

بنام خدا

گزارش تمرین سری سوم درس یارگیری عمیق

سارا بارونی

۹۸۴۴۳۰۳۱

### الگوریتم Adagrad

این الگوریتم مقدار learning rate ثابتی که در حالت معمول داشتیم را تقسیم میکند بر جذر مجموع توان دو گرادیان‌ها به طوری که هیستوری را در نظر میگیرد و فلسفه‌ی این کار این است که آداپته شود با تغییرات شیب به گونه‌ای که اگر شیب کم باشد چون گرادیان‌ها در مخرج قرار دارند در نتیجه حاصل learning rate افزایش خواهد یافت و این کمک میکند در قسمت‌هایی که تابع هزینه flat است بتواند سریع تر پیش برود و این مشکل برطرف بشود ضمن اینکه اگر شیب زیاد باشد مقدار learning rate کم میشود و با گام‌های کوچک‌تری پیش میرود و این کمک میکند تا اگر نزدیک نقطه‌ی اپتیمم باشد از آن عبور نکند و محتاط‌تر پیش برود. این الگوریتم در تابع هزینه‌ای که convex است خوب عمل میکند. اما ایرادی که میتوان از این روش گرفت درواقع آن است که کل هیستوری را در نظر میگیرد مثلاً یکی از مشکلاتی که میتواند ایجاد شود آن است که شیب همچنان زیاد شود و این باعث میشود learning rate کم و کمتر شود و گاهی ممکن است آنقدر کم بشود که اثری در آپدیت شدن وزن‌ها نداشته باشد.

در شکل ۱ نیز شمای کلی از این الگوریتم مشاهده میشود.

در پیاده سازی الگوریتم مشاهده میشود که نمودار دقت داده‌ها صعودی است و صعودی بودن این نمودار در test به معنای آن است که overfit نشده است.

لازم به ذکر است در تمامی شکل‌ها نمودار قرمز مربوط به دقت داده‌های test و آبی دقت داده‌های train و سبز دقت داده‌های validation را نشان می‌دهد.

پارامترهای این الگوریتم شامل epsilon و delta می‌باشد که معمولاً epsilon را 0.001 در نظر میگیرند و delta را  $1e-8$  در نظر میگیرند. این پارامترها در شکل ۱ که این الگوریتم را نشان میدهد مشاهده میشوند که به چه شکل در این الگوریتم به کار برده شده اند.

### الگوریتم Adadelta

این الگوریتم مشابه الگوریتم Adagrad می‌باشد اما به نوعی اثر هیستوری‌های دورتر را کم میکند. الگوریتم در قالب فرمول‌های زیر بیان شده است.

$$W(t+1) = W(t) - \frac{\sqrt{(D(t-1) + \epsilon)}}{\sqrt{(v(t) + \epsilon)}} \frac{\partial L}{\partial W(t)} \quad (1)$$

$$D(t) = \beta * D(t-1) + (1 - \beta) \Delta w(t) \quad (2)$$

$$v(t) = \beta * v(t-1) + (1 - \beta) \left( \frac{\partial L}{\partial W(t)} \right)^2 \quad (3)$$

$$\Delta w(t) = W(t) - W(t-1) \quad (4)$$

در این الگوریتم از پارامترهای  $\beta$  که معمولاً 0.95 آن را در نظر میگیرند و epsilon که معمولاً آن را  $1e-8$  در نظر میگیرند استفاده می‌شود.

در شکل ۱۲ خروجی مدل با این الگوریتم و همچنین الگوریتم adagrad باهم رسم شده است که همانطور که مشاهده میشود به طور کلی در الگوریتم adadelta افزایش یافته اما مشکلی که دارد از قسمتی به بعد نزولی

شده است. **الگوریتم RMSprob**

این الگوریتم همان الگوریتم Adagrad است با این تفاوت که بجای مجموع توان دو گرادیان‌ها میانگین وزن‌دار نمایی آن‌ها را حساب میکند. این الگوریتم زمانی که تابع هزینه nonconvex باشد میتواند بهتر عمل کند و در مواردی که تابع هزینه convex باشد هم سریع تر همگرا میشود باتوجه به اینکه مشکلی که در الگوریتم Adagrad وجود داشت (هیستوری را درنظر میگرفت و ممکن بود قبل از رسیدن به نقطه‌ی اپتیمم وزن‌ها آپدیت نشوند) را حل کرده است به این شکل که اثر هیستوری گرادیان‌ها را هرچه دورتر باشد کم میکند. پارامترهایی که در این الگوریتم وجود دارد شامل beta و epsilon میباشد که به طور معمول beta را 0.95 و epsilon را 1e-8 درنظر میگیرند.

**الگوریتم Adam**

روش Adaptive momentum روش RMSprob را با روش momentum ترکیب میکند درواقع -momentum کمک میکند نوسانات موجود در جواب را کم کنیم و RMSprob کمک میکند تا نرخ یادگیری را تنظیم کنیم

این الگوریتم به طور کامل در شکل ۲ مشاهده میشود. که پارامترهایی که به کاربرده شده اند در آن شامل beta و delta و epsilon می‌باشند که معمولاً beta را 0.999 و 1e-7 و epsilon را 0.001 درنظر میگیرند.

