

# CM20315 - Machine Learning

Prof. Simon Prince  
7b Initialization



# 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

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

- 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:

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

- How do we initialize the biases and weights?
- Equivalent to choosing starting point in Gabor/Linear regression models

# Initialization

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

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

- 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:

$$\begin{aligned}\frac{\partial\ell_i}{\partial\boldsymbol{\beta}_k} &= \frac{\partial\ell_i}{\partial\mathbf{f}_k} & k \in \{K, K-1, \dots, 1\} \\ \frac{\partial\ell_i}{\partial\boldsymbol{\Omega}_k} &= \frac{\partial\ell_i}{\partial\mathbf{f}_k} \mathbf{h}_k^T & k \in \{K, K-1, \dots, 1\} \\ \frac{\partial\ell_i}{\partial\mathbf{f}_{k-1}} &= \mathbb{I}[\mathbf{f}_{k-1} > 0] \odot \left( \boldsymbol{\Omega}_k^T \frac{\partial\ell_i}{\partial\mathbf{f}_k} \right), & k \in \{K, K-1, \dots, 1\}\end{aligned}\tag{7.13}$$

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:

$$\begin{aligned}\frac{\partial\ell_i}{\partial\boldsymbol{\beta}_0} &= \frac{\partial\ell_i}{\partial\mathbf{f}_0} \\ \frac{\partial\ell_i}{\partial\boldsymbol{\Omega}_0} &= \frac{\partial\ell_i}{\partial\mathbf{f}_0} \mathbf{x}_i^T\end{aligned}$$

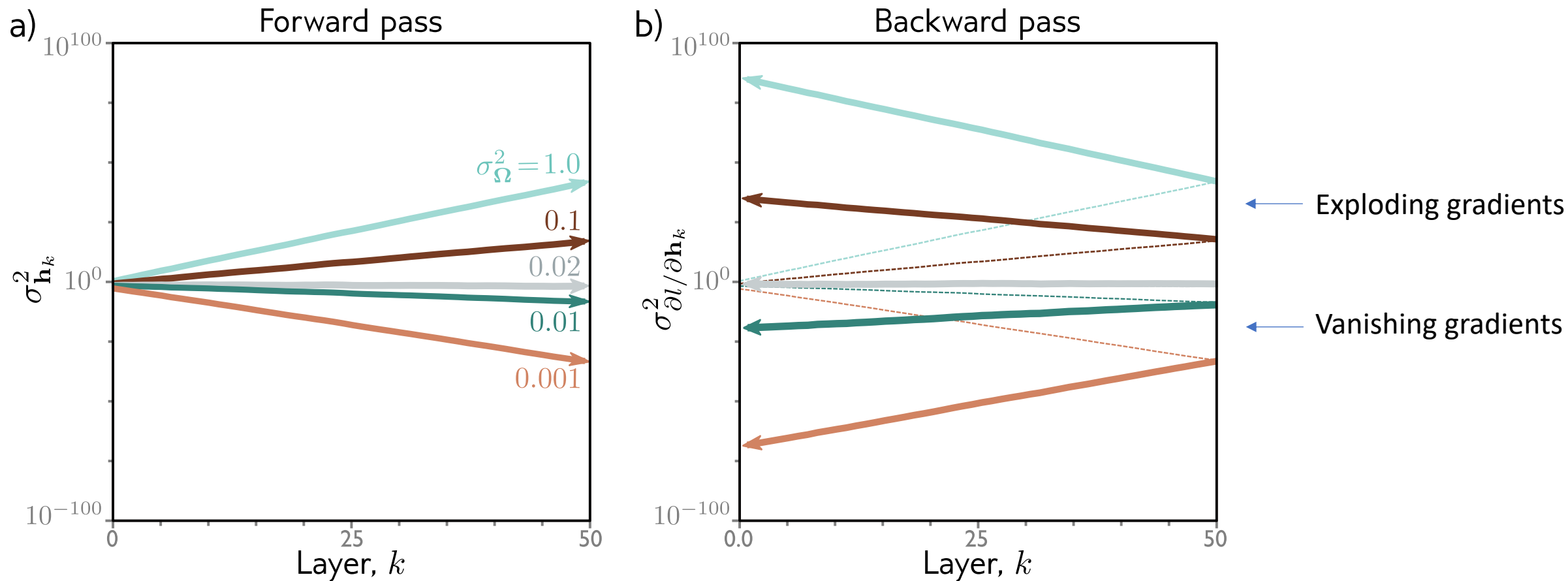
# Initialization

- 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

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

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

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



**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  $\beta_k$  are initialized to zero and the weight matrices  $\Omega_k$  are initialized with a normal distribution with mean zero and five different variances  $\sigma_{\Omega}^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$ . a)

