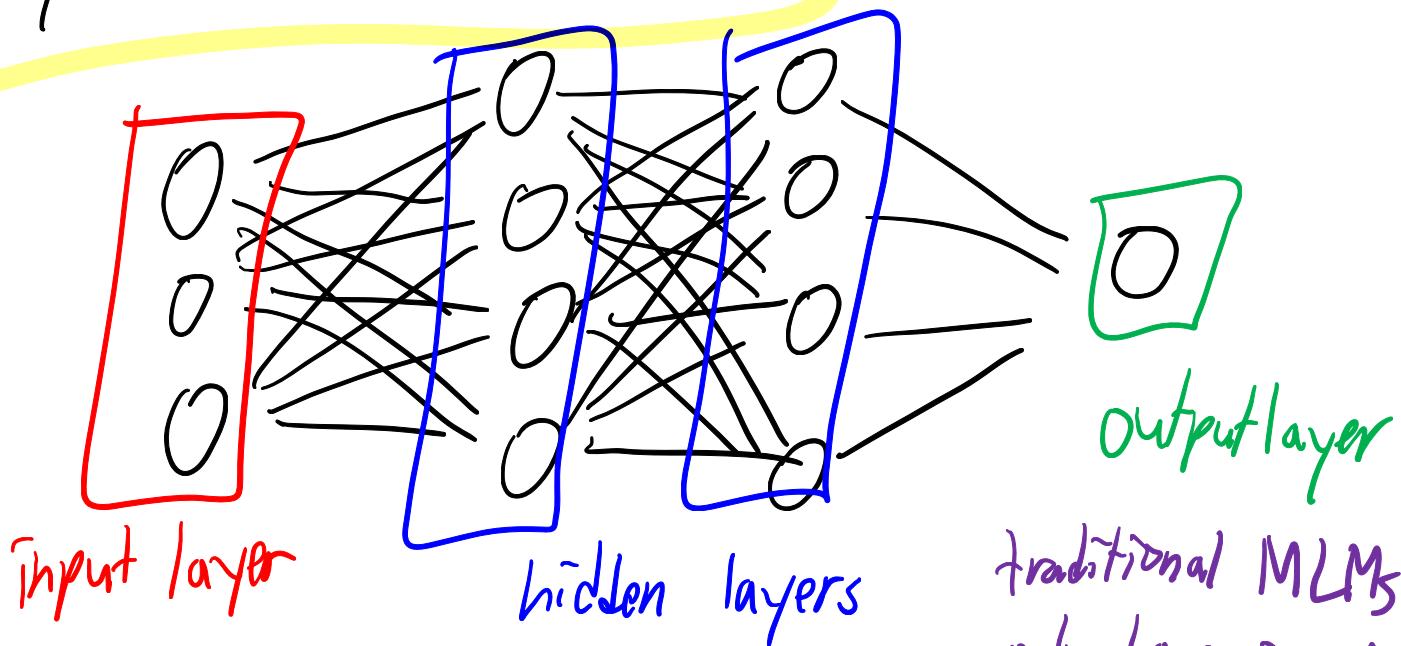


Multi-Layer Perceptron (MLP)

4/4

Fully-Connected Neural Network

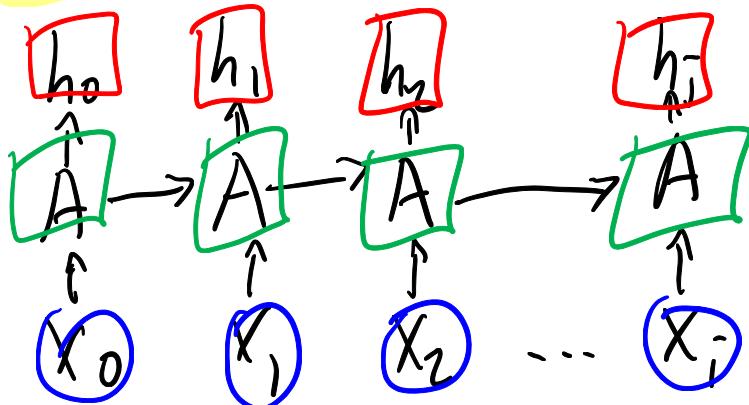


traditional MLPs
only have input
and output

Convolutional Neural Network

- each neuron is connected to only part of the previous neuron
- often used for Image data