**الگوریتم Gradient descent**

در این روش وزن‌ها با مقدار یک learning rate ثابت آپدیت میشوند.

**الگوریتم L2 Regularization**

در رگرسیون یک بخش جدید یعنی تابعی از بردار وزن‌ها را به تابع هزینه اضافه میکنیم

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N L(y(x, \theta)) + R(\theta) \quad (5)$$

که این بخشی که اضافه میکنیم در L2 regularization به این شکل تعریف میشود

$$R_{L2}(\theta) = \frac{\lambda}{2} \sum_{i=1}^N W_i^2 \quad (6)$$

به طور معمول از L2 regularization بیشتر از L1 regularization استفاده می‌شود و L2 regularization اثربخشی بیشتری در شبکه‌های عصبی دارد اثر این روش این است که وزن‌ها نسبت به قبل مقادیر کوچکتری داشته باشند و این جلو overfit شدن را بگیرد و مدل ساده‌تری خواهیم داشت. در ادامه روابط ریاضی آپدیت شدن وزن‌ها در این روش را داریم:

$$W = W - \alpha \left( \frac{\partial J}{\partial W} + \lambda \frac{\partial R}{\partial W} \right) \quad (7)$$

$$W = W - \alpha \left( \frac{\partial J}{\partial W} + \lambda W \right) \quad (8)$$

$$W = (1 - \alpha\lambda)W - \alpha \frac{\partial J}{\partial W} \quad (9)$$

در پیاده‌سازی این الگوریتم اینگونه عمل شد که به تابع `function lrcost` مجموع توان دو وزن‌ها را اضافه کردم و برای مشتق گرفتن از آن اثر آن را در الگوریتم `backpropagation` اعمال کردم به این شکل که به ازای هر لایه که بخواهم مشتق تابع هزینه را نسبت به وزن‌های مربوط به آن بدست آورم چون دارم مشتق را نسبت به آن وزن‌ها بدست می‌آورم فقط آن وزن‌ها باقی میمانند زیرا مشتق وزن‌های دیگر نسبت به آن وزن صفر میشود به این شکل مشتق آن را در الگوریتم `backpropagation` قرار دادم. همانطور که در شکل ۱۴ هم مشاهده میشود نشده `overfit` است زیرا هر دو نمودار تست و ترین و همچنین ولیدیشن هم صعودی هستند.

مقدار لاندا در این الگوریتم 0.001 در نظر گرفته شد که نتایج بهتری نسبت به مقادیر دیگر میداد. **Regu- L1 larization**

در این روش مشابه روش `L2 regularization` است با این تفاوت که بجای به توان دو رساندن وزن‌ها قدرمطلق آنها را قرار میدهم. در این روش هم هدف کاهش دادن وزن‌ها است اما در عمل در این روش بردار وزن‌ها تراکم کمتری دارد و این یک سری از فیچرها را حذف میکند درواقع فیچرهایی که اهمیت کمتری در مدل دارند را حذف میکند به همین دلیل این روش را به عنوان یک مکانیزم برای `feature selection` هم در نظر میگیرند.

در این الگوریتم در صورتی که مقدار لاندا بزرگ انتخاب میشد حتی در حدود 0.1 مقدار خطا زیاد میشد در حد ۷ زیاد میشد و مقدار آن را خیلی کوچک در نظر گرفتیم. در مقدار لاندا 0.000001 خطا به حدود 0.6 رسید خروجی مدل با در نظر گرفتن `L1 regularization` در شکل ۱۶ و ۱۷ مشاهده میشود نمودار دقت داده ی تست همچنان افزایش پیدا کرده اما به طور کلی دقت کم شده است. همچنین نمودار خطا روندی نزولی دارد و این نکته مثبتی است.

---

**Algorithm 8.4** The AdaGrad algorithm

---

**Require:** Global learning rate  $\epsilon$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ .

    Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---

شکل ۱: الگوریتم AdaGrad برگرفته از کتاب deep learning

---

---

**Algorithm 8.5** The RMSProp algorithm

---

**Require:** Global learning rate  $\epsilon$ , decay rate  $\rho$

**Require:** Initial parameter  $\theta$

**Require:** Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers

Initialize accumulation variables  $\mathbf{r} = \mathbf{0}$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ .

    Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$ .

