

Unfavorable SNR in Forward Gradients

Charles Frye

February 2022

The forward gradient method of [1] allows the computation of approximate gradients without the need for a separate “backwards pass”.

The backwards pass has the same computational complexity, in the big- \mathcal{O} sense, as the forwards pass, and so doesn’t change the complexity class of optimization algorithms that require forwards passes. But the backwards pass typically incurs a fixed (low, i.e. 2-4x) multiplicative cost relative to forward computation and occasionally requires a large amount of memory, so there can be meaningful practical benefits to avoiding it.

There is, however, a catch to the forward gradient method that is not clearly stated in [1]: the signal-to-noise ratio (SNR) using the forward gradient method is at most $\mathcal{O}(1/n)$, where n is the number of parameters.

High SNR is good, and so this bound is bad news. Noise increases the effective Lipschitz constant and so decreases the highest stable learning rate. Lower learning rates mean slower convergence (in iteration count), which can dwarf any gains from faster iterations. Specifically, this bound implies that for any model, once n is large enough, descent by forward gradients will converge more slowly in wall time than will backward gradients¹.

We therefore should expect (and do observe) poorer performance when directly substituting a forward gradient oracle in for a backwards gradient oracle. Perhaps additional tricks, like gradient clipping or a special optimizer, might make the forward gradient method more applicable in more settings.

Specifically, denoting the (full) gradient ∇ and its forward estimate g , we have that each entry of g has variance approximately equal to² the gradient norm:

$$\mathbb{V}[g]_i \approx \|\nabla\|^2 \tag{1}$$

Assuming that the gradient ∇ has $\mathcal{O}(1)$ norm and homogeneous entries, we can obtain an asymptotic, in parameter count, expression for the SNR:

$$\frac{\nabla_i^2}{\mathbb{V}[g]_i} \approx \frac{1/n}{\|\nabla\|^2} = \mathcal{O}(1/n) \tag{2}$$

¹Space constraints might still favor forward gradients.

²In the exact expression it is slightly higher.

Forward Gradient Descent

Let ∇ be the full-batch gradient. We calculate an estimator g as in the forward gradient method [1]:

$$g \stackrel{\text{def}}{=} (\nabla^\top v) \cdot v \quad (3)$$

where $v \sim \mathcal{N}(0, I)$ is a random Gaussian vector of length n . Effectively, the forward gradient method produces a biased random walk, taking larger steps when the randomly chosen vector v happens to align with the gradient. $\nabla^\top v$ can be calculated efficiently as part of the same forward pass that computes the loss, as it is a Jacobian-vector product.

It is proven in [1] that the expectation of g is ∇ . This is sufficient to prove convergence of the algorithm by any number of established means, treating the forward gradient as the (noisy) gradient oracle.

But what about the variance?

We can break down the variance of the i th entry of g as

$$\begin{aligned} \mathbb{V}[g_i] &= \mathbb{V}[(\nabla^\top v) \cdot v_i] && \text{(Def'n of } g) \\ &= \mathbb{V}\left[\sum_j v_j v_i \nabla_j\right] && \text{(Def'n of } \nabla) \\ &= \mathbb{V}[v_i^2 \nabla_i] + \mathbb{V}\left[\sum_{j \neq i} v_j v_i \nabla_j\right] + 2 \cdot \text{Cov}\left(v_i^2 \nabla_i, \sum_{j \neq i} v_j v_i \nabla_j\right) \\ &&& \text{(Variance of sum formula)} \end{aligned}$$

Notice that the variance of a single entry g_i depends on a (random) sum over all of the entries.

