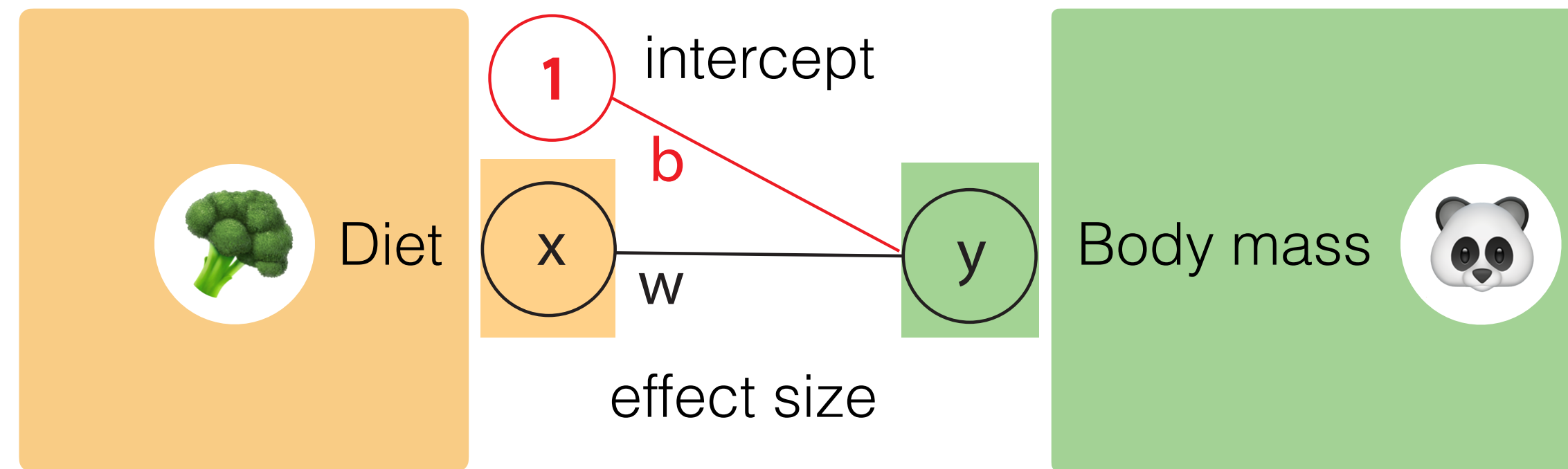


Neural networks: (very) short version...

