

Topic: Optimization for Neural Networks

- Background:

- ① Neural Networks (NN) are a broad class of predictive models that learns the relationship between the features and the outcomes.
- ② Training a NN involves updating the parameters describing the connections in order to minimize some loss function over the training data.
- ③ NN have been applied to broad range of fields and reached the state-of-the-art performance in many tasks.
- ④ NN require large datasets and sufficient computing power.
- ⑤ Optimizing NN means reducing overfitting, having a lower tendency to get stuck in local saddle points, or converging in a smaller number of epochs.

- Goal:

- ① Code Deep Neural Network from scratch. Apply regularization and momentum techniques. Main Resource: [*A Computational Approach to Statistical Learning*](#)
- ② Implement two variants of SGD as the optimizers for the NN and experiment with simulated data. Compare their performances.

Review of the Basics

Architecture

- Layers ($l \in 0 : L$)
 - ▶ Input Layer x
 - ▶ Output Layer \hat{y}
 - ▶ Hidden Layer(s)
- Nodes (Neurons) a^l
 - ▶ Weights w^l
 - ▶ Biases b^l
- Activation Function σ
 - ▶ ReLU
 - ▶ Softmax (before output layer)
- Relationship (Forward propagation)
 - ▶ $a^0 = x$
 - ▶ $z^l = w^l a^{l-1} + b^l$
 - ▶ $a^l = \sigma(z^l)$

Backward Propagation

- Goal: Adjusting the weights to reach the minimized loss function
- Gradient descent updates:
$$w^{l+1} = w^l - \eta \nabla_w f(w)$$
$$b^{l+1} = b^l - \eta \nabla_b f(b)$$
- NNs generally need many thousands of iterations to converge to a reasonable minimizer of the loss function. With large datasets and models, while still feasible, gradient descent can become quite slow. To avoid working with the entire dataset at each step, we use SGD, which incrementally updates the weights in the model.

Backward Propagation

- Notation: For each $l \in (0, L)$, $w_{j,k}^l$ denotes the weights on the k th neuron in the $(l-1)$ -st layer within the j th neuron in the l st layer; b_j^l denotes the bias term for the j th neuron in the l st layer

- SGD updates:

$$w^{l+1} = w^l - (\eta/n) \nabla_{w^l} f(w^l)$$

$$b^{l+1} = b^l - (\eta/n) \nabla_{b^l} f(b^l)$$

- ▶ Loss: $f = (a^L - y)^2$
(regression)

$$f = -\sum_k y_k \cdot \log(a_k^L)$$

(classification)

- ▶ $\partial_{b_j^l} f = \frac{\partial f}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = \frac{\partial f}{\partial z_j^l}$
- ▶ $\partial_{w_{j,k}^l} f = \frac{\partial f}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{j,k}^l} = \frac{\partial f}{\partial z_j^l} a_k^{l-1}$

$$\triangleright z_k^{l+1} = \sum_m w_{k,m}^{l+1} \sigma(z_m^l) + b_k^{l+1}$$

$$\triangleright \frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{k,j}^{l+1} \sigma'(z_j^l)$$

$$\triangleright \frac{\partial f}{\partial z_j^l} = \sum_k \frac{\partial f}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \frac{\partial f}{\partial z_k^{l+1}} w_{k,j}^{l+1} \sigma'(z_j^l)$$

$$\triangleright \Rightarrow \frac{\partial f}{\partial z^l} = \left[(w^{l+1})^T \left(\frac{\partial f}{\partial z^{l+1}} \right) \right] \odot \sigma'(z^l)$$

$$\triangleright (\text{Regression}) \quad \frac{\partial f}{\partial z_j^l} = \sum_k \frac{\partial f}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^l} =$$

$$\frac{\partial f}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^l} = \frac{\partial f}{\partial a_j^L} \sigma'(z_j^l)$$

$$\triangleright (\text{Classification}) \quad \frac{\partial f}{\partial z_j^l} = a_j^L - y_i$$