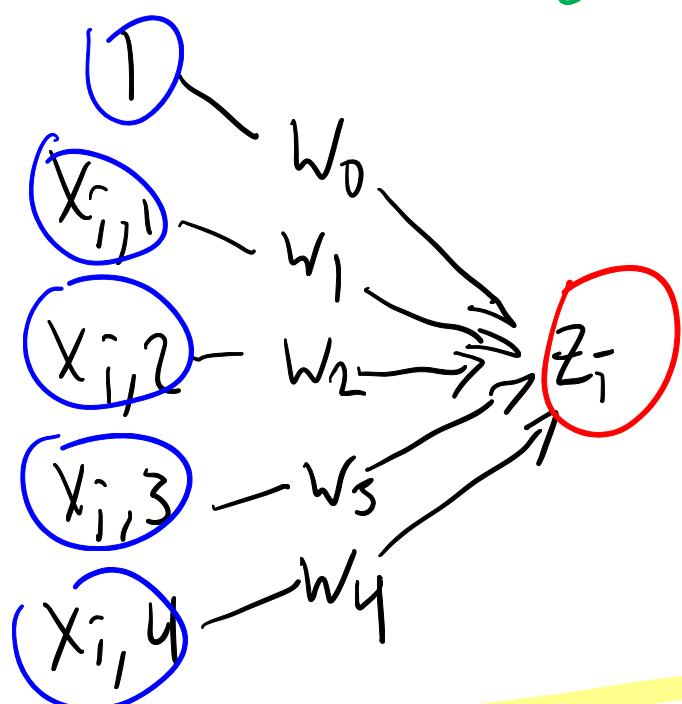
Recurrent Neural Network



- used when there are dependencies in sequential data

Logistic Regression

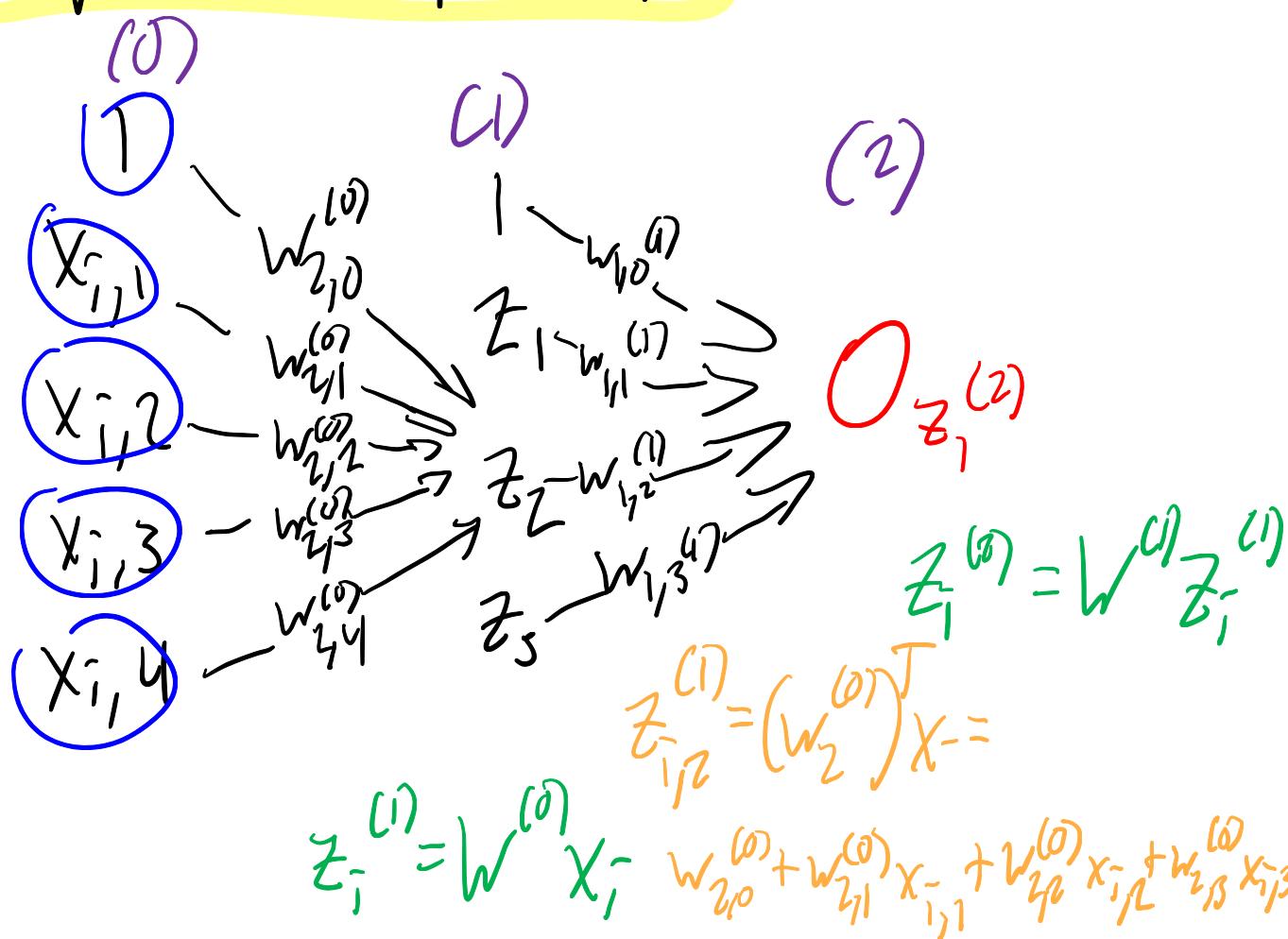
Given n samples: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$



