌دهد و گرادیان ۰ نیز ۰ است، در فرآیند BP گرادیان‌ها نمیتوانند به عقب بازگردند و وزن‌ها آپدیت نمی‌شوند.

## ۵.۱ Leaky ReLU

این تابع شرایط ReLU را دارد و برای همگرایی سریع است و دچار dying ReLU problem نمی‌شود چون برای مقادیر منفی نیز مقداری شیب منفی و در نتیجه گرادیان منفی دارد.

## وزن‌دهی اولیه پارمترهای شبکه

آموزش شبکه عصبی به میزان قابل توجهی به وزن‌دهی اولیه پارمترها بستگی دارد. این نقطه اولیه میتواند حتی در اینکه شبکه همگرا بشود یا نه هم تاثیر داشته باشد. به عنوان مثال مقدار دهی اولیه می‌تواند به گونه‌ای unstable انجام شود که مشکلات عددی به وجود بیاید. حتی در زمانی که یادگیری همگرا می‌شود، مقدار دهی اولیه مناسب میتواند در اینکه با چه سرعتی همگرایی به دست می‌آید یا به نقطه ای با چه هزینه‌ای همگرا می‌شود تاثیر داشته باشد. علاوه بر این، وزن‌دهی اولیه بر تعمیم پذیری نیز اثر دارد.

## ۱.۲ وزن دهی تمام صفر

در این روش تمام پارمترهای اولیه با صفر مقداردهی اولیه می‌شوند. در این روش یکی از اصول مهم در وزن‌دهی اولی که نیاز به از بین بردن تقارن یا break symmetry است انجام نمی‌شود. دو نورون مخفی که با تایع فعالیت یکسان که به ورودی‌های مشابهی متصل هستند، باید پارامترهای اولیه متفاوتی داشته باشند و اگر اینگونه نباشد یک الگوریتم یادگیری قطعی (deterministic) که بر روی یک تابع هزینه قطعی اعمال شده است، به صورت دائم هر دو نورون را به یک صورت آپدیت خواهد کرد.

حتی در صورت آموزش مدل با الگوریتمی که می‌تواند از ناپیچینی (stochasticity) برای محاسبه آپدیت‌های نورون‌ها استفاده کند (به طور مثال با استفاده از dropout) باز هم بهتر است وزن‌ها هر نورون به گونه‌ای تعیین شوند که هر نورون تابعی متفاوت از نورون‌های دیگر محاسبه کند.