Regularization

- Problem: Overfitting
- Solution: Includes a penalty term directly into the loss function. For example, adding an l_2 penalty yields

$$f_{\lambda}(w, b) = f(w, b) + \frac{\lambda}{2} \|w\|_2^2$$
$$\nabla_w f_{\lambda} = \nabla_w f(w, b) + \lambda \cdot w$$

- The updating rule becomes

$$\begin{aligned} w^{l+1} &= w^l - \eta \nabla_w f(w^{l-1}) \\ &= w^l - \eta [\nabla_w f(w^{l-1}) + \lambda w^{l-1}] \\ &= [1 - \eta \lambda] w^{l-1} - \eta \nabla_w f(w^{l-1}) \end{aligned}$$

SGD with Momentum

- Problem: Relying only on the gradient leads to models getting stuck at saddle points and being very sensitive to the learning rate.
- Solution: Incorporating the momentum vector v into the algorithm
- Advantages:
 - 1 if the gradient remains relatively unchanged over several steps it will 'pick-up speed' (momentum) and effectively use a larger learning rate;
 - 2 if the gradient is changing rapidly, the step-sizes will shrink accordingly;
 - 3 when passing through a saddle point, the built up momentum from prior steps helps propel the algorithm past the point.
- Updating Rule:

$$\begin{aligned}v^l &= v^{l-1} \cdot \mu - \eta \cdot \nabla_w f(w^{l-1}) \\w^l &= (1 - \eta\lambda)w^{l-1} + v^l\end{aligned}$$

where $\mu \in [0, 1]$

Stochastic Average Gradient (SAG)

- Roux et al. 2012
- SAG incorporates a memory of previous gradient values in order to achieve a linear convergence rate.
- Updating Rule

$$w^{l+1} = w^l - \frac{\eta}{n} \sum_{i=1}^n y_i^l$$

where at each iteration a random training example i_l is selected and we set

$$y_i^l = \begin{cases} f'_i(w^l) & \text{if } i = i_l \\ y_i^{l-1} & \text{otherwise} \end{cases}$$

Stochastic Variance Reduced Gradient (SVRG)

- Johnson et al. 2013
- SVRG was developed to reduce the variance of the gradients
- Start with \tilde{w} which is close to the optimal w , and calculate the average gradient by one pass over the data using \tilde{w}

$$\tilde{\mu} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\tilde{w})$$

- Updating Rule:
Draw i_t randomly from $\{1, \dots, n\}$

$$w^{(l+1)} = w^{(l)} - \eta[\nabla f_{i_t}(w^{(l)}) - \nabla f_{i_t}(\tilde{w}) + \tilde{\mu}]$$

Update \tilde{w} and $\tilde{\mu}$ after m SGD iterations.

SAGA

- Defazio et al. 2014
- SAGA is inspired by SVRG. Instead of approximating the parameters w , SAGA stores the gradient vectors for each sample.

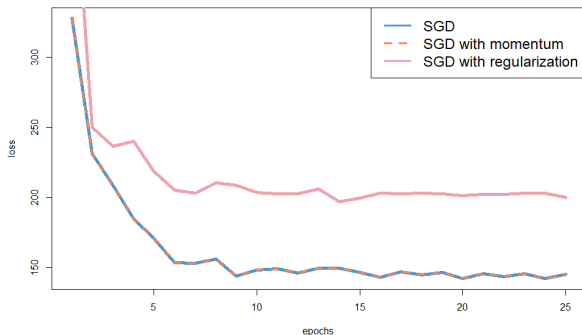
```
1: Choose an initial iterate  $w_1 \in \mathbb{R}^d$  and stepsize  $\alpha > 0$ .
2: for  $i = 1, \dots, n$  do
3:   Compute  $\nabla f_i(w_1)$ .
4:   Store  $\nabla f_i(w_{[i]}) \leftarrow \nabla f_i(w_1)$ .
5: for  $k = 1, 2, \dots$  do
6:   Choose  $j$  uniformly in  $\{1, \dots, n\}$ .
7:   Compute  $\nabla f_j(w_k)$ .
8:   Set  $g_k \leftarrow \nabla f_j(w_k) - \nabla f_j(w_{[j]}) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_{[i]})$ .
9:   Store  $\nabla f_j(w_{[j]}) \leftarrow \nabla f_j(w_k)$ .
10:  Set  $w_{k+1} \leftarrow w_k - \alpha g_k$ .
```