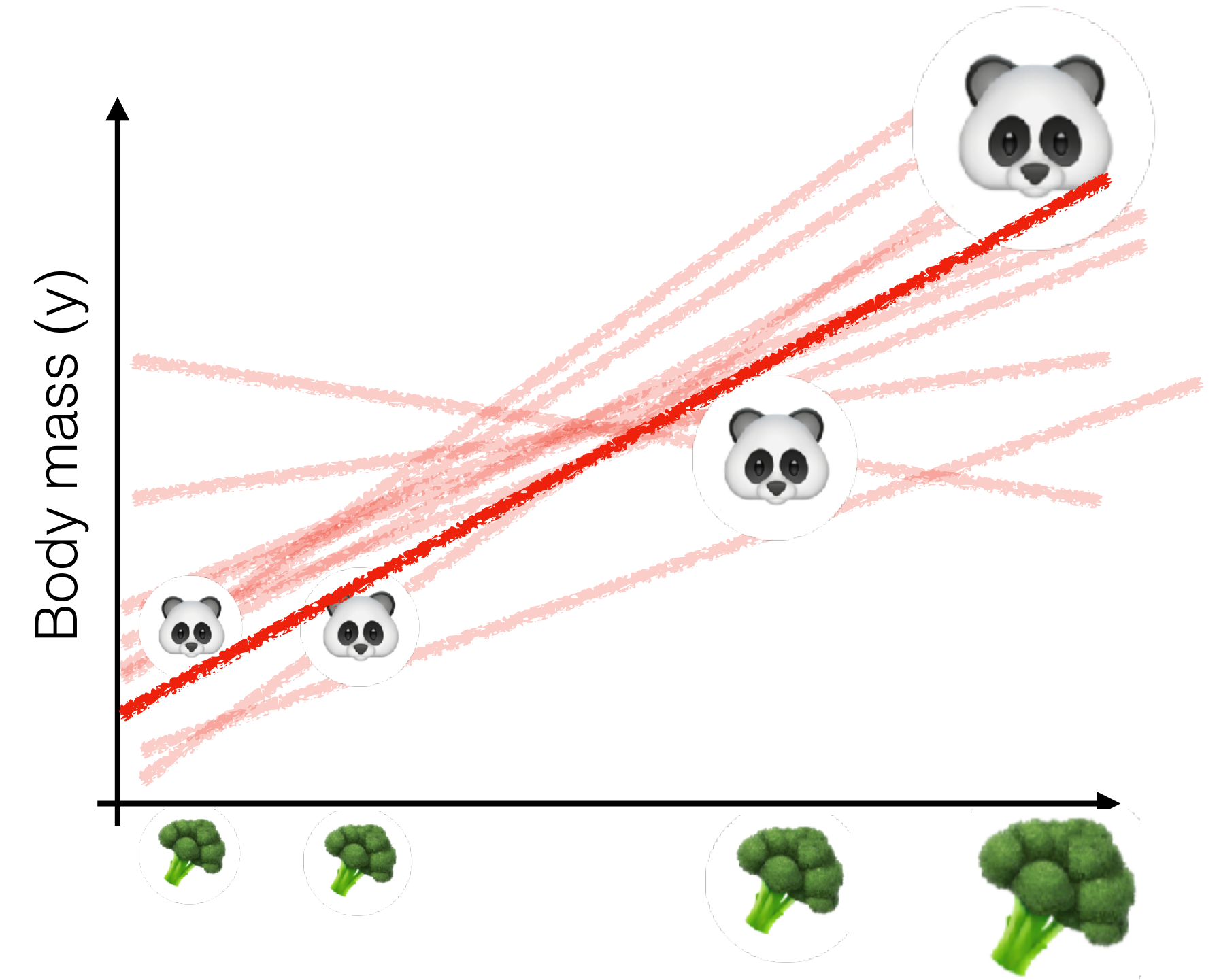
Linear regression



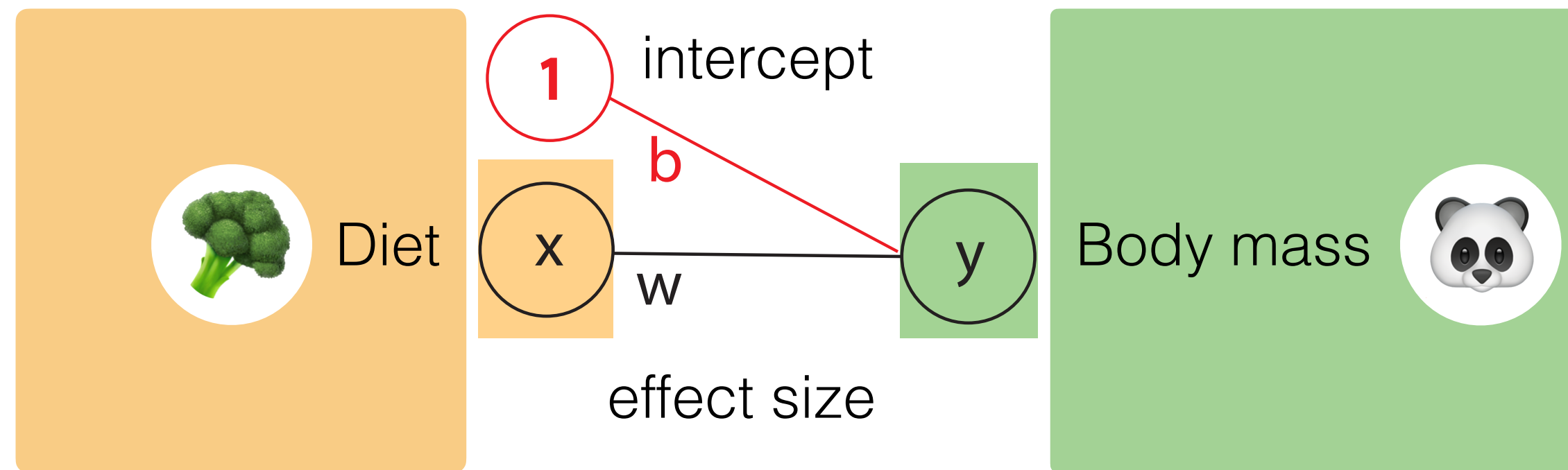
$$y = wx + b$$

prediction effect size intercept

Optimizing ('training') a model



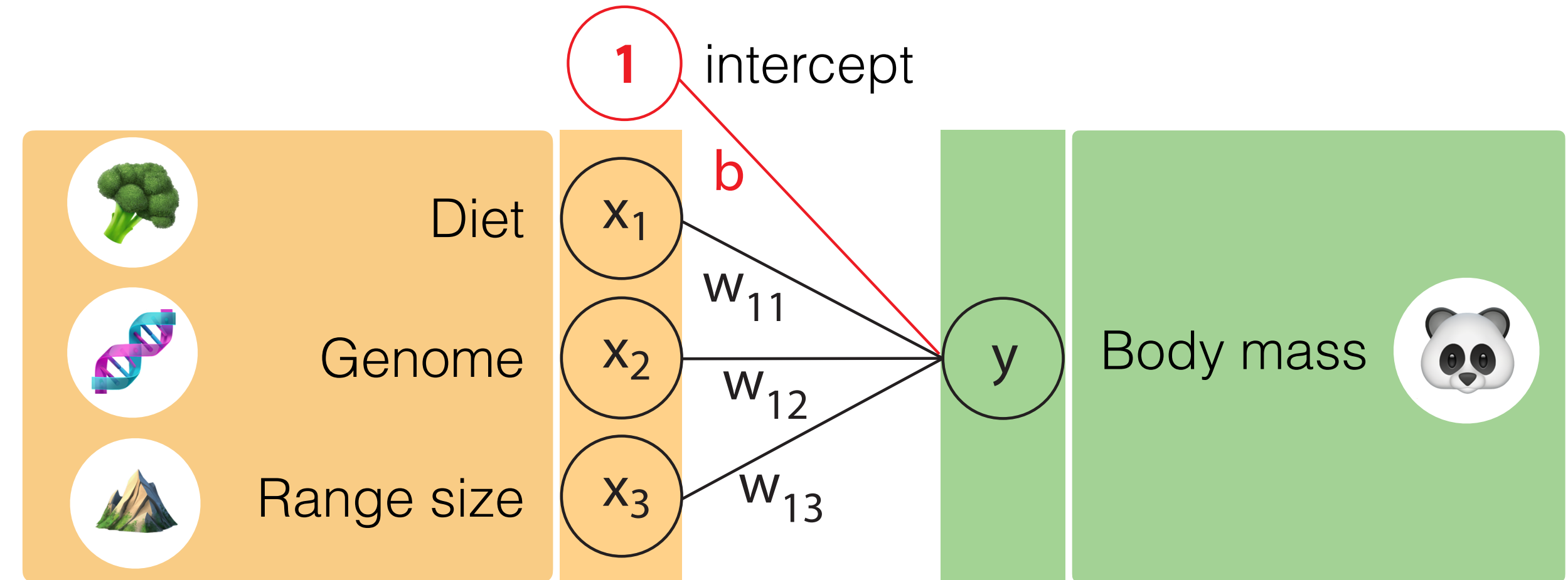
Linear regression



$$y = wx + b$$

prediction effect size intercept

Multiple linear regression



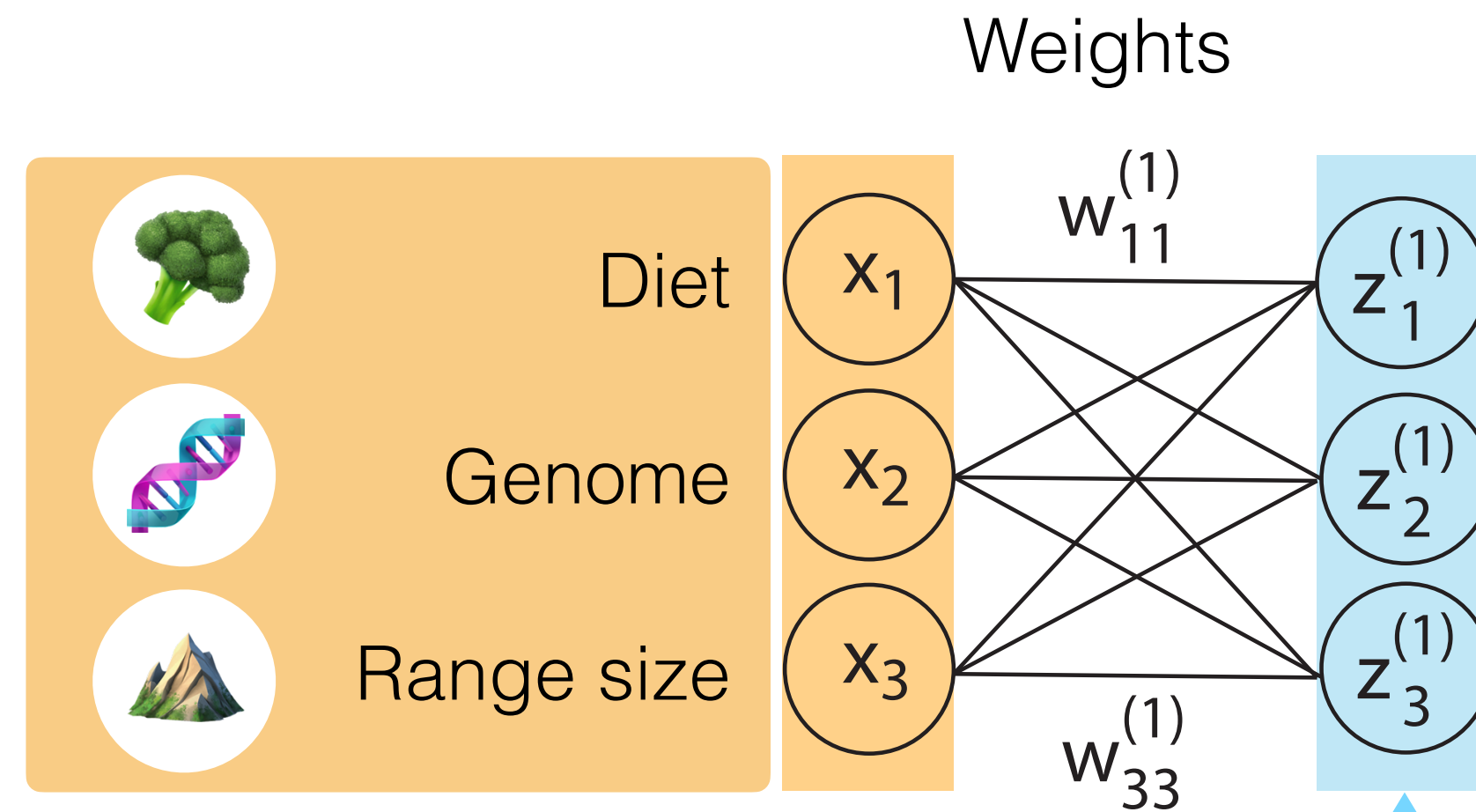
$$y = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b$$

prediction effect sizes intercept

$$y = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} (w_{11} \quad w_{12} \quad w_{13}) + b$$

We can re-write it as a matrix multiplication

Neural network



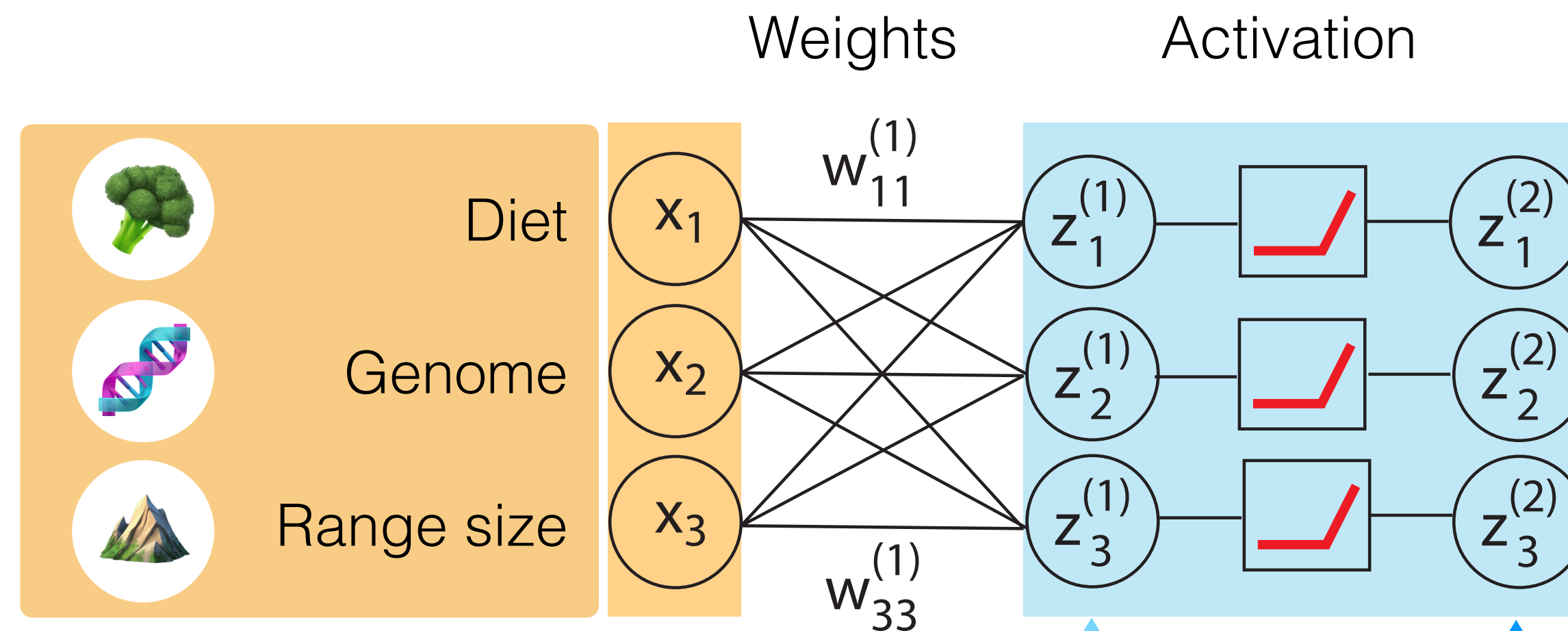
$$z_i^{(1)} = w_{i1}^{(1)} x_1 + w_{i2}^{(1)} x_2 + w_{i3}^{(1)} x_3$$

