

Applied Artificial Intelligence

Session 20: DNN for Regression, Advanced Optimization, CNNs, and Motivation for RNNs

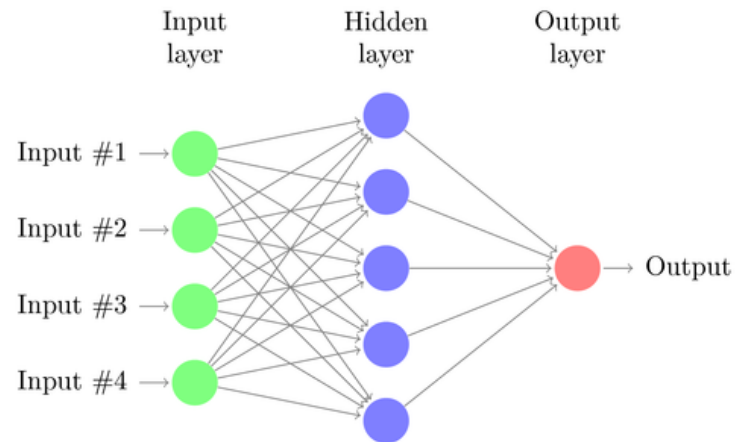
Fall 2018

NC State University

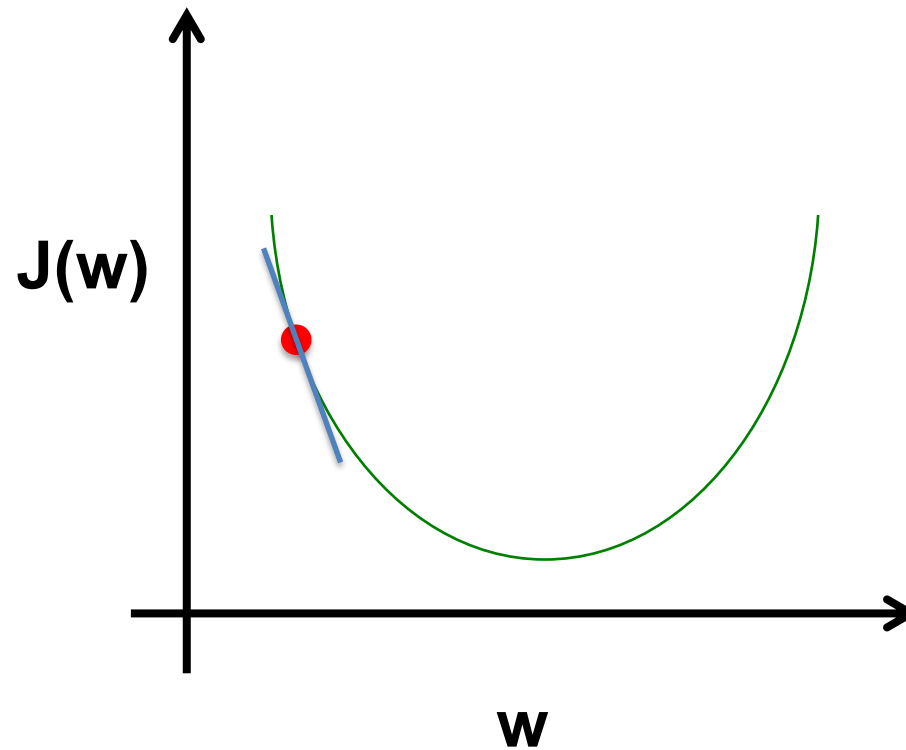