Simulation Studies

- Simple case:

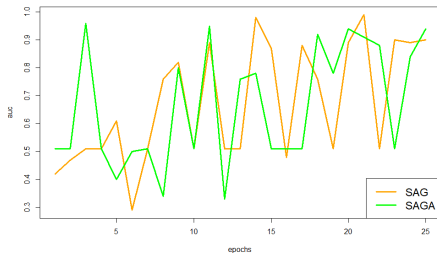
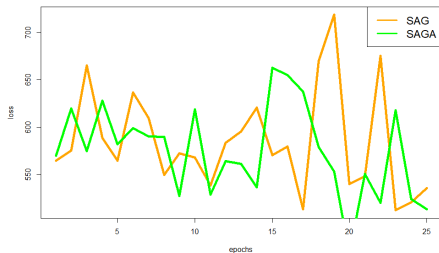
- ▶ Data: $X \sim 1000 \times 1$, where $x_i \sim \text{Unif}(-1, 1)$
 $y_i = x_i^2 + \epsilon$ where $\epsilon \sim N(0, 0.1)$
 $y_i = I_{(y_i > 0.5)}$
- ▶ Total epochs: 25
Learning rate: 0.01

SGD, SGD with momentum (0.09), and SGD with regularization (0.01)



AUC: 0.98 ~ 0.99 with 0.5 as the discriminating threshold

SAG vs SAGA



2021 Phenology Data (Blooming)

- Real data from *rnpn* package
- training data 500×5

	latitude	longitude	elevation_in_meters	day_of_year	phenophase_status
1	36.44084	-84.32359	384	1	0
2	36.44084	-84.32359	384	1	0
3	36.44084	-84.32359	384	1	0
4	36.44084	-84.32359	384	1	0
5	36.44084	-84.32359	384	1	0
6	36.44084	-84.32359	384	1	0
7	32.39724	-111.00216	801	1	0
8	32.39724	-111.00216	801	1	0
9	32.39724	-111.00216	801	1	0
10	32.39724	-111.00216	801	1	0
11	32.37423	-110.96838	800	1	0
12	32.37423	-110.96838	800	1	0
13	32.37423	-110.96838	800	1	0
14	32.37423	-110.96838	800	1	0

Training NN only using 1 feature (*day of year*)

	epoch	Hidden Layers (25 nodes)	loss (median)
SGD	50	10	538.87
SAGA	50	10	453.21
SAG	50	10	893.72

Training NN using all 4 features

	epoch	Hidden Layers (25 nodes)	loss (median)
SGD	50	10	1336.4
SGD *	50	10	940.38
SAGA	50	10	782.92
SAGA*	50	10	658.24
SAG	50	10	1557
SAG*	50	10	1763.28

Conclusions

- Summary

- ▶ For the simple case where X is uniform and Y is X^2 plus a normal noise, all types of NNs are working great, showing promising AUC values. While for the real data, where the relationship between features and responses is unclear, SAGA displays better performances. SAG is performing the worst. Also by updating the gradients selectively, SAG and SAGA require less cpu time than SGD.
- ▶ If we exclude the cases of human error, we may reach the conclusion that SAG and SAGA have the potential of complex machine learning tasks where tens of thousand feature may exist.

- Limits

- ▶ learning rate