Intermediate (hidden) values $z^{(1)}$: linear function of all predictors

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{pmatrix} = \begin{pmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \end{pmatrix}$$

x
w⁽¹⁾
z⁽¹⁾

Neural network



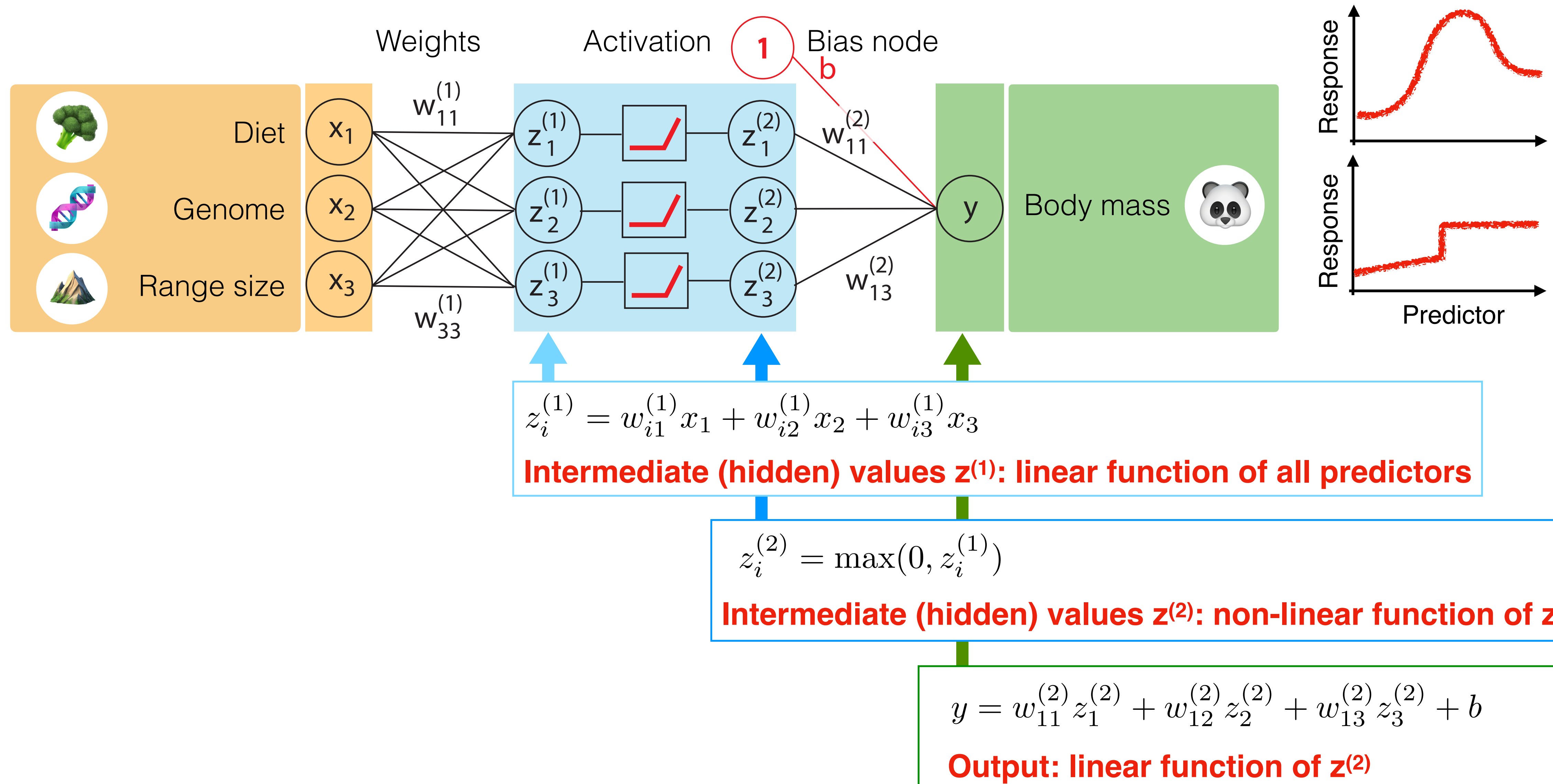
$$z_i^{(1)} = w_{i1}^{(1)} x_1 + w_{i2}^{(1)} x_2 + w_{i3}^{(1)} x_3$$

Intermediate (hidden) values $z^{(1)}$: linear function of all predictors

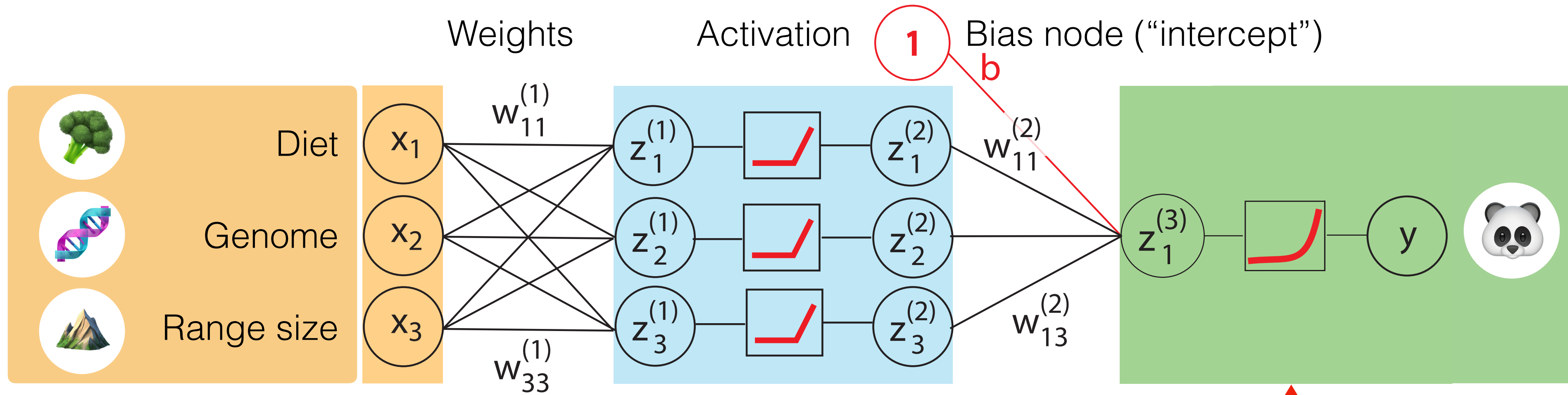
$$z_i^{(2)} = \max(0, z_i^{(1)})$$

Intermediate (hidden) values $z^{(2)}$: non-linear function of $z^{(1)}$

Neural network

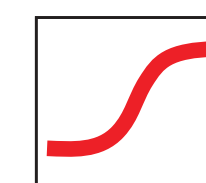


Neural network



$$y = e^z \leftarrow \text{Must be positive}$$

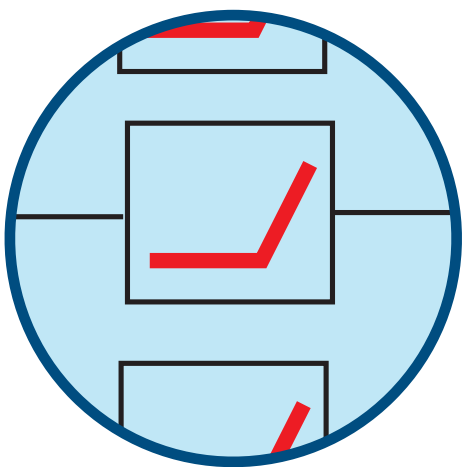
**Add an output activation function
to match expected range**



