# Initialization

- Consider standard building block of NN in terms of *preactivations*:

$$\mathbf{f}_k = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{h}_k$$
$$= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathrm{a}[\mathbf{f}_{k-1}]$$

- Set all the biases to 0

$$\boldsymbol{\beta}_k = \mathbf{0}$$

- Weights normally distributed
  - mean 0
  - variance $\sigma_\Omega^2$

- What will happen as we move through the network if $\sigma_\Omega^2$ is very small?
- What will happen as we move through the network if $\sigma_\Omega^2$ is very large?

# Backprop summary

**Backward pass:** We start with the derivative $\partial \ell_i / \partial \mathbf{f}_K$ of the loss function $\ell_i$ with respect to the network output $\mathbf{f}_K$ and work backward through the network:

$$\frac{\partial \ell_i}{\partial \boldsymbol{\beta}_k} = \frac{\partial \ell_i}{\partial \mathbf{f}_k} \qquad\qquad k \in \{K, K-1, \ldots 1\}$$

$$\frac{\partial \ell_i}{\partial \boldsymbol{\Omega}_k} = \frac{\partial \ell_i}{\partial \mathbf{f}_k} \mathbf{h}_k^T \qquad\qquad k \in \{K, K-1, \ldots 1\}$$

$$\frac{\partial \ell_i}{\partial \mathbf{f}_{k-1}} = \mathbb{I}[\mathbf{f}_{k-1} > 0] \odot \left( \boldsymbol{\Omega}_k^T \frac{\partial \ell_i}{\partial \mathbf{f}_k} \right), \qquad k \in \{K, K-1, \ldots 1\} \qquad (7.13)$$

where $\odot$ denotes pointwise multiplication and $\mathbb{I}[\mathbf{f}_{k-1} > 0]$ is a vector containing ones where $\mathbf{f}_{k-1}$ is greater than zero and zeros elsewhere. Finally, we compute the derivatives with respect to the first set of biases and weights:

$$\frac{\partial \ell_i}{\partial \boldsymbol{\beta}_0} = \frac{\partial \ell_i}{\partial \mathbf{f}_0}$$

$$\frac{\partial \ell_i}{\partial \boldsymbol{\Omega}_0} = \frac{\partial \ell_i}{\partial \mathbf{f}_0} \mathbf{x}_i^T$$

# Initialization

- Need for initialization
- He initialization
- Interlude: Expectations
- Show that $\mathbb{E}[f'_i] = 0$
- Write variance of pre-activations $f'$ in terms of activations $h$ in previous layer

$$\sigma^2_{f'} = \sigma^2_\Omega \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

- Write variance of pre-activations f' in terms of pre-activations f in previous layer