$$z_i = w^T x_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_4 x_{i,4}$$

There are no hidden layers.
Shallow learning

Single \rightarrow Multiple Layers



$$\text{Layer 1: } z^{(0)} = w^{(0)} x_i$$

$$\text{Layer 2: } z^{(1)} = w^{(1)} x_i$$

$$\text{Layer } L: z^{(L)} = w^{(L-1)} x_i$$

$$w^{(0)} \in R^{d_1 \times d_0}$$

$$w^{(1)} \in R^{d_2 \times d_1}$$

$$w^{(L-1)} \in R^{d_L \times d_{L-1}}$$

$d_L = \# \text{ of nodes in } L$

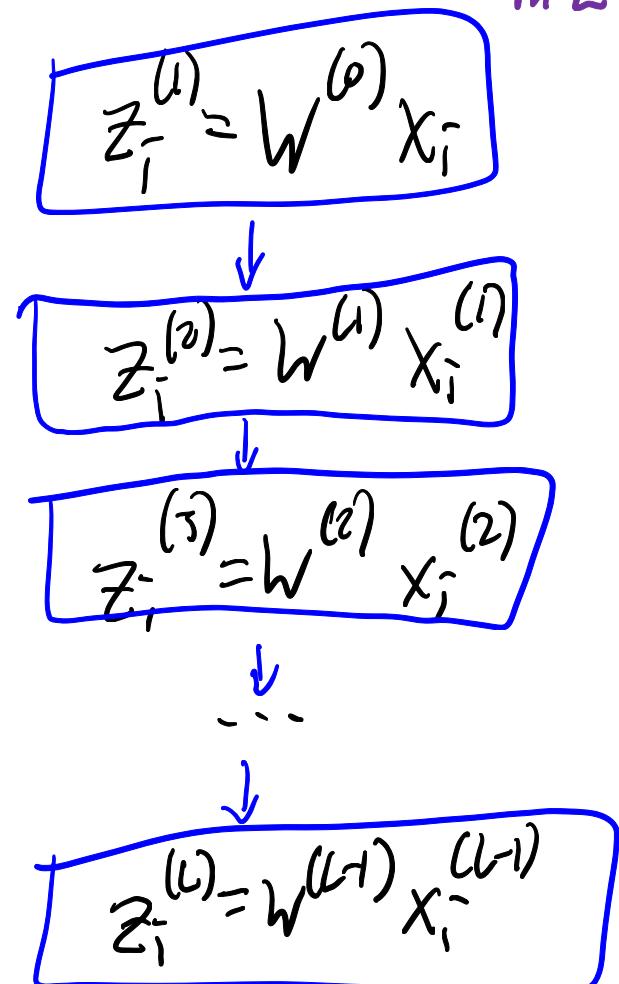
- Regression Task

$$f(x_i) = z_i^{(L)}$$

- Classification Task

$$f(x_i) = \sigma(z_i^{(L)})$$

σ for binary classification,
use softmax for multi-class



From Linear to Non-linear Model

- stacking multiple linear models is still a linear model

$$z_i^{(L)} = w^{(L-1)} \dots w^{(2)} w^{(1)} w^{(0)} x_i$$

deep linear networks are no more expressive than linear regression

Activation Function

- add (non-linear) activation functions to hidden layers
- multilayer fully-connected neural nets with nonlinear activation functions are (universal) approximators;
they can approximate any function arbitrarily well

- sigmoid function

Linear Model

$$\bar{z}_i^{(1)} = w^{(0)} \bar{x}_i$$

$$\bar{z}_i^{(2)} = w^{(1)} h_i^{(1)}$$

$$\bar{z}_i^{(3)} = w^{(2)} h_i^{(2)}$$

...

$$\bar{z}_i^{(L)} = w^{(L-1)} h_i^{(L-1)}$$

Non-linear model

$$\bar{z}_i^{(1)} = w^{(0)} \bar{x}_i, h_i^{(1)} = \sigma(\bar{z}_i^{(1)})$$

$$\bar{z}_i^{(2)} = w^{(1)} h_i^{(1)}, h_i^{(2)} = \sigma(\bar{z}_i^{(2)})$$

$$\bar{z}_i^{(3)} = w^{(2)} h_i^{(2)}, h_i^{(3)} = \sigma(\bar{z}_i^{(3)})$$

$$\bar{z}_i^{(L)} = w^{(L-1)} h_i^{(L-1)}, h_i^{(L)} = \sigma(\bar{z}_i^{(L)})$$

- tanh function
- ReLu function if the neuron is negative, it turns to 0;
some other wise