**end while**

---

شکل ۲: الگوریتم RMSprop برگرفته از کتاب deep learning

---

---

**Algorithm 8.7** The Adam algorithm

---

**Require:** Step size  $\epsilon$  (Suggested default: 0.001)

**Require:** Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ . (Suggested defaults: 0.9 and 0.999 respectively)

**Require:** Small constant  $\delta$  used for numerical stabilization (Suggested default:  $10^{-8}$ )

**Require:** Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

Initialize time step  $t = 0$

**while** stopping criterion not met **do**

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

    Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

    Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

    Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

    Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

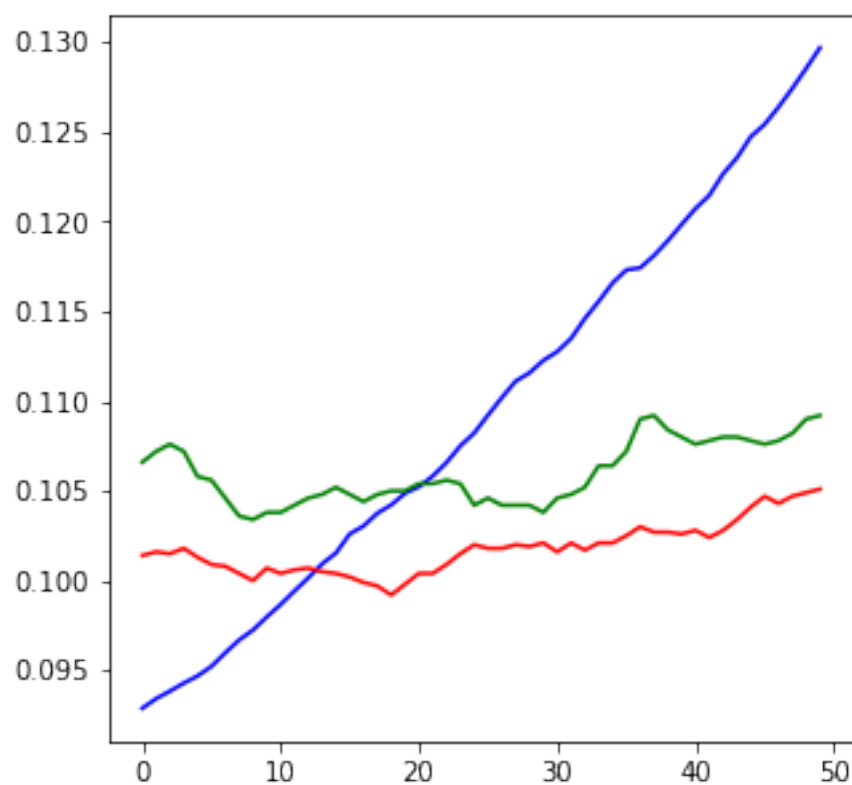
    Compute update:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\delta + \hat{\mathbf{r}}}}$  (operations applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$

