# Computational Graph for Multi-layer Perceptron (MLP)

definitions:
input: $x \in \mathbb{R}^n$; $y \in \mathbb{R}^m$

Our MLP: $f(x; W)$ — weight matrix parameterizing our MLP

loss function: $L(f(x, W), y)$

$\sigma(\cdot)$: activation function, e.g., tanh, sigmoid
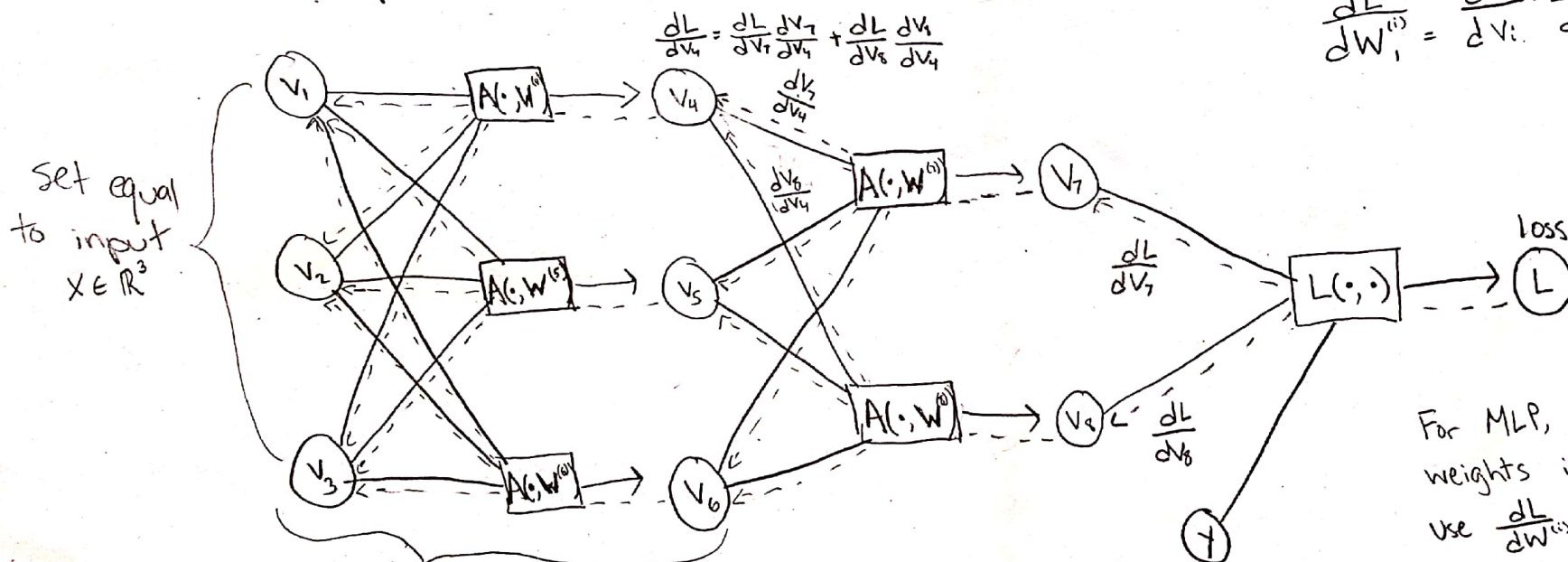
$*$ For simplicity, we consider $A(x, W) = \sigma(x^T W)$ as a primitive

How to calculate $\dfrac{dL}{dW^{(i)}}$?

$$v_i = A(\cdot, W^{(i)}) = \sigma(x^T W^{(i)}) = \sigma(x_1 W_1^{(i)} + \dots + x_n W_n^{(i)})$$

$$\frac{dL}{dW_i^{(i)}} = \frac{dL}{dv_i} \cdot \frac{dv_i}{dW_i^{(i)}} = \frac{dL}{dv_i} \sigma'(x^T W^{(i)}) x_i^{(i)}$$

$*$ here we use $x = \langle v_{pa(i)} \rangle$, the set of nodes coming into $v_i$



$$\frac{dL}{dv_4} = \frac{dL}{dv_7}\frac{dv_7}{dv_4} + \frac{dL}{dv_8}\frac{dv_8}{dv_4}$$

$\dfrac{dv_7}{dv_4}$

$\dfrac{dv_8}{dv_4}$

$\dfrac{dL}{dv_7}$

$\dfrac{dL}{dv_8}$

Set equal to input $x \in \mathbb{R}^3$

loss

close up

For MLP, we optimize over $W$, all weights in our MLP. We can then use $\dfrac{dL}{dW^{(i)}}$ in e.g., gradient descent updates

$$W^{(i)} \leftarrow W^{(i)} - \alpha \frac{dL}{dW^{(i)}}$$

learning rate