# He initialization (assumes ReLU)

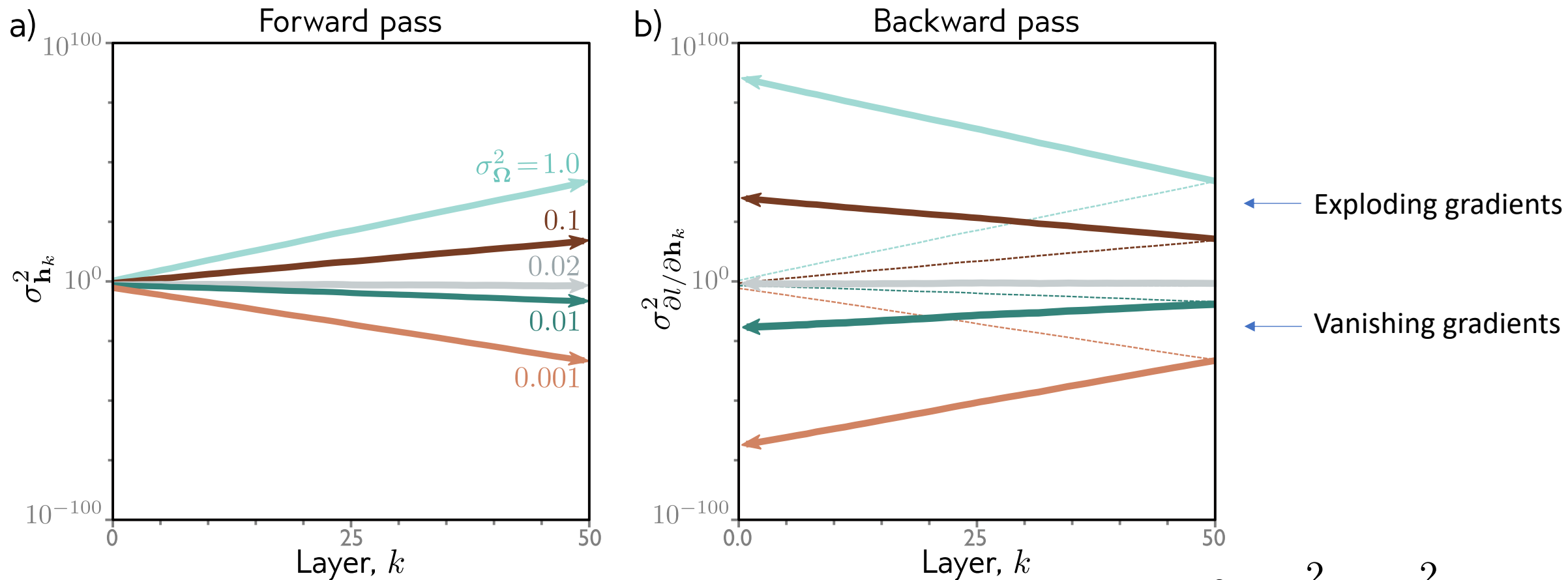
- Forward pass: want the variance of hidden unit activations in layer k+1 to be the same as variance of activations in layer k:

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

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

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





$$\sigma_{\Omega}^2 = \frac{2}{D_h} = \frac{2}{100} = 0.02$$

**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  $\beta_k$  are initialized to zero and the weight matrices  $\Omega_k$  are initialized with a normal distribution with mean zero and five different variances  $\sigma_{\Omega}^2 \in \{0.001, 0.01, 0.02, 0.1, 1.0\}$ . a)

# 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

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

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

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

# Expectations

$$\mathbb{E}[g[x]] = \int g[x] Pr(x) dx,$$

Interpretation: what is the average value of  $g[x]$  when taking into account the probability of  $x$ ?