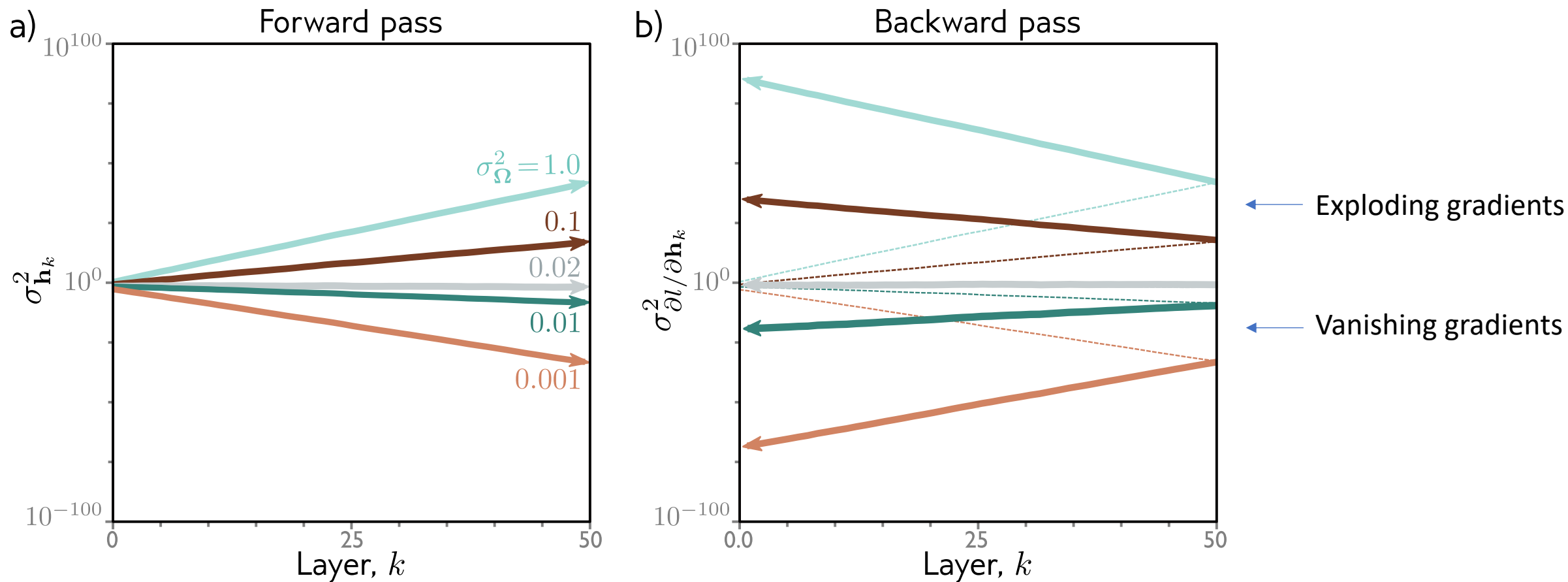$$\sigma^2_{f'} = \frac{D_h \sigma^2_\Omega \sigma^2_f}{2}$$

**Figure 7.4** Weight initialization. Consider a deep network with 50 hidden layers and $D_h = 100$ hidden units per layer. The network has a 100 dimensional input $\mathbf{x}$ initialized with values from a standard normal distribution, a single output fixed at $y = 0$, and a least squares loss function. The bias vectors $\boldsymbol{\beta}_k$ are initialized to zero and the weight matrices $\boldsymbol{\Omega}_k$ are initialized with a normal distribution with mean zero and five different variances $\sigma_{\boldsymbol{\Omega}}^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$. a)
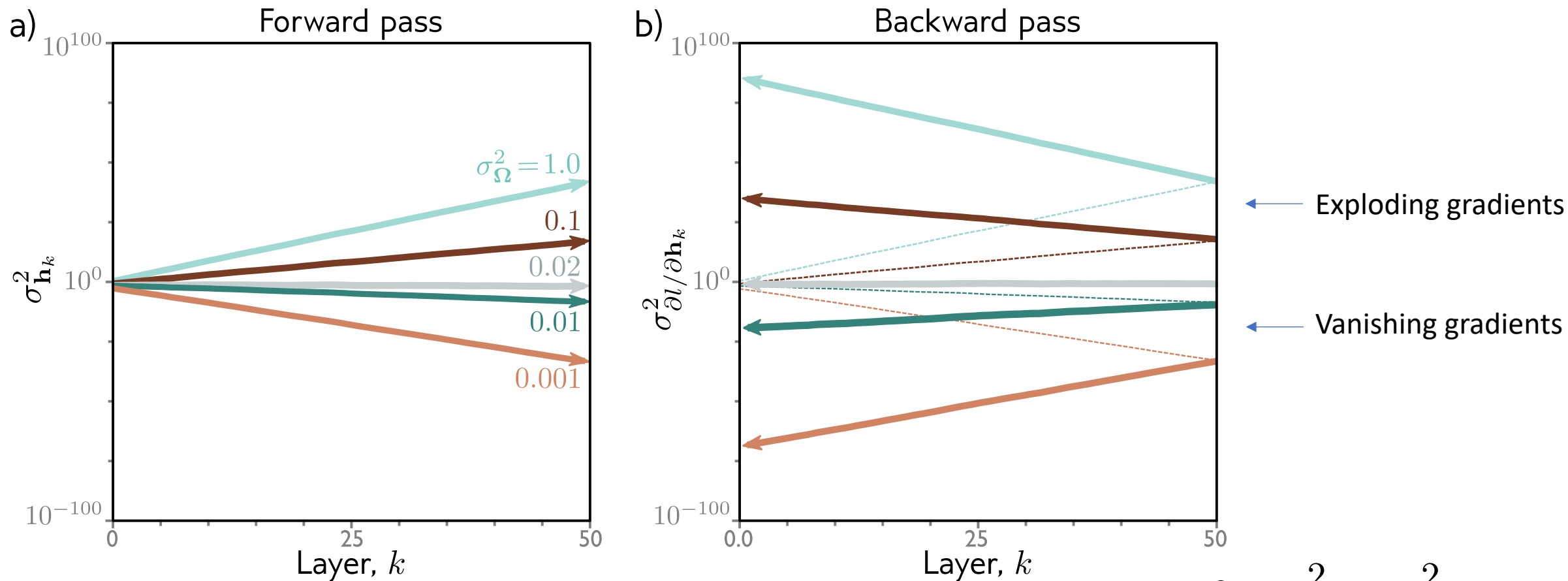
# He initialization (assumes ReLU)