**end while**

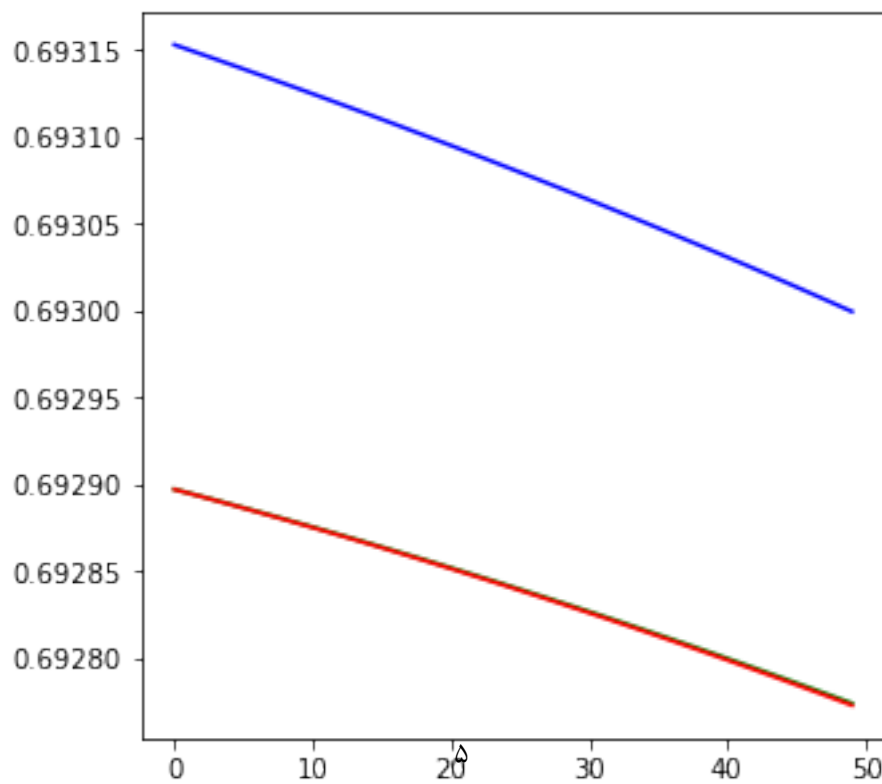
---

شکل ۳: الگوریتم Adam برگرفته از کتاب deep learning



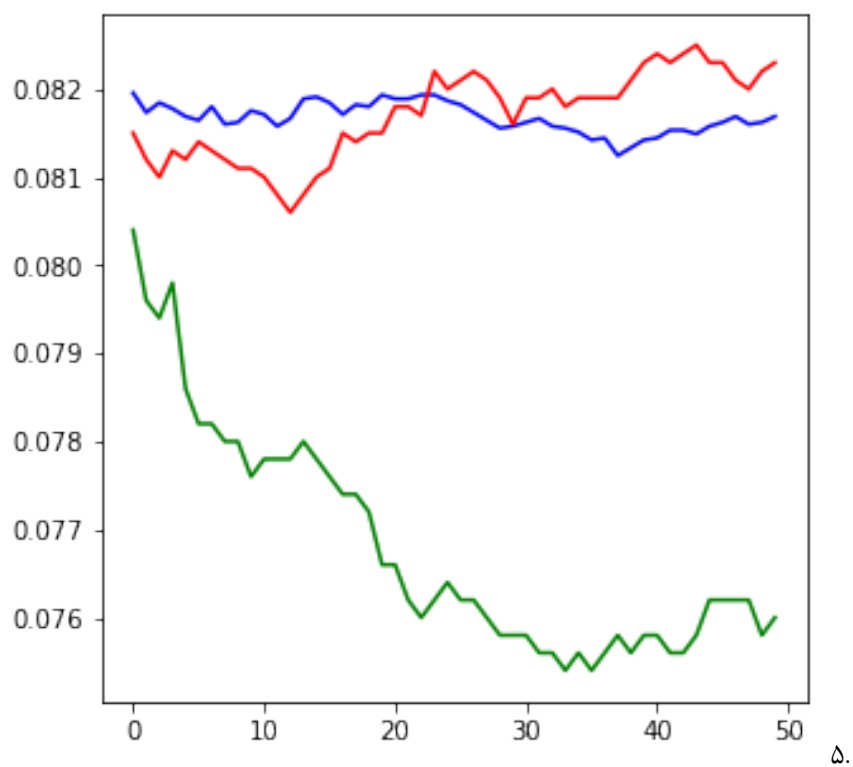
۵.

شکل ۴: نمودار دقت Gradient descent الگوریتم

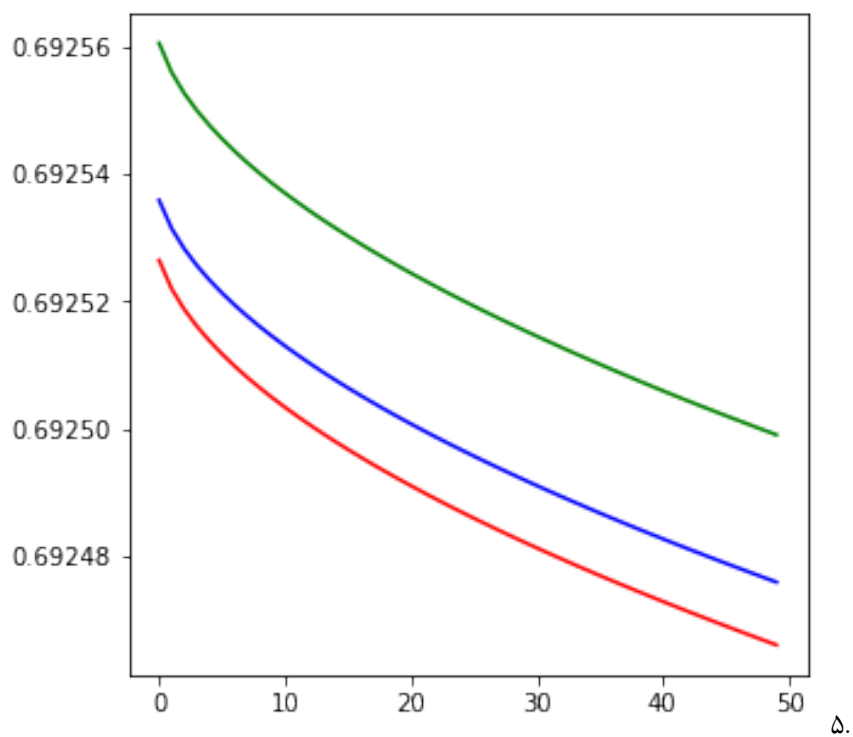


۵.

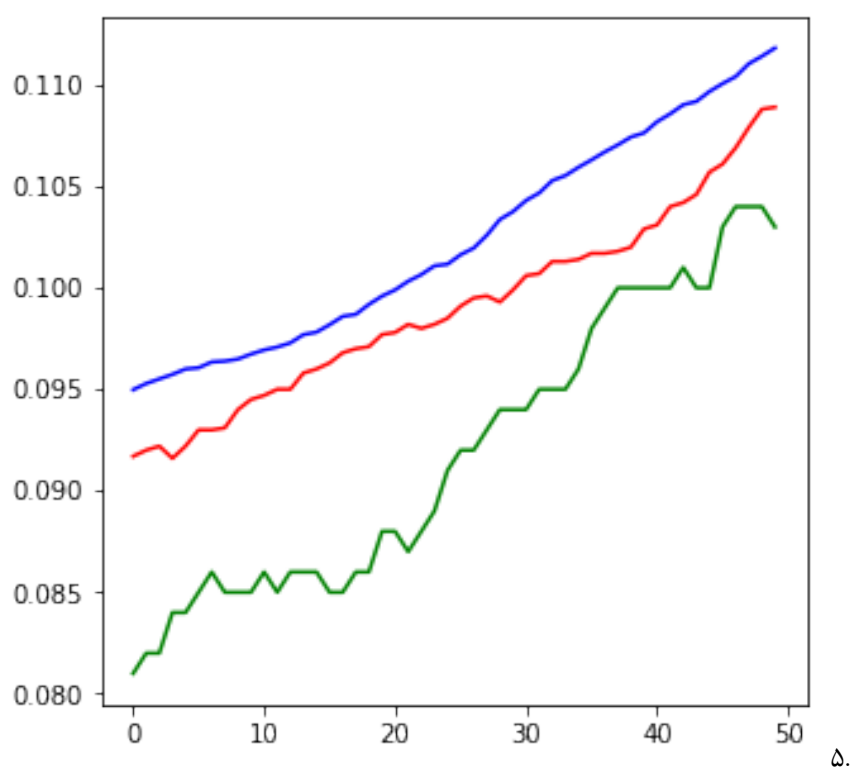
شکل ۵: نمودار خطا Gradient descent الگوریتم



