

<https://arxiv.org/pdf/1502.01852.pdf>

<https://stats.stackexchange.com/questions/15978/variance-of-product-of-dependent-variables>

What are we trying to initialize?

<https://datascience.stackexchange.com/questions/82917/why-do-we-want-the-variance-of-the-layers-to-remain-the-same-throughout-a-deep-n>

What exactly is helped by the variance being stable in this way?

Following the notation of the article, let's compute the gradient of the cost function w.r.t. the parameters in two consecutive layers (that we are going to call: layer i and layer $i + 1$). In this layers, the quantity used to update their respective weight matrices, is given by:

$$\text{Layer 1} \rightarrow \frac{\partial \text{Cost}}{\partial W^i} = -\frac{\partial \text{Cost}}{\partial s^i} \frac{\partial s^i}{\partial W^i} = -\frac{\partial \text{Cost}}{\partial s^i} (z^{i-1})^T$$

$$\text{Layer 2} \rightarrow \frac{\partial \text{Cost}}{\partial W^{i+1}} = -\frac{\partial \text{Cost}}{\partial s^{i+1}} \frac{\partial s^{i+1}}{\partial W^{i+1}} = -\frac{\partial \text{Cost}}{\partial s^{i+1}} (z^i)^T$$

There we can see that if we consider:

$$\text{Var} \left(\frac{\partial \text{Cost}}{\partial s^i} \right) = \text{Var} \left(\frac{\partial \text{Cost}}{\partial s^{i+1}} \right) \leftrightarrow \text{Var}(z^i) = \text{Var}(z^{i-1})$$

Then we would have:

$$\text{Var} \left(\frac{\partial \text{Cost}}{\partial W^i} \right) = \text{Var} \left(\frac{\partial \text{Cost}}{\partial W^{i+1}} \right)$$

This is a good thing because having the same variance in the updates of both layers means that the updates are globally spread in the same way, so assuming that the mean value of $\partial \text{Cost} / \partial W$ in both layers is the same, then this would mean that globally these **layers are learning at the same rythm**.

Whatever the answer to (2) is, what is the proof or evidence that that's the case?

A good advantage of the previous reasoning is that if we could achieve this happening in the whole neural network, then all the layers in the NN would be learning at the same rythm! So problems like vanishing or exploding gradient would be avoided.

$$\begin{aligned} \text{var}(XY) &= E((XY - E(XY))^2) \\ &= E(X^2Y^2) - E(XY)^2 \end{aligned}$$

by definition

$$\begin{aligned} \text{cov}(X, Y) &= E((X - E(X))(Y - E(Y))) \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

and

$$\text{cov}(X^2, Y^2) = E(X^2Y^2) - E(X^2)E(Y^2)$$

$$\text{so } \text{var}(XY) = \text{cov}(X^2, Y^2) + E(X^2)E(Y^2) - (\text{cov}(X, Y) + E(X)E(Y))^2$$

$$\text{if } \text{cov}(X, Y) = \text{cov}(X^2, Y^2) = 0,$$

$$\text{var}(XY) = E(X^2)E(Y^2) - (E(X)E(Y))^2$$

$$\begin{aligned}
&= (\text{var}(X) + E(X)^2)(\text{var}(Y) + E(Y)^2) - (E(X)E(Y))^2 \\
&= \text{var}(X)\text{var}(Y) + \text{var}(X)E(Y)^2 + \text{var}(Y)E(X)^2
\end{aligned}$$

we want to keep the variance of gradient updates $\frac{\partial L}{\partial w} = X^t \frac{\partial L}{\partial y}$ constant across every layer, so that the variance can be kept the same across every layer

assuming X^t and $\frac{\partial L}{\partial y}$ at each level are independent,

$$\text{var}\left(\frac{\partial L}{\partial w}\right) = \text{var}\left(X^t \frac{\partial L}{\partial y}\right) = \text{var}(X^t)\text{var}\left(\frac{\partial L}{\partial y}\right) + \text{var}(X^t)E\left(\frac{\partial L}{\partial y}\right)^2 + \text{var}\left(\frac{\partial L}{\partial y}\right)E(X^t)^2$$

if we can keep $\text{var}(X)$ and $\text{var}\left(\frac{\partial L}{\partial y}\right)$ constant across all the layers, then $\text{var}\left(\frac{\partial L}{\partial w}\right)$ will also be constant!

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$\text{var}(\mathbf{y}) = \text{trace of covariance matrix of } \mathbf{y} = \text{var}(\sum w_i x_i) = n \text{var}(w_i x_i)$

assuming \mathbf{x} and \mathbf{W} are mutually independent

and $E(w) = 0$

$$\begin{aligned}
\text{var}(w_i x_i) &= \text{var}(w_i)\text{var}(x_i) + \text{var}(w_i)E(x_i)^2 + \text{var}(x_i)E(w_i)^2 \\
&= \text{var}(w_i)(\text{var}(x_i) + E(x_i)^2) = \text{var}(w_i)E(x_i^2)
\end{aligned}$$

$$\text{so } \text{var}(\mathbf{y}) = n_l \text{var}(w_i)E(x_i^2)$$

if w^{l-1} has a symmetric distribution around 0, then $y^{l-1} \equiv f^{-1}(x^l)$ has a symmetric distribution around 0 and mean 0, and if f is RELU

$$\begin{aligned}
\text{then } E(x^2) &= E(f(y)^2) \\
&= E(0.5y^2) = 0.5E(y^2) = \frac{1}{2}\text{var}(y)
\end{aligned}$$

$$\text{so } \text{var}(\mathbf{y}^l) = \frac{1}{2}n_l \text{var}(w_i) \text{var}(y^{l-1})$$

this yields the recursive relationship

$$\text{var}(\mathbf{y}^l) = \text{var}(\mathbf{y}^1) \prod_{l=2}^L \frac{1}{2}n_l \text{var}(w^l)$$

A proper initialization method should avoid reducing or magnifying the magnitudes of input signals exponentially.

to keep $\text{var}(\mathbf{y}^l) \equiv \text{var}(\mathbf{y}^1)$, we need to make $\prod_{l=2}^L \frac{1}{2} n_l \text{var}(w^l) = 1 \rightarrow$ so if we set $\text{var}(w^l) = 2/n_l$ for every level, we might be able to!

For the backpropagation, something similar, see the paper for details