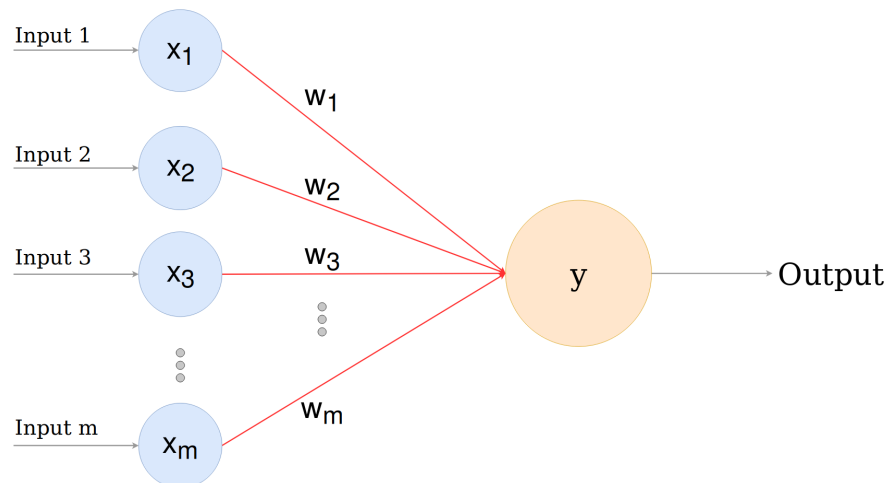
## ۲.۲ وزن‌دهی رندم

روش دیگر وزن‌دهی اولیه بامقادیر رندوم از توزیع یکنواخت یا گوسی می‌باشد. انتخاب بین این دو نوع توزیع اهمیت چندانی ندارد و به طور جدی نیز بررسی نشده است. هر بار که وزن‌های شبکه مقداردهی اولیه می‌شوند باعث می‌شود یک مجموعه متفاوت از وزن‌ها به معنای یک نقطه شروع جدید در فرآیند بهینه‌سازی و احتمالاً مقادیر نهایی متفاوت وزن‌ها (با ویژگی‌ها و عملکرد متفاوت) ایجاد می‌شود.

روش‌های همخوان‌شده تری با انواع مختلف توابع فعالیت برای مقداردهی اولیه وزن‌ها وجود دارد که در ادامه به بررسی برخی خواهیم پرداخت. همچنین در مقداردهی باید به این نکته توجه داشت که مقادیر اولیه متفاوت برای این وزن‌ها می‌تواند به این کمک کند که هیچ الگو ورودی (Input pattern) در فضای پوچ پیمایش رو به جلو (forward propagation) گم نشود و همچنین هیچ الگو گرادیانی در فضای پوچ (back-propagation) از دست نرود.

## ۳.۲ مقداردهی اولیه وزن‌ها برای توابع فعالیت Tanh و Sigmoid

یک نورون پرسپترون را در نظر بگیرید. اگر داشته باشیم:



شکل ۲: یک نورون perceptron

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n, \quad y = g(z)$$

هر چه قدر تعداد ورودی‌ها بیشتر باشد، برای اینکه  $z$  بزرگ نشود (به منظور جلوگیری از اشباع تابع فعالیت)، باید  $w$  کوچکتر باشد. اگر  $W$  ماتریس وزن این نرون باشد، به دلیل نگاه داشتن این رابطه معکوس می‌توان برای واریانس این ماتریس نوشت:

$$Var(W) = \frac{1}{n}$$

بر همین اساس برای مقداردهی اولیه وزن‌ها در شبکه‌ای با توابع فعالیت Sigmoid و Tanh از روشی به نام Xavier استفاده می‌شود که پیاده‌سازی آن در پایتون می‌تواند به صورت زیر انجام بگیرد:

```
import numpy as np

def initialize_weights(shape):
    n = shape[0] # Number of input units
    return np.random.randn(*shape) * np.sqrt(1/n)

# Example usage:
shape = (10, 5) # Example shape for the weight matrix (input_size, output_size)
W1 = initialize_weights(shape)
```

شکل ۳: Xavier Initialization

این مقداردهی وزن‌ها را بر اساس یک توزیع گوسی با واریانس  $\frac{1}{n}$  مقدار می‌دهد. (از توزیع یکنواخت نیز می‌توان استفاده کرد)

## ۴.۲ مقداردهی اولیه وزن‌ها برای تابع فعالیت ReLU

روش Xavier برای تابع فعالیت ReLU مشکلاتی دارد به همین منظور از روش دیگری به نام He Initialization استفاده می‌شود. این مقداردهی وزن‌ها را بر اساس یک توزیع گوسی با واریانس  $\frac{2}{n}$  مقدار می‌دهد. (از توزیع یکنواخت نیز می‌توان استفاده کرد)

## MiniBatch Gradient descent

الگوریتم Gradient descent که به آن Batch Gradient descent هم گفته می‌شود، برای هر گام آپدیت گرادیانی یک بار همه داده‌های دیتاست آموزشی را طی می‌کند و بعد وزن‌ها را یک بار آپدیت می‌کند. برای این الگوریتم اگر نمودار تابع هزینه بر حسب شماره پیمایش دیتاست آموزشی رسم کنیم مشاهده می‌شود که همواره نمودار کاهشی است و نویزی نمی‌باشد. در این الگوریتم می‌توان نرخ یادگیری را بزرگ‌تر اختیار کرد و گام‌های آپدیت بلند با نوبز کم برمی‌دارد.