Could approximate, by sampling many values of  $x$  from the distribution, calculating  $g[x]$ , and taking average:

$$\mathbb{E}[g[x]] \approx \frac{1}{N} \sum_{n=1}^N g[x_n^*] \quad \text{where} \quad x_n^* \sim Pr(x)$$

# Expectations

Function $g[\bullet]$	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

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

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

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

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

# Rule 1

$$\mathbb{E}[g[x]] = \int g[x]Pr(x)dx,$$

---

$$\begin{aligned}\mathbb{E}[\kappa] &= \int \kappa Pr(x)dx \\ &= \kappa \int Pr(x)dx \\ &= \kappa.\end{aligned}$$

# Rules for manipulating expectation

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

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

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

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

## Rule 2

$$\mathbb{E}[g[x]] = \int g[x]Pr(x)dx,$$

---

$$\begin{aligned}\mathbb{E}[\kappa \cdot g[x]] &= \int \kappa \cdot g[x]Pr(x)dx \\ &= \kappa \cdot \int g[x]Pr(x)dx \\ &= \kappa \cdot \mathbb{E}[g[x]]\end{aligned}$$



# Rules for manipulating expectation

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

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

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

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

# Rule 3

$$\mathbb{E}[g[x]] = \int g[x]Pr(x)dx,$$

---

$$\begin{aligned}\mathbb{E}[f[x] + g[x]] &= \int (f[x] + g[x])Pr(x)dx \\ &= \int (f[x]Pr(x) + g[x]Pr(x)) dx \\ &= \int f[x]Pr(x)dx + \int g[x]Pr(x)dx \\ &= \mathbb{E}[f[x]] + \mathbb{E}[g[x]]\end{aligned}$$