Lecturer: Dr. Behnam Kia

Course Website: <https://appliedai.wordpress.ncsu.edu/>

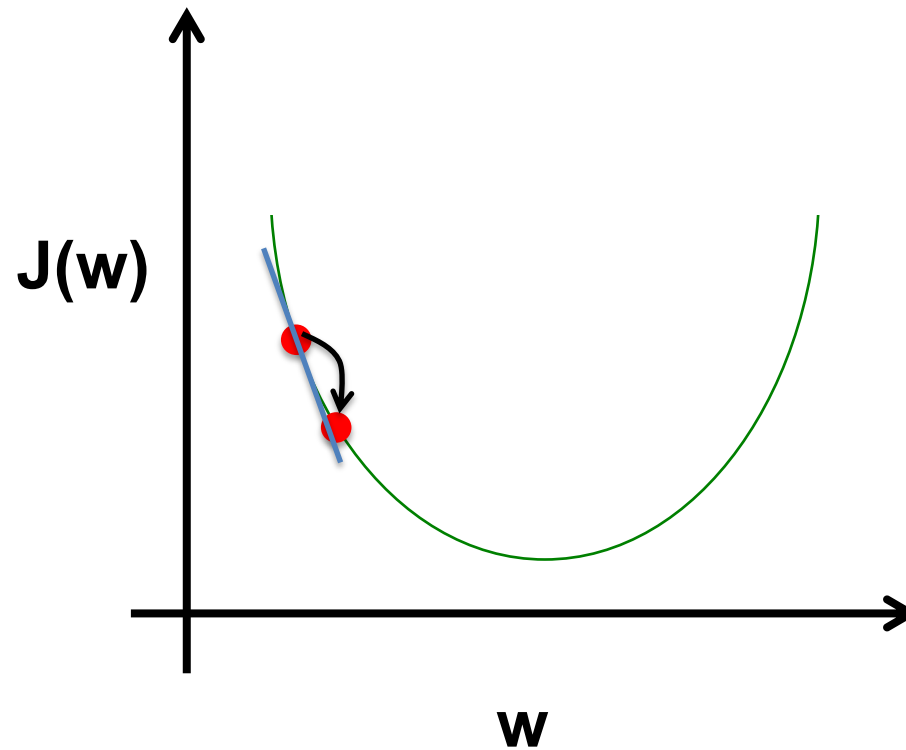
Feedforward networks for Regression

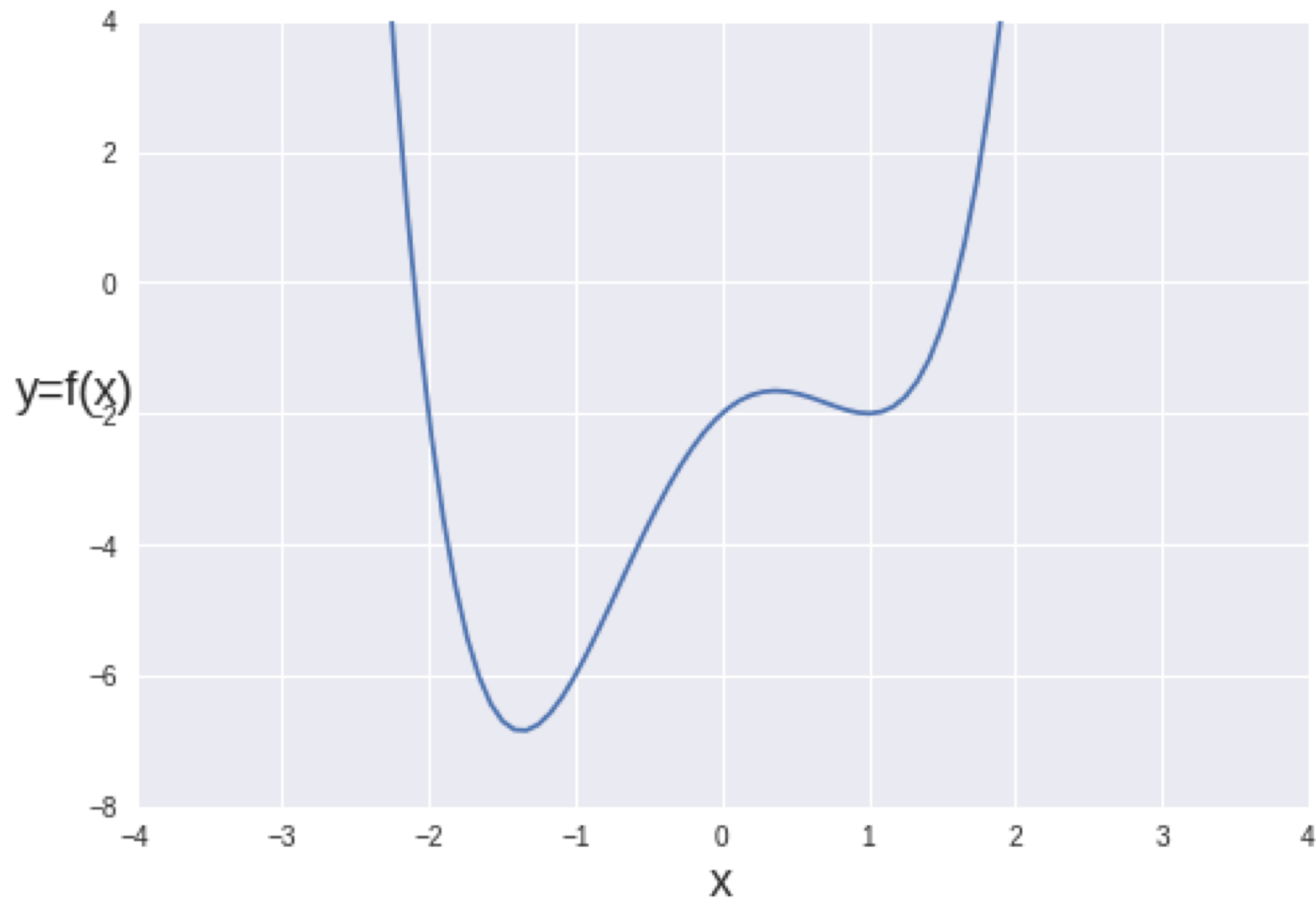