The covariance term is 0:

$$\begin{aligned} \text{Cov}\left(v_i^2 \nabla_i, \sum_{j \neq i} v_j v_i \nabla_j\right) &= \mathbb{E}\left[v_i^2 \nabla_i \cdot \sum_{j \neq i} v_j v_i \nabla_j\right] - \mathbb{E}[v_i^2 \nabla_i] \mathbb{E}\left[\sum_{j \neq i} v_j v_i \nabla_j\right] \\ &&& \text{(Def'n of covariance)} \\ &= \mathbb{E}\left[v_i^2 \nabla_i \cdot \sum_{j \neq i} v_j v_i \nabla_j\right] && (\mathbb{E}[v_i^2] = 0) \\ &= \nabla^\top \mathbb{E}[v^2] && \text{(Def'n of } \nabla, \text{ Linearity of } \mathbb{E}) \\ &= 0 && (\mathbb{E}[v_j^2] = 0) \end{aligned}$$

and so we have:

$$\mathbb{V}[g_i] = \mathbb{V}[v_i^2 \nabla_i] + \mathbb{V}\left[\sum_{j \neq i} v_j v_i \nabla_j\right] \quad (4)$$

We tackle the simpler term first:

$$\mathbb{V} [v_i^2 \nabla_i] = \nabla_i^2 \mathbb{V} [v_i^2] \quad (\text{Scaling rule for } \mathbb{V})$$

and notice that v_i is a unit normal random variable. Products of unit normals follow distributions from the χ^2 family. These distributions are indexed by the degree-of-freedom parameter k and their variances are given by $2k$. In this case, $k = 1$ and so the variance is 2.

$$\mathbb{V} [v_i^2 \nabla_i] = 2 \cdot \nabla_i^2 \quad (\mathbb{V} [\chi^2(1)] = 2)$$

Now we tackle the sum term:

$$\mathbb{V} \left[\sum_{j \neq i} v_j v_i \nabla_j \right] = \sum_{j \neq i} \mathbb{V} [v_j v_i \nabla_j] + 2 \cdot \sum_{j \neq i \neq k, j < k} \text{Cov} (v_j v_i \nabla_j, v_k v_i \nabla_k) \quad (\text{Variance of sum formula})$$

Switching things up this time, we try the harder term first and peek into the covariances in the sum:

$$\begin{aligned} \text{Cov} (v_j v_i \nabla_j, v_k v_i \nabla_k) &= \nabla_j \nabla_k \text{Cov} (v_j v_i, v_k v_i) \quad (\text{Bilinearity of the covariance}) \\ \text{Cov} (v_j v_i, v_k v_i) &= \mathbb{E} [v_j v_i v_i v_k] + \mathbb{E} [v_j v_i] \cdot \mathbb{E} [v_k v_i] \quad (\text{Def'n of the covariance}) \\ &= \mathbb{E} [v_j] \mathbb{E} [v_i v_i] \mathbb{E} [v_k] + \mathbb{E} [v_j] \mathbb{E} [v_i] \cdot \mathbb{E} [v_k] \mathbb{E} [v_i] \quad (\text{Product rule for } \mathbb{E}) \\ &= 0 \quad (\text{Expectations of } v_i, v_j, \text{ and } v_k) \end{aligned}$$

That is, the covariance terms are all identically 0,³ so we can return our focus to the simpler term in the variance of the sum of products:

$$\begin{aligned} \mathbb{V} \left[\sum_{j \neq i} v_j v_i \nabla_j \right] &= \sum_{j \neq i} \mathbb{V} [v_j v_i \nabla_j] + 0 \quad (\text{From above}) \\ &= \sum_{j \neq i} \nabla_j^2 \mathbb{V} [v_j v_i] \quad (\text{Scaling rule for } \mathbb{V}) \end{aligned} \quad (5)$$

We again have a product of normals, but it's not quite a χ^2 . Instead we use a trick⁴ to express $v_i v_j$ as the difference of two independent χ^2 variables, which is again χ^2 - *distributed*

³Intuition: knowing the value of $v_j v_i$, e.g. that it is far from its mean, does help you guess the value of $v_k v_i$, e.g. that it will be far from its mean, so the variables are dependent. But knowing which direction $v_j v_i$ has deviated from its mean doesn't help you guess which direction $v_k v_i$ has deviated from its mean, since the signs of v_j and v_k could differ or match with equal probability, and therefore the variables are uncorrelated.

