

FNN

Feed forward neural network.

- (1980- ANN

2006 Auto-encoder

2012 Relu.

Softmax layer: output a probability for each class.

- multiple classes & one-hot encoding for labels.

→ preprocessing = < learn faster: SGD converges faster.

• min-max normalization:  $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$   
(when we know the max)

• standardization:  $z_i = \frac{x_i - \mu}{\sigma} \sim N(0, 1)$

- overfitting: • Dropout  
• Regularization

- parameters = - Depth: < 1000

- Width: 5-30

- Connectivity: usually fully connected

- Activation function: - provide non-linearity.

- Relu (overcome Vanishing gradient).

- Loss function = - (Classification: cross entropy)  $= - \sum_{i=1}^K y_i \log \hat{y}_i$

- Regression: squared loss:  $\sum_{i=1}^K (y_i - \hat{y}_i)^2$

## - Training of FNN.

- Mini-batch: divide data into minibatches (iterations)
  - Go through the whole sample (epoch)
  - Update weights after each iteration.

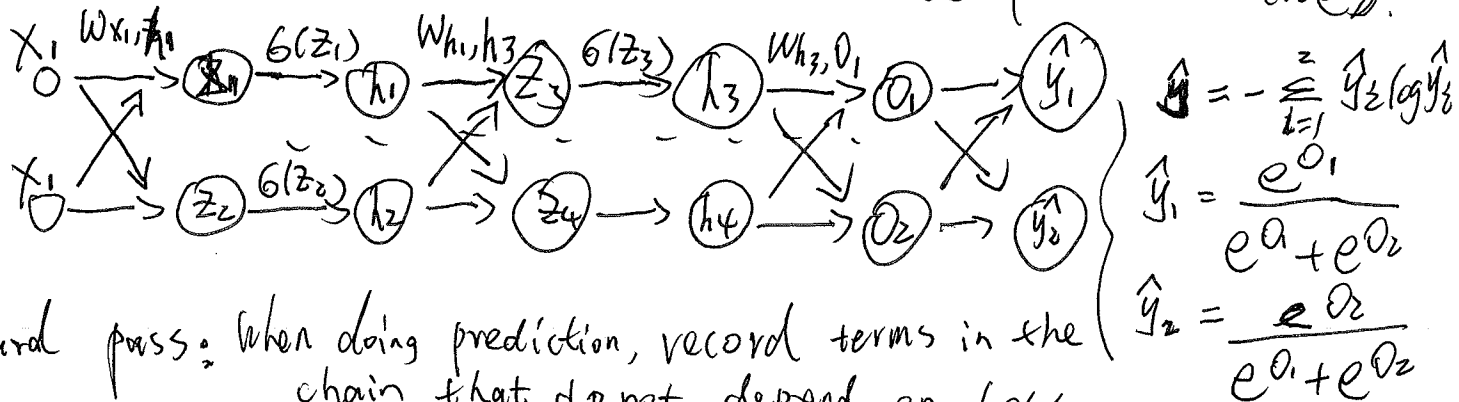
• SGD:  $W_i \leftarrow W_i - \eta \cdot \frac{\partial L}{\partial W_i}$

- learning rate:
  - Constant
  - $\eta \leftarrow \eta / \sqrt{i+1}$  (AdaGrad),  $i$  - epoch.

• weight-based  $\eta_w = \frac{\eta}{\sum \frac{\partial L}{\partial W_i}^2}$   $g_i = \frac{\partial L}{\partial W_i}$

- Dropout:
  - only in training ( $p$ )
  - in testing, reweight  $\times (1-p)$
  - "ensemble" models.

- Backward propagating (update the weights using chain rule for derivative).



\* Vanishing gradient ....

sequence data.  
RNN - LSTM

NN with memory

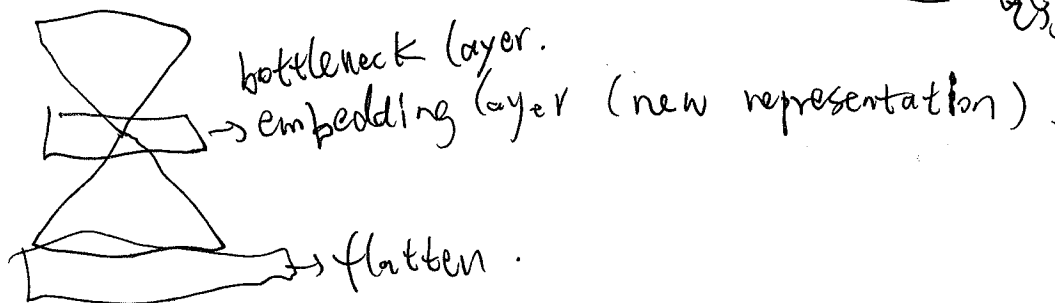
pca, SVD (linear)

- Auto-encoder.