الگوریتم MiniBatch Gradient descent داده را به قسمت‌هایی (mini-batch) تقسیم می‌کند به این صورت که برای داده‌های هر قسمت از داده‌های آموزش سمپل گرفته می‌شود و اپدیت وزن‌ها بعد از پیمایش هر mini-batch

صورت می‌گیرد. این فرآیند باعث ایجاد stochasticity در فرآیند بهینه‌سازی و ایجاد نویز در تخمین گرادیان‌ها می‌شود به علت استفاده از زیرمجموعه‌هایی از داده‌های آموزش به جای کل دیتاست آموزشی. هر چه قدر اندازه mini-batch کوچکتر باشد نویز هم بیشتر می‌شود. برای این الگوریتم اگر نمودار تابع هزینه بر حسب شماره پیمایش دیتاست آموزشی رسم کنیم مشاهده می‌شود که نمودار نویزی است اما روند کلی کاهشی دارد.

از آنجایی که در مسائل امروزه شبکه عصبی سائز دیتاست های آموزشی بسیار بالا می‌باشد اسفاده از روش full batch که با پیمایش کل دیتاست تنها یک بار وزن‌ها را آپدیت می‌کند در مسائل با سائز زیاد داده ممکن نیست و به همین دلیل از روش mini-batch استفاده می‌شود. معمولا سائز mini-batch را توانی از ۲ استفاده می‌کنند و استفاده از سائزهای ۶۴، ۱۲۸، ۲۵۶ و ۵۱۲ مرسوم است.

### **Batch Gradient descent**

for epoch in number of epochs:

- for all the training instances in the dataset compute the derivative of the cost function
- update the weights

### **Mini-batch Gradient descent**

for epoch in number of epochs:

for batch in num of batches:

- for all the training instances in the batch sample compute the derivative of the cost function
- update the weights

## **Stochastic Gradient descent (SGD)**

در الگوریتم MiniBatch Gradient descent یک حالت extreme این است که سائز mini-batch برابر با سائز دیتاست آموزشی باشد که در این حالت به همان الگوریتم گرادیان کاهشی فول بچ می‌رسیم. حالت دیگر extreme این است که سائز هر mini-batch برابر ۱ باشد و به عبارتی هر داده خود یک mini-batch باشد. در این حالت به الگوریتم Stochastic Gradient descent (SGD) می‌رسیم. این الگوریتم در هر Iteration تنها یک sample از داده را می‌بیند. در برخی از اپدیت ها گام‌ها در مسیر global optimum پیش می‌روند اما در برخی از آپدیت ها نیز دور می‌شوند و به همین علت SGD بسیار نویزی است. نکته دیگر این است که SGD هیچ گاه به خود مینیمم نمی‌رسد و در اطراف آن نوسان می‌کند. نویزی بودن این الگوریتم را می‌توان با انتخاب نرخ یادگیری مناسب (کوچک) حل کرد.

## Stochastic Gradient descent (SGD)

for epoch in number of epochs:

for instance in total dataset:

- for the current instance compute the derivative of the cost function
- update the weights

## SGD with Momentum

به دلیل نوسانات آپدیت ها در الگوریتم SGD باید از نرخ یادگیری کوچکی استفاده کرد تا واگرایی رخ ندهد و این موضوع باعث کند شدن فرآیند بهینه سازی می شود.