4/7

Image Classification

- 60,000 training samples $\{(x_i, y_i)\}$

with Logistic Regression:

- vectorize each 28×28 image to 784-dim vector, $x_i \in \mathbb{R}^{784}$
- add constant 1 to each x_i (bias), $x_i \in \mathbb{R}^{785}$
- denote model parameter $W \in \mathbb{R}^{10 \times 785}$
- $Z_i = W x_i$
- output prediction with softmax

with MLP:

input image $x_i \in \mathbb{R}^{785}$

$$z_i^{(1)} = w^{(0)} x_i \in \mathbb{R}^{256}$$

hidden layer 1

$$h_i^{(1)} = \text{relu}(z_i^{(1)}) \in \mathbb{R}^{256}$$

$$z_i^{(2)} = w^{(1)} h_i^{(1)} \in \mathbb{R}^{128}$$

hidden layer 2

$$h_i^{(2)} = \text{relu}(z_i^{(2)}) \in \mathbb{R}^{128}$$

$$z_i^{(3)} = w^{(2)} h_i^{(2)} \in \mathbb{R}^{10}$$

multiclassification \rightarrow softmax

$$\hat{y}_i = \text{Softmax}(z_i^{(3)}) \in \mathbb{R}^{10}$$

Output layer

$$\text{loss function: } L = - \sum_{i=1}^{60000} \sum_{j=1}^{10} y_{ij} \log(\hat{y}_{ij})$$

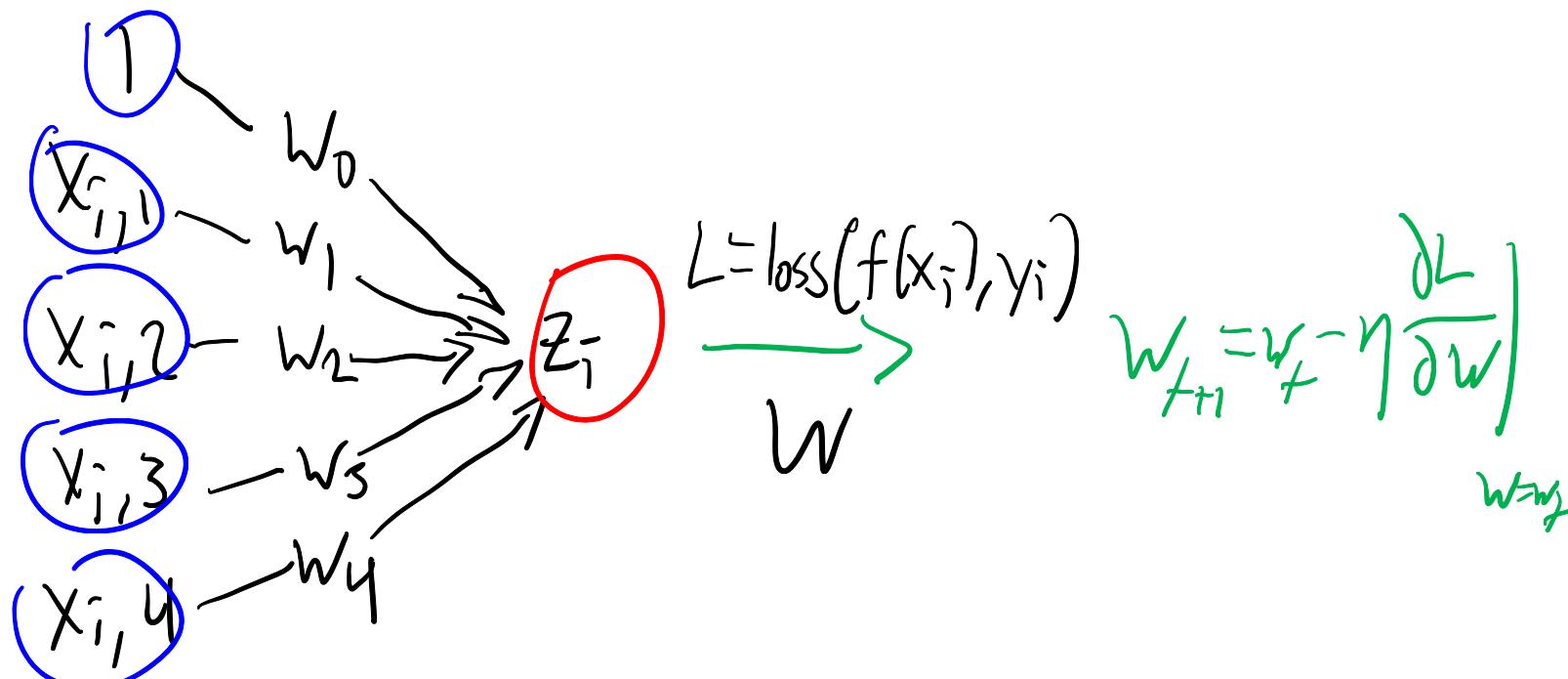
$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$w^{(0)} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1d} \\ w_{21} & & & w_{2d} \\ \vdots & & & \vdots \\ w_{d'1} & w_{d'2} & \dots & w_{d'd'} \end{bmatrix}$$