# Rules for manipulating expectation

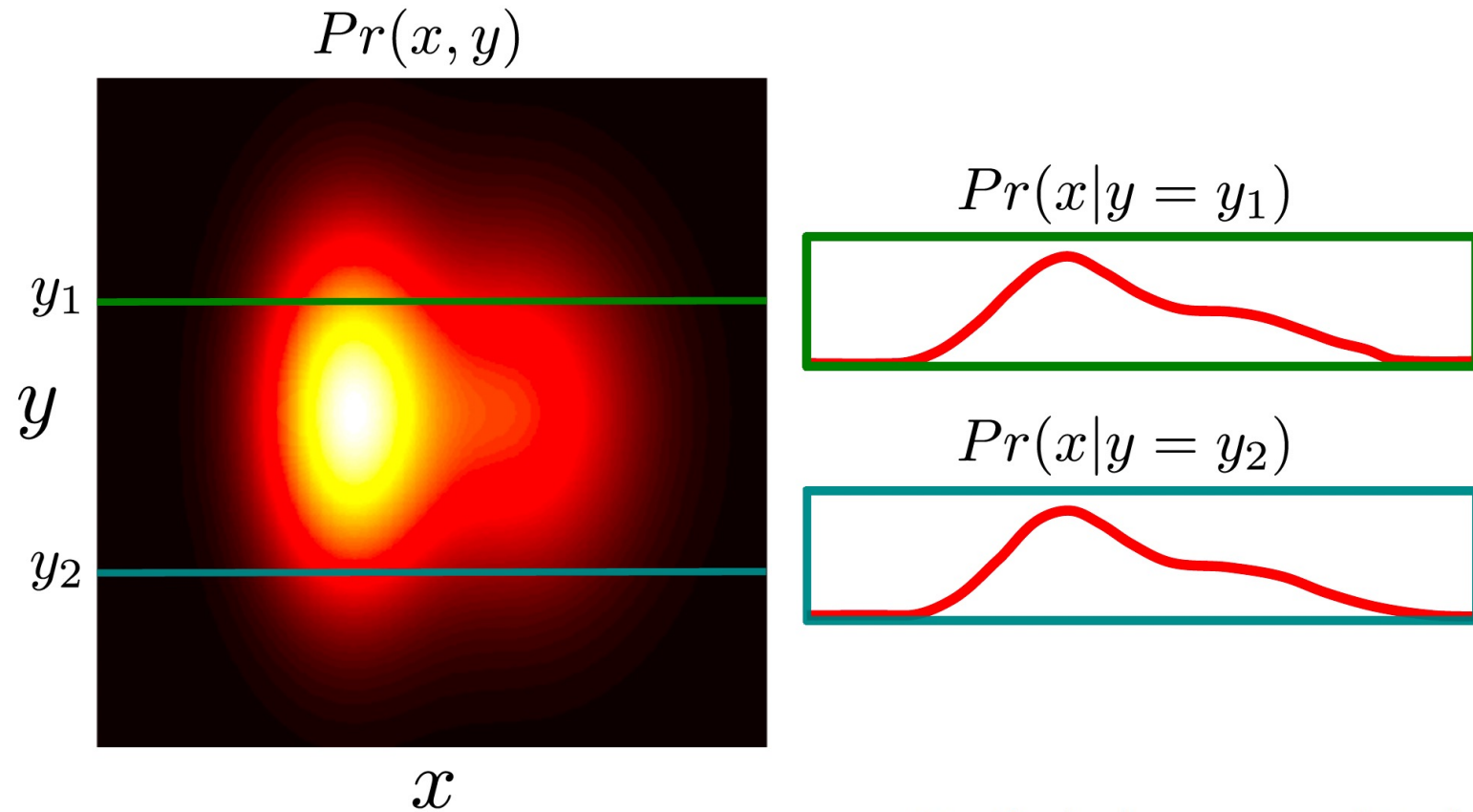
$$\mathbb{E}[k] = k$$

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

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

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

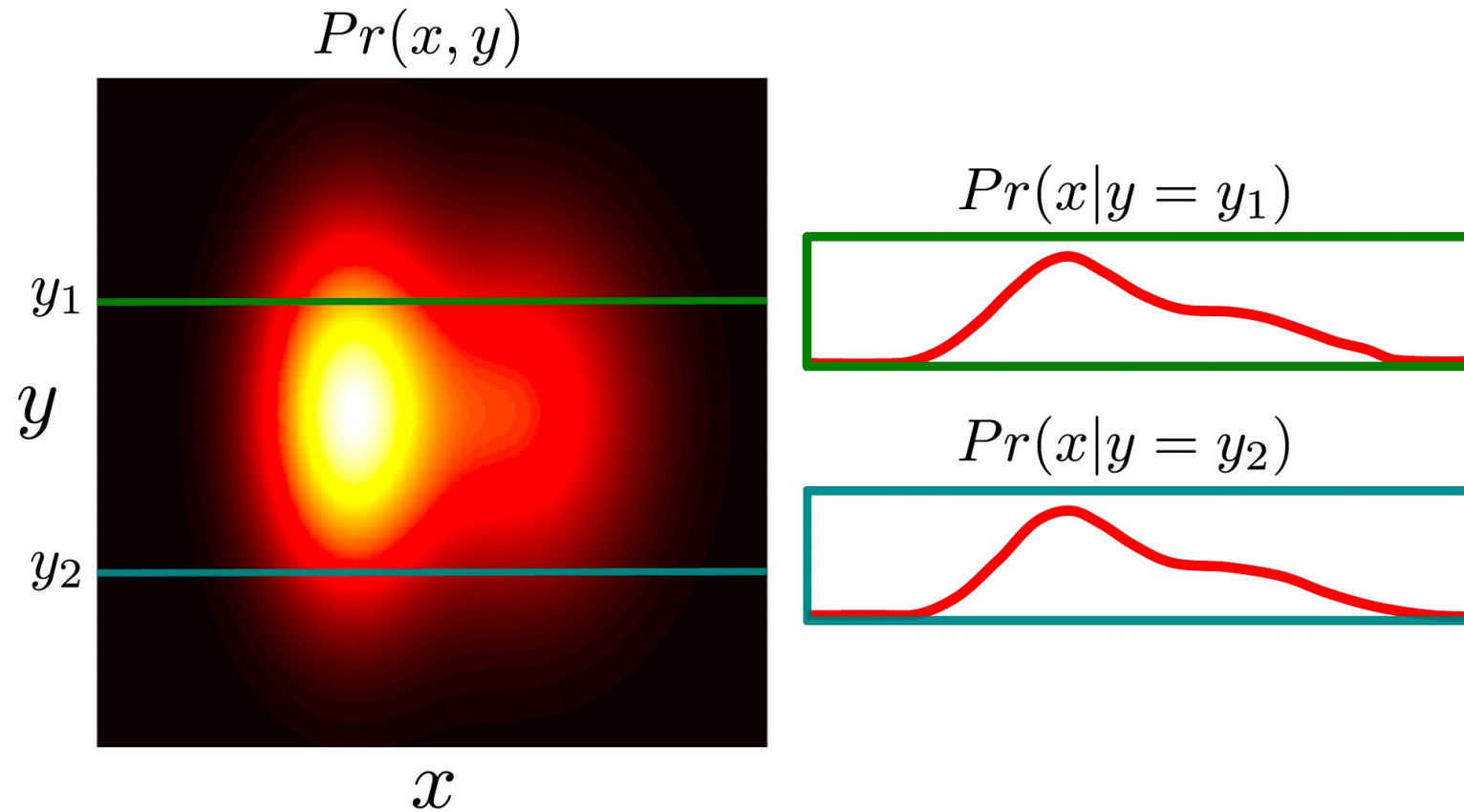
# Independence



Probability of x and y

$$\begin{aligned} Pr(x|y) &= Pr(x) \\ Pr(y|x) &= Pr(y) \end{aligned}$$

# Independence



$$Pr(x, y) = Pr(x)Pr(y)$$

Probability of x and y

# Rule 4

$$\mathbb{E}[g[x]] = \int g[x] Pr(x) dx,$$

---

$$\begin{aligned}\mathbb{E}[f[x] \cdot g[y]] &= \int \int f[x] \cdot g[y] Pr(x, y) dx dy \\ &= \int \int f[x] \cdot g[y] Pr(x) Pr(y) dx dy \quad \leftarrow \text{Because independent} \\ &= \int f[x] Pr(x) dx \int g[y] Pr(y) dy \\ &= \mathbb{E}[f[x]] \mathbb{E}[g[y]] \quad \text{if } x, y \text{ independent}\end{aligned}$$

Now let's prove:

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

Keeping in mind:

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

Now let's prove:

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

Keeping in mind:

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



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

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

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

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

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

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

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

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