---

### Algorithm 1 SGD with Momentum

---

Require:  $\alpha$  and  $\beta$

On iteration t:

Compute gradients  $d_w$  and  $d_b$  for the sample

$$V_{d_w} = \beta V_{d_w} + (1 - \beta) d_w$$

$$V_{d_b} = \beta V_{d_b} + (1 - \beta) d_b$$

$$W = W - \alpha V_{d_w}$$

$$b = b - \alpha V_{d_b}$$

---

این الگوریتم moving average برای گرادیان ها را محاسبه می کند و با استفاده از آن آپدیت وزن ها را انجام می دهد. این الگوریتم گام های الگوریتم گرادیان کاهشی را ملایم تر می کند. این میانگین گیری باعث می شود که در بعد نوسانات، میانگین این نوسانات به حدودا صفر برسد و مسیر رسیدن به نقطه بهینه هموار تر شود. برای توضیح استفاده از جمله مومنتوم از analogy پرت کردن یک توپ در یک دره استفاده می شود که در باره آن می توانید از [اینجا](#) بخوانید.

## Root Mean Squared prop(RMSprop)

---

**Algorithm 2** Root Mean Squared prop(RMSprop)

---

Require:  $\alpha$  and  $\beta$

On iteration  $t$ :

Compute gradients  $d_w$  and  $d_b$  for the mini-batch

$$S_{d_w} = \beta V_{d_w} + (1 - \beta) d_w^2$$

$$S_{d_b} = \beta V_{d_b} + (1 - \beta) d_b^2$$

$$W = W - \alpha \frac{d_w}{\sqrt{S_{d_w}}}$$

$$b = b - \alpha \frac{d_b}{\sqrt{S_{d_b}}}$$

---

در این الگوریتم برای بعدی که در آن سرعت بهینه سازی بیشتری می‌خواهیم مقدار میانگین کمتر خواهد بود به این صورت که در بعدی که نوسان بیشتری دارد چون مقدار گرادیان ها و در نتیجه توان ۲ آن ها بیشتر می‌شود مقدار میانگین گرادیانی آنها بیشتر شده و در نتیجه در آپدیت مقدار کمتری آپدیت می‌شوند (چون جذر میانگین در مخرج قرار گرفته است در فرمول‌های آپدیت) و سرعت آپدیت در بعد با نوسان زیاد کمتر می‌شود. اثر این الگوریتم این است که می‌توان بدون واگرایی از نرخ یادگیری بزرگتری استفاده کرد که این به معنای افزایش سرعت یادگیری است. (منظور از ابعاد پارمترهایی هستند که در حال آپدیت آن ها هستیم)

### Adaptive Moment Estimation(Adam)

---

**Algorithm 3** Adaptive Moment Estimation(Adam)

---

Require:  $\alpha$  and  $\beta_1$  and  $\beta_2$

On iteration  $t$ :

Compute gradients  $d_w$  and  $d_b$  for the mini-batch

$$V_{d_w} = \beta_1 V_{d_w} + (1 - \beta_1) d_w \text{ (Update biased first moment estimate)}$$

$$V_{d_b} = \beta_1 V_{d_b} + (1 - \beta_1) d_b \text{ (Update biased first moment estimate)}$$

$$S_{d_w} = \beta_2 V_{d_w} + (1 - \beta_2) d_w^2 \text{ (Update biased second moment estimate)}$$

$$S_{d_b} = \beta_2 V_{d_b} + (1 - \beta_2) d_b^2 \text{ (Update biased second moment estimate)}$$

$$V_{d_w}^{corrected} = \frac{V_{d_w}}{1 - \beta_1^t} \text{ (Bias-corrected)}$$

$$V_{d_b}^{corrected} = \frac{V_{d_b}}{1 - \beta_1^t} \text{ (Bias-corrected)}$$

$$S_{d_w}^{corrected} = \frac{S_{d_w}}{1 - \beta_2^t} \text{ (Bias-corrected)}$$

$$S_{d_b}^{corrected} = \frac{S_{d_b}}{1 - \beta_2^t} \text{ (Bias-corrected)}$$

$$W = W - \alpha \frac{V_{d_w}^{corrected}}{\sqrt{S_{d_w}^{corrected} + \epsilon}} \text{ (Update parameters)}$$

$$b = b - \alpha \frac{V_{d_b}^{corrected}}{\sqrt{S_{d_b}^{corrected} + \epsilon}} \text{ (Update parameters)}$$

