

Perceptron

ក្នុងអត្ថបទនេះយើងនឹងលើកយក Perceptron Algorithm មកបង្ហាញ ។ Perceptron គឺជា Algorithm មួយដែលត្រូវបានណែនាំជាដំបូងដោយអ្នកស្រាវជ្រាវអាមេរិកគឺលោក Alexander Rosenblatt នៅឆ្នាំ១៩៥៧ ។ ទោះបីជា Perceptron Algorithm ជាវិធីសាស្ត្រចាស់ក្តី ប៉ុន្តែជាត្រូវបានគេស្គាល់ថាជាប្រភពដើមនៃវិធីសាស្ត្រគណនាក្នុងម៉ូដែល Neural Network ឬ Deep Learning ដែលកំពុងរីកដុះដាលយ៉ាងសកម្មនាសម័យនេះ ។ ហេតុនេះ ការសិក្សាអំពី Perceptron Algorithm អាចជួយឱ្យយើងងាយស្រួលក្នុងការឈានទៅសិក្សាអំពី Neural Network ឬ Deep Learning ។

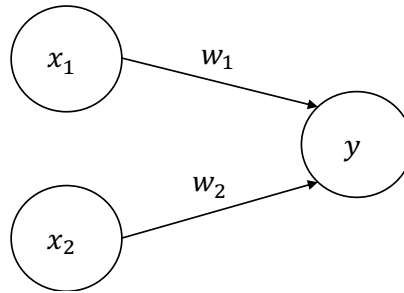
1. អំពី Perceptron

Perceptron ឬហៅថាម៉ូដែលណឺរ៉ូនសិប្បនិម្មិត (artificial neuron) ជាម៉ូដែលគណិតវិទ្យាមួយដែលទទួលសញ្ញាណ (signal) ឬ ធាតុចូល (input) ច្រើន និងផ្តល់នូវលទ្ធផល (output) មួយ ។ សញ្ញាណដែលទទួលនៅទីនេះអាចប្រៀបបានជាចរន្តអគ្គិសនីឬសញ្ញាណព័ត៌មានដូចដែលណឺរ៉ូននៃប្រព័ន្ធប្រសាទរបស់ការវាស់ទទួលដែរ ។ ប៉ុន្តែសញ្ញាណឬធាតុចូលក្នុង Perceptron កត់យកតម្លៃធម្មតាពេលគឺបញ្ជូនបន្តឬមិនបញ្ជូនបន្តដោយតម្លៃ (1 ឬ 0) ។

រូបទី១ខាងក្រោមបង្ហាញអំពី Perceptron ដែលទទួលសញ្ញាណឬធាតុចូល (input) 2 ។ រង្វង់ដែលមានក្នុងរូបហៅថាណឺរ៉ូន (neuron) ឬ node ។ x_1, x_2 ជាសញ្ញាណឬធាតុចូល ឯ y គឺជាលទ្ធផលនៃណឺរ៉ូននោះ ។ w_1, w_2 ជាទម្ងន់ផ្ទាល់នៃណឺរ៉ូនចំពោះធាតុចូលនីមួយៗ ពេលគឺតម្លៃដែលបង្ហាញនូវកម្រិតឥទ្ធិពលនៃធាតុចូលនីមួយៗទៅលើលទ្ធផល ។ តម្លៃនៃទម្ងន់កាន់តែធំបង្ហាញពីកម្រិតសំខាន់នៃធាតុចូលនោះ ។

នៅពេលដែលណឺរ៉ូនមួយទទួលបាននូវសញ្ញាណឬធាតុចូល នោះផលគុណរវាងតម្លៃនៃធាតុចូលនោះនិងតម្លៃនៃទម្ងន់ផ្ទាល់របស់ណឺរ៉ូននោះ (w_1x_1, w_2x_2) ត្រូវបានគណនា ។ លទ្ធផលដែលណឺរ៉ូននោះត្រូវផ្តល់គឺអាស្រ័យនឹងផលបូកនៃគ្រប់ធាតុចូលទាំងអស់ ដោយកំណត់តាមលក្ខខណ្ឌខាងក្រោម ។ លក្ខខណ្ឌនេះគឺ ណឺរ៉ូននឹងបញ្ចេញលទ្ធផល 1 បើផលបូកនៃផលគុណធាតុចូលនិងទម្ងន់របស់វាធំជាងតម្លៃនៃកម្រិតកំណត់របស់ណឺរ៉ូន ឬ បញ្ចេញលទ្ធផល 0 បើតូចជាង ។ តម្លៃនៃកម្រិតកំណត់របស់ណឺរ៉ូនហៅថា Threshold ។ នៅទីនេះ θ ជាតម្លៃ Threshold ។