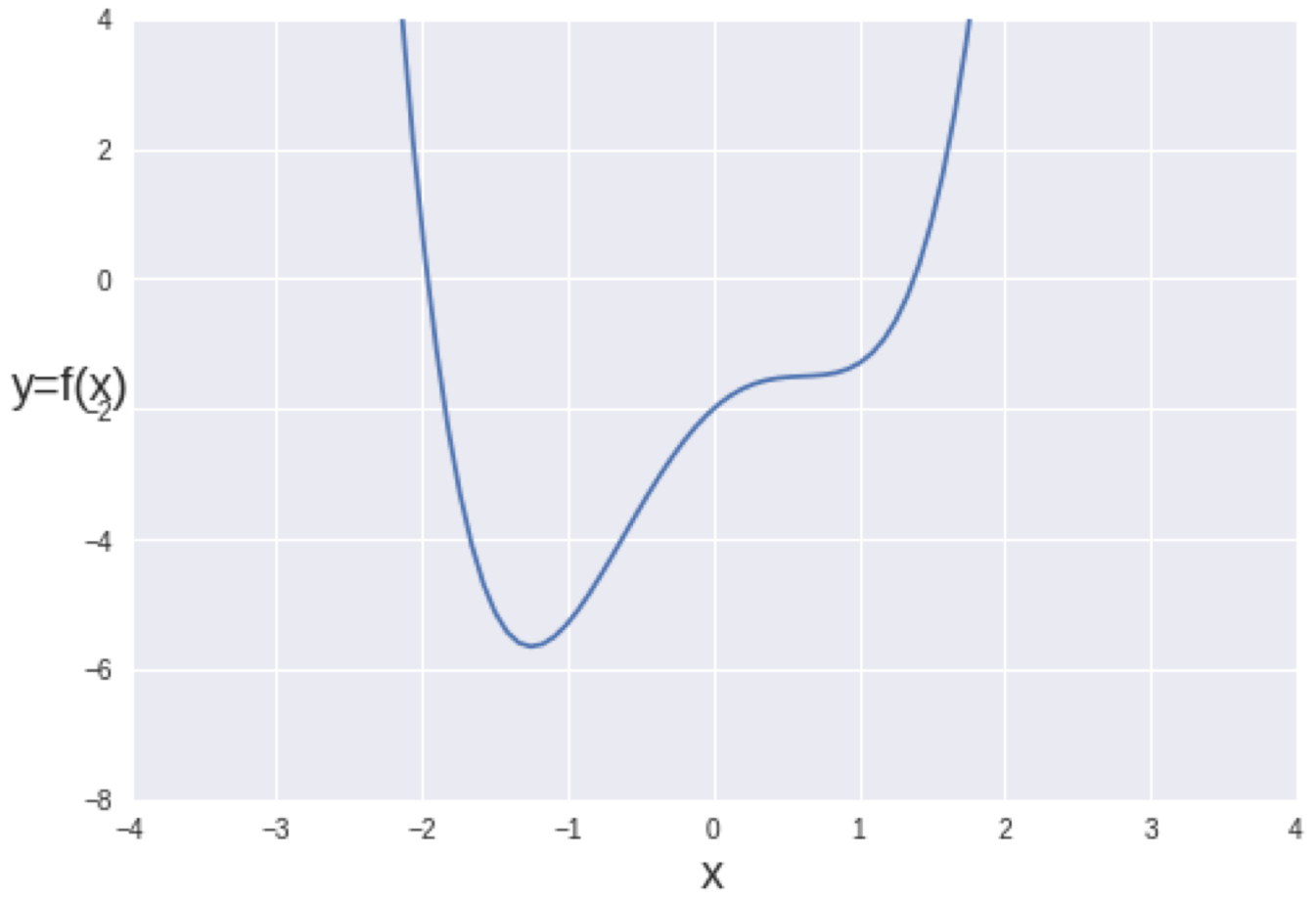
Gradient Descent



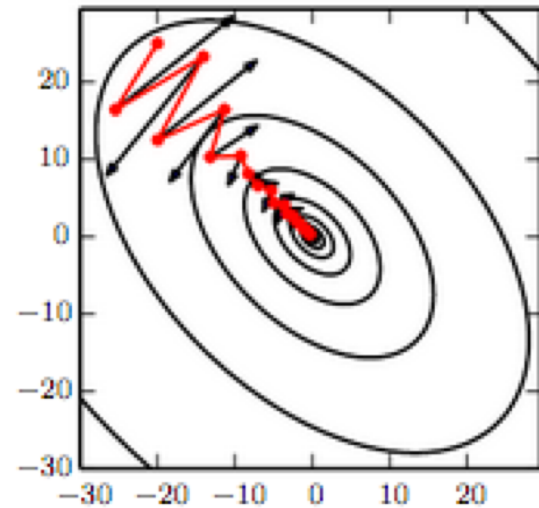
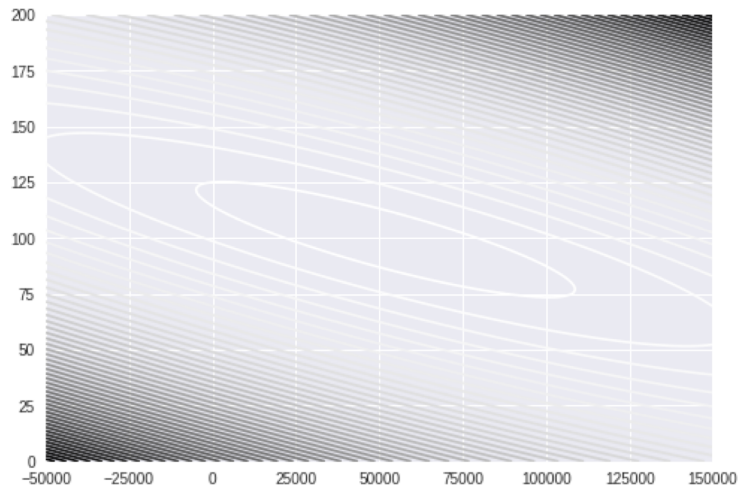
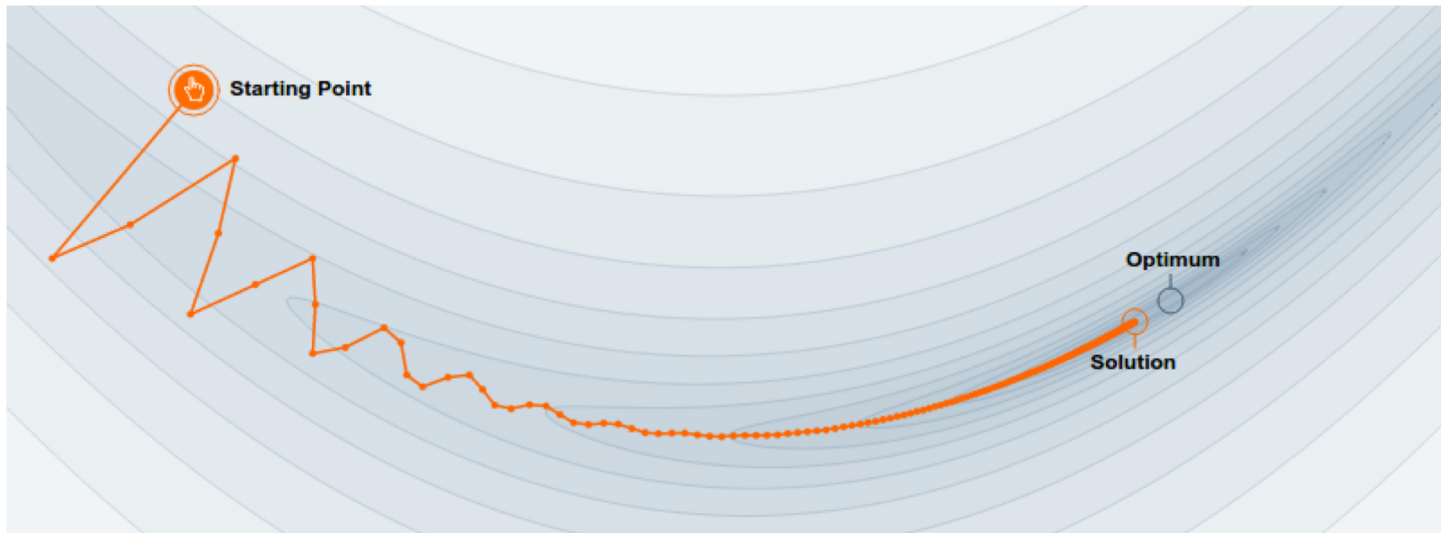
Gradient Descent







Momentum



Momentum

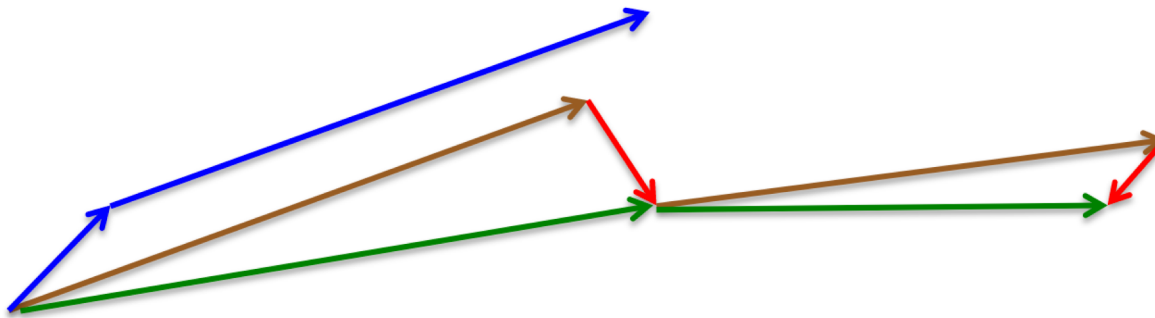
- Exponentially decaying moving average of the gradients
- *Numerical Example:*

$$v = -[0.9^1 \times g_{t-1} + 0.9^2 \times g_{t-2} + \dots + 0.9^1 \times g_{t-10} + (1 - 0.9) \times g_t]$$

- g_i is gradient at step i
- Window size is 10 here, it could be infinity. With exponentially decaying weights it won't differ too much.
- $v_t = 0.9 \times v_{t-1} - (1 - 0.9) \times g_t$
- Instead of 0.9 we can take other numbers between 0 and 1.
- The closer to 1, the higher the momentum.

A picture of the Nesterov method

- **First** make a big jump in the direction of the previous accumulated gradient.
- **Then** measure the gradient where you end up and make a correction.



brown vector = jump, red vector = correction, green vector = accumulated gradient

blue vectors = standard momentum

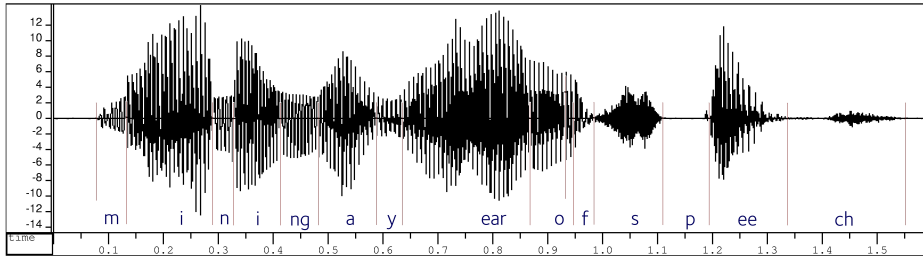
RMSProp

- **RMSProp**: Adapting step size separately for each parameter.
 - For each direction (parameter) divide the gradient by a running average of its recent magnitude

Adam

- Adam \approx RMSProp + Momentum

Sequence Processing



Amazon Movie Reviews @AmznMovieRevws · Jul 24

Life of Pi.

3.1415HOLY

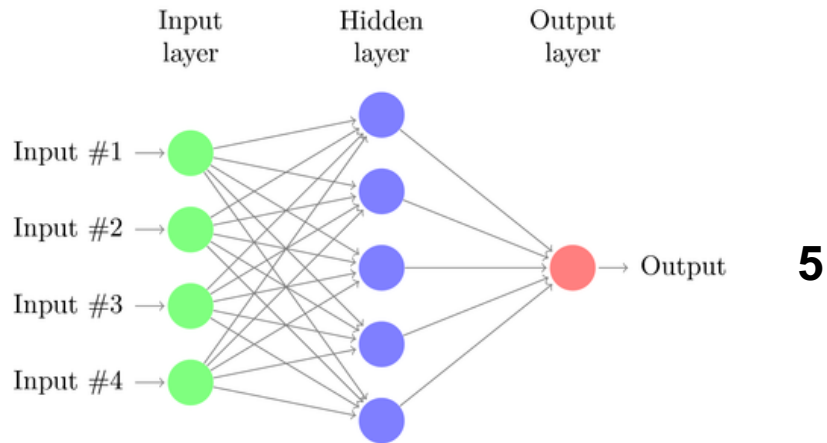
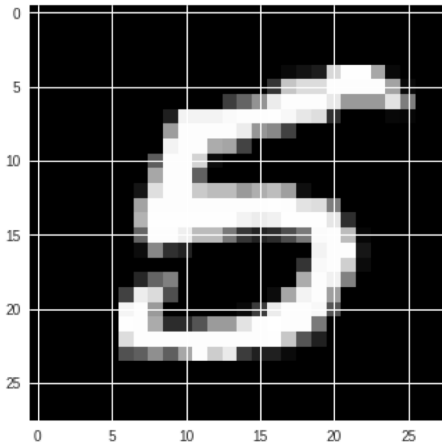


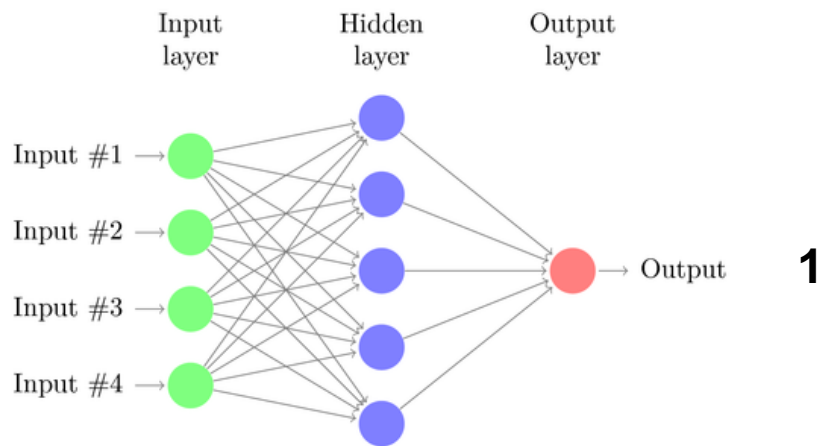
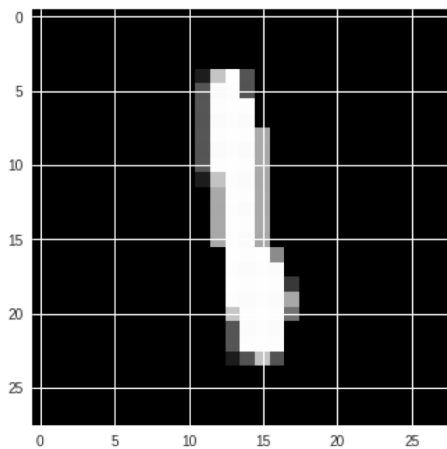
Misleading

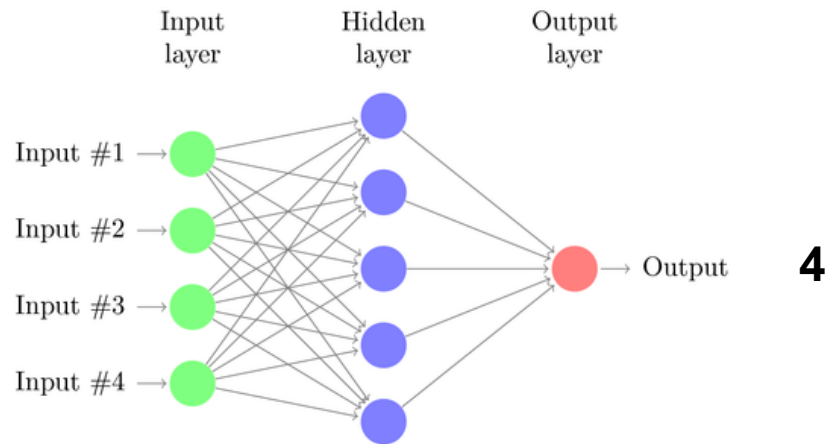
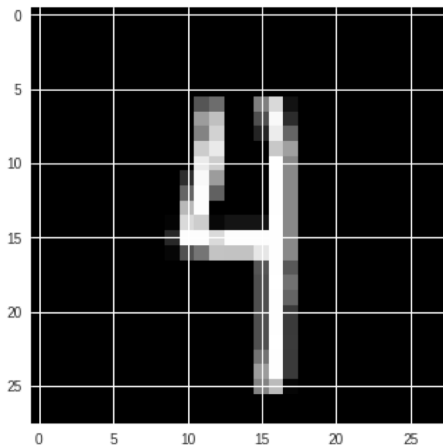
By Michael M. - May 8, 2015

I screened this film for the students in my 9th grade algebra class and was very upset.

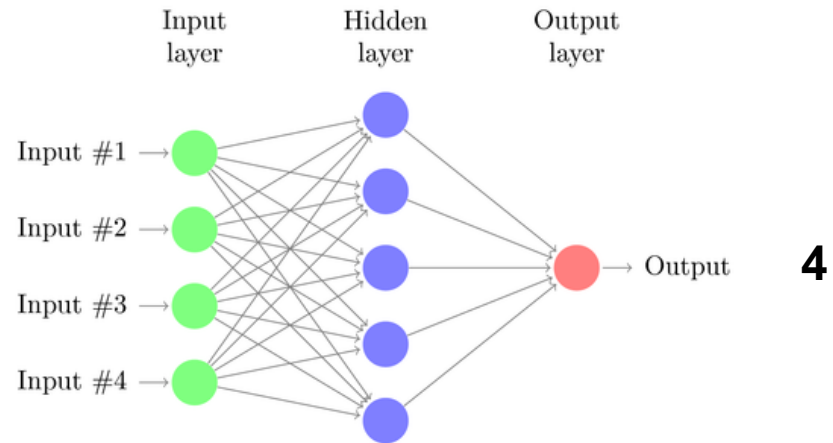
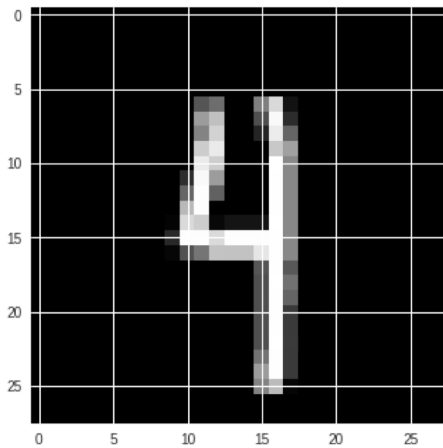


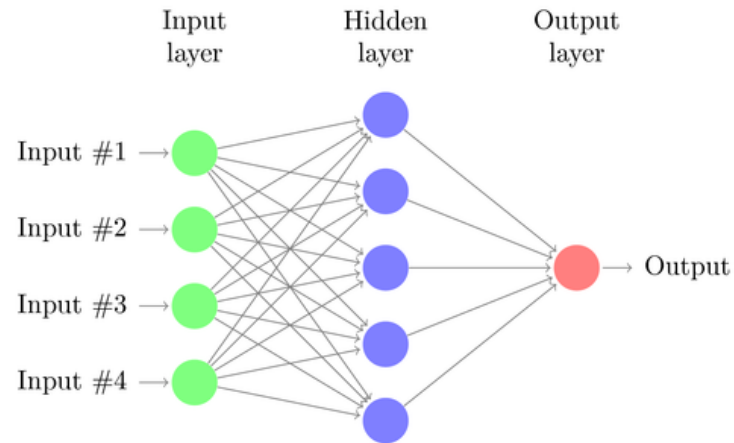
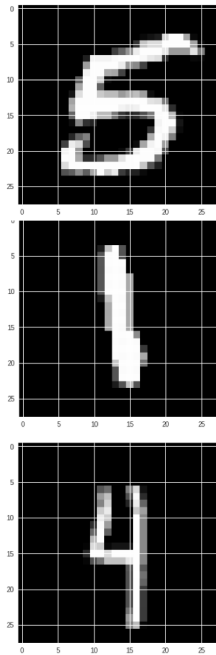






Feedforward networks have no memory

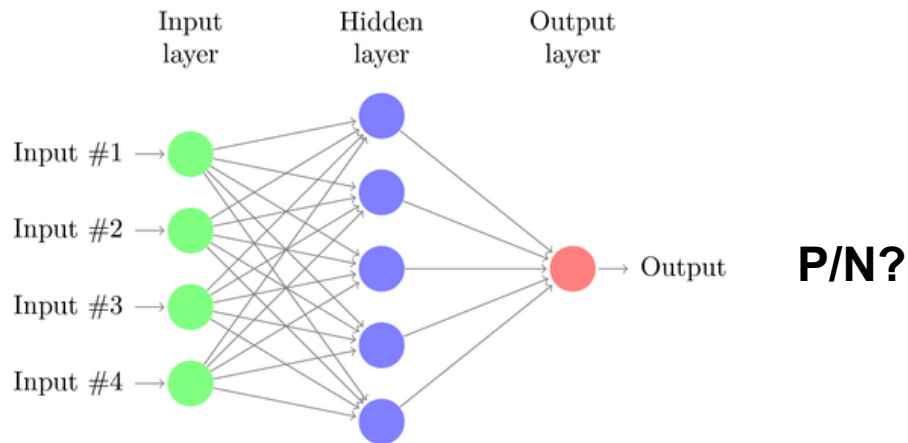
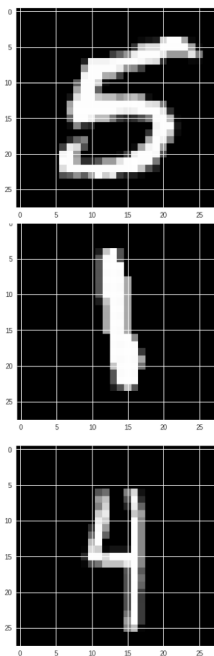




P/N?

-
-
-

Not feasible for very large sequences

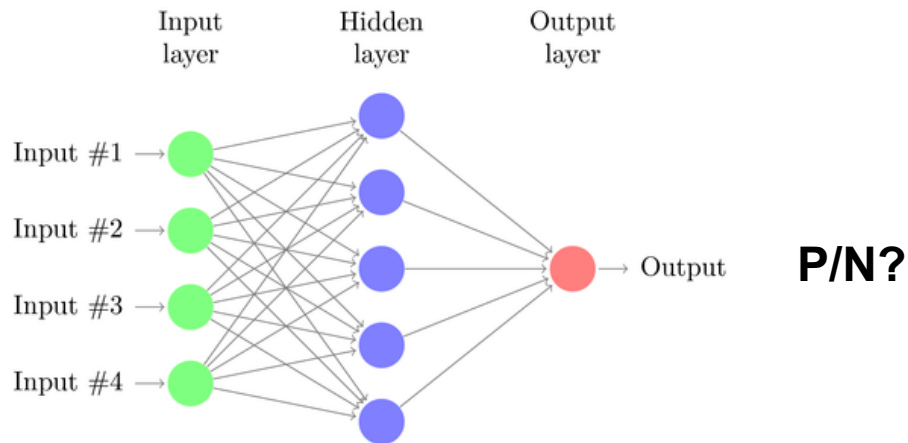
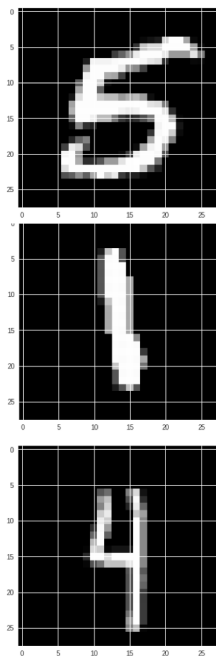


⋮

I went to Nepal in 2009.

In 2009, I went to Nepal.

Separate parameters for each time step.

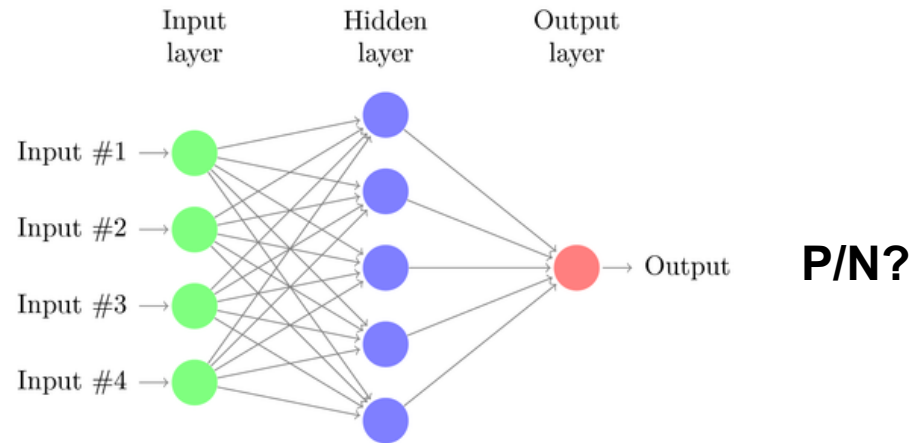
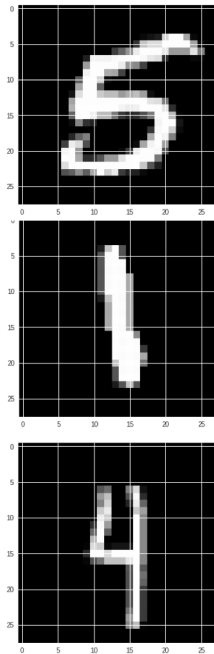


⋮

I went to Nepal in 2009.

In 2009, I went to Nepal.

No/poor generalization, especially for sequence sizes not observed during the training.

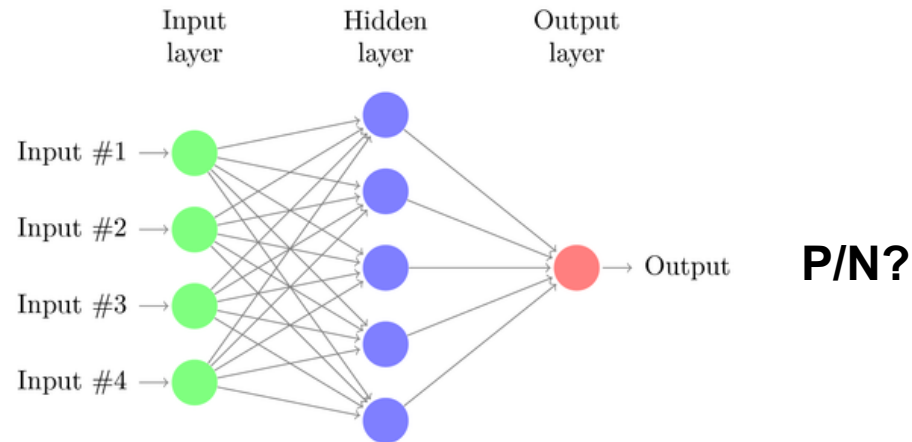
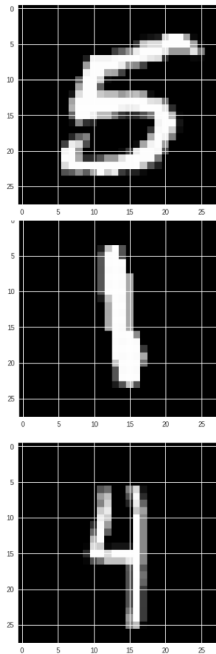


⋮

I went to Nepal in 2009.

