

A comparative evaluation of gradient-based optimization algorithms for training Terrain GAN

Kai Qin

Dept. of Electrical and Computer Engineering
University of Texas at Austin
Austin, TX
kai.qin@utexas.edu

Yi Han

Dept. of Electrical and Computer Engineering
University of Texas at Austin
Austin, TX
yh5598@utexas.edu

I. INTRODUCTION

Generative Adversarial Networks (GAN) [1] have been widely used in generative applications such as anime generation, human-face generation, and video generation. In this study, we focus on evaluating the performance of three gradient-based optimization algorithms in training GANs for procedural terrain generation [2]. "Procedural terrain generation for video games has been traditionally done with smartly designed but handcrafted algorithms that generate hightmaps" [2] *{DESCRIBE dataset}* The first goal of this project is for us as a term project to extend what we have learned in class about gradient-based optimization algorithms and to dive deeper in this topic by experimenting with three state-of-the-art gradient based algorithm for training neural networks. The second goal is to explore the practical differences that each optimizer can make for TerrianGAN. The third goal is to understands the results of fine-tuning parameters on the different optimizers and their effects on TerrianGAN.

II. GRADIENT-BASED OPTIMIZATION ALGORITHMS REVIEW

Before we dive into the evaluation of three gradient-based optimizers we compared in this study, SGD with momentum, RMSProp, and ADAM, we want to provide a review of these three algorithms. We will start from the classic gradient descent algorithm. In one sentence, the general idea of the classic gradient descent algorithm is that at every iteration, the desired parameters are moving in the direction of the negtive gradient of the objective function based on the entire training dataset to decrease the loss function. The vanilla gradient descent may be accelerated considerably by using stochastic gradient descent (SGD) which follows the gradient of randomly selected minibatches. Even though SGD ingtroduced this very important idea of training with minibatches, learning with it can sometimes be very slow. For example, let's say that you're trying to optimize a cost function which has a contour in image 1. The red dot denotes the position of the minimum. And we can see this up and own oscillations in black lines slows down the gradient descent. Therefore, it's rear to see large scale neural network training with classic SGD. One improvement made to SGD is by adding momentum. The basic idea of SGD with momentum is to compute an exponentially

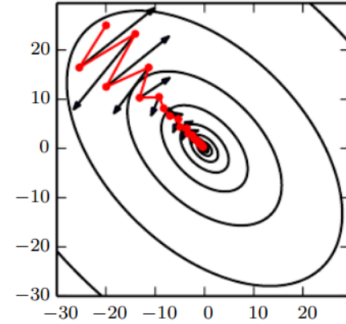


Fig. 1. Visualization of the effect of momentum [?]

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

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

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

Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

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

end while

Fig. 2. SGD with momentum [?]

decaying moving average of the past gradients and continues to move in that direction [?]. The velocity update step in algorithm 2 shows how to compute the exponentially weighted averages of the last N iterations, where N is determined by the hyperparameter α . For example, when α equals 0.99, which is commonly used in practice, we are approximately averaging over the last 100 iterations. This method can only approximate the average value over the last N days, but it takes very little memory. Image 1 illustrates the effect of momentum. If you average out these gradients of the black arrows, you can find that the oscillations in the forward diagonal direction will tend to average out to something closer to zero. In the mean time, it will keep the direvative on the backward diagonal direction. So this allows the optimizer to take a more straight foward path or to damp out the oscillations in the path to the minimum.

Root mean squared prop (RMSProp) is another algorithm

Require: Global learning rate ϵ , decay rate ρ .
Require: Initial parameter θ
Require: Small constant δ , 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

Fig. 3. RMSProp [?]

that can speed up the classic gradient descent. This algorithm (see figure in 3) is very similar to SGD with momentum, but instead of accumulating the gradient, we are taking the average of the squared gradient and divide the gradient by the square root of the average squared gradient [?]. By dividing the average magnitude of the derivative in each dimension, the net effect is that the forward diagonal derivative in image 1 is divided by a relatively larger number, and it therefore helps damp out the oscillations. Comparing to SGD with momentum, we can use a larger learning rate, and get faster learning without diverging in the forward diagonal direction, where in SGD with momentum, the learning rate is applying to all each dimension with same step size.

By combining both ideas of momentum and RMSProp, we get Adam optimization algorithm which has been shown to work well across a wider range of deep learning architectures. In Adam, in each iteration, in the first highlighted equation in figure 4 we first update the first moment estimate (the mean of the derivatives), which is exactly what we had when we're implementing SGD with momentum. And similarly, we do the rms prop to update the second moment estimate in the second equation. Then we do bias corrections. Finally, we compute the update by dividing the bias corrected first moment estimate with the square root of the second moment estimate and reversing the sign. Common choices for ρ_1 , ρ_2 , σ are 0.9, 0.99 and 10^{-8} respectively. The common practice for Adam hyperparameter tuning is to keep ρ_1 , ρ_2 , and σ default and try a range of values of the learning rate to see which works the best.