$$y = \begin{cases} 1 & (w_1x_1 + w_2x_2 > \theta) \\ 0 & (w_1x_1 + w_2x_2 \leq \theta) \end{cases}$$



រូបទី១ Perceptronដែលទទួលសញ្ញាណប្រជាតិចូល (input) 2

ក្នុងការអនុវត្តភាគច្រើន គេច្រើនប្រើប្រាស់មធ្យមនៃខ្លឹមសារដែលមានកម្រិតកំណត់នៃ លីរ៉ូនស្មើ០ ដោយហៅតម្លៃ θ ដែលត្រូវបានបញ្ជូនទៅអង្គទី១ ($b = -\theta$) ដោយ bias ។

$$y = \begin{cases} 1 & (b + w_1x_1 + w_2x_2 > 0) \\ 0 & (b + w_1x_1 + w_2x_2 \leq 0) \end{cases}$$

2. បង្ហាញសៀគ្វីឡូស៊ីកងាយៗដោយម៉ូដែលPerceptron

ដើម្បីស្វែងយល់អំពីដំណើរការរបស់Perceptron នៅទីនេះយើងលើកយកសៀគ្វីឡូស៊ីក (Logic circuits) ងាយៗដូចជា AND gate, OR gate មកបង្ហាញដោយប្រើម៉ូដែលPerceptron ។

2.1. AND gate

ដូចដែលអ្នកបានដឹង AND gate ផ្តល់នូវលទ្ធផលអាស្រ័យនឹងតម្លៃភាពពិតនៃធាតុចូលរបស់ វាដូចក្នុងតារាងខាងក្រោម។

តារាងទី១ តម្លៃភាពពិតនៃAND gate

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1

ក្នុងករណីនេះ បើយើងយកPerceptronមួយដែលទទួលធាតុចូលពីរនិងមានតម្លៃទម្ងន់និង កម្រិតកំណត់ (threshold): $(w_1, w_2, \theta) = (0.5, 0.5, 0.8)$ នោះយើងអាចបង្ហាញAND gateដោយ Perceptronបាន។ ឧទាហរណ៍ ករណី $(x_1, x_2) = (1, 0)$ នោះ $w_1x_1 + w_2x_2 = 0.5 < 0.8$ ហេតុ នេះ $y = 0$ ។ ករណី $(x_1, x_2) = (1, 1)$ នោះ $w_1x_1 + w_2x_2 = 1 > 0.8$ ហេតុនេះ $y = 1$ ។

2.2. OR gate

OR gate ផ្តល់នូវលទ្ធផលអាស្រ័យនឹងតម្លៃភាពពិតនៃធាតុចូលរបស់វាដូចក្នុងតារាងខាងក្រោម។

តារាងទី២ តម្លៃភាពពិតនៃOR gate

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

ក្នុងករណីនេះ បើយើងយកPerceptronមួយដែលទទួលធាតុចូលពីរនិងមានតម្លៃទម្ងន់និងកម្រិតកំណត់ (threshold) : $(w_1, w_2, \theta) = (0.5, 0.5, 0.2)$ នោះយើងអាចបង្ហាញOR gateដោយPerceptronបាន។ ឧទាហរណ៍ ករណី $(x_1, x_2) = (1, 0)$ នោះ $w_1x_1 + w_2x_2 = 0.5 > 0.2$ ហេតុនេះ $y = 1$ ។ ករណី $(x_1, x_2) = (0, 0)$ នោះ $w_1x_1 + w_2x_2 = 0 < 0.2$ ហេតុនេះ $y = 0$ ។