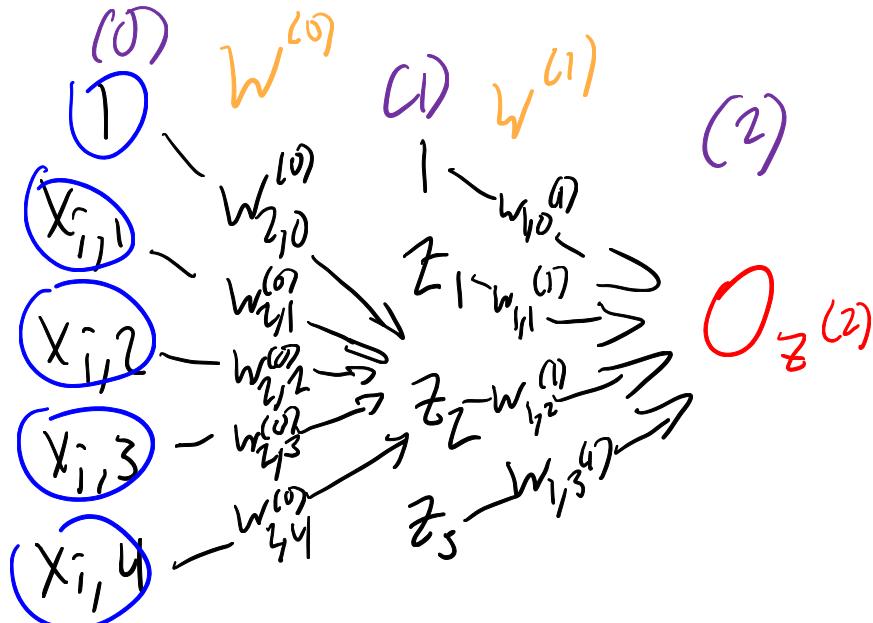
Optimization

- Stochastic Gradient Descent to learn model parameters

Single layer:



multiple layer:



It's easy to calculate $w^{(1)}$.

Calculate $w^{(0)}$ through backpropagation

Calculate features going forward.
Calculate gradients going backward.

chain rule:

$$h(x) = f(g(x))$$

gradient: $\frac{\partial h(x)}{\partial x} = \frac{\partial f(y)}{\partial y} \frac{\partial g(x)}{\partial x}$

$$f(y) = y^2$$

$$\frac{\partial f(y)}{\partial y} = 2y$$

$$g(x) = 3x + 5$$

$$\frac{\partial g(x)}{\partial x} = 3$$

$$h(x) = f(g(x))$$

$$\frac{\partial h(x)}{\partial x} = (2(3x + 5)) \times 3$$

Back propagation

D Compute gradients of last layer

Compute gradients as regular models

$$\frac{\partial L}{\partial w^{(l)}} \quad \frac{\partial L}{\partial z^{(l)}}$$

Loss function is function of $w^{(l)}$ and $z^{(l)}$

② Compute gradients of hidden layers
based on chain rule

$z^{(1)}$ is function of $w^{(0)}$: $z^{(1)} = w^{(0)}x$

$$\frac{\partial L}{\partial w^{(0)}} = \frac{\partial L}{\partial z^{(1)}} \boxed{\frac{\partial z^{(1)}}{\partial w^{(0)}}} = x$$

Forward Pass

$$z_i^{(1)} = w^{(0)}x_i \in \mathbb{R}^{256}$$

$$z_i^{(2)} = w^{(1)}x_i \in \mathbb{R}^{128}$$

$$z_i^{(3)} = w^{(2)} \text{relu}(z_i^{(2)}) \in \mathbb{R}^{10}$$

Backward Pass

$$\frac{\partial L}{\partial w^{(0)}} = \frac{\partial L}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(0)}}$$

$$\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial w^{(1)}}$$

$$\frac{\partial L}{\partial z^{(0)}} = \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial z^{(0)}}$$