$$y = 1/(1 + e^{-z}) \leftarrow \text{Must be in } [0, 1]$$

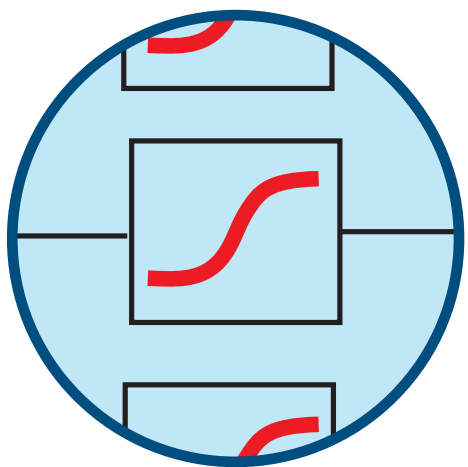
Parameterization of a neural network

Activation functions



ReLU: rectified linear unit

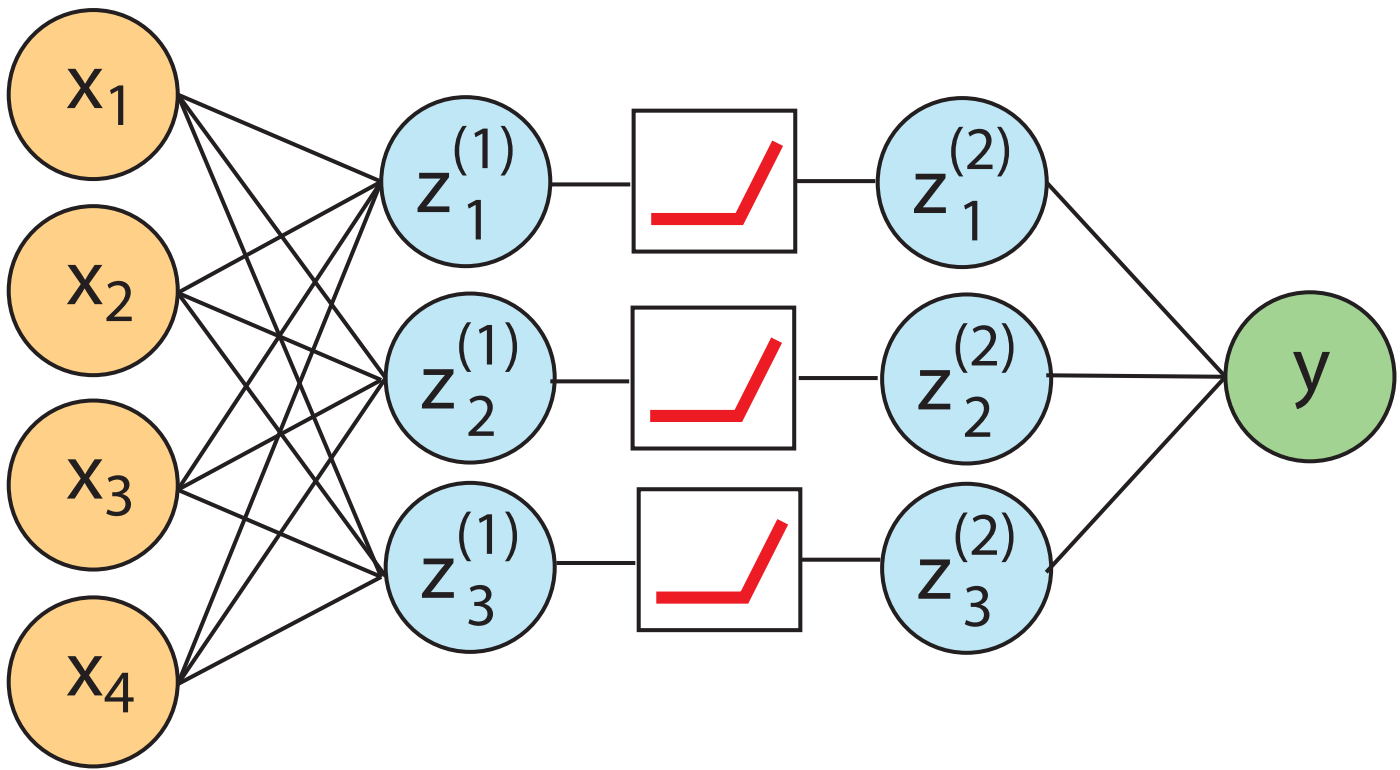
$$\text{ReLU}(x) = \max(0, x)$$



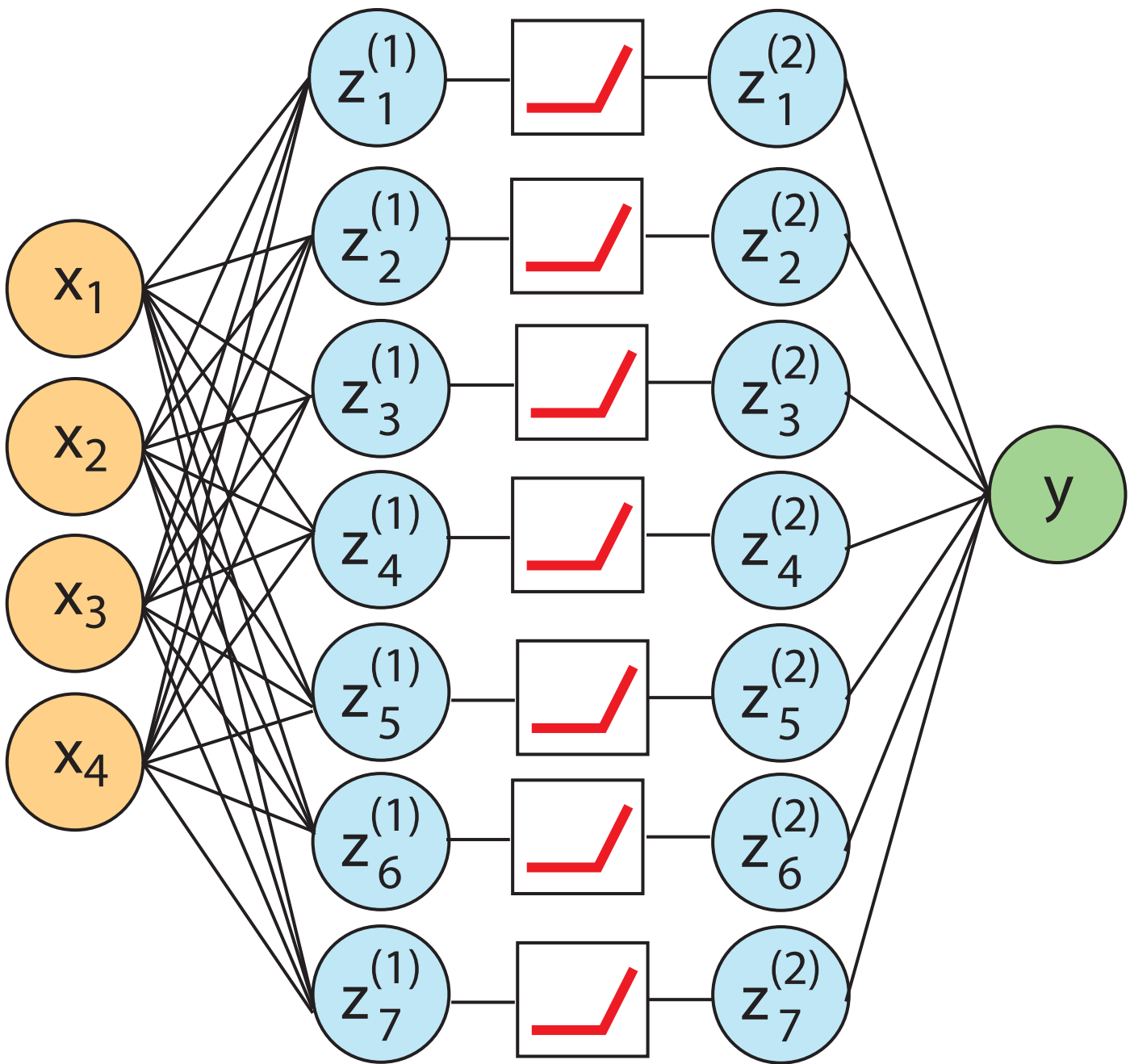
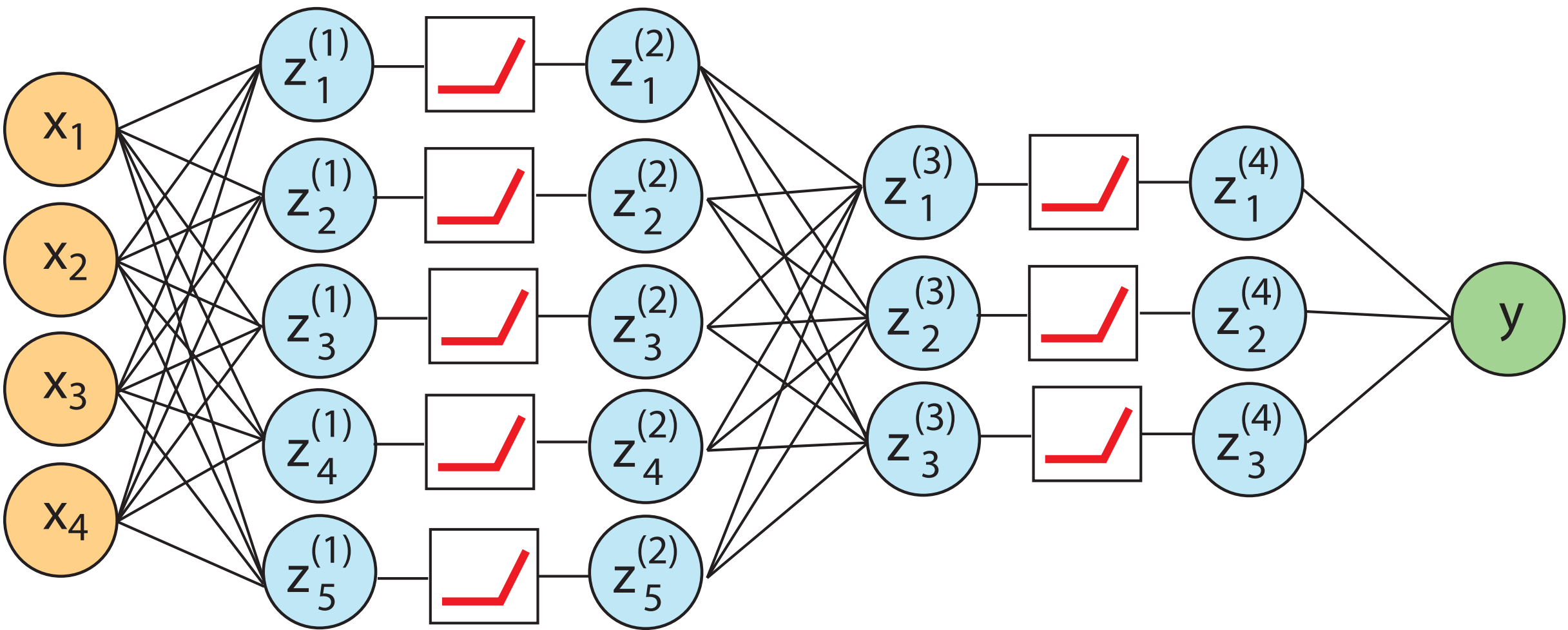
Sigmoid