- Forward pass:  want the variance of hidden unit activations in layer k+1 to be the same as variance of activations in layer k:

$$\sigma_\Omega^2 = \frac{2}{D_h}$$ ← Number of units at layer k

- Backward pass:  want the variance of gradients at layer k to be the same as variance of gradient in layer k+1:

$$\sigma_\Omega^2 = \frac{2}{D_{h'}}$$ ← Number of units at layer k+1

**Figure 7.4** Weight initialization. Consider a deep network with 50 hidden layers and $D_h = 100$ hidden units per layer. The network has a 100 dimensional input $\mathbf{x}$ initialized with values from a standard normal distribution, a single output fixed at $y = 0$, and a least squares loss function. The bias vectors $\boldsymbol{\beta}_k$ are initialized to zero and the weight matrices $\boldsymbol{\Omega}_k$ are initialized with a normal distribution with mean zero and five different variances $\sigma_{\boldsymbol{\Omega}}^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$. a)

# Expectations

| Function g[●] | Expectation |
|:---:|:---:|
| $x$ | mean, $\mu$ |
| $x^k$ | $k$th moment about zero |
| $(x - \mu)^k$ | $k$th moment about the mean |
| $(x - \mu)^2$ | variance |
| $(x - \mu)^3$ | skew |
| $(x - \mu)^4$ | kurtosis |

**Table B.1** Special cases of expectation. For some functions g[$x$], the expectation $\mathbb{E}[g[x]]$ is given a special name. Here we use the notation $\mu_x$ to represent the mean with respect to random variable $x$.

# Rules for manipulating expectation

$$\mathbb{E}\Big[k\Big] = k$$

$$\mathbb{E}\Big[k \cdot \mathrm{g}[x]\Big] = k \cdot \mathbb{E}\Big[\mathrm{g}[x]\Big]$$

$$\mathbb{E}\Big[\mathrm{f}[x] + \mathrm{g}[x]\Big] = \mathbb{E}\Big[\mathrm{f}[x]\Big] + \mathbb{E}\Big[\mathrm{g}[x]\Big]$$

$$\mathbb{E}\Big[\mathrm{f}[x]g[y]\Big] = \mathbb{E}\Big[\mathrm{f}[x]\Big]\mathbb{E}\Big[\mathrm{g}[y]\Big] \qquad \text{if} \quad x, y \quad \text{independent}$$