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



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

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

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

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

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

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

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

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

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

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



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

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

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

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

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

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

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

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

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

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



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

# 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

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

- 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*:

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

- 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

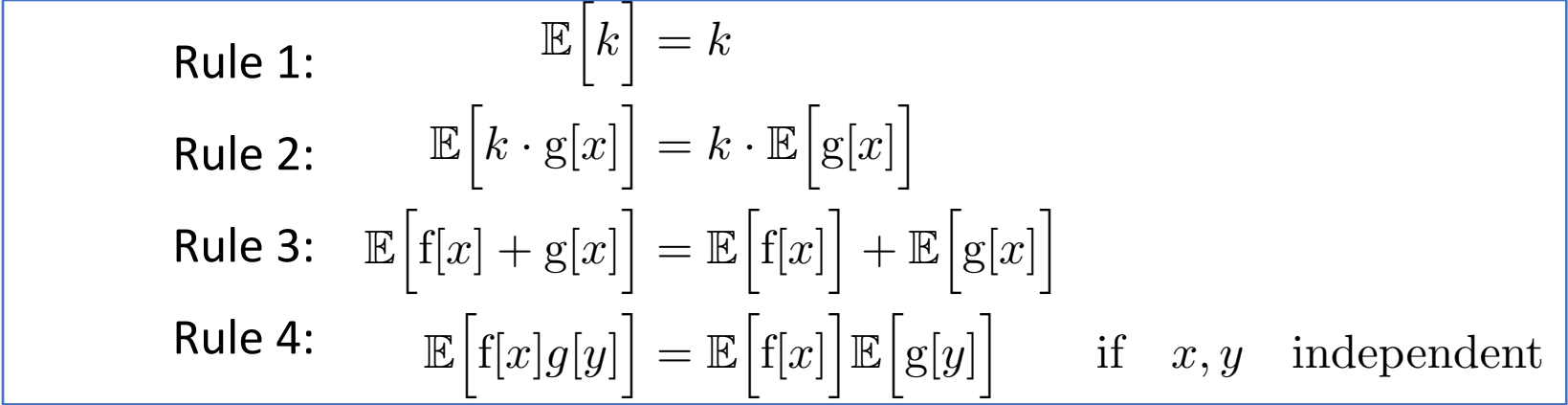
$$\mathbf{f}' = \boldsymbol{\beta} + \boldsymbol{\Omega}\mathbf{h}$$

