

Neural Nets: Initialisation Strategies

David Barber

Initialisation

- Particularly for saturating transfer functions that have regions of zero gradient, sensible initialisation of the weights is essential.
- Consider a simple regression net with error

$$\sum_n (y^n - f(\mathbf{w}^T \mathbf{x}^n))^2$$

- Firstly, it makes sense to scale the error so that it remains roughly order 1, even as the number of datapoints increases. For example

$$E(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (y^n - f(\mathbf{w}^T \mathbf{x}^n))^2$$

- Gradient is

$$E(\mathbf{w}) = -\frac{2}{N} \sum_{n=1}^N (y^n - f(\mathbf{w}^T \mathbf{x}^n)) f'(\mathbf{w}^T \mathbf{x}^n) \mathbf{x}^n$$

- For $f(x) = 1/(1 + e^{-x})$, for large $|x|$, this function saturates to 1 or 0 and the gradient becomes zero. We need therefore to ensure we don't immediately get trapped in these zero gradient regions.

Initialisation

- Let's assume that we have scaled the inputs x_i so that they have zero mean and unit variance

$$x_i^n \rightarrow \frac{x_i^n - \mu_i}{\sigma_i}$$

where

$$\mu_i = \sum_{n=1}^N x_i^n, \quad \sigma_i^2 = \sum_{n=1}^N (x_i^n - \mu_i)^2$$

After this rescaling,

$$\frac{1}{N} \sum_{n=1}^N x_i^n = 0, \quad \frac{1}{N} \sum_{n=1}^N (x_i^n)^2 = 1$$

Initialisation

- Define the activation (input of the transfer function) as

$$z_n \equiv \mathbf{w}^\top \mathbf{x}^n$$

Let's assume that we sample each w_i identically and independently from a zero mean distribution $\mathbb{E}(w_i) = 0$.

$$\mathbb{E}(z_n) = \sum_i \mathbb{E}(w_i) x_i^n = 0$$

$$\begin{aligned}\mathbb{E}(z_n^2) &= \sum_{i \neq j} \mathbb{E}(w_i) \mathbb{E}(w_j) x_i^n x_j^n + \sum_i \mathbb{E}(w_i^2) (x_i^n)^2 \\ &= \mathbb{E}(w^2) \sum_{i=1}^D (x_i^n)^2\end{aligned}$$

- Hence if we sample the w_i from a distribution with zero mean and variance $1/D$, then the activation will be zero mean with variance

$$\mathbb{E}(z_n^2) = \frac{1}{D} \sum_{i=1}^D (x_i^n)^2$$

Initialisation

- Since each x_i^n has values roughly in the range -2 to 2 (2 standard deviations from the zero mean), then the activation will have values roughly in the range -2 to 2 as well.
- A reasonable initialisation then of the weights of a network with transfer function $f(x)$ that has a non-saturating region when x is between -2 and 2, is: Draw each w_i from a distribution with zero mean and unit variance and then rescale. For example:

$$w_i \sim \mathcal{N}(w_i|0, 1), \quad w_i \rightarrow \frac{w_i}{\sqrt{D}}$$

$$w_i \sim \text{sign}(\mathcal{N}(w_i|0, 1)), \quad w_i \rightarrow \frac{w_i}{\sqrt{D}}$$

$$w_i \sim \mathcal{U}(w_i|-\sqrt{3}, \sqrt{3}), \quad w_i \rightarrow \frac{w_i}{\sqrt{D}}$$

- For nets with multiple hidden layers, the previous initialisation strategy is also commonly used. or weights in layer l , we set (for example)

$$w_i^l \sim \mathcal{N}(w_i^l|0, 1), \quad w_i^l \rightarrow \frac{w_i^l}{\sqrt{D_{l-1}}}$$

where D_l is the number of units in layer l .

Decorrelating Inputs

For an $\mathbf{x} \rightarrow y$ net the Hessian is

$$\mathbf{H} = -\frac{2}{N} \sum_n \underbrace{\left[(y^n - f(\mathbf{w}^T \mathbf{x}_n)) f''(\mathbf{w}^T \mathbf{x}_n) - (f'(\mathbf{w}^T \mathbf{x}_n))^2 \right]}_{\equiv \gamma_n} \mathbf{x}_n \mathbf{x}_n^T$$

As a very crude approximation, we can assume that all the γ_n are roughly equal to a value γ . In that case, the approximate Hessian is

$$\mathbf{H} \approx -\frac{2\gamma}{N} \sum_n \mathbf{x}_n \mathbf{x}_n^T$$

If we perform SVD of the matrix $\mathbf{X} = [\mathbf{x}^1, \dots, \mathbf{x}^N]$ then

$$\sum_n \mathbf{x}_n \mathbf{x}_n^T = \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{S} \mathbf{V}^T \mathbf{V} \mathbf{S}^T \mathbf{U}^T = \mathbf{U} \mathbf{S}^2 \mathbf{U}^T$$

Hence applying the 'whitening' transformation

$$\mathbf{x}^n \rightarrow \mathbf{S}^{-1} \mathbf{U}^T \mathbf{x}^n$$

will make the approximate Hessian diagonal and can be a useful initial preprocessing step to make optimisation easier.

A more complex discussion

- Let's consider what happens for the second hidden layer. This will have inputs

$$f(z_{nj}), \quad z_{nj} \equiv \mathbf{w}_j^T \mathbf{x}^n$$

where j is the index of the neuron in the first layer.

- And the activation to a unit in the second layer will be

$$\tilde{z}_n \equiv \sum_j u_j f(z_{nj})$$

- If we assume that $z_{nj} \sim \mathcal{N}(z_{nj}|0,1)$ and the u_j are independently sampled from a zero mean distribution then the activation of the second layer has zero mean and variance

$$\mathbb{E}(u^2) \sum_{j=1}^{D_1} (f(z_{nj}))^2$$

- If we therefore draw each u_j^l from a distribution with zero mean and variance

$$\frac{1}{D_{l-1} \mathbb{E}(f(z)^2)_{\mathcal{N}(z|0,1)}}$$

the activation of each unit j in layer l will be roughly zero mean with unit variance distributed.