$$\text{sigmoid}(x) = 1 / (1 + \exp(-x))$$

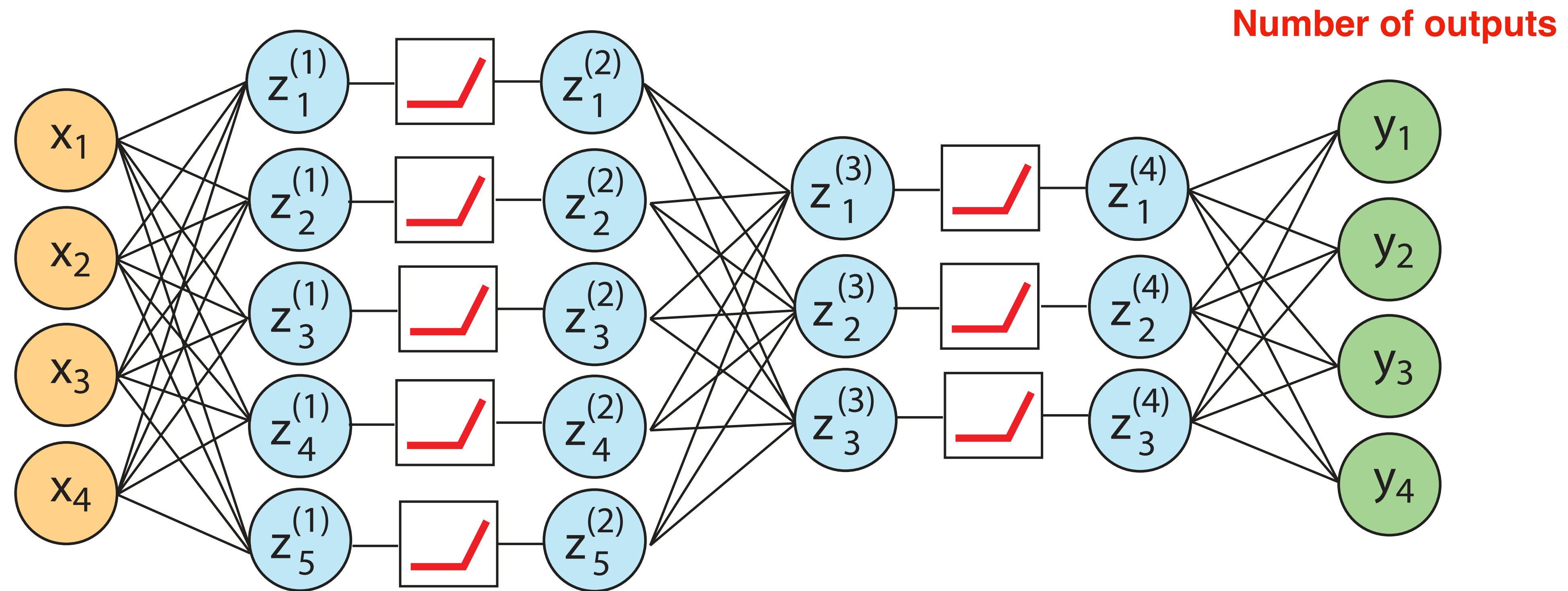
Number of nodes



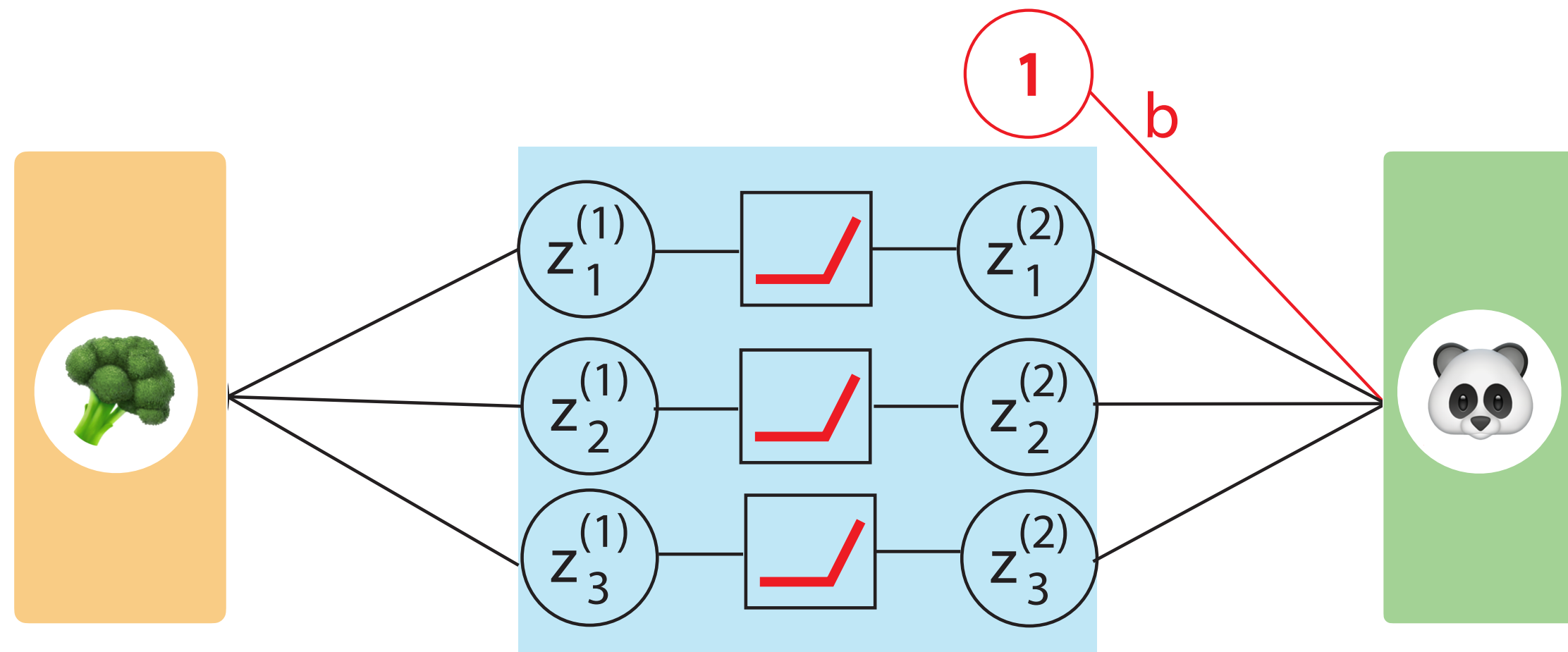
Number of hidden layers (deep NNs)



Parameterization of a neural network



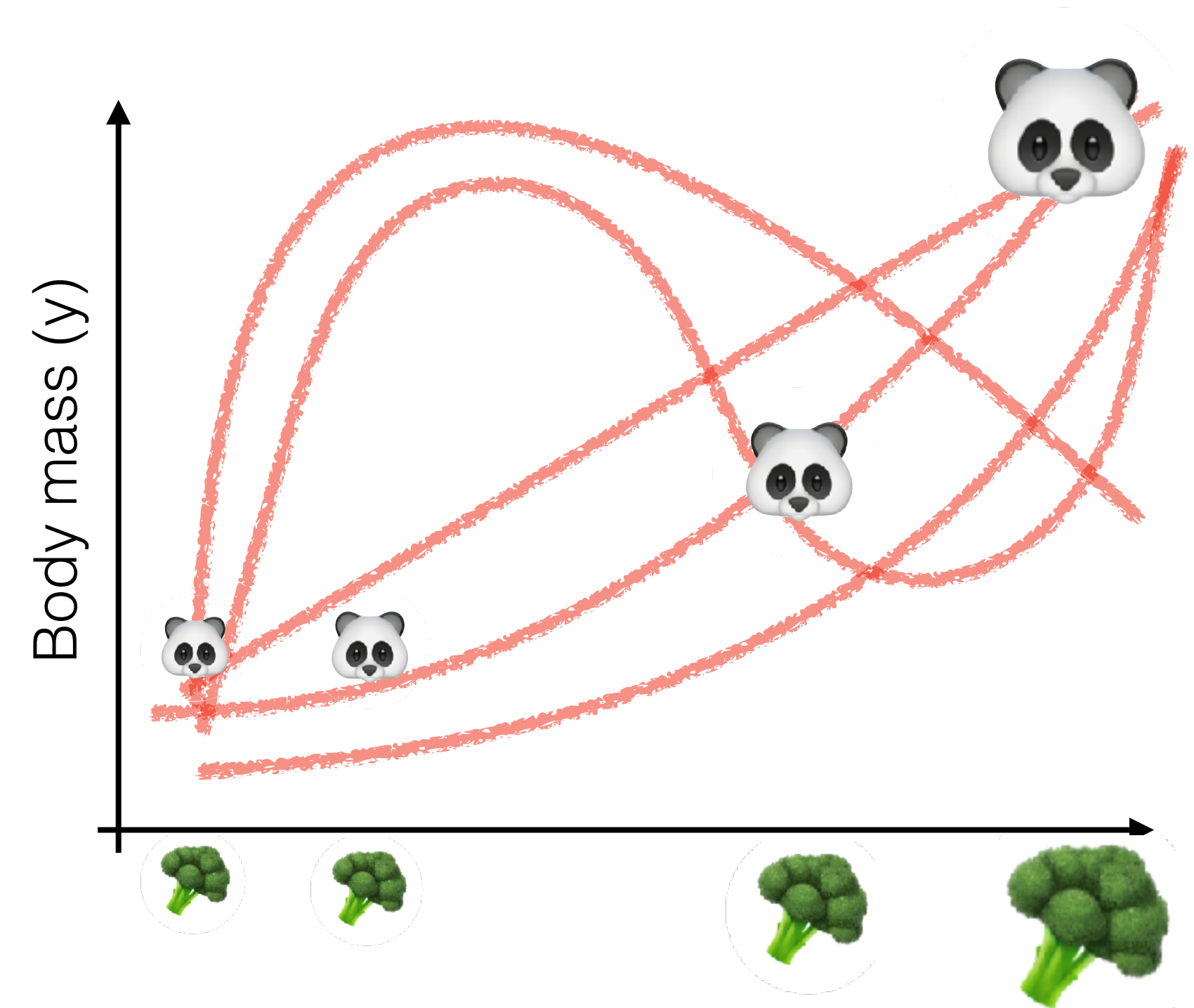
NN regression



Data ('training set')