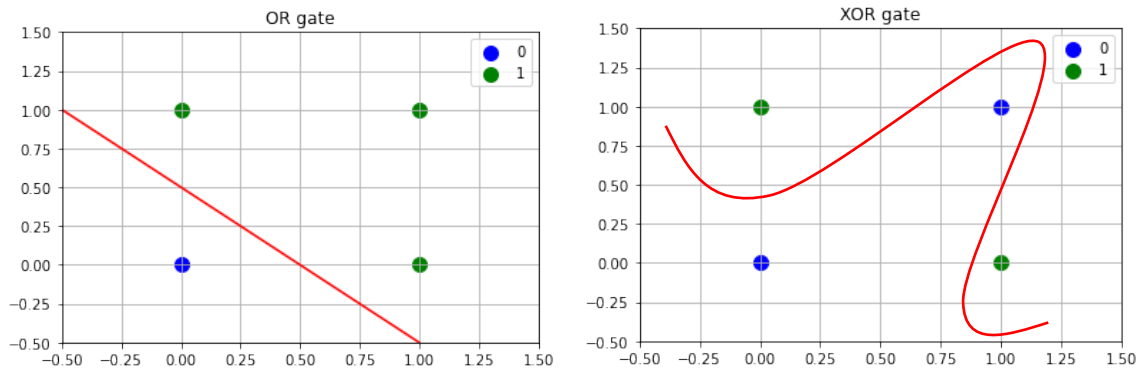
ដូចដែលបង្ហាញខាងលើ ដោយប្រើPerceptronយើងអាចបង្ហាញសៀគ្វីឡូស៊ីកងាយៗបាន។ ចំណុចសំខាន់នៅទីនេះគឺ ចំពោះAND gate, OR gate ដែលមានទម្រង់សៀគ្វីឡូស៊ីកផ្សេងៗគ្នា ក៏Perceptronដែលយើងប្រើមានគោលគំនិតឬទម្រង់តែមួយមិនប្រែប្រួលឡើយ។ អ្វីដែលខុសគ្នាគឺតម្លៃនៃប៉ារ៉ាម៉ែត្រ (ទម្ងន់និងកម្រិតកំណត់របស់វា) តែប៉ុណ្ណោះ។ ពោលគឺដោយប្រើទម្រង់នៃម៉ូដែលតែមួយយើងអាចបង្ហាញទម្រង់នៃសៀគ្វីឡូស៊ីកដែលជាគ្រឹះនៃសៀគ្វីអេឡិចត្រូនិចនានាបានដោយគ្រាន់តែកែសម្រួលតម្លៃនៃប៉ារ៉ាម៉ែត្ររបស់វាតែប៉ុណ្ណោះ។

3. ព្រំដែនសមត្ថភាពនៃPerceptron

យើងឃើញថាPerceptron អាចបង្ហាញ AND gate, OR gateបានយ៉ាងងាយ។ បន្តទៅនេះយើងនឹងពិនិត្យលើករណីនៃ XOR gate។

តារាងទី៣ តម្លៃភាពពិតនៃXOR gate

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0



រូបទី២ ប្រៀបធៀបករណី OR gate និង XOR gate

មុននឹងឈានទៅមើលករណី XOR gate យើងបង្ហាញព្រមគ្នាជាមួយករណី OR gate ដោយប្រើក្រាបដូចរូបទី២ ។ ក្នុងករណី OR gate យើងអាចកំណត់តម្លៃប៉ារ៉ាម៉ែត្ររបស់ Perceptron បានដូចឧទាហរណ៍ក្នុងចំណុច 2. ដែលលទ្ធផល 1 ឬ 0 អាចបែងចែកដាច់ពីគ្នាបានដោយបន្ទាត់ត្រង់មួយបាន។ ផ្ទុយពីនេះ ករណី XOR gate យើងមិនអាចកំណត់បន្ទាត់ត្រង់ដើម្បីបែងចែកករណីលទ្ធផល 1 ឬ 0 បានឡើយ។ ពោលគឺមិនអាចកំណត់ប៉ារ៉ាម៉ែត្រណាដែលអាចឱ្យ Perceptron បង្ហាញ XOR gate បានទេ។