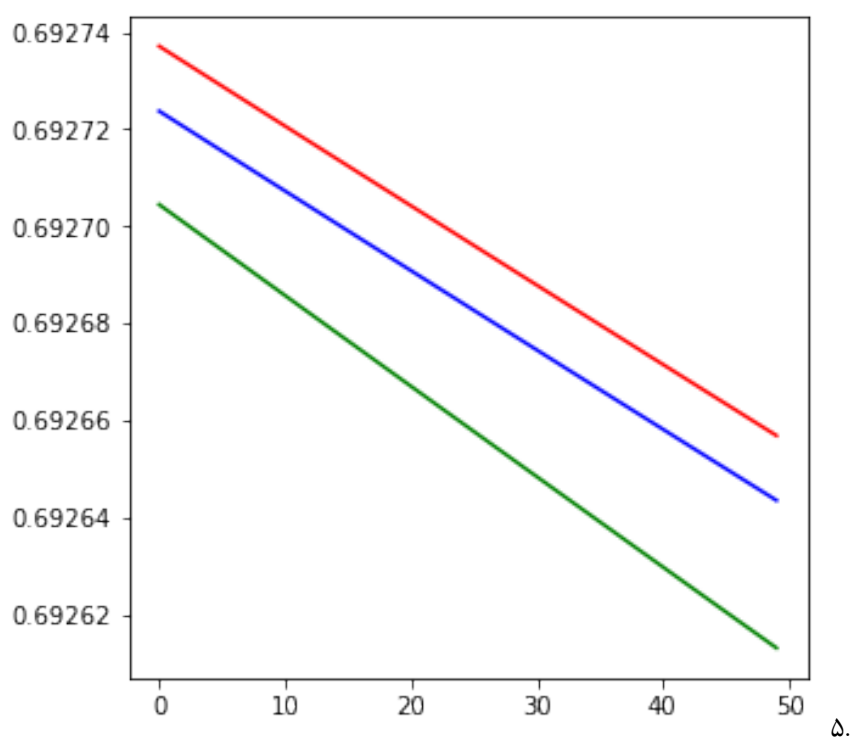
شکل ۶: نمودار دقت با learning rate مربوط به الگوریتم adograd



