



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده ی ریاضی و علوم کامپیوتر

پایان نامه کارشناسی ارشد  
گرایش علوم کامپیوتر

مقایسه روشهای فراابتکاری برای بهینه سازی ضرایب  
شبکه عصبی

گزارش هشتم

نگارش

پویا پارسا

استاد راهنما

دکتر قطعی

خرداد ۱۴۰۰

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

# فصل اول

## دیتاست

## ۱-۱ مقدمه

در این گزارش سعی بر بررسی ساز و کار بهینه ساز های متفاوت و هم چنین استفاده از آن ها در آموزش یک شبکه عصبی ساده پرداخته ایم. دو بهینه ساز مد نظر PSO که جزء روش های شاخص بدون گرادیان و از آن طرف بهینه ساز Adam که در زمره ی بهترین روش گرادیان محور قرار می گیرند بوده اند.

## ۲-۱ اطلاعات کلی

به دلیل تمرکز این گزارش بر روی بهینه ساز ها ، از استفاده ی دیتاست پیچیده پرهیز کرده ام و از دیتاست معروف و قدیمی Iris استفاده کرده ام. این دیتاست اطلاعات سه نوع گل را دربردارد و شامل ۱۵۰ رکورد می باشد و برای هر رکورد ۴ فیچر شامل اطلاعاتی در مورد طول گلبرگ و کاسبرگ و ... را داراست.

## ۳-۱ معیار استفاده شده

معیار استفاده شده برای ارزیابی نتایج برای هر دو بهینه سازی Accuracy بوده است. که در گزارش قبل به دلیل انتخاب آن اشاره شده است.

# فصل دوم

## پیاده سازی

## ۱-۲ بهینه سازی ازدحام ذرات یا PSO

کدهای استفاده شده در این گزارش از طریق **این لینک** در دسترس است. Particle swarm optimization یا PSO [۵] یک روش جمعیت محور تصادفی برای بهینه سازی است در این روش از گرادیان استفاده نمی شود و در زمره ی روش های Swarm Intelligence قرار می گیرد. مزایای این روش به طور موردی به شرح زیر است :

- لزومی ندارد تابع هدف مشتق پذیر باشد.
  - تعداد هایپر پارامتر های خیلی کمی دارد
  - به راحتی می توان آن را درک کرد
  - و بر روی طیف وسیعی از مسائل می توان آن را به کار برد.
- به نظر می رسد این روش الهام گرفته از رفتار اجتماعی پرندگان و دسته های بزرگ ماهی برای پیدا کردن غذا است.
- کلیت این روش به این شکل است که تعدادی ذره به طور تصادفی در فضای جستجو پراکنده می شوند؛ در ابتدا جهت حرکت و سرعت آن ها نیز به طور تصادفی تعیین می شود سپس در هر مرحله جهت و سرعت ذرات با توجه به بهترین جواب پیدا شده توسط خود ذره تا آن لحظه و همچنین بهترین جواب پیدا شده توسط کل جمعیت به روز می شود. به بیان دقیق تر:

$$P_i^{t+1} = P_i^t + V_i^{t+1}$$

$$V_i^{t+1} = \underbrace{wV_i^t}_{\text{Inertia}} + \underbrace{c_1r_1(P_{best(i)}^t - P_i^t)}_{\text{Cognitive (Personal)}} + \underbrace{c_2r_2(P_{bestglobal}^t - P_i^t)}_{\text{Social (Global)}}$$

شکل ۱-۲: نحوه تعیین سرعت ذره  $i$  در  $t+1$  iteration.

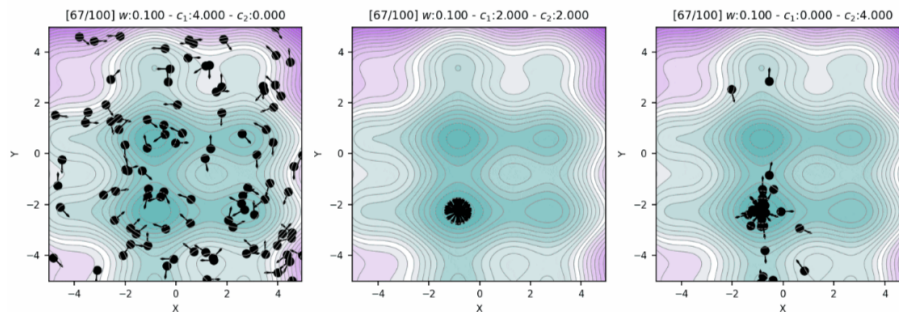
در فرمول بالا  $r_1, r_2$  ضرایبی تصادفی هستند که رفتار ذره ما را تعیین می کنند در واقع چه میزان به ندای درون و چه میزان به پیام های جمعی توجه کند! دقت داشته باشید که  $r_1, r_2$  هایپر پارامتر نیستند و برای هر ذره و در هر مرحله متفاوتند

### ۱-۱-۲ هایپر پارامتر ها

۱. تمایل به تغییر سرعت یا inertia weight در واقع قدرت تغییر مسیر ذره ها را مشخص می کند هر چه این پارامتر بزرگتر باشد گردش و اکتشاف ذرات بیشتر و هر چه کمتر باشد همگرایی قوی تر و بهره وری از نقاط یافت بیشتر است.

