

An example of how we learn XOR using FFNN with hidden layer  
=> Hidden layer mapping our data to a separable space.

## Gradient-Based Learning

- cost functions

- Learning conditional distribution with Maximum Likelihood

Relationship of likelihood  $P(X|\theta)$ , with the cost function  $-\ln P(X|\theta)$

- Learning Conditional Statistics

which means we often learn a statistics instead of a distribution of  $P(y|x, \theta)$

- Output Units

Discussion of different type of output units

- Linear output for Gaussian Output Distribution

- Sigmoid units for Bernoulli Distribution

Discussion of some interesting property of sigmoid function when we combine it with logarithm.

- Softmax units for Multinoulli Distribution

- other types of output units

Here discussion focus on we can also model mixture Gaussian with FFNN or predict/inference variance using FFNN - by modifying the output of the neural network.

## Hidden Units

Talking about how we design Hidden Units (Considering Gradients and non-linearity)

- Rectified Units and their generation

Keep the linearity of positive values, some generation (max out) also

keep linearity.

- Logistic Sigmoid and Hyperbolic Tangent (Tanh)

For FFNN networks, sigmoid function may be not a good choice, some times tanh is more preferable cause it has more steep shape.

- Other Hidden Units

## - Architecture design

- Universal Approximation properties and depth

Universal Approximation Theorem: MLP can theoretically simulate all functions but deeper neural network saves the parameters.

- Other considerations

## - Back-propagation and other differentiation algorithms

- Computational Graphs

Fundamental data structure of backward and forward calculation.

- Chain Rule of Calculus

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} \quad \text{for scalar}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i} \quad \text{for vector}$$

$$\nabla_x z = \sum_j (\nabla_x y_j) \frac{\partial z}{\partial y_j} \quad \text{for tensor}$$

- Recursively applying chain rule to get Back prop.

Defining computational process  $\underbrace{u^1, u^2, \dots, u^i, \dots, u^n}_{\text{input node}} \quad \frac{u^{(n)}}{\text{output}}$

$$\text{Forward: } A^{(i)} \leftarrow \{u^{(j)} \mid j \in \text{Pa}(u^{(i)})\}$$

$$u^{(i)} \leftarrow f^{(i)}(A^{(i)})$$

$$\text{Backward: } \text{grad-table}[u^{(n)}] = 1$$

For  $j=n-1$ , down to 1 do

$$\frac{\partial u^{(n)}}{\partial u^{(i)}} = \sum_{j: j \in \text{Pa}(u^{(i)})} \boxed{\frac{\partial u^{(n)}}{\partial u^{(j)}}} \cdot \frac{\partial u^{(j)}}{\partial u^{(i)}}$$

using grad-table[ $u^{(i)}$ ]

- Back-Propagation Computation in fully connected MLP

Forward:  $a^{(k)} = b^{(k)} + w^{(k)} h^{(k-1)}$

$h^{(k)} = f(a^{(k)})$ . until we get  $h^{(L)}$ ,  $J = L(\hat{y}, y) + \lambda \Omega$

Backward:  $g \leftarrow \nabla g J = \nabla g L(\hat{y}, y)$

for  $k=L, \dots, 1$  do:

$g \leftarrow \nabla a^{(k)} J = g \odot f'(a^{(k)})$

compute activation gradient

Compute gradient for parameters  $\left\{ \begin{array}{l} \nabla b^{(k)} J = g + \lambda \nabla b^{(k)} \Omega(\theta) \\ \nabla w^{(k)} J = g h^{(k-1)} + \lambda \nabla w^{(k)} \Omega(\theta) \end{array} \right.$

$g \leftarrow \nabla h^{(k-1)} J = w^{(k)T} \cdot g$

- Symbol-to-Symbol Derivatives

Two different computing approach (Symbol-to-number / Symbol)

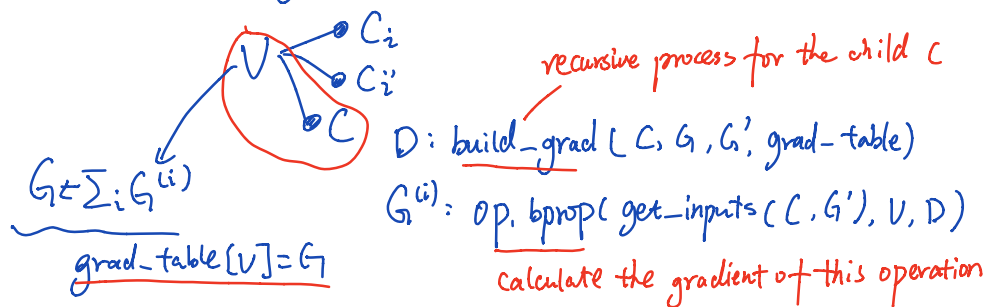
- General Back-Propagation

General framework: updating grad-table

$\text{grad-table}[z] \leftarrow 1$

for  $v$  in  $T$  do: build\_grad( $v, G, G', \text{grad-table}$ )

zoom into build\_grad:



- An example of back-prop MLP

- Complications

- Differentiate outside deep learning community.
  - Back-Prop is a kind of automatic differentiation
  - Consideration of order of vector-matrix multiplication
  - Back-prop cannot utilize (generally) math tricks due to its generality.