شکل ۷: نمودار خطا با learning rate مربوط به الگوریتم adograd

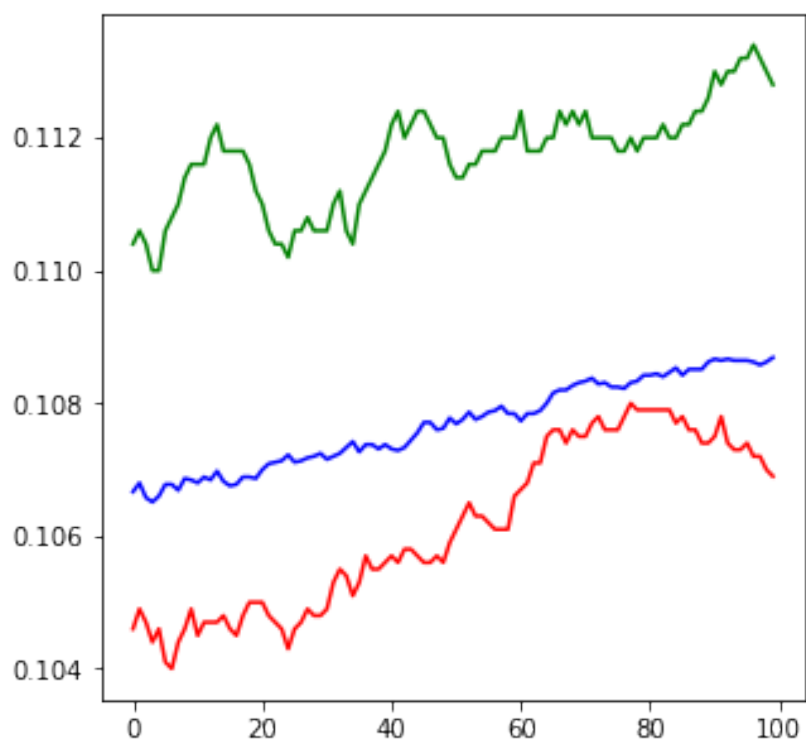


شکل ۸: نمودار دقت با learning rate مربوط به الگوریتم adam



شکل ۹: نمودار خطا با learning rate مربوط به الگوریتم adam

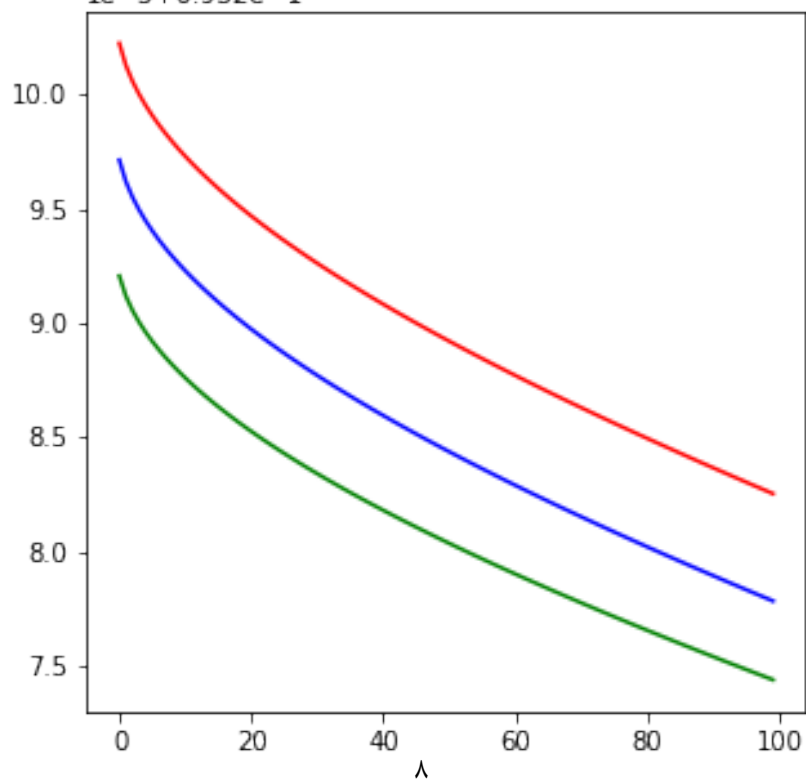




۵.

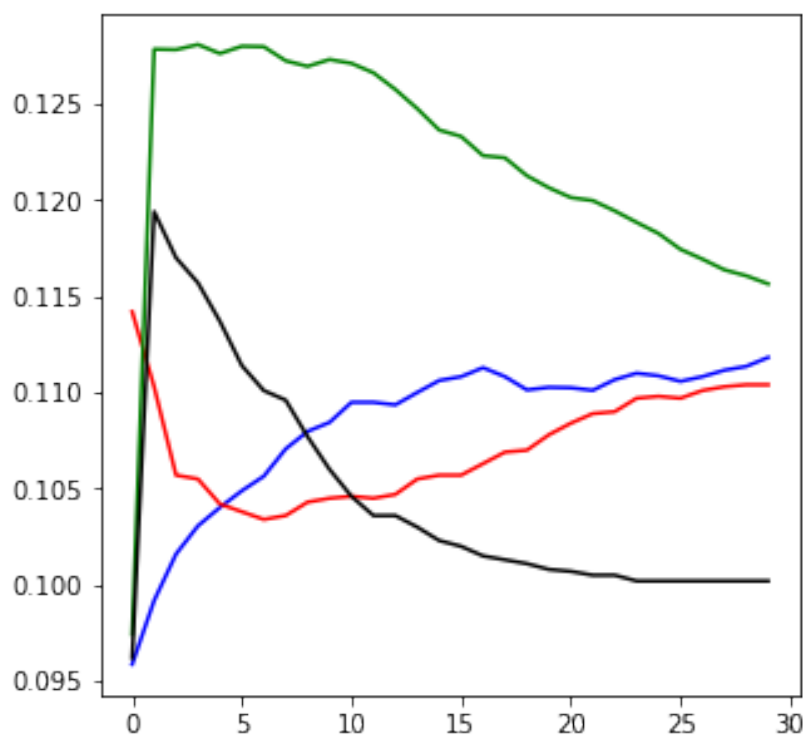
شکل ۱۰: نمودار دقت با learning rate مربوط به الگوریتم rmsprob

$1e-5+6.932e-1$



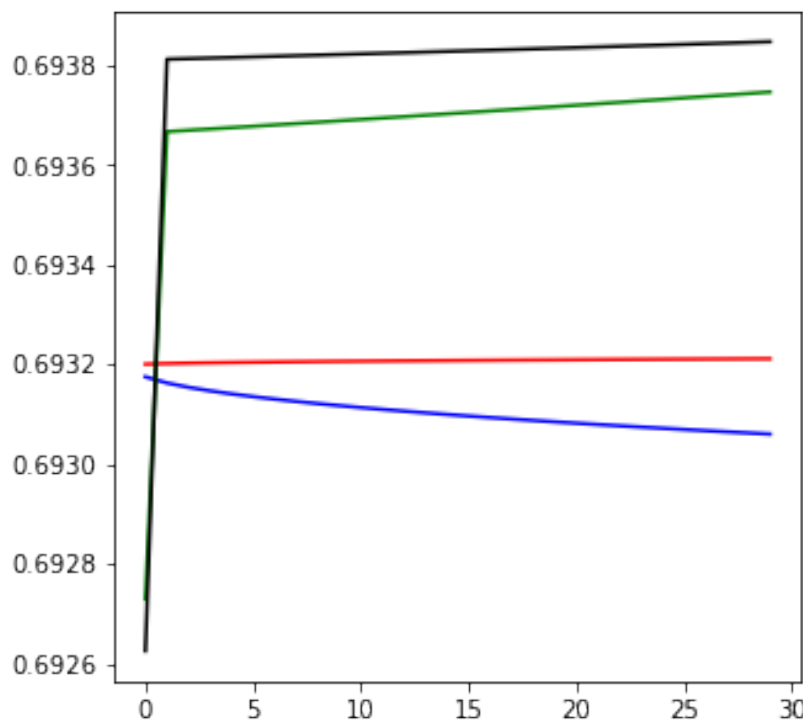
۵.