$$\frac{\partial L}{\partial w^{(2)}}$$

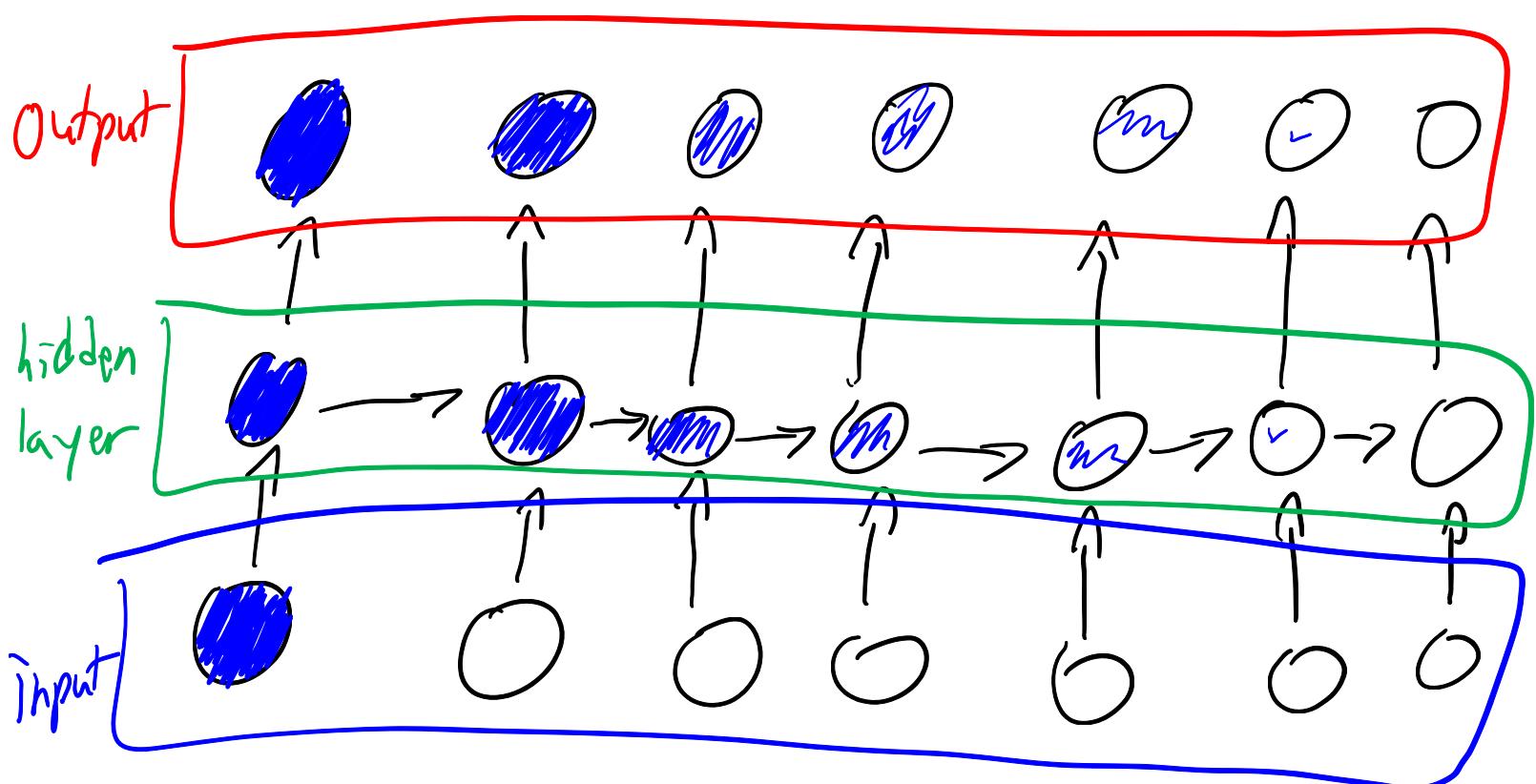
$$\frac{\partial L}{\partial z^{(0)}}$$

Recurrent Neural Networks

4/14

- an "internal state" that is updated as a sequence is processed

