# Weight initialization: why it matters

$$z = Wx + b, \quad X \sim P(0, \sigma_x^2) \quad, \quad X \text{ is } n \text{ samples} \times k \text{ features}$$

$$E[z] = E[Wx + b] \quad, \quad \text{assume } X \text{ are } \sim N(0,1)$$
$$b = 0 \text{ for simplicity}$$

for a single neuron $i$:

$$E[z_i] = E[w_i x + \cancel{b_i}] = E\left[\sum_{j=1}^{k} w_{ij} x_j\right]$$

$$= E\left[\sum^{k} w_{ij}\right] E[x_j] \quad \longleftarrow \quad \textcolor{red}{\text{assuming } W \text{ and } X \text{ are independent}}$$

$$= \qquad 0 \qquad\qquad = 0$$

$$\sigma_{z_i} = E[z_i^2] - \cancel{(E[z_i])^2}$$

$$= E\left[\left(\sum^{k} w_{ij} x_j\right)^2\right] = \sum_{1}^{k} E[w_{ij}^2] \underbrace{E[x_j^2]}_{\textcolor{red}{\sigma_x = 1}}$$

$$= k \, \sigma_w^2 \, \sigma_x^2 \; /\!/$$

$$\textcolor{red}{\uparrow \text{ at each layer, variance increases by}}$$
$$\textcolor{red}{\# \text{ of inputs}}$$