شکل ۱۱: نمودار خطا با learning rate مربوط به الگوریتم rmsprob



۵.

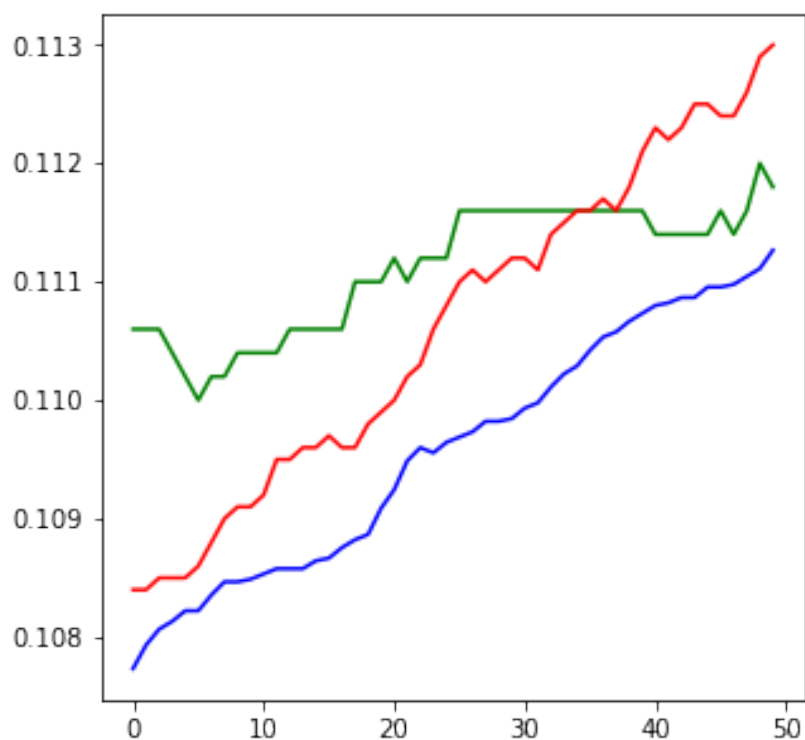
شکل ۱۲: نمودار دقت با learning rate مربوط به الگوریتم adograd و adadelta میباید که نمودار سبز مربوط به داده های آموزش و مشکی داده های تست adagrad و آبی آموزش و قرمز تست adadelta میباید



۵.

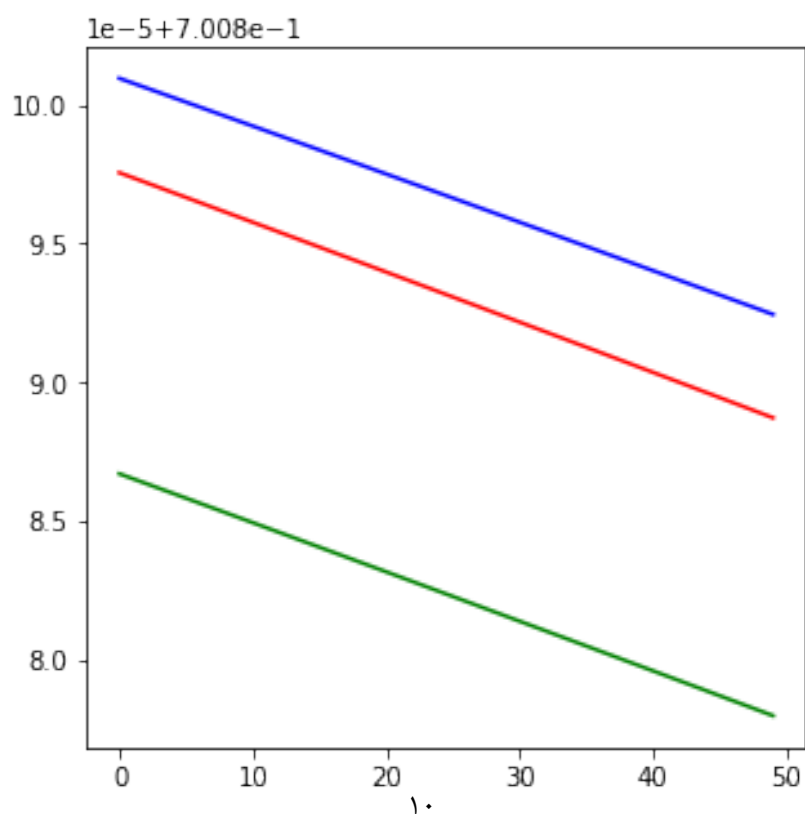
۹

شکل ۱۳: نمودار خطا با learning rate مربوط به الگوریتم adograd و adadelta میباید که نمودار سبز مربوط به داده های آموزش و مشکی داده های تست adagrad و آبی آموزش و قرمز تست adadelta میباید