۲. تاثیر پذیری از جمع یا خود  $c_1$ ,  $c_2$

$c_1$  میزان تاثیر پذیری از پاسخ یافت شده توسط خود ذره و در مقابل آن  $c_2$  میزان تاثیر گرفتن ذرات از بهترین پاسخ یافت شده را مشخص می کنند. یک تعادل مناسب بین این دو مقدار باعث می شود هم ذرات فضای اطراف خود را جستجو کنند و هم به سمت بهترین جواب های پیدا شده همگرا شوند.



شکل ۲-۲: بررسی تاثیر  $c_1$  و  $c_2$  های متفاوت

## ۲-۱-۲ کد

کد های استفاده شده از [۳] هستند و در این بخش بیشتر به توصیف آن ها می پردازم.

```
1 def logits_function(p):
2     """ Calculate roll-back the weights and biases
3
4     Inputs
5     -----
6     p: np.ndarray
7         The dimensions should include an unrolled version of the
8         weights and biases.
9
10    Returns
11    -----
12    numpy.ndarray of logits for layer 2
13
14    """
15    # Roll-back the weights and biases
16    W1 = p[0:80].reshape((n_inputs,n_hidden))
17    b1 = p[80:100].reshape((n_hidden,))
18    W2 = p[100:160].reshape((n_hidden,n_classes))
19    b2 = p[160:163].reshape((n_classes,))
20
21    # Perform forward propagation
22    z1 = X.dot(W1) + b1 # Pre-activation in Layer 1
23    a1 = np.tanh(z1)    # Activation in Layer 1
24    logits = a1.dot(W2) + b2 # Pre-activation in Layer 2
25    return logits      # Logits for Layer 2
```

شکل ۲-۳: نحوه ی بازگردان ذره از حالت برداری به وزن های شبکه

برای ساخت ذره در واقع ما وزن های بین لایه ها ( به همراه بایاس ) را در یک بردار تک بعدی می ریزیم و به نوعی ساختار شبکه ی عصبی را صاف می کنیم. هر ذره برداری است با ابعاد وزن های شبکه عصبی. و در هر مرحله برای محاسبه ی fitness ذره بردار متناظر با آن را به ساختار شبکه ی عصبی آن برمی گردانیم. این عملیات در تکه کد بالا انجام می شود:

و با انجام یک forward propagation و محاسبه ی خطا میزان بهینگی ذره را به دست می آوریم. همانطور که می بینید ابتدا با استفاده از Softmax خروجی را تبدیل به یک توزیع احتمالی می کنیم و سپس از طریق تابع خطایمان که negative likelihood است و ضابطه ی این چینی دارد

$$L_i = -\log(p_{y_i})$$

که در آن  $p_{y_i}$  مقدار احتمال حساب شده توسط مدل برای کلاس صحیح است.

```

1  # Forward propagation
2  def forward_prop(params):
3      """Forward propagation as objective function
4
5      This computes for the forward propagation of the neural network, as
6      well as the loss.
7
8      Inputs
9      -----
10     params: np.ndarray
11         The dimensions should include an unrolled version of the
12         weights and biases.
13
14     Returns
15     -----
16     float
17         The computed negative log-likelihood loss given the parameters
18     """
19
20     logits = logits_function(params)
21
22     # Compute for the softmax of the logits
23     exp_scores = np.exp(logits)
24     probs = exp_scores / np.sum(exp_scores, axis=1, keepdims=True)
25
26     # Compute for the negative log likelihood
27
28     correct_logprobs = -np.log(probs[range(num_samples), y])
29     loss = np.sum(correct_logprobs) / num_samples
30
31     return loss
32

```

شکل ۲-۴: محاسبه ی خطا برای یک ذره.