ទិន្នន័យដូចក្នុងករណី OR gate ហៅថាទិន្នន័យដែលអាចបែងចែកលីនេអ៊ែរបាន (linear seperable) ឯទិន្នន័យដូចក្នុងករណី XOR gate ហៅថាទិន្នន័យដែលមិនអាចបែងចែកលីនេអ៊ែរបាន (linear non-seperable) ។ ហេតុនេះ យើងអាចនិយាយបានថា Perceptron ដែលមានទម្រង់ដូចណែនាំក្នុងចំណុច 1. មិនអាចប្រើជាម៉ូដែលសម្រាប់ទិន្នន័យមិនអាចបែងចែកលីនេអ៊ែរបានឡើយ។

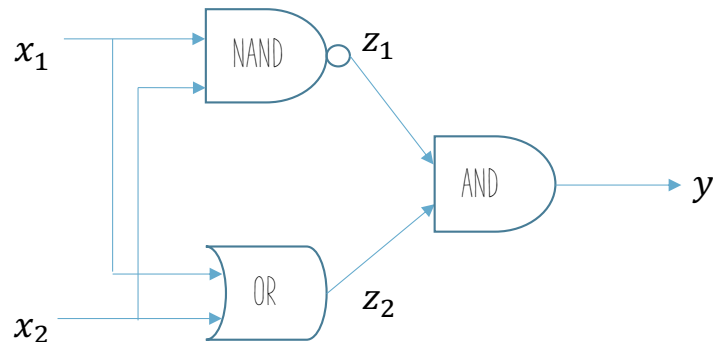
4. Perceptron ច្រើនថ្នាក់ (Multilayer Perceptron)

ដូចដែលបានពិនិត្យខាងលើ Perceptron ទម្រង់ធម្មតាមិនអាចធ្វើការពណ៌នា XOR gate បានប្រសើរឡើយ។ ដើម្បីស្វែងយល់ពីដំណោះស្រាយតាមរយៈការបង្កើនចំនួនថ្នាក់នៃ Perceptron យើងពិនិត្យលើការបង្ហាញទម្រង់ XOR gate ដោយប្រើបង្គុំនៃសៀគ្វីឡូស៊ីកមូលដ្ឋាន AND gate, OR gate, NAND gate ។

4.1. បង្ហាញ XOR gate ដោយប្រើ AND gate, OR gate, NAND gate

នៅទីនេះយើងនឹងមិនធ្វើការបកស្រាយលំអិតអំពីលក្ខណៈនៃសៀគ្វីឡូស៊ីកមូលឡើយ ប៉ុន្តែតាមពិតទៅ XOR gate អាចបង្ហាញដោយប្រើ AND gate, OR gate, NAND gate បានដោយធ្វើតំណភ្ជាប់ដូចរូបទី៣ ។ ដោយសារតែ AND gate, OR gate, NAND gate អាចបង្ហាញដោយប្រើ

Perceptron ទម្រង់ធម្មតាដូចក្នុងចំណុច1. 2.ខាងលើបាន ហេតុនេះ គំនិតសំខាន់ដែលយើងអាចសិក្សាពីចំណុចនេះគឺថា យើងអាចផ្គុំPerceptronធម្មតាជាច្រើនថ្នាក់ដើម្បីបង្ហាញXOR gateបាន។



រូបទី៣ ការបង្ហាញXOR gate ដោយប្រើAND gate, OR gate, NAND gate

4.2. ការបង្កើតPerceptronងាយៗនិងច្រើនថ្នាក់ដោយប្រើPython

```
import numpy as np
```

4.2.1. ករណីទម្រង់ធម្មតា(មិនប្រើទម្រង់Bias)

```
def AND(x1, x2):
    (w1,w2,theta) = (0.5, 0.5, 0.8)
    t = w1*x1 + w2*x2
    y = 1 if t > theta else 0
    return y
```

```
def OR(x1, x2):
    (w1,w2,theta) = (0.5, 0.5, 0.2)
    t = w1*x1 + w2*x2
    y = 1 if t > theta else 0
    return y
```

```
def NAND(x1, x2):
    (w1,w2,theta) = (-0.5, -0.5, -0.8)
    t = w1*x1 + w2*x2
    y = 1 if t > theta else 0
    return y
```