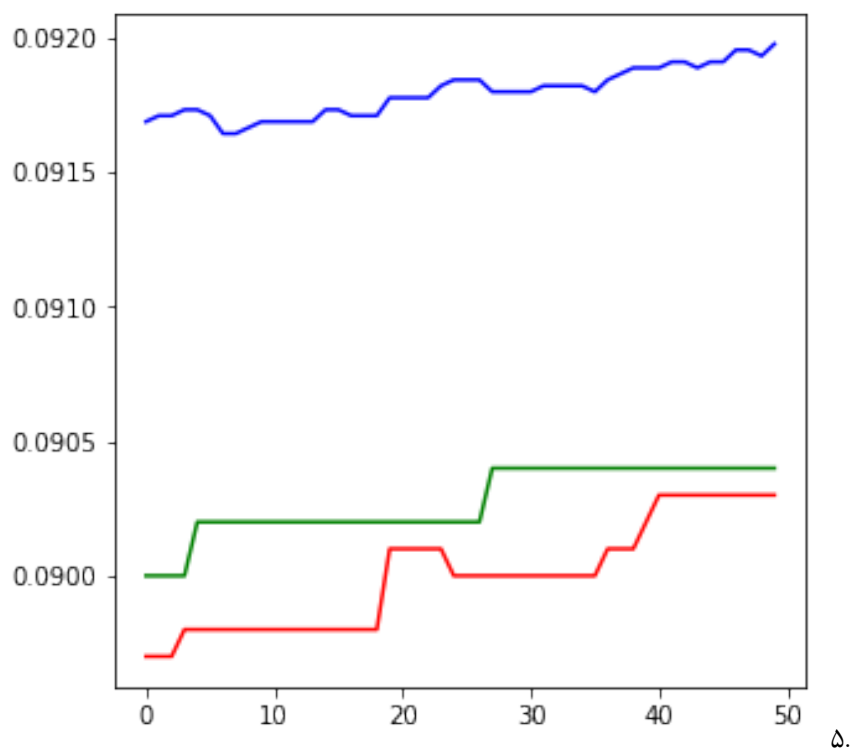
۵.

شکل ۱۴: نمودار دقت مدلی که  $L_2$  Regularization در آن استفاده شده است.

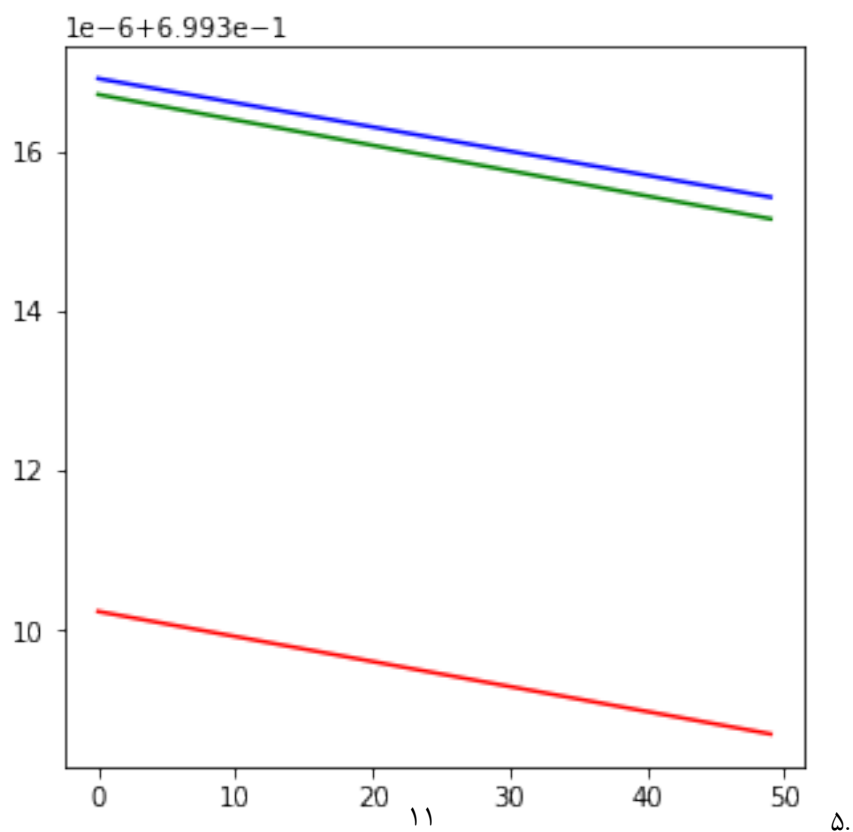


۵.

شکل ۱۵: نمودار خطا مدلی که  $L_2$  Regularization در آن استفاده شده است.



شکل ۱۶: نمودار دقت مدلی که L1 Regularization در آن استفاده شده است.



شکل ۱۷: نمودار خطا مدلی که L1 Regularization در آن استفاده شده است.