Optimizing ('training') a model



Likelihood function: the normal density function

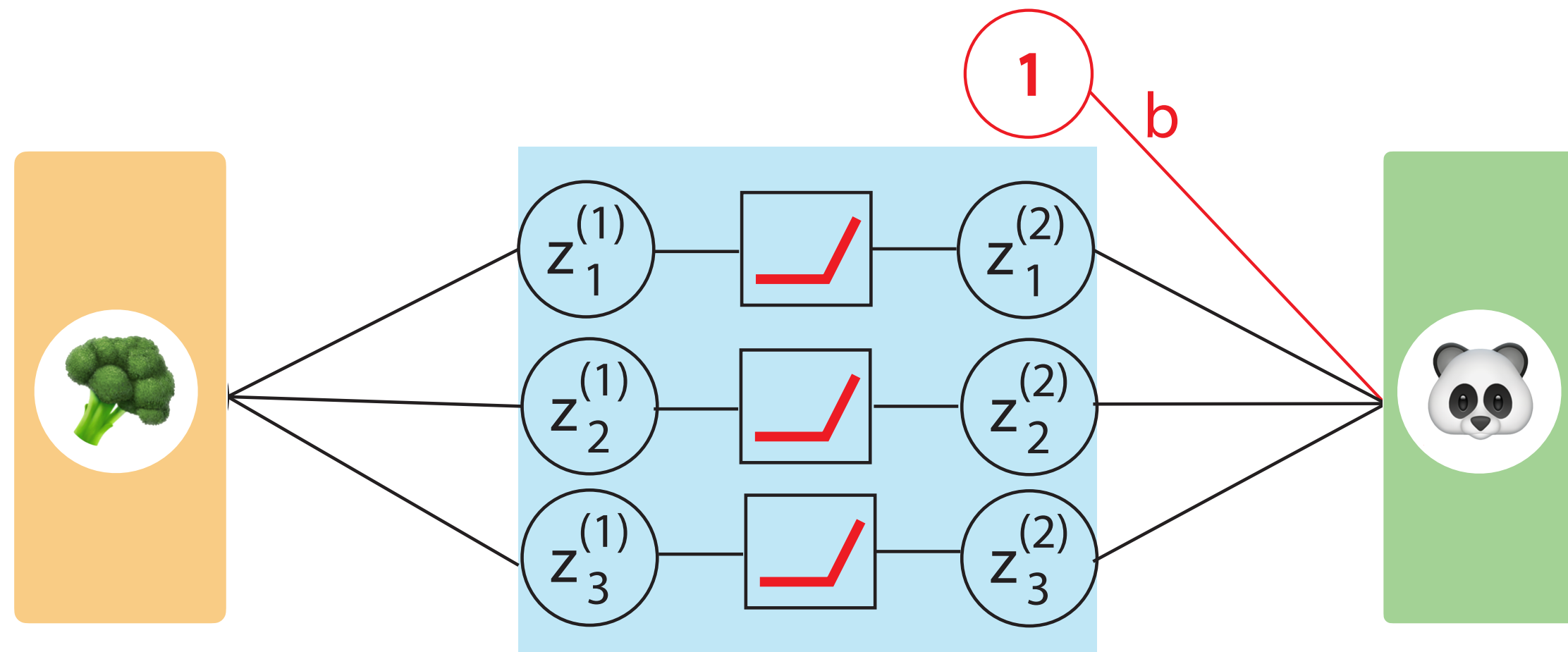
$$P(\text{panda} | \text{broccoli}, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\text{panda} - \text{broccoli})^2}{2\sigma^2}}$$

Loss function: mean squared error

$$-\log P(\text{panda} | \text{broccoli}) \propto (\text{panda} - \text{broccoli})^2$$

Squared difference between truth and prediction

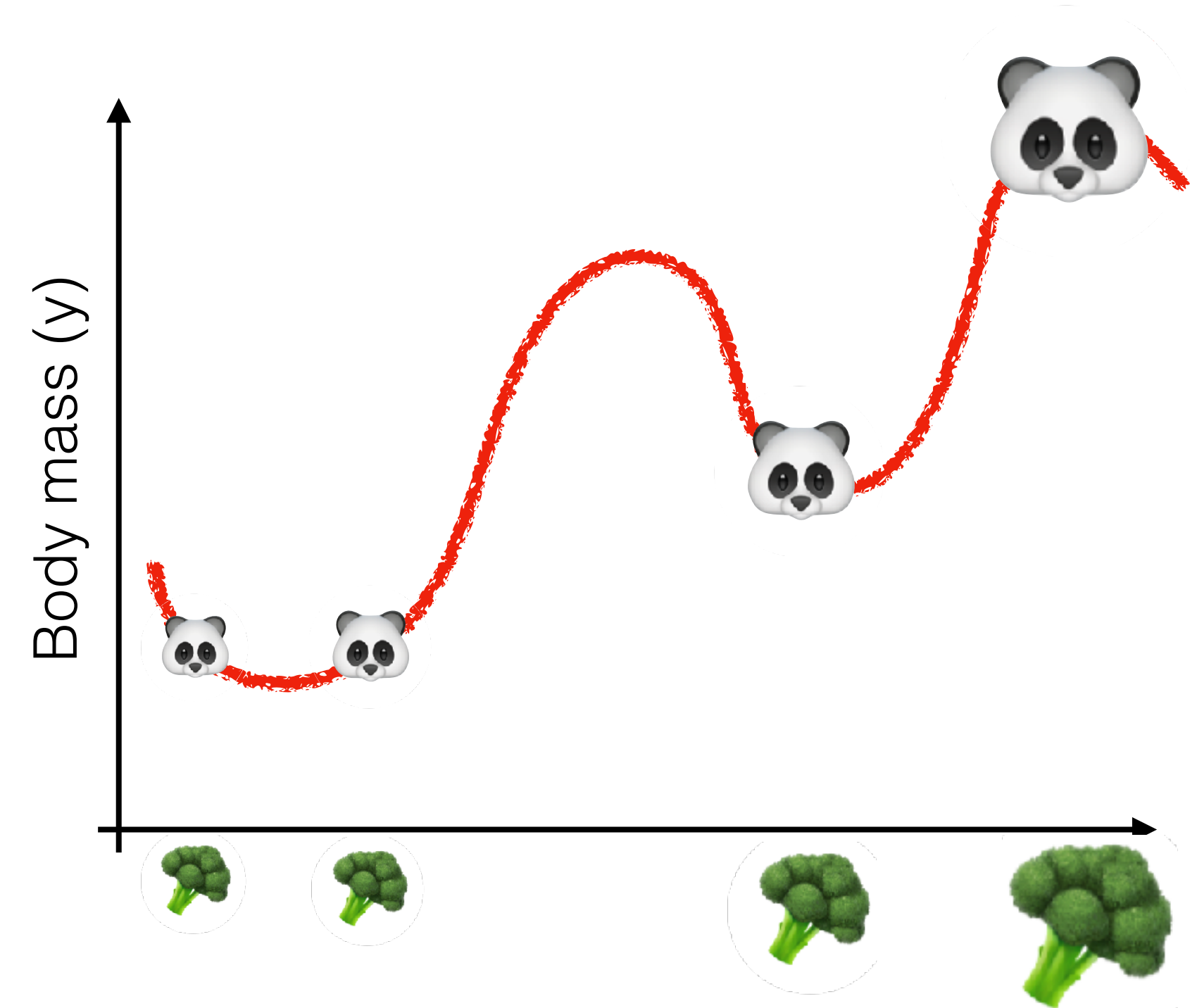
NN regression



Data ('training set')

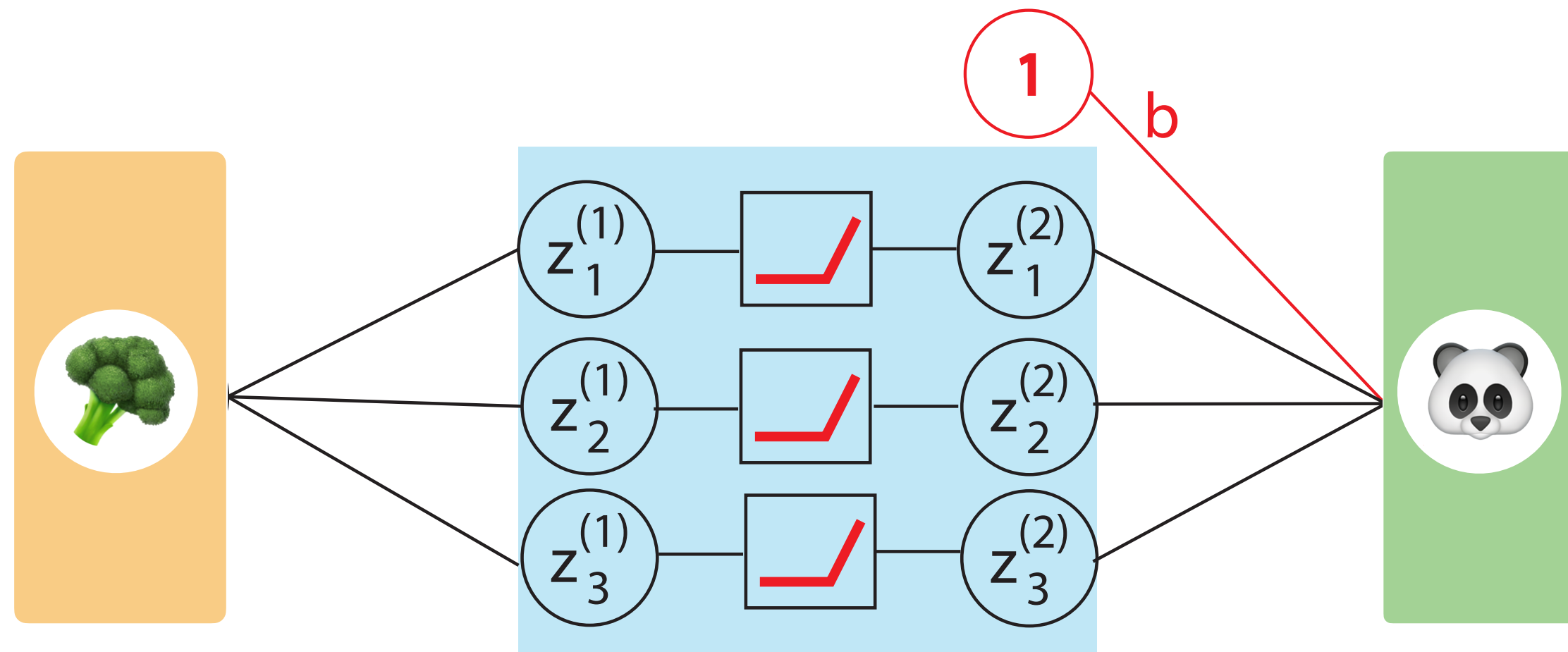


Best-fitting model



Neural networks are over-parameterized models and will overfit with a maximum likelihood approach