4.2.2. ករណីទម្រង់ប្រើBias

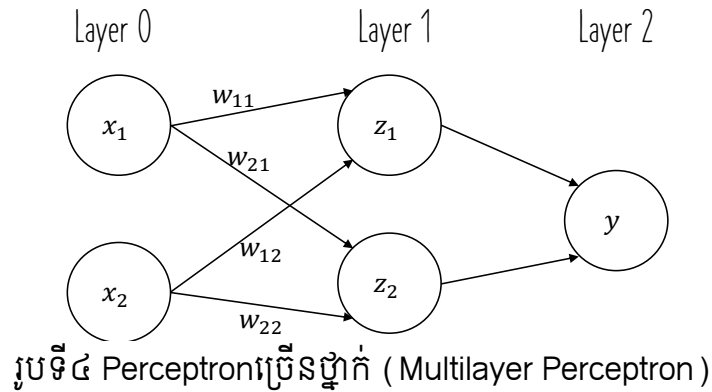
```
def AND(x1, x2):  
    x_vec = np.array([x1,x2])  
    w = np.array([0.5,0.5])  
    b = -0.8  
    t = np.sum(w*x_vec) + b  
    y = 1 if t > 0 else 0  
    return y
```

```
def OR(x1, x2):  
    x_vec = np.array([x1,x2])  
    w = np.array([0.5,0.5])  
    b = -0.2  
    t = np.sum(w*x_vec) + b  
    y = 1 if t > 0 else 0  
    return y
```

```
def NAND(x1, x2):  
    x_vec = np.array([x1,x2])  
    w = np.array([-0.5,-0.5])  
    b = 0.8  
    t = np.sum(w*x_vec) + b  
    y = 1 if t > 0 else 0  
    return y
```

4.2.3. XOR gate

```
def XOR(x1,x2):  
    s1 = NAND(x1,x2)  
    s2 = OR(x1,x2)  
    y = AND(s1,s2)  
    return y
```



ទម្រង់នៃXOR gateដោយបង្កើននៃAND gate, OR gate, NAND gateដូចក្នុងរូបទី៣ អាចប្រដូចបានជាករណីមួយនៃករណីទូទៅរបស់បង្កើននៃPerceptronទម្រង់ធម្មតាច្រើនបញ្ចូលគ្នាដូចក្នុងរូបទី៤ ។ ទម្រង់បែបនេះហៅថា Multilayer Perceptron ដែលក្នុងអត្ថបទនេះនិងបន្តបន្ទាប់យើងកំណត់ហៅថា Perceptronច្រើនថ្នាក់ ។

ដូចក្នុងករណីXOR gate ដែរ ការប្រើMultilayer Perceptronអាចឱ្យយើងបង្ហាញទិន្នន័យដែលមិនអាចបែងចែកលីនេអ៊ែរដោយម៉ូដែលPerceptronបាន ។ ការបង្កើនចំនួនថ្នាក់ (Layer) ក្នុង Perceptronនឹងជួយបង្កើនសមត្ថភាពរបស់វាក្នុងការពណ៌នាលក្ខណៈរបស់ទិន្នន័យដែលកាន់តែស្មុគស្មាញបានដែលនេះជាគោលគំនិតគ្រឹះក្នុង Artificial Neural Network ឬ Deep Learningដែលយើងនឹងលើកមកសិក្សាក្នុងអត្ថបទក្រោយៗ ។

5. ការកំណត់ប៉ារ៉ាម៉ែត្រនៃPerceptron ១ថ្នាក់ (Perceptron Algorithm)

ចំពោះPaceptronទម្រង់ធម្មតា១ថ្នាក់ យើងអាចកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រពីទិន្នន័យដែលមានបានដោយអនុវត្តតាមវិធីសាស្ត្រខាងក្រោម ។ នៅទីនេះ $\mathbf{w}^{(t)}$ សម្គាល់តម្លៃនៃប៉ារ៉ាម៉ែត្រនៅដំណាក់កាលផ្លាស់ប្តូរទី t , η សម្គាល់តម្លៃនៃកម្រិតផ្លាស់ប្តូរប៉ារ៉ាម៉ែត្រដែលហៅថា learning rate ។