در انتها با استفاده از کتاب خانه ی PySwarms به جستجوی وزن های بهینه می پردازیم و با استفاده از معیار Accuracy دقت مدل را می سنجیم همان طور که می بینید بر روی دیتاست ساده ی Iris به دقت ۹۹ درصد رسیدیم که نشان دهنده کارا بودن این الگوریتم است.



```

1 # Initialize swarm
2 options = {'c1': 0.5, 'c2': 0.3, 'w': 0.9}
3
4 # Call instance of PSO
5 dimensions = (n_inputs * n_hidden) + (n_hidden * n_classes) + n_hidden + n_classes
6 optimizer = ps.Single.GlobalBestPSO(n_particles=100, dimensions=dimensions, options=options)
7
8 # Perform optimization
9 cost, pos = optimizer.optimize(f, iters=1000)

```

```

1 def predict(pos):
2     """
3     Use the trained weights to perform class predictions.
4
5     Inputs
6     -----
7     pos: numpy.ndarray
8         Position matrix found by the swarm. Will be rolled
9         into weights and biases.
10    """
11    logits = logits_function(pos)
12    y_pred = np.argmax(logits, axis=1)
13    return y_pred

```

```

1 (predict(pos) == y).mean()

```

0.9933333333333333

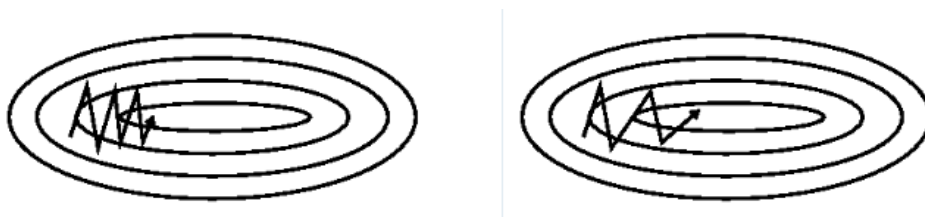
شکل ۲-۵: استفاده از کتاب خانه ی PySwarms و محاسبه ی دقت.

## ۲-۲ بهینه ساز Adam [۲]

Adam از جمله روش های تدریجی (Iterative) که در واقع بسط داده شده ی روش Stochastic Gradient Descent است. SGD روشی بر پایه گرادیان است که به شدت در کاربر های بینایی ماشین و NLP به کار گرفته شده است. می دانیم که شبکه های عصبی را می توان تابعی از ورودی ها به خروجی ها دید روش SGD با مشتق گیری از تابع خطا نسبت به وزن های شبکه، با استفاده از مشتق زنجیره ای اقدام به آپدیت کردن وزن های شبکه می کند.

تا اینجا همه ی توضیحات درباره ی SGD بود اما Adam دقیقاً چگونه کار می کند ؟ Adam که از کلمات adaptive moment estimation برگرفته شده است بر خلاف SGD که learning rate ثابتی را برای کل روند train در نظر می گیرد؛ learning rate متمایزی را بر هر یک از پارامتر های شبکه نگهداری می کند و هر یک از این پارامتر ها را با توجه به تغییرات گرادیان سازگار و آپدیت می کند. در واقع شاید بتوان گفت روش Adam مزیت های روش AdaGrad و RMSProp را با هم داراست

در واقع به جای آن که مانند RMSProp تنها از Momentum مرتبه اول استفاده کند از Momentum مرحله ی دوم نیز استفاده می کند. در توضیح Momentum نیز باید بگوییم به جای آنکه تنها گرادیان در مرحله ی آخر را در نظر بگیریم می توانیم گرادیان های مرحله ی قبل را نیز دخیل کنیم که این کار منجر به حذف نویز از گرادیان می شود و تشکیل یک میانگین هندسی از گرادیان ها می شود.



شکل ۲-۶: استفاده از momentum و تاثیر آن در جهت گرادیان [۴].

## هایپر پارامترها

۱. آلفا

همان میزان تاثیر گرادیان جدید در محاسبه آپدیت وزن هاست

۲. beta1

ضریب کاهش گشتاور مرتبه ی اول در میانگین هندسی

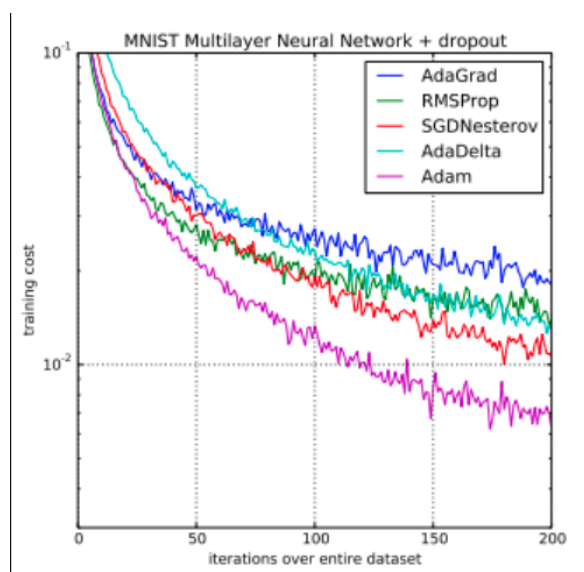
۳. beta2

ضریب کاهش گشتاور مرتبه ی دوم در میانگین هندسی

۴. epsilon

یک عدد بسیار کوچک برای جلوگیری از تقسیم بر صفر در طی عملیات.

