

Non-convex cost function optimization

Adagrad, RMSProp, Adam

Alex-Antoine Fortin

American Family Insurance

March 31, 2017



What role is playing your optimizer?

- Computing the cost of your intermediate model
- Compute $\nabla_{\theta} Cost$
- Tuning the learning rate
- Computing the parameter update to your model
- and more ...

delta-bar-delta (Jacobs, 1988)

Let

$$\nabla_{\theta} Cost = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$$

then, update:

$$\theta^* \leftarrow \theta + \epsilon \nabla_{\theta} Cost$$

Intuition:

if $sign(\nabla Cost; \theta)_i = sign(\nabla Cost; \theta^*)_i$; then ϵ should be bigger

if $sign(\nabla Cost; \theta)_i \neq sign(\nabla Cost; \theta^*)_i$; then ϵ should be smaller

Idea

Parameters are scaled down proportionally to the square root of their historical value

Net effect

More progress in the gently sloped directions of the parameter space

Algorithm 1 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ for numerical stability

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

while stopping criterion not met **do**

Get minibatch: $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$

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: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$

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

end while

Idea

AdaGrad is really good for convex optimization. In non-convex optimization, the learning trajectory may pass through many structures and eventually arrive at a locally convex region.

RMSProp use an EMA to discard updates from the distant past so that it can converge rapidly once it found a locally convex region.

Thus, RMSProp introduces a parameter ρ that controls the length scale of the MA.

Algorithm 2 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ

Require: Initial parameter θ

Require: Small constant δ for numerical stability

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

while stopping criterion not met **do**

 Get minibatch: $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$

 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 update: $\Delta \theta \leftarrow -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$

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

end while

Adam stands for *Adaptative moments*.

Idea

RMSProp has the capacity to tune ϵ for each θ_i individually as well as to forget distant past. However, the squared gradient accumulator, \mathbf{r} , initialized to $\mathbf{0}$ introduces a bias toward $\mathbf{0}$ in its calculation.

Adam introduces a bias correction term to both the calculated gradient and its second-moment, \mathbf{r} . It also adds momentum to the estimation of the gradient.

Algorithm 3 The Adam algorithm

Require: Global learning rate ϵ , decay rate ρ_1, ρ_2 , param θ , const δ

Initialize accumulation variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$ and $t = 0$

while stopping criterion not met **do**

Get minibatch: $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$

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: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

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

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: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\hat{\mathbf{r}}}} \odot \hat{\mathbf{s}}$

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

end while