(ជំហានទី១) កំណត់តម្លៃដើមនៃប៉ារ៉ាម៉ែត្រ \mathbf{w} ដោយ ០ ឬតម្លៃបំណែងចែកចៃដន្យ

(ជំហានទី២) ចំពោះទិន្នន័យ (training data) (x_i, y_i) អនុវត្តជំហានខាងក្រោមរហូតដល់គ្មាន

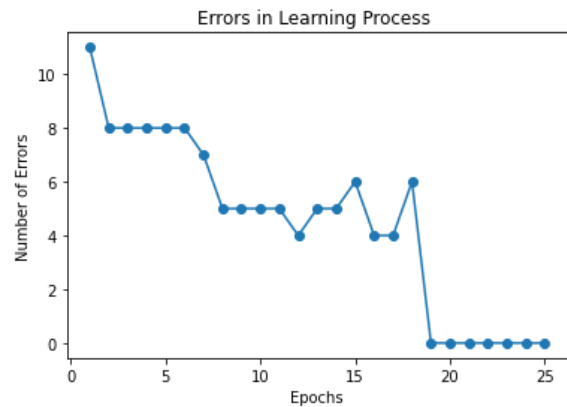
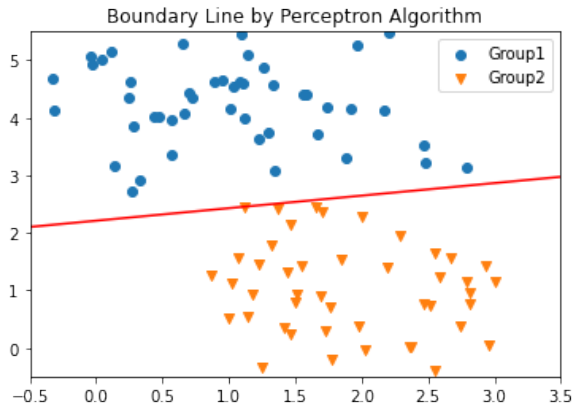
កំហុសក្នុងការប៉ាន់ស្មានតម្លៃ ពោលគឺ $\Delta \mathbf{w}^{(t)} = \mathbf{0}$

ក. គណនាលទ្ធផលនៃណឺរ៉ូន $\hat{y}_i = x_i^T \mathbf{w}^{(t)}$ ដោយប្រើតម្លៃប៉ារ៉ាម៉ែត្របច្ចុប្បន្ន

ខ. ផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្រ

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \Delta \mathbf{w}^{(t)}$$

$$\Delta \mathbf{w}^{(t)} = \eta (y_i - \hat{y}_i) x_i$$



រូបទី៥ ការធ្វើចំណាត់ថ្នាក់ក្រុមទិន្នន័យ២ក្រុមដោយPerceptron Algorithm

```
import numpy as np

class Perceptron(object):

    def __init__(self, learning_rate=0.01, iteration_number=100, random_state=1):

        self.eta = learning_rate

        self.n_iteration = iteration_number

        self.random_state = random_state

    def fit(self, X, y):

        rand_gen = np.random.RandomState(self.random_state)

        self.w = rand_gen.normal(loc=0.0, scale=0.1, size=1+X.shape[1])

        self.errors = []

        XP = np.ones((X.shape[0],X.shape[1]+1))

        XP[:, :-1]=X

        for t in range(self.n_iteration):

            error = 0

            for xi,yi in zip(XP,y):

                delta = self.eta * (yi - self.predict(xi))

                self.w += delta * xi

                error += int(delta != 0.0)

            self.errors.append(error)

        return self

    def predict(self, X):

        z = X@self.w

        return np.where(z >= 0.0, 1, 0)
```

Feedforward Neural Network (FNN)