# Now let's prove:

$$\mathbb{E}\left[(x - \mu)^2\right] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

Keeping in mind:

$$\mathbb{E}[x] = \mu$$

Rule 1: $$\mathbb{E}\left[k\right] = k$$

Rule 2: $$\mathbb{E}\left[k \cdot \mathrm{g}[x]\right] = k \cdot \mathbb{E}\left[\mathrm{g}[x]\right]$$

Rule 3: $$\mathbb{E}\left[\mathrm{f}[x] + \mathrm{g}[x]\right] = \mathbb{E}\left[\mathrm{f}[x]\right] + \mathbb{E}\left[\mathrm{g}[x]\right]$$

Def'n $$\mathbb{E}[x] = \mu$$

$$\mathbb{E}[(x - \mu^2)] = \mathbb{E}[x^2 - 2x\mu + \mu^2]$$

$$= \mathbb{E}[x^2] - \mathbb{E}[2x\mu] + \mathbb{E}[\mu^2]$$

$$= \mathbb{E}[x^2] - 2\mu\mathbb{E}[x] + \mu^2$$

$$= \mathbb{E}[x^2] - 2\mu^2 + \mu^2$$

$$= \mathbb{E}[x^2] - \mu^2$$

$$= \mathbb{E}[x^2] - E[x]^2$$

# Initialization

- Need for initialization

- He initialization

- Interlude: Expectations

- Show that $\quad \mathbb{E}[f_i'] = 0$

- Write variance of pre-activations *f'* in terms of activations *h* in previous layer

$$\sigma_{f'}^2 = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

- Write variance of pre-activations f' in terms of pre-activations f in previous layer

$$\sigma_{f'}^2 = \frac{D_h \sigma_\Omega^2 \sigma_f^2}{2}$$

# Initialization

- Consider standard building block of NN in terms of *preactivations*:

$$\mathbf{f}_k = \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathbf{h}_k$$
$$= \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \mathrm{a}[\mathbf{f}_{k-1}]$$

- Set all the biases to 0

$$\boldsymbol{\beta}_k = \mathbf{0}$$

- Weights normally distributed
  - mean 0
  - variance $\sigma_\Omega^2$

- What will happen as we move through the network if $\sigma_\Omega^2$ is very small?
- What will happen as we move through the network if $\sigma_\Omega^2$ is very large?

# Aim: keep variance same between two layers

$$\mathbf{f}' = \boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{h}$$

Consider the mean of the pre-activations:

$$\mathbb{E}[f_i'] = \mathbb{E}\left[\beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j\right]$$

Rule 1: $\mathbb{E}\left[k\right] = k$

Rule 2: $\mathbb{E}\left[k \cdot \mathrm{g}[x]\right] = k \cdot \mathbb{E}\left[\mathrm{g}[x]\right]$

Rule 3: $\mathbb{E}\left[\mathrm{f}[x] + \mathrm{g}[x]\right] = \mathbb{E}\left[\mathrm{f}[x]\right] + \mathbb{E}\left[\mathrm{g}[x]\right]$

Rule 4: $\mathbb{E}\left[\mathrm{f}[x]\mathrm{g}[y]\right] = \mathbb{E}\left[\mathrm{f}[x]\right]\mathbb{E}\left[\mathrm{g}[y]\right]$     if    $x, y$    independent

$$\mathbb{E}[f_i'] = \mathbb{E}\left[\beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j\right]$$

$$= \mathbb{E}\left[\beta_i\right] + \sum_{j=1}^{D_h} \mathbb{E}\left[\Omega_{ij} h_j\right]$$

Set all the biases to 0

$$= \mathbb{E}\left[\beta_i\right] + \sum_{j=1}^{D_h} \mathbb{E}\left[\Omega_{ij}\right] \mathbb{E}\left[h_j\right]$$