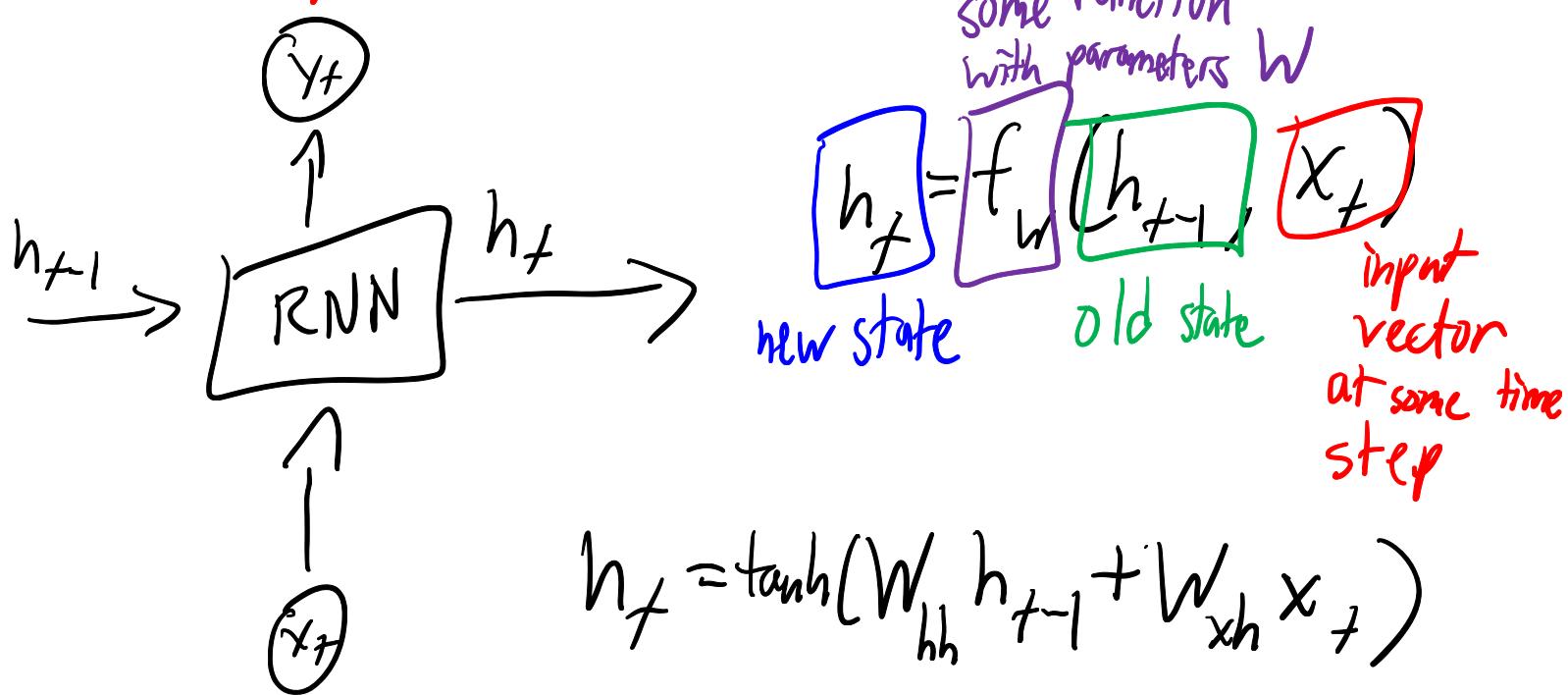
The last input includes all of the previous ones.



hidden layer captures dependence across words

- process a sequence of vectors x by applying a recurrence formula at every time step