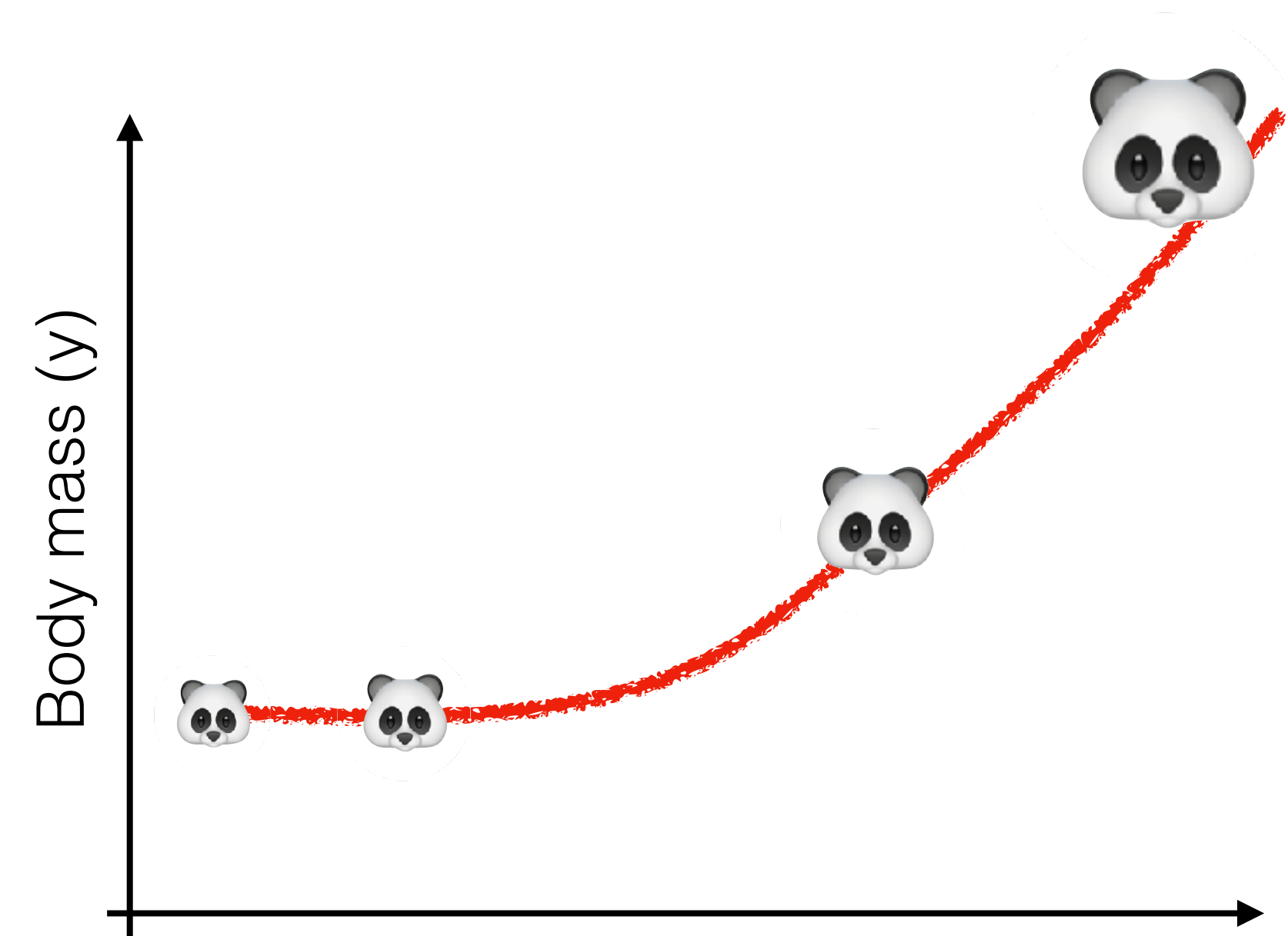
NN regression



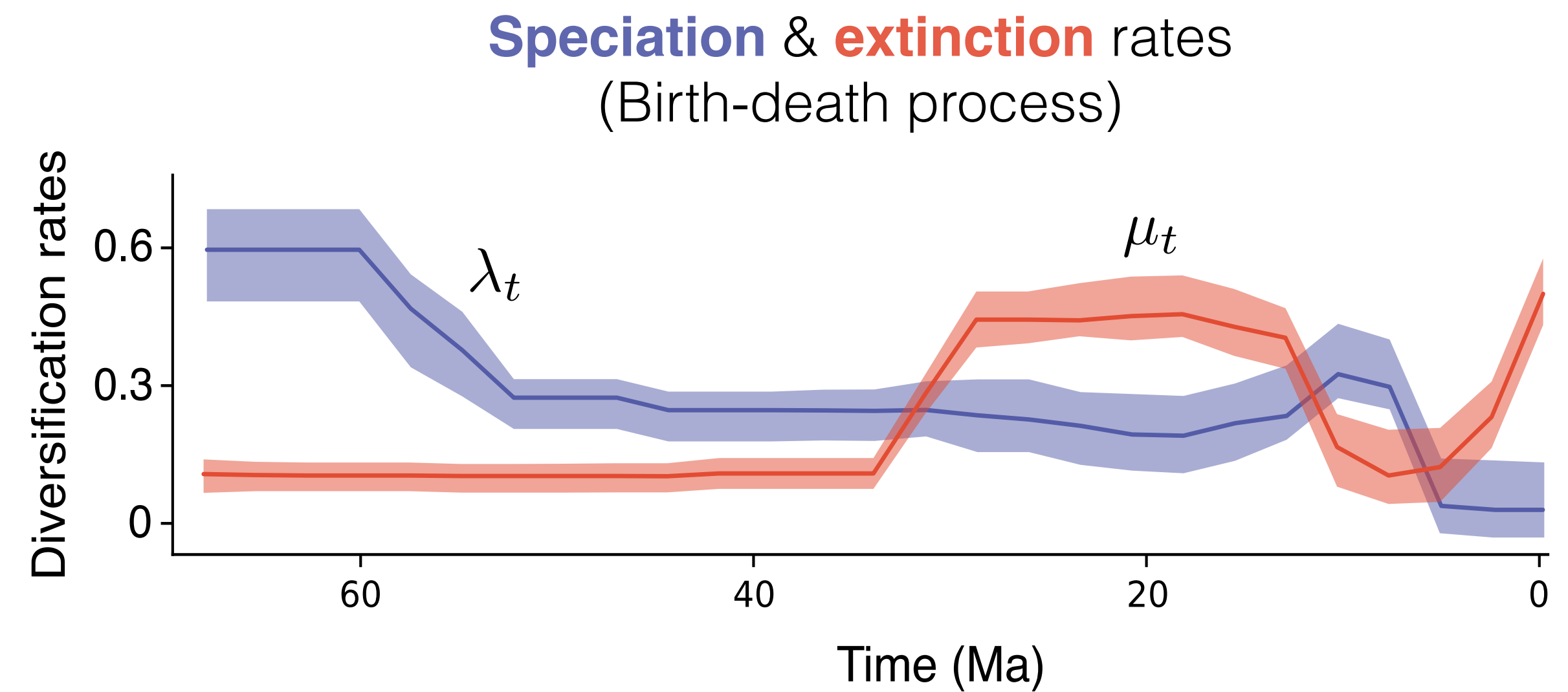
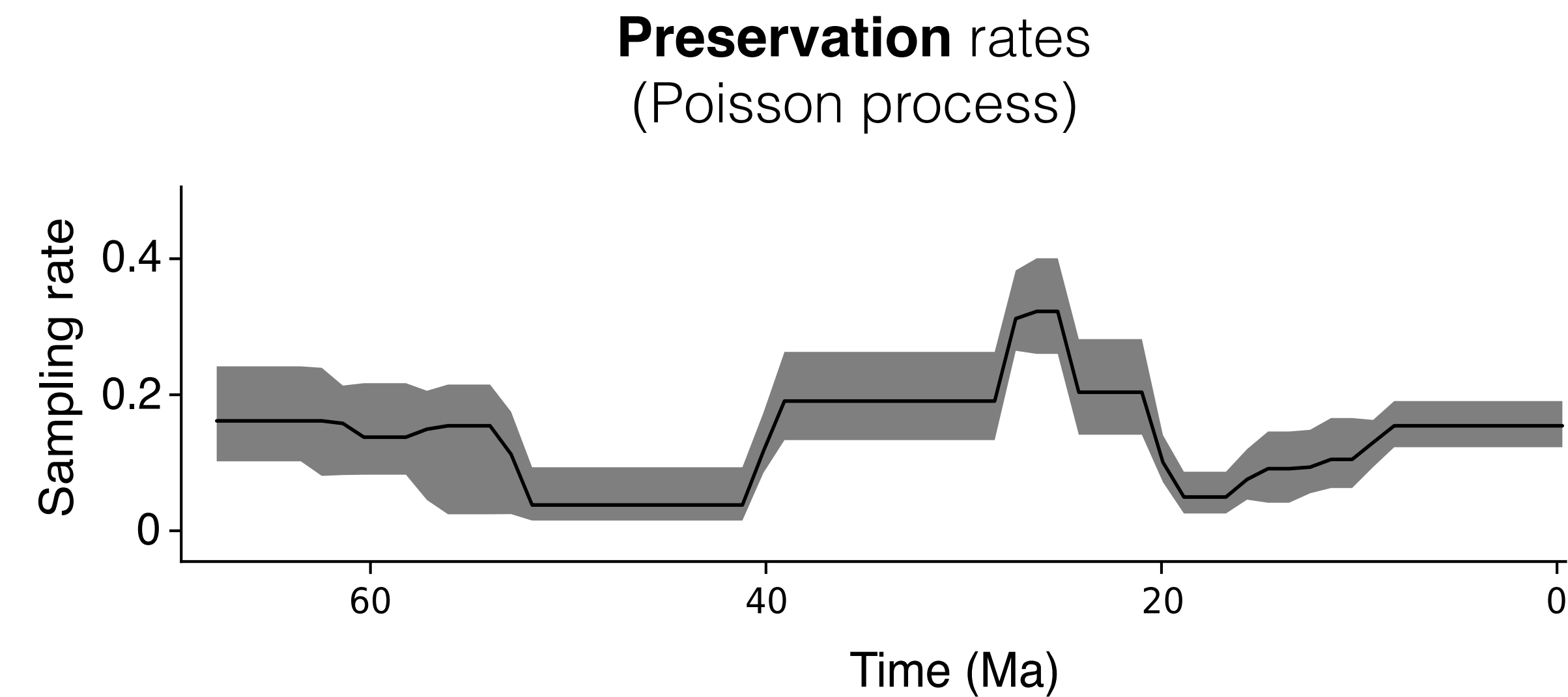
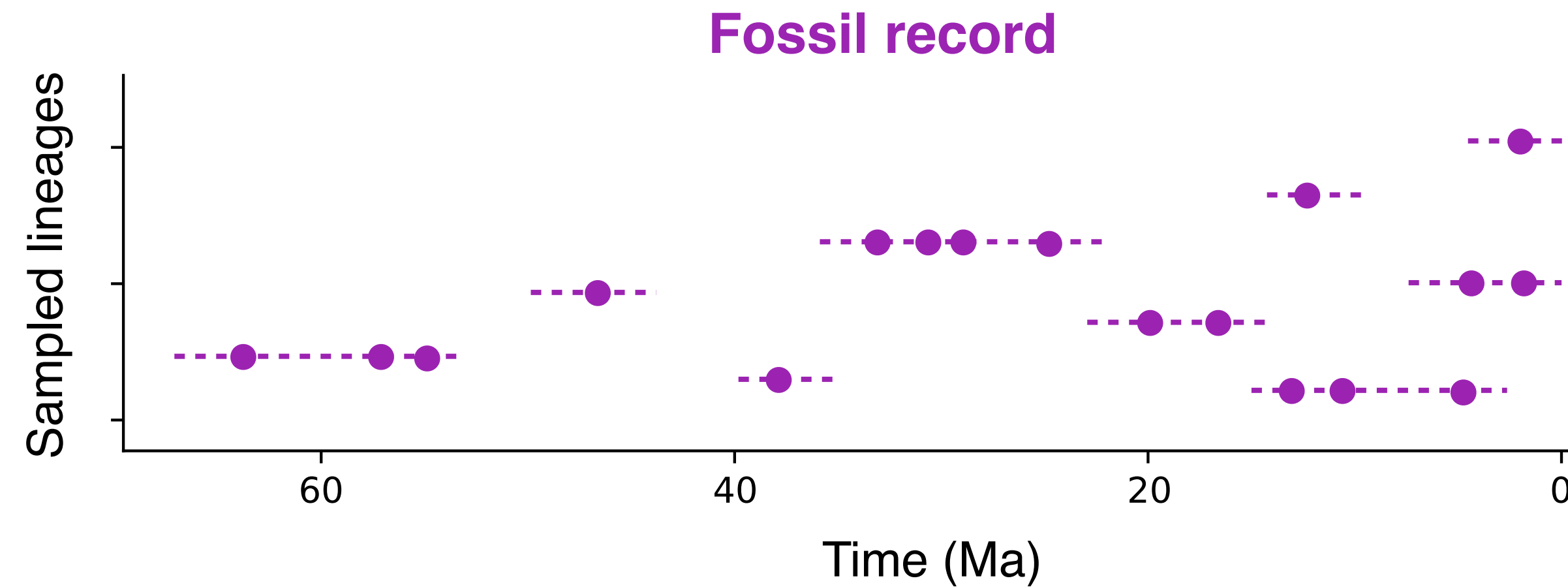
Data ('training set')