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:  $\mathbb{E}[k] = k$
- Rule 2:  $\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
- Rule 3:  $\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
- Rule 4:  $\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$  if  $x, y$  independent
- 

$$\begin{aligned}\mathbb{E}[f'_i] &= \mathbb{E} \left[ \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right] \\ &= \mathbb{E} [\beta_i] + \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij} h_j]\end{aligned}$$

- Rule 1:  $\mathbb{E}[k] = k$
- Rule 2:  $\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
- Rule 3:  $\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
- Rule 4:  $\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$  if  $x, y$  independent
- 

$$\begin{aligned}\mathbb{E}[f'_i] &= \mathbb{E} \left[ \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right] \\ &= \mathbb{E} [\beta_i] + \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij} h_j] \\ &= \mathbb{E} [\beta_i] + \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij}] \mathbb{E} [h_j]\end{aligned}$$

Rule 1:	$\mathbb{E}[k] = k$
Rule 2:	$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
Rule 3:	$\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
Rule 4:	$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$ if $x, y$ independent

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

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

Set all the biases to 0

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

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

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

# 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

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

- 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

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

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

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

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

Rule 1:	$\mathbb{E}[k] = k$
Rule 2:	$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
Rule 3:	$\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
Rule 4:	$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$ if $x, y$ independent

$$\begin{aligned}\sigma_{f'}^2 &= \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2 \\ &= \mathbb{E} \left[ \left( \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] - 0\end{aligned}$$

Set all the biases to 0

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

Rule 1:	$\mathbb{E}[k] = k$
Rule 2:	$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
Rule 3:	$\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
Rule 4:	$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$ if $x, y$ independent


$$\begin{aligned}
 \sigma_{f'}^2 &= \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2 \\
 &= \mathbb{E} \left[ \left( \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] - 0 \\
 &= \mathbb{E} \left[ \left( \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right]
 \end{aligned}$$

Set all the biases to 0 

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



Rule 1:	$\mathbb{E}[k] = k$
Rule 2:	$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
Rule 3:	$\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
Rule 4:	$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$ if $x, y$ independent



$$\begin{aligned}
 \sigma_{f'}^2 &= \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2 \\
 &= \mathbb{E} \left[ \left( \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] - 0 \\
 &= \mathbb{E} \left[ \left( \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] \\
 &= \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij}^2] \mathbb{E} [h_j^2]
 \end{aligned}$$

Set all the biases to 0

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

Rule 1:	$\mathbb{E}[k] = k$
Rule 2:	$\mathbb{E}[k \cdot g[x]] = k \cdot \mathbb{E}[g[x]]$
Rule 3:	$\mathbb{E}[f[x] + g[x]] = \mathbb{E}[f[x]] + \mathbb{E}[g[x]]$
Rule 4:	$\mathbb{E}[f[x]g[y]] = \mathbb{E}[f[x]]\mathbb{E}[g[y]]$ if $x, y$ independent

$$\begin{aligned}\sigma_{f'}^2 &= \mathbb{E}[f_i'^2] - \mathbb{E}[f_i']^2 \\ &= \mathbb{E} \left[ \left( \beta_i + \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right] - 0 \\ &= \mathbb{E} \left[ \left( \sum_{j=1}^{D_h} \Omega_{ij} h_j \right)^2 \right]\end{aligned}$$

Set all the biases to 0

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

$$\begin{aligned}&= \sum_{j=1}^{D_h} \mathbb{E} [\Omega_{ij}^2] \mathbb{E} [h_j^2] \\ &= \sum_{j=1}^{D_h} \sigma_{\Omega}^2 \mathbb{E} [h_j^2] = \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \mathbb{E} [h_j^2]\end{aligned}$$

# 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

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

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

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

$$\begin{aligned}
\sigma_{f'}^2 &= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \mathbb{E} [h_j^2] \\
&= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \mathbb{E} [\text{ReLU}[f_j]^2] \\
&= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \int_{-\infty}^{\infty} \text{ReLU}[f_j]^2 Pr(f_j) df_j \\
&= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \int_{-\infty}^{\infty} (\mathbb{I}[f_j > 0] f_j)^2 Pr(f_j) df_j \\
&= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \int_0^{\infty} f_j^2 Pr(f_j) df_j \\
&= \sigma_{\Omega}^2 \sum_{j=1}^{D_h} \frac{\sigma_f^2}{2} = \frac{D_h \sigma_{\Omega}^2 \sigma_f^2}{2}
\end{aligned}$$

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**.



Feedback