- 1) update hidden state



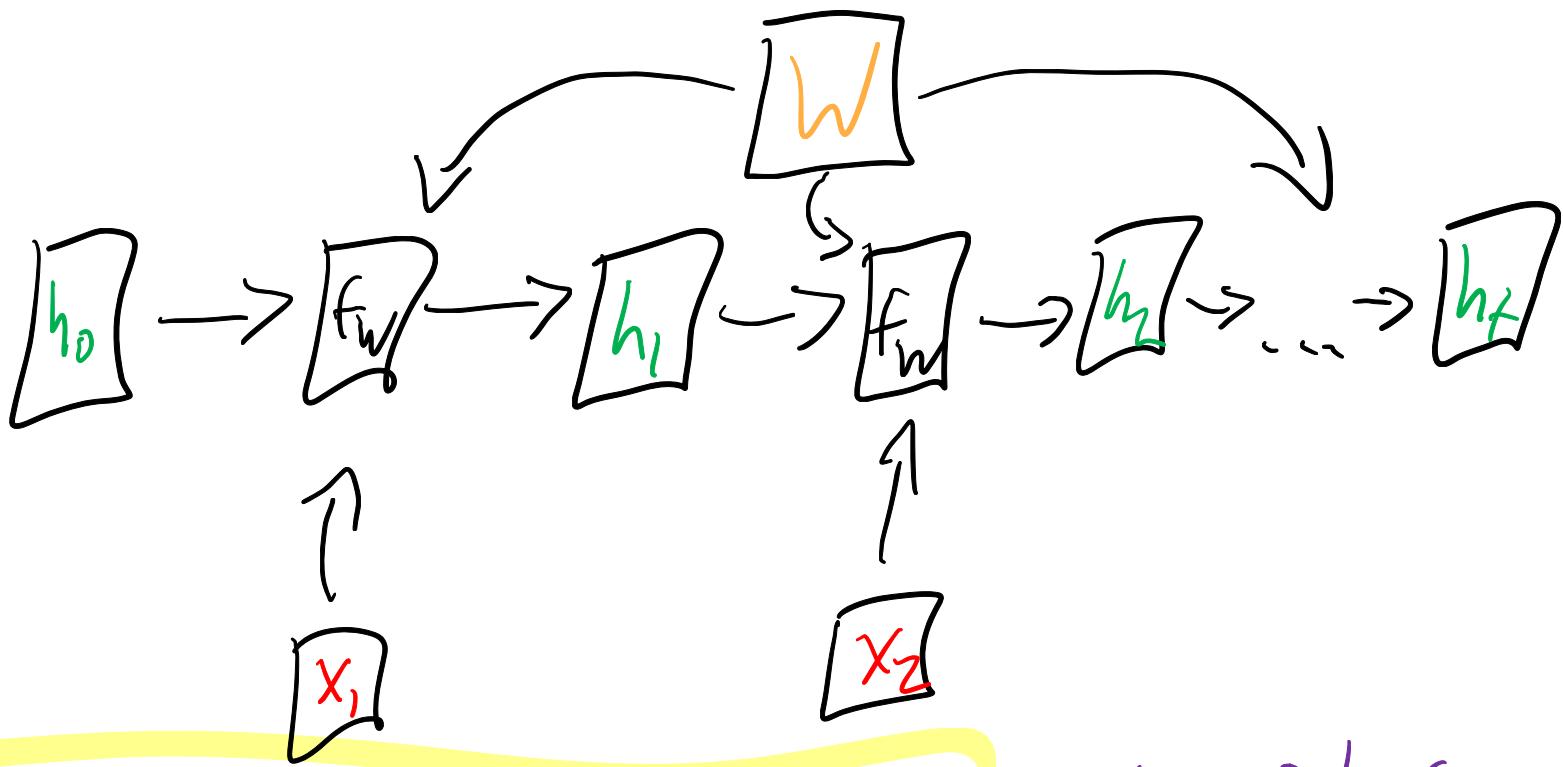
- update output

$$y_t = W_{hy} h_t$$

- Forward pass

$$h_f = \tanh(W_{hh} h_{t-1} + W_{xh} x_t + b)$$

- weight sharing: same for all states

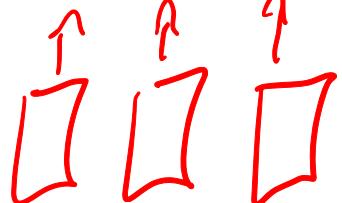


Example - Sentiment Classification

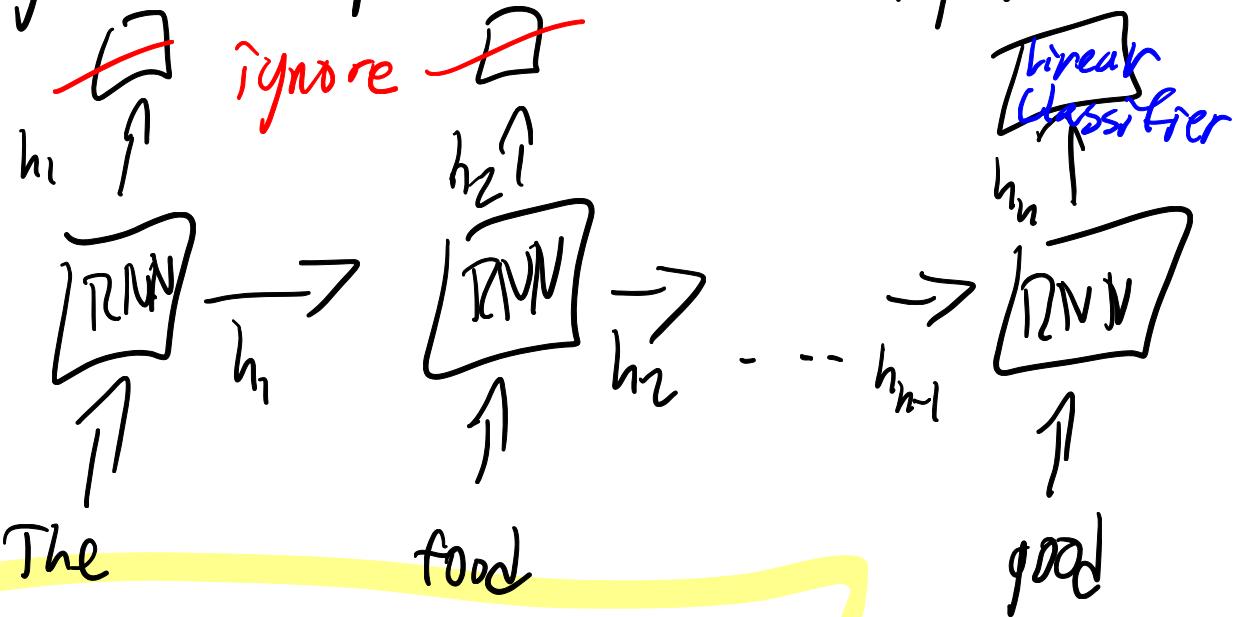
- Many to one

- input: one or more sentences
- output: positive/negative class

We only care about the output at the last time step $\boxed{\quad} \rightarrow \boxed{\quad} \rightarrow \boxed{\quad}$



- Ignore output of internal input



Example - Image Captioning

- One to many

- input: image features (from CNN)
- output: sentence describing image