- ~~unsupervised~~ supervised learning
- non-linear compression. (embedding).

• used to initialize parameters for NN (Early time).

• Trained using RBM:  $E(V, h) = - \sum_i b_i V_i - \sum_j b_j h_j - \sum_{i,j} V_i h_j W_{ij}$

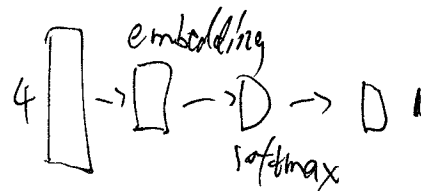


- Word embedding. (vector)

- New representation of the original data with certain meaning (semantics).

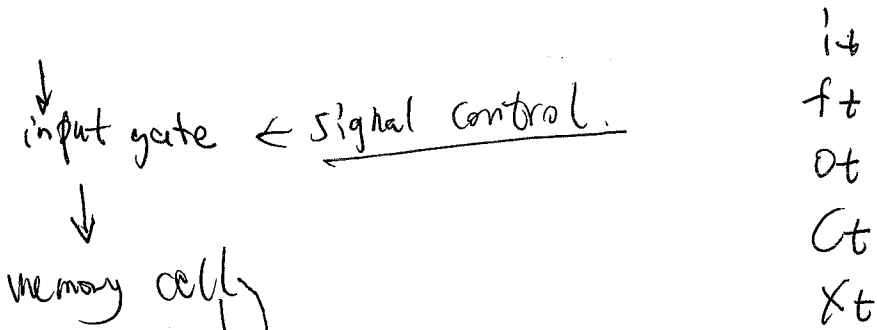
• CBOW : 5 words, with the middle as label.

(more popular) SKIP-gram : 5 words, with the middle word as feature, others as labels



LSTM.

— LSTM. memory to store previous NN output.



signal control, & output gate      forget gate ← signal control

$$\begin{cases} i_t = \sigma(W_i[h_{t-1}, x_t] + b_i) \\ f_t = \sigma(W_f[h_{t-1}, x_t] + b_f) \\ o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \end{cases}$$

$$\tilde{C}_t = \tanh[W_c[h_{t-1}, x_t] + b_c], \quad C_t = i_t \cdot \tilde{C}_t + f_t \cdot C_{t-1}$$

$$\boxed{h_t = o_t \cdot \tanh(C_t)} \rightarrow \text{final output.}$$

- CNN**
- capture local patterns - subsampling
  - less parameters (weight sharing).

= Structure: data → convolution → max-pooling → flatten → FFNN

↓  
 repeat  
 ↓  
 local patterns are smaller  
 ↓  
 2nd convolution patterns in different regions. ↓ subsampling  
 does not change the object.

- Data augmentation (generate more samples, ~~avoid~~ <sup>reduce</sup> overfitting).

- rotation
- random cropping
- add random noises.

objects are still there.

- Convolutional layer:

filter ("matrix dot product") → feature map

• Zero padding: ~~the~~ feature map size = original data size.

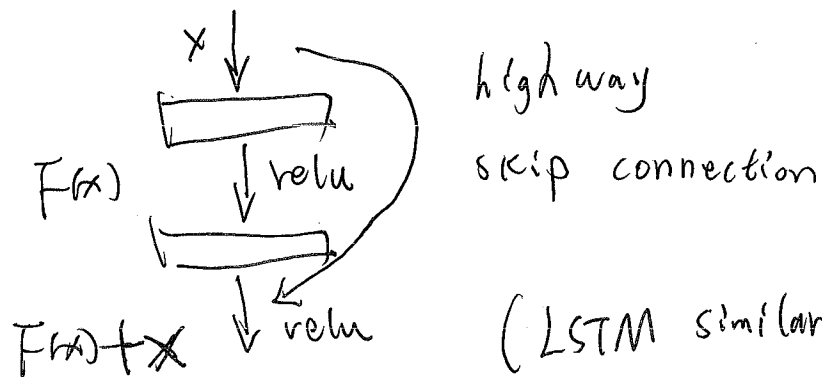
= Max-pooling: new representation, smaller input data size.

{ filter size: 5-13  
 # : sample/para.  
 stride: 1

CNN variants (~~original~~ normal is called AlexNet).

- Inception Net: irregular <sup>filter sizes</sup>  
branches

- Residual Net: deeper (overcome vanishing gradient).



(LSTM similar if  $C_t = 1 \cdot C_{t-1} + 0 \cdot C_t$ )

- Style transferring (styles are also patterns)

$$L_{total} = \alpha L_{content} + \beta L_{style}.$$

- Devise (Multi-model model).

source word  $\xrightarrow{\text{skip-gram}}$  embedding vector

image  $\xrightarrow{\text{CNN} \rightarrow \text{embed}}$  embedding vector  
repeat { convolution  
          maxpooling

> combine using similar metric.