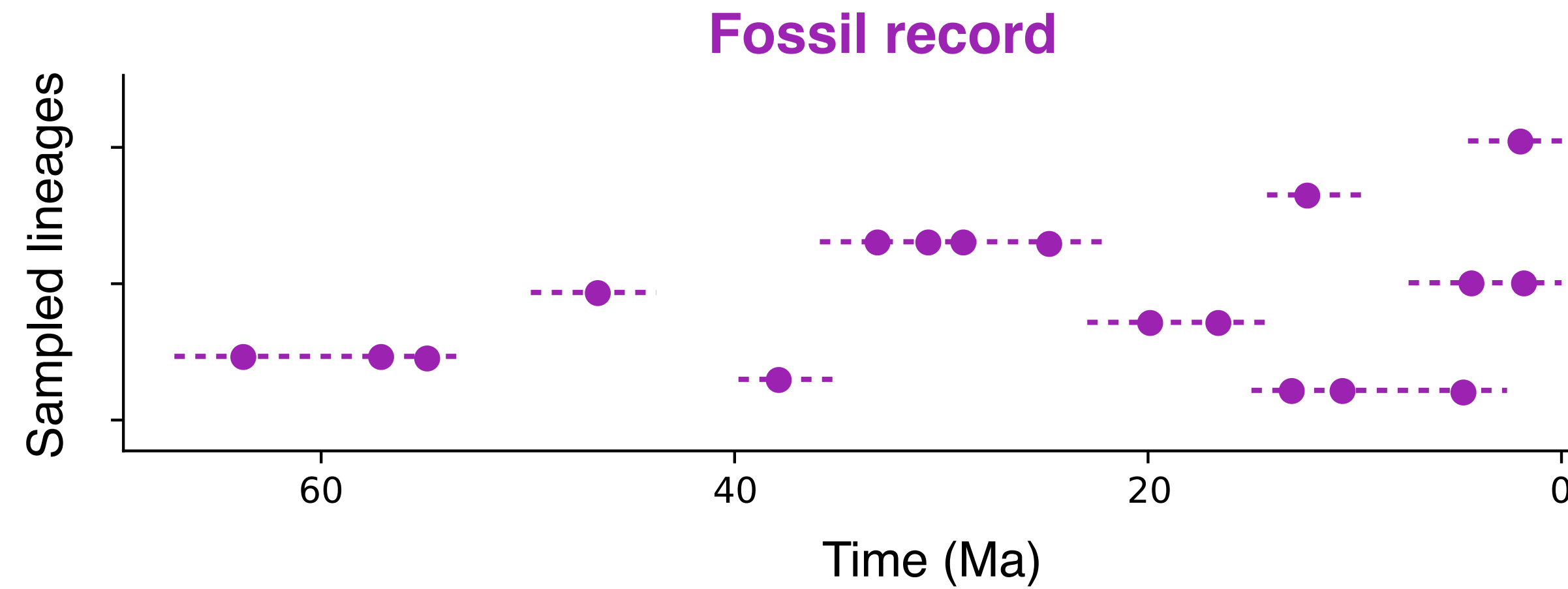
Trained model



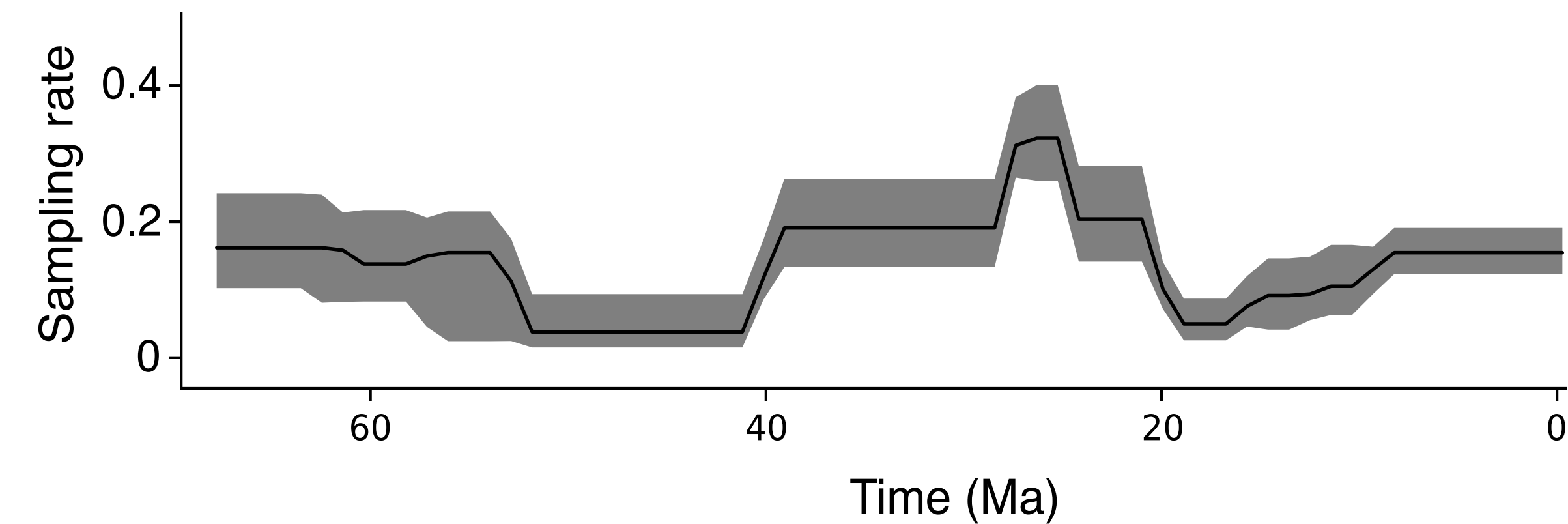
Bayesian (unsupervised) estimation of speciation and extinction rates from fossils



A birth-death neural network model of speciation and extinction



Preservation rates



Time-varying **Speciation** & **extinction** rates modeled by a NN