⁴Read more at <https://math.stackexchange.com/questions/101062/is-the-product-of-two-gaussian-random-variables-also-a-gaussian>

$$\begin{aligned}
\mathbb{V}[v_j v_i] &= \mathbb{V}\left[\frac{1}{4}(v_i + v_j)^2 - \frac{1}{4}(v_i - v_j)^2\right] && \text{(Clever rewrite trick)} \\
&= \frac{1}{16} (\mathbb{V}[(v_i + v_j)^2] + \mathbb{V}[(v_i - v_j)^2]) && \text{(Sum of variance for independent r.v.s)} \\
&= \frac{1}{4} (\mathbb{V}[(v_i + v_j)^2/2] + \mathbb{V}[(v_i - v_j)^2/2]) && \text{(Normalize to unit variance, pull out factor of 4)} \\
&= \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 = 1 && (\mathbb{V}[\chi^2(1)] = 2)
\end{aligned}$$

We can now combine our two expressions for the variance terms in the sum-of-variances decomposition

$$\mathbb{V}[g_i] = \mathbb{V}[v_i^2 \nabla_i] + \sum_{j \neq i} \nabla_j^2 \mathbb{V}[v_j v_i] \quad (6)$$

$$= 2 \cdot \nabla_i^2 + \sum_{j \neq i} \nabla_j^2 \quad (7)$$

When n is large, the $i \neq j$ sum predominates, and so we can introduce the approximation

$$\mathbb{V}[g_i] = 2 \cdot \nabla_i^2 + \sum_{j \neq i} \nabla_j^2 \quad (8)$$

$$\approx \sum_j \nabla_j^2 = \|\nabla\|^2 \quad (9)$$

Is this so bad? At first glance, things look fine – the variance is finite and it’s on the scale of the gradient norm, which is often manipulated to be $\mathcal{O}(1)$ to obtain favorable numerics in optimization.

But consider what this means for the signal-to-noise ratio: if the gradient norm is of fixed size, the entries must be getting smaller as more parameters are added to the model. If they are homogeneous, each ∇_i should be of absolute value around $\frac{1}{\sqrt{n}}$, and so we obtain

$$\frac{\nabla_i^2}{\mathbb{V}[g]_i} \approx \frac{1/n}{\|\nabla\|^2} = \mathcal{O}(1/n) \quad (10)$$

as desired.

Commentary

The proof above, like [1], assumes full-batch gradients, rather than mini-batch gradients. Experiments and rough calculations suggest that the problem of forward gradient variance becomes worse in the stochastic setting.

We note that, because it looks like a product of normals, g has unfavorable characteristics in its higher moments, not just in its variance (or second moment). In simulation, we observed highly kurtotic (heavy-tailed) distributions with prominent peaks at 0.

There are two intuitions for why this is the case: one at the level of individual entries and one at the level of the entire vector. When the randomly drawn value of v_i is small, the gradient value is multiplied by the much smaller value v_i^2 and vice versa when the randomly drawn value of v_i is large. Alternatively, as the dimension of the space increases, two randomly chosen vectors become increasingly orthogonal. The inner product $\nabla^\top v$ thus has a very sparse distribution and the vectors $\nabla^\top v \cdot v$ have norms with a very right-skew distribution.

Occasional large parameter updates produced by heavy-tailed noise can cause numerical issues for gradient descent algorithms, e.g. by throwing the parameters into an area of the function where the local Lipschitz constant is much higher, turning a convergent learning rate into a divergent one.

These distributional issues suggest that an alternative distribution for the entries of v might help speed up forward gradient. Without a priori information about the gradient vector, the isotropy of the unit normal is a helpful property; an anisotropic distribution might be even less likely to produce vectors v that . Note also that the proof of unbiasedness in [1] requires that the components of the vector be independent, and the unit normal is the only isotropic distribution with independent components.

One useful form of anisotropic distribution would be to produce probing vectors v that are more likely to point in the direction of the gradient. Indeed, if the distribution of v is a Dirac delta at ∇ , the variance of g drops to 0 and forward gradients align perfectly with true gradients. Of course, in that case there is no need to compute a gradient estimate, because the gradient is already at hand, but perhaps there is a happy medium between complete ignorance and complete knowledge of the gradient where the forward gradient becomes optimal.

References

- [1] Atılım Güneş Baydin, Barak A. Pearlmutter, Don Syme, Frank Wood, and Philip Torr. Gradients without backpropagation, 2022.