ក្នុងអត្ថបទមុន យើងបានសិក្សាអំពី Perceptron ដែលជាម៉ូដែលអាចបែងចែកទិន្នន័យបំណែងចែកលីនេអ៊ែរបានយ៉ាងងាយដោយការកំណត់តម្លៃប៉ារ៉ាម៉ែត្រសមស្រប។ លើសពីនេះ ក្នុងករណីទិន្នន័យមិនអាចបែងចែកលីនេអ៊ែរ ការបង្កើនចំនួនថ្នាក់នៃ Perceptron ត្រូវបានប្រើប្រាស់។ ម៉ូដែលបែបនេះហៅថា Multilayer Perceptron ។ ដោយការភ្ជាប់ណឺរ៉ូន (node) ច្រើនបន្តគ្នាជាច្រើនថ្នាក់ដែលស្រដៀងគ្នានឹងទម្រង់នៃប្រព័ន្ធប្រសាទរបស់ការ៉ាវីសផងនោះម៉ូដែលបែបនេះក៏ត្រូវបានគេហៅថា Artificial Neural Network ផងដែរ។

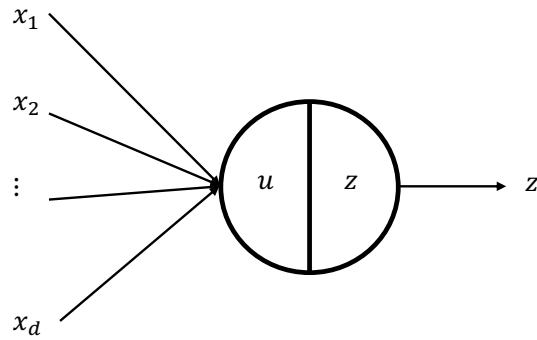
1. លទ្ធផលបញ្ជូនបន្តនៃណឺរ៉ូន

Feedforward Neural Network (FNN) គឺជាទម្រង់មួយនៃ Artificial Neural Network ដែលមានណឺរ៉ូន (node) ច្រើនតម្រៀបគ្នាជាថ្នាក់និងភ្ជាប់គ្នា និងគ្នារវាងថ្នាក់នៅជាប់បន្តបន្ទាប់គ្នាដោយទម្ងន់ផ្ទាល់ខ្លួន។ សញ្ញាណប្បធាតុចូលនៃ FNN ត្រូវបានបញ្ជូនពីផ្នែកថ្នាក់ធាតុចូល (input layer) ទៅកាន់ផ្នែកនៃថ្នាក់លទ្ធផល (output layer) តាមទិសតែមួយ។ លទ្ធផលដែលបញ្ចេញដោយណឺរ៉ូននៅថ្នាក់លទ្ធផលត្រូវបានគណនាដូចក្នុងករណី Perceptron ទម្រង់ធម្មតាដែរ ប៉ុន្តែនៅទីនេះលទ្ធផលមិនគ្រាន់តែប្រៀបធៀបផលបូកនៃធាតុចូលនិងកម្រិតកំណត់ (threshold) នៃណឺរ៉ូនប៉ុណ្ណោះទេ តែអនុគមន៍មិនលីនេអ៊ែរត្រូវបានអនុវត្តលើលទ្ធផលនៃផលបូកនោះដើម្បីកំណត់នូវលទ្ធផលដែលត្រូវបញ្ជូនបន្ត។ អនុគមន៍ដែលអនុវត្តលើលទ្ធផលនៃផលបូកធាតុចូលនេះហៅថា អនុគមន៍សកម្ម (activation function) ។ យើងនឹងធ្វើការបកស្រាយលម្អិតអំពីអនុគមន៍សកម្មនៅចំណុចបន្ទាប់។

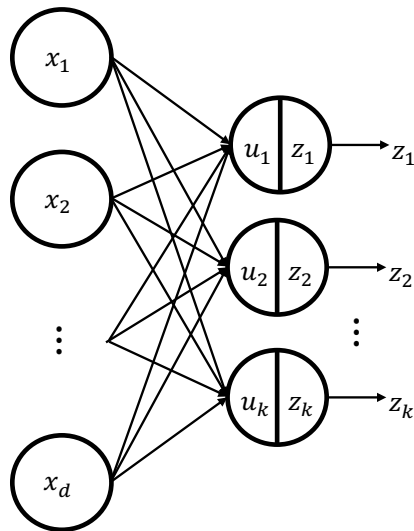
នៅពេលដែលធាតុចូល x_1, x_2, \dots, