

Machine Learning for Finance (FIN 570)

Neural Network (Deep Learning)

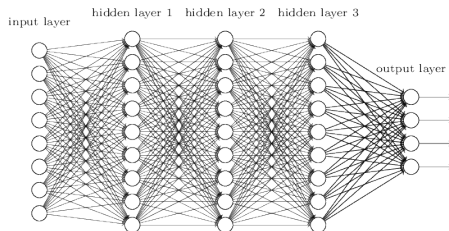
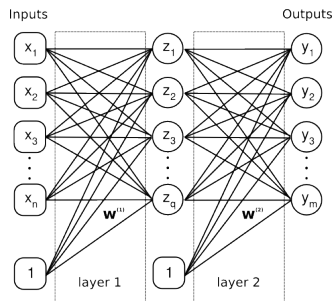
Instructor: Jaehyuk Choi

Peking University HSBC Business School, Shenzhen, China

2019-20 Module 3 (Spring 2020)

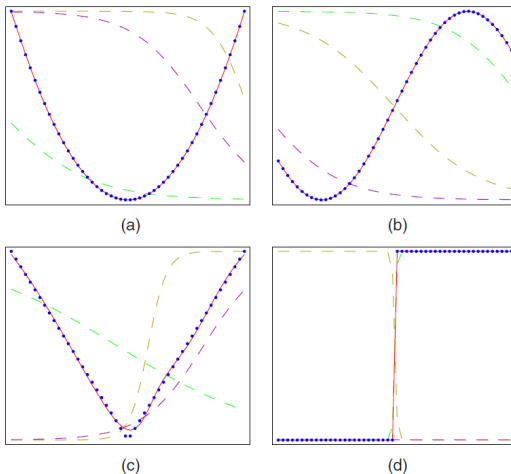
Neural Network

- Multi-layer 'perceptron' (MLP): logistic regression
- **Deep learning**: a set of methods to efficiently train multi-layer NN
 - backpropagation, activation function, dropout, etc



Neural Network

Figure 5.3 Illustration of the capability of a multilayer perceptron to approximate four different functions comprising (a) $f(x) = x^2$, (b) $f(x) = \sin(x)$, (c) $f(x) = |x|$, and (d) $f(x) = H(x)$ where $H(x)$ is the Heaviside step function. In each case, $N = 50$ data points, shown as blue dots, have been sampled uniformly in x over the interval $(-1, 1)$ and the corresponding values of $f(x)$ evaluated. These data points are then used to train a two-layer network having 3 hidden units with 'tanh' activation functions and linear output units. The resulting network functions are shown by the red curves, and the outputs of the three hidden units are shown by the three dashed curves.



NN Notations

The conventions are **different** from other books.

Single layer

$$\mathbf{y} = \phi(\mathbf{x}\mathbf{W}) \Rightarrow [\cdots y_j \cdots] = [\cdots x_i \cdots] \begin{bmatrix} \cdots \\ \cdots W_{ij} \cdots \\ \cdots \end{bmatrix} \quad (x_0 = 1)$$

W_{ij} connects i -th input x_i to j -th output y_j . $\phi(x)$ is the sigmoid function.

Multi-layer: L layers

$$\mathbf{x}^{(1)} = \phi(\mathbf{a}^{(1)} = \mathbf{x}^{(0)}\mathbf{W}^{(1)})$$

$$\mathbf{x}^{(2)} = \phi(\mathbf{a}^{(2)} = \mathbf{x}^{(1)}\mathbf{W}^{(2)})$$

\dots

$$\mathbf{x}^{(L)} = \phi(\mathbf{a}^{(L)} = \mathbf{x}^{(L-1)}\mathbf{W}^{(L)})$$

$$\mathbf{y} = \phi(\cdots \phi(\phi(\mathbf{x}^{(0)}\mathbf{W}^{(1)})\mathbf{W}^{(2)})\cdots)$$

Input: $\mathbf{x} = \mathbf{x}^{(0)}$, Output: $\mathbf{y} = \mathbf{x}^{(L)}$

$$[\cdots y_k \cdots] = [\cdots x_i \cdots]$$

$$\times \begin{bmatrix} \cdots \\ \cdots W_{ij}^{(1)} \cdots \\ \cdots \end{bmatrix} \cdots \begin{bmatrix} \cdots \\ \cdots \times W_{jk}^{(L)} \cdots \\ \cdots \end{bmatrix}$$