III. GAN, DCGAN, AND LSGAN REVIEW

GAN was originally proposed by Ian Goodfellow, et al. in [1]. DCGAN [?] showed us that GAN is capable of generating perceptually good samples (see figure 5 and interpolations with fine-tuned architecture details. Figure 6 shows the pseudocode for GAN [1]. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$ in all of our experiments in the next section. The basic architecture for LSGAN [?] is the same as DCGAN, but instead of using the cross entropy as the loss function, it adopted the least squares lossfunction since the original loss function may lead to the vanishing gradients problem. Also in the original Terrain GAN

Require: Step size ϵ (Suggested default: 0.001)
Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$. (Suggested defaults: 0.9 and 0.999 respectively)
Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})
Require: Initial parameters θ
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: $\Delta \theta = -\frac{\epsilon \hat{\mathbf{s}}}{\sqrt{\delta + \hat{\mathbf{r}}}}$ (operations applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta \theta$
end while

Fig. 4. ADAM [?]



Fig. 5. DCGAN samples on faces [?]

paper [2], the author claims that LSGAN shows better results than DCGAN.

IV. EXPERIMENTS AND RESULTS

As we can see in the graph on the right, the loss function of both G and D oscillate and don't converge. That's one

for number of training iterations **do**
 for k steps **do**
 • Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
 • Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
 • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)})))]$$

 end for
 • Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
 • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)})))$$

end for
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Fig. 6. Minibatch stochastic gradient descent training of GAN [1]

of the big challenge of GAN training is that you don't get this clean monotonic improvement that you are used to with training supervised models. In GANs, the objective function for the generator and the discriminator usually measures how well they are doing relative to the opponent. For example, we measure how well the generator is fooling the discriminator.

Sometime it oscillates and doesn't converge. That's one of the big challenges of GAN training is that you don't get this clean monotonic improvement that you are used to with training supervised models.

"It is still an open problem to have good metrics. If you had a really really good metric, you can probably use it as an optimization criteria and optimize against it" "Assuming you don't do optimization directly on the Frechet Inception Distance, it can be a nice independent measure of the amount of variation and crispness of the image generated ; when you do optimize against it, you will find results not exactly what you want" Adam 0.0002 lr and 16 batch size, we think based on the 3d model generated by this software looks realistic enough to meet our experience. We manually swept the lr rate and batch size using ADAM and collected the loss function over training iterations. Next step will collect results using other optimizers and figure out a metric to quantify the results we see. The discriminator can always

V. DISCUSSION ABOUT METRICS

As we have seen in loss graph, the oscillation of the loss makes it harder to measure the training progress. Currently, it is still an open problem for gan training to have good metrics. On the other hand, If you had a really really good metric, you can probably use it as an optimization criteria and optimize against it. Inception score and FID are two state of the art metrics people use to measure the performance. However, both metrics don't apply to the data we. Inception score requires categorical labeled data. The Frechet Inception Distance score (FID) requires a pre-trained classification neural network. Since the goal of gan is to learn the probability distribution from the training data, and implicitly represent it using neural network. Thus we have proposed a new metric to measure the Terrain GAN performance, the KL-divergence between normalized intensity histogram of training images and generated images.

This is how way calculate this metric, first: Image intensity histogram, which is a widely used technique in image processing. The definition is that Histograms plots how many times (frequency) each intensity value in an image occurs. Here we have an example of intensity histogram of a generated height map.

However we want the distribution over a large set of images, the original intensity histogram doesn't work here, so we calculate the average occurrences of each intensity value and divide it by the total number of pixels over all images, and define it the normalized histogram. The top image shows the normalized intensity histogram of the training images. The bottom image shows the normalized intensity histogram of 50 generated image using the neural net trained by Adam with lr

.0002 and batchsize 16 for 150 epochs. We can see that they both have mean around 30 gray levels.

Finally, Here is the KL divergence equation. In our case the $P(x)$ is the real dist and Q_X is the generated height map. From the normalized histogram we can tell that the ADAM has more overlap than the SGD. This also matches what we see from the generated images. And the KL divergence value confirmed our observations. The SGD with this set of hyperparameter has no overlap with the training data thus the KL div is $+\infty$.

Here is a bar plot of the KL divergence of different hyperparameter settings of adam after 150 epochs. Visually, we can see that the higher the metric, the harder you can provide some geographic meaning to the generated images.

Again, we have all the results of the lowest KL divergence model.

1 pixel is 1 square km; Original image pixel 21600x10800, and with 128x128 there are 14000, we excluded the ocean by excluding images with 90 percent pixels with max min value less than 0.05 (total range is from 0 to 1). And we end up with 4305. We have excluded flat land and waters, so we expect no such images from the output. We propose a metric of $\max - \min / \min$ as a metric.

Compare performance of each optimizer: we need a metric to tell if the generated image is good or bad. Based on our goal, we want noticeable features like mountains or lakes, valley, coasts. Compare min-max diff histogram of the training data and the output data.

Distribution distances: KL-divergence

VI. CONCLUSION

VII. FUTURE WORKS

REFERENCES

- [1] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," 2014.
- [2] C. Beckham and C. Pal, "A step towards procedural terrain generation with gans," 2017.