با وجود الگوریتم های مشابه که نام آنان در بالا ذکر شد توصیه ی بسیاری از متخصصان بر استفاده از Adam در وهله ی اول می باشد [۲] در تصویر پایین عملکرد بهینه ساز های مختلف را در آموزش یک شبکه ی پرسپترون چند لایه می بینید.



شکل ۲-۷: عملکرد بهینه ساز های متفاوت در آموزش.

## کد ۱-۲-۲

کد استفاده شده از [۱] می باشد. با استفاده از یک ساختار شبکه عصبی مشابه اما این بار در pytorch به استفاده از بهینه ساز Adam می پردازیم

```

1 class Model(nn.Module):
2     def __init__(self):
3         super(Model, self).__init__()
4         self.layer1 = nn.Linear(n_inputs, n_hidden)
5         self.layer2 = nn.Linear(n_hidden, n_classes)
6
7     def forward(self, x):
8         x = F.relu(self.layer1(x))
9         x = F.softmax(self.layer2(x), dim=1)
10        return x

```

```

1 model = Model()
2 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
3 loss_fn = nn.CrossEntropyLoss()
4 model

```

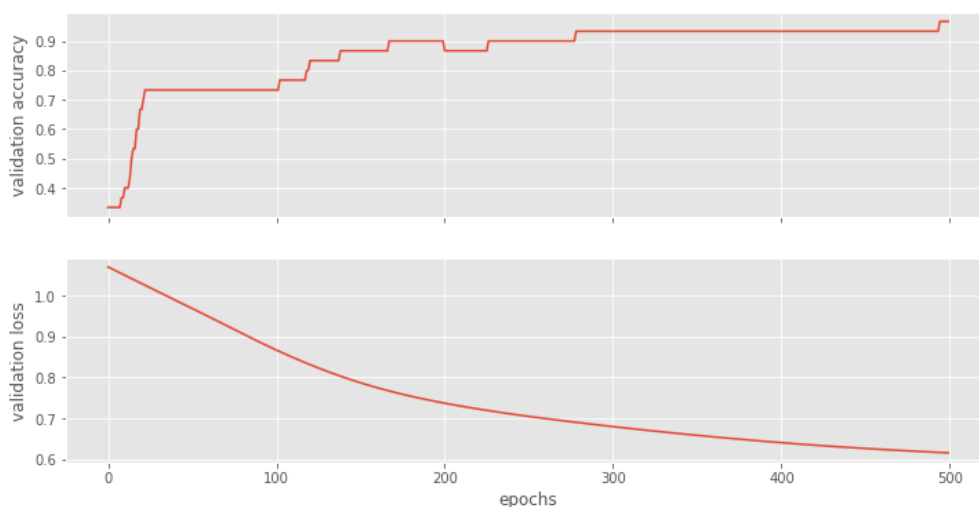
```

Model(
  (layer1): Linear(in_features=4, out_features=20, bias=True)
  (layer2): Linear(in_features=20, out_features=3, bias=True)
)

```

شکل ۲-۸: پیاده سازی شبکه و استفاده از بهینه ساز adam

همچنین روند یادگیری مدل را در نمودار زیر می توانیم مشاهده کنیم. همان طور که دیده می شود دقت تقریباً مشابه ای را بر روی داده های تست به دست می دهد.



شکل ۲-۹: روند بهبود در آموزش توسط Adam

## ۳-۲ جمع بندی

در آزمایش های انجام شده در این گزارش روش PSO هم سرعت همگرایی بیشتری و هم دقت بالاتری را به نمایش گذاشت البته به دلیل کوچک بودن مدل و دیتاست باید آزمایش های تکمیلی انجام داد.

## منابع و مراجع

- [1] Brownlee, Jason. Classifying the iris data set with pytorch. <https://janakiev.com/blog/pytorch-iris/>.
- [2] Brownlee, Jason. Gentle introduction to the adam optimization algorithm for deep learning. <https://machinelearningmastery.com/adam-optimization-algorithm-for-deep-learning/>: :text=Adam
- [3] pyswarms. Pyswarms training a neural network. [https://pyswarms.readthedocs.io/en/latest/examples/usecases/train\\_neural\\_network.html](https://pyswarms.readthedocs.io/en/latest/examples/usecases/train_neural_network.html).
- [4] RUDER, SEBASTIAN. An overview of gradient descent optimization algorithms. <https://ruder.io/optimizing-gradient-descent/index.html#adam>.
- [5] Thevenot, Axel. Particle swarm optimization (pso) visually explained. <https://towardsdatascience.com/particle-swarm-optimization-visually-explained-46289eeb2e14>.