Training NN

- Multi-variate response variable (n : samples, j : responses)

$$\mathbf{Y} = \begin{bmatrix} \cdots \\ \cdots Y_{nj} \cdots \\ \cdots \end{bmatrix} = \mathbf{X}_{n*} \mathbf{W}_{*j}$$

- Error (loss) function for the n -th sample:

$$J_n(\mathbf{W}) = - \sum_j Y_{nj} \log \phi(X_{nj}) + (1 - Y_{nj}) \log (1 - \phi(X_{nj}))$$

- Gradient: (i : input, j : output)

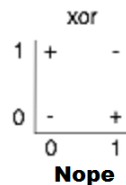
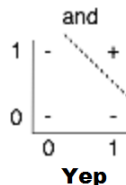
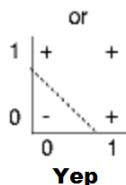
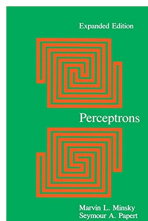
$$\frac{\partial J_n(\mathbf{W})}{\partial W_{ij}} = -(Y_{nj} - \phi(X_{nj})) X_{ni}$$

- Total loss function with regularization:

$$J(\mathbf{W}) = \sum_n J_n(\mathbf{W}) + \lambda_1 \sum_k \sum_{ij} |W_{ij}^{(k)}| + \frac{\lambda_2}{2} \sum_k \sum_{ij} (W_{ij}^{(k)})^2$$

The 1st AI winter (1969~1986)

- *Perceptrons* by Minsky and Papert (1969) mathematically proved that the single layer can not learn exclusive OR (XOR) gate due to its nonlinear feature. MLP is needed but it's too complicated to train the parameters.
- The criticism from the leading AI figure effectively killed all AI research.
- The difficulty is resolved by backpropagation by Werbos (1974,1982), Hinton(1986).



Backpropagation: chain rule

It's a fancy name of chain rule in calculus. Imagine input x_0 goes through functions f_1 , f_2 and f_3 to reach the output $y = f_3(f_2(f_1(x_0)))$:

$$x_0 \rightarrow f_1(x_0) = x_1 \rightarrow f_2(x_1) = x_2 \rightarrow f_3(x_2) = x_3 = y$$

The chain rule is to multiply the derivative of each function:

$$\frac{\partial y}{\partial x_0} = \frac{\partial x_3}{\partial x_2} \frac{\partial x_2}{\partial x_1} \frac{\partial x_1}{\partial x_0} = f'_3(x_2) f'_2(x_1) f'_1(x_0),$$

Now imagine the weight w is considered, $y = f_3(f_2(f_1(x_0 w_1) w_2) w_3)$

$$x_0 \rightarrow f_1(a_1 = x_0 w_1) = x_1 \rightarrow f_2(a_2 = x_1 w_2) = x_2 \rightarrow f_3(a_3 = x_2 w_3) = x_3 = y$$

$$\frac{\partial y}{\partial w_3} = \frac{\partial}{\partial w_3} f_3(x_2 w_3) = f'_3(x_2 w_3) x_2 = f'_3(a_3) x_2 = \frac{\partial y}{\partial a_3} x_2$$

$$\frac{\partial y}{\partial w_2} = \frac{\partial}{\partial w_2} f_3(f_2(x_1 w_2) w_3) = f'_3(a_3) w_3 \cdot f'_2(a_2) x_1 = \frac{\partial y}{\partial a_2} x_1$$

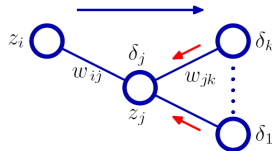
$$\frac{\partial y}{\partial w_1} = \frac{\partial}{\partial w_1} f_3(f_2(f_1(x_0 w_1) w_2) w_3) = f'_3(a_3) w_3 \cdot f'_2(a_2) w_2 \cdot f'_1(a_1) x_0 = \frac{\partial y}{\partial a_1} x_0$$

Backpropagation: error term δ

Imagine W_{ij} connects unit x_i (layer $m - 1$) to x_j (layer m) and the unit x_j is connected to the units x_k in the next layer $m + 1$. Remind that

$$a_j = \sum_i x_i W_{ij}, \quad x_j = \phi(a_j)$$

$$\frac{\partial J_n}{\partial W_{ij}} = \frac{\partial J_n}{\partial a_j} \frac{\partial a_j}{\partial W_{ij}} = \delta_j x_i, \quad \text{where} \quad \delta_j := \frac{\partial J_n}{\partial a_j}$$



The term δ_j is referred to as error because $\delta_j = \hat{y}_j - y_j$ for the output units. Now apply chain rule to obtain the backpropagation,

$$\delta_j := \frac{\partial J_n}{\partial a_j} = \sum_k \frac{\partial J_n}{\partial a_k} \cdot \frac{\partial a_k}{\partial x_j} \cdot \frac{\partial x_j}{\partial a_j} = \phi'(a_j) \sum_k W_{jk} \delta_k$$

$$\text{Back-propagation: } \boldsymbol{\delta}^{(m)} = \phi'(\mathbf{a}^{(m)}) * \left(\boldsymbol{\delta}^{(m+1)} (\mathbf{W}^{(m+1)})^T \right)$$

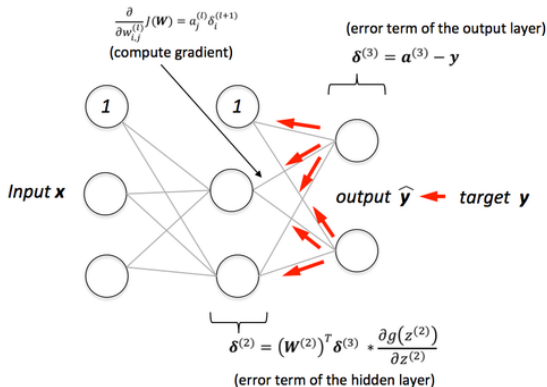
$$\text{v.s. Forward-propagation: } \mathbf{x}^{(m+1)} = \phi(\mathbf{a}^{(m)}) = \mathbf{x}^{(m)} \mathbf{W}^{(m+1)}$$

where $*$ is the element-wise multiplication.

Backpropagation: error term δ (cont)

$$x = \phi(a) = \frac{1}{1 + e^{-a}} \rightarrow \frac{\partial x}{\partial a} = \phi'(a) = \frac{e^{-a}}{(1 + e^{-a})^2} = \phi(a)(1 - \phi(a)) = x(1 - x)$$

$$x = \phi(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} \rightarrow \frac{\partial x}{\partial a} = \phi'(a) = (1 - \phi^2(a)) = (1 - x^2)$$



Backpropagation: Summary and benefit

Procedure

- 1 Forward-propagate \mathbf{x} 's and \mathbf{a} 's: $\mathbf{x}^{(L+1)} = \phi(\mathbf{a}^{(L)} = \mathbf{x}^{(L)}\mathbf{W})$
- 2 Evaluate $\delta^{(L)} = \hat{\mathbf{y}} - \mathbf{y}$ at the output units
- 3 Backpropagate δ 's: $\delta^{(m)} = \phi'(\mathbf{a}^{(m)}) * (\delta^{(m+1)}(\mathbf{W}^{(m+1)})^T)$
- 4 Obtain derivatives: $\frac{\partial J_n}{\partial W_{ij}^{(m)}} = \delta_j^{(m)} X_{ni}^{(m-1)}, \quad \frac{\partial J}{\partial W_{ij}^{(m)}} = \sum \delta_j^{(m)} X_{ni}^{(m-1)}$

Benefits

Numerical differentiation:

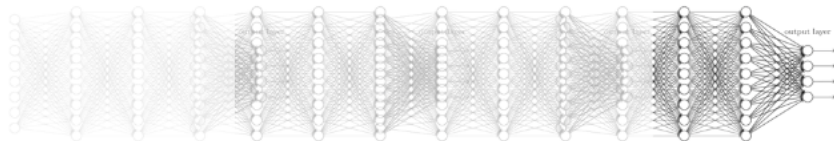
- $\frac{\partial J}{\partial W_{ij}} = \frac{J(W_{ij}+\varepsilon) - J(W_{ij}-\varepsilon)}{2\varepsilon}$
- require $O(N_w^2)$ operations for the total number of weights N_w .
- can be used to validate the backpropagation

Back-propagation:

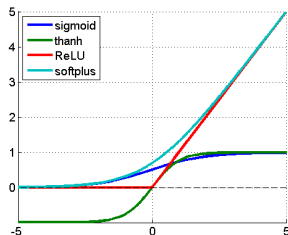
- require $O(N_w)$ operations.

The 2nd AI winter (1986-2006)

- Even the backpropagation is not enough for 'deep' layers as the gradient vanishes.



- Other 'better' activation functions used: $\phi(x) =$



Sigmoid (logistic) $\frac{1}{1+e^{-x}}$

Tanh $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

ReLU (Rectified linear unit)
 $\max(0, x)$

Softplus $\log(1 + e^x)$

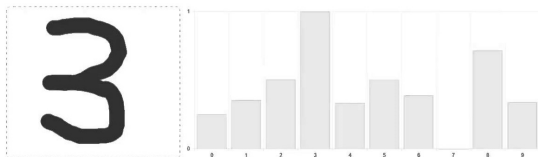
ReLU or **Softplus** activations have non-vanishing gradient.

Soft-max classification function

- In multi-class classification, we need to *normalize* the outcome for the last unit to the probability on each class for a given sample.
- Given the output \mathbf{a} , we generalize the logistic function:

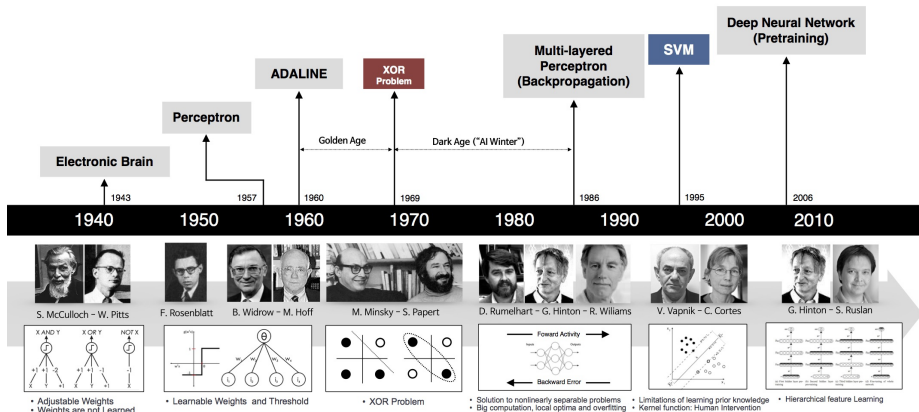
$$P(y = i|\mathbf{a}) = \phi_{\text{softmax}}(\mathbf{a}) = \frac{e^{a_i}}{\sum_k e^{a_k}} \quad \left(\sum_i P(y = i|\mathbf{a}) = 1 \right)$$

- The probability is largest for the largest value of a_i (vs argmax)
- The logistic function is a special case with $a_0 = 0$: $\phi(a) = e^a / (e^0 + e^a)$.
- In NN for multi-class classification, soft-max is applied to the last units to match the 'one-hot' encoding of the response variable.



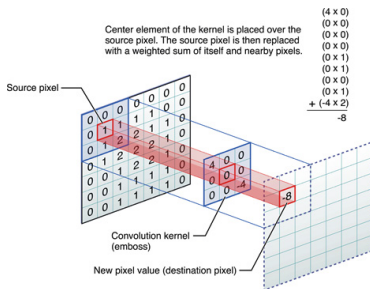
Timeline of NN

From Andrew Beam's Deep Learning 101



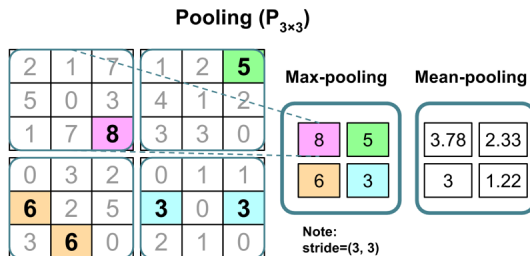
Convolutional Neural Network (CNN)

- Apply kernel/convolution matrix to image to create feature maps.
- For various image processing kernels, see [Wikipedia](#)
- The original image is “padded” with zero to avoid boundary loss.
- The layer is locally connected (vs fully connected), so the weight matrix is sparse.



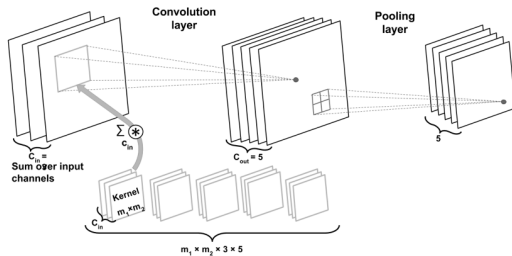
CNN: Pooling/Subsampling

- Either **max** or **mean** applied in non-overlapping way.
- Captures local invariance, robust to noise.
- Reduces feature sizes (e.g. pixel), so avoid overfitting.
- Predetermined operation. No need for learning weights.

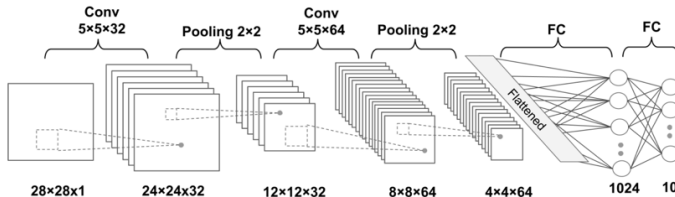


CNN: PML example

- Convolution + pooling unit



- Tensor dimensions



Dropout (in words)

- Deep neural nets with a large number of parameters are very powerful machine learning systems. However, overfitting is a serious problem in such networks. Large networks are also slow to use, making it difficult to deal with overfitting by combining the predictions of many different large neural nets at test time.
- Dropout is a technique for addressing this problem. The key idea is to **randomly drop units (along with their connections)** from the neural network during training. This prevents units from **co-adapting** too much.
- During training, dropout samples from an exponential number of different “thinned” networks. At test time, it is easy to approximate the effect of averaging the predictions of all these thinned networks by simply using a single unthinned network that has smaller weights. This significantly reduces overfitting and gives **major improvements over other regularization methods**.
...

Dropout (in graph)

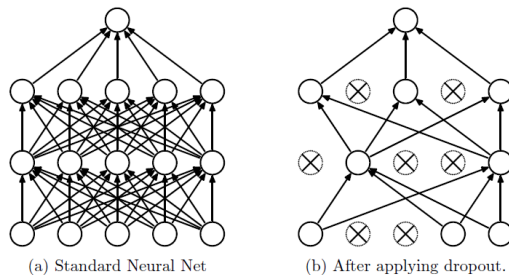


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

Co-adaptation

- In $3 \times 2 \times 1$ network, suppose

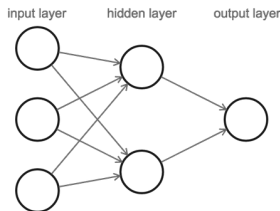
$$W^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -2 & 2 \end{bmatrix}, \quad W^{(2)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

- Assuming linear activation $\phi(x) = x$,

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} W^{(1)} W^{(2)} = 2x_1$$

Only x_1 matters (x_2 and x_3 cancel each other.)

- A typical type of over-fitting in deep learning.



Training

- If p is the probability of keeping unit ($1 - p$: dropout),

$$a = (e_1 x_1 w_1 + \cdots + e_n x_n w_n) / p,$$

where $e_i = 1$ or 0 with probability p or $1 - p$.

- e_i is fixed during one epoch of weight update.
- $p = 1/2$ works best.

Test/Prediction

- Use all units without dropout:

$$a = x_1 w_1 + \cdots + x_n w_n,$$

which is understood as the average over the all possible (2^n) cases.