---



این الگوریتم ترکیب دو الگوریتم قبلی می‌باشد. مشاهده می‌شود که چهار هایپرپارامتر دارد که نرخ یادگیری باید توسط تست مقدار مناسبی برایش اختیار شود. مقدار  $\beta_1$  را معمولا ۹۰٪ و مقدار  $\beta_2$  ۹۹۹۰٪ استفاده می‌شود. مقدار  $\epsilon$  نیز برای جلوگیری از مشکلات عددی مانند مخرج صفر است و مقدار آن اهمیت چندانی ندارد و تنها باید مقدار کوچکی باشد.

## L2 Regularization

در درس با مشکل overfitting آشنا شدید و مشاهده کردید که یک روش برای جلوگیری از آن، استفاده از روش های regularization است. در این بخش، شما باید L2 regularization که یکی از روش های regularization است را پیاده سازی کنید. در تمرین شماره ۳، با روش های regularization بیشتری آشنا خواهید شد. در L2 regularization یک ترم اضافی مربعات وزن به تابع هزینه اضافه می‌شود. این کار باعث می‌شود تا وزن های شبکه خیلی بزرگ نشوند و در نتیجه مدل خیلی پیچیده نمی‌شود. این کار باعث جلوگیری از overfitting می‌شود.

فرمول تابع هزینه بعد از اضافه کردن ترم L2 regularization به صورت زیر است:

$$loss(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \gamma \sum_{j=1}^n \theta_j^2 \right]$$

که در آن  $h_{\theta}$  خروجی شبکه، و  $\theta$  وزن های شبکه می‌باشد.

## آزمایشات

شما باید ابتدا تمامی توابع فعالیت و بهینه ساز های معرفی شده به همراه L2 Regularization را پیاده سازی کنید. در این تمرین برای انجام آزمایشات باید از شبکه که در تمرین ۱ پیاده سازی کردید و دیتاست Fashion MNIST استفاده کنید. و فایل های کد و گزارش خود را در سامانه VU آپلود کنید.

• توابع فعالیت داکيومنت را بر روی شبکه MLP با تعداد لایه ها و نوروں های متفاوت تست کنید. (دقت کنید که حتما توابع فعالیت مختلف را بر روی شبکه ای با ۴ یا ۵ لایه تست کنید تا تاثیر متفاوت آن ها بر Vanishing Gradient رامشاهده کنید). برای هر تابع فعالیت بر روی مدل های متفاوتی که تست می‌کنید نمودار تابع هزینه و اکیورسی بر حسب ایپاک روی داده آموزش و تست را رسم کنید و در گزارشات خود درمورد سرعت یادگیری و مقدار بهینه و نتیجه که مدل های متفاوت با هم داشته اند بررسی انجام دهید.

• وزن دهی اولیه تمام صفر، رندوم اعداد کوچک، رندوم اعداد بزرگ، Xavier و He را با بکدیگر مقایسه کنید. بر روی یک مدل ۳ لایه تست کنید و نتایج هر روش را با رسم توابع هزینه و اکیورسی روی داده های تست و آموزش مقایسه کنید.

• شبکه ۴ لایه ای را با بهینه سازهای متفاوت آموزش دهید و learning curve را برای هر کدام رسم کنید و نتایج را بررسی کنید.

- این بار به این شبکه ۴ لایه L2 Regularization نیز اضافه کنید و دوباره این شبکه را به یک بهینه ساز دلخواه آموزش دهید و تابع هزینه بر روی داده آموزش و تست بر حسب ایپاک را با هم پلات کنید.  
و با حالت قبلی (شبکه ۴ لایه با همان بهینه ساز اما بدون L2 Regularization ) مقایسه کنید.