Weights normally distributed
    mean 0
    variance $\sigma_\Omega^2$

$$= 0 + \sum_{j=1}^{D_h} 0 \cdot \mathbb{E}\left[h_j\right] = 0$$

# Initialization

- Need for initialization

- He initialization

- Interlude: Expectations

- Show that $\mathbb{E}[f_i'] = 0$

- Write variance of pre-activations *f'* in terms of activations *h* in previous layer

$$\sigma_{f'}^2 = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

- Write variance of pre-activations f' in terms of pre-activations f in previous layer

$$\sigma_{f'}^2 = \frac{D_h \sigma_\Omega^2 \sigma_f^2}{2}$$

# Aim: keep variance same between two layers

$$\mathbf{f}' = \boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{h}$$

$$\mathbf{h} = \mathbf{a}[\mathbf{f}],$$

$$\sigma_{f'}^2 = \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2$$

$$\longrightarrow \quad \mathbb{E}\left[(x - \mu)^2\right] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

Rule 1: $\mathbb{E}\left[k\right] = k$

Rule 2: $\mathbb{E}\left[k \cdot g[x]\right] = k \cdot \mathbb{E}\left[g[x]\right]$

Rule 3: $\mathbb{E}\left[f[x] + g[x]\right] = \mathbb{E}\left[f[x]\right] + \mathbb{E}\left[g[x]\right]$

Rule 4: $\mathbb{E}\left[f[x]g[y]\right] = \mathbb{E}\left[f[x]\right]\mathbb{E}\left[g[y]\right]$    if  $x, y$  independent

$$\sigma_{f'}^2 = \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2$$

$$= \mathbb{E}\left[\left(\beta_i + \sum_{j=1}^{D_h} \Omega_{ij}h_j\right)^2\right] - 0$$

$$= \mathbb{E}\left[\left(\sum_{j=1}^{D_h} \Omega_{ij}h_j\right)^2\right]$$

$$= \sum_{j=1}^{D_h} \mathbb{E}\left[\Omega_{ij}^2\right]\mathbb{E}\left[h_j^2\right]$$

$$= \sum_{j=1}^{D_h} \sigma_\Omega^2 \mathbb{E}\left[h_j^2\right] = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

Set all the biases to 0

Weights normally distributed
    mean 0
    variance $\sigma_\Omega^2$

# Initialization

- Need for initialization

- He initialization

- Interlude: Expectations

- Show that $\mathbb{E}[f_i'] = 0$

- Write variance of pre-activations $f'$ in terms of activations $h$ in previous layer

$$\sigma_{f'}^2 = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

- Write variance of pre-activations f' in terms of pre-activations f in previous layer

$$\sigma_{f'}^2 = \frac{D_h \sigma_\Omega^2 \sigma_f^2}{2}$$

$$\sigma_{f'}^2 = \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[h_j^2\right]$$

$$= \sigma_\Omega^2 \sum_{j=1}^{D_h} \mathbb{E}\left[\text{ReLU}[f_j]^2\right]$$

$$= \sigma_\Omega^2 \sum_{j=1}^{D_h} \int_{-\infty}^{\infty} \text{ReLU}[f_j]^2 Pr(f_j) df_j$$

$$= \sigma_\Omega^2 \sum_{j=1}^{D_h} \int_{-\infty}^{\infty} (\mathbb{I}[f_j > 0]f_j)^2 Pr(f_j) df_j$$

$$= \sigma_\Omega^2 \sum_{j=1}^{D_h} \int_{0}^{\infty} f_j^2 Pr(f_j) df_j$$

$$= \sigma_\Omega^2 \sum_{j=1}^{D_h} \frac{\sigma_f^2}{2} = \frac{D_h \sigma_\Omega^2 \sigma_f^2}{2}$$

# Aim: keep variance same between two layers

$$\sigma_{f'}^2 = \frac{D_h \sigma_{\Omega}^2 \sigma_f^2}{2}$$

Should choose:

$$\sigma_{\Omega}^2 = \frac{2}{D_h}$$

This is called He initialization.