In 2009, I went to Nepal.

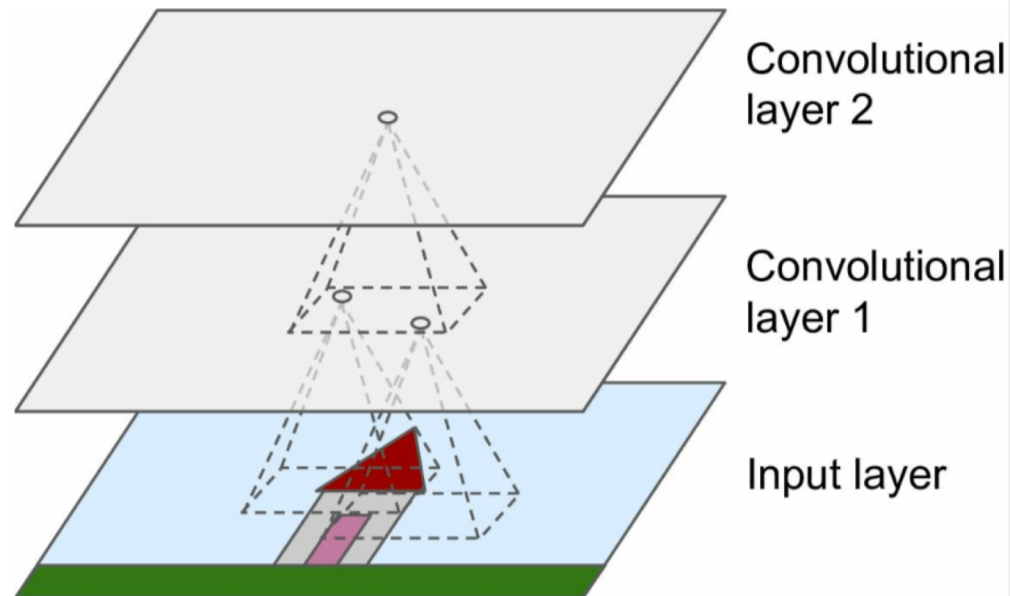
No/poor generalization, especially for sequence sizes not observed during the training.



When you have separate parameters for each location, all rules need to be separately learned for each position.

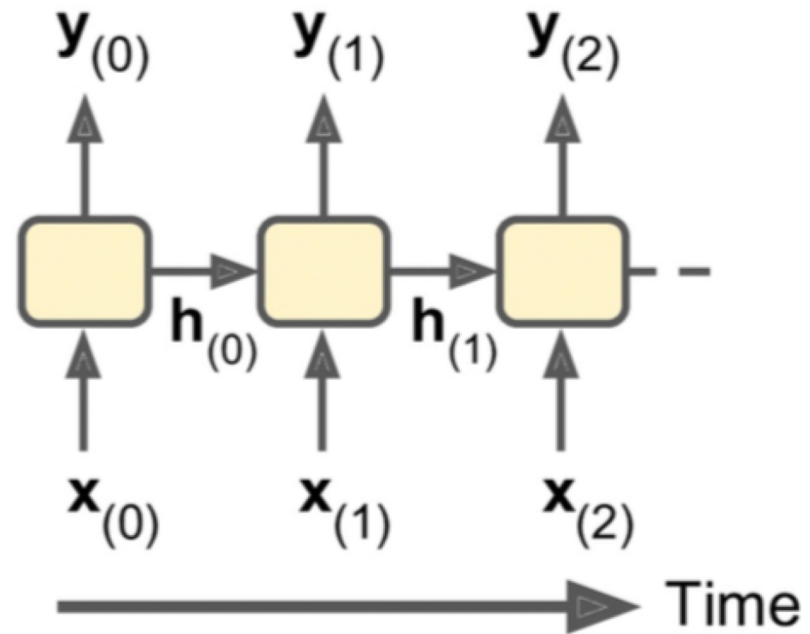
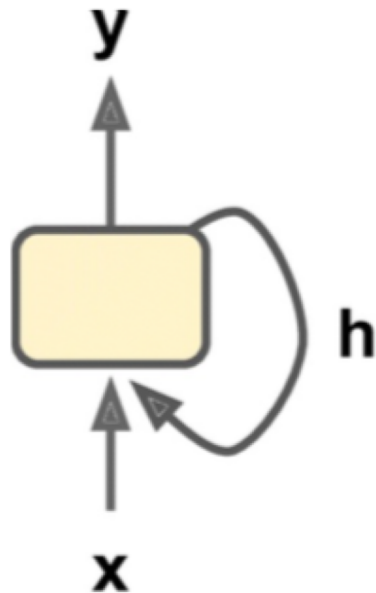
CNNs

- CNNs are specialized to handle a grid of data, e.g. an images.
- And they can scale easily for larger images.
- Parameter sharing through application of kernels.



Recurrent Neural Networks

(opposed to feedforward Neural Networks)



$$\hat{y} = \sigma(X.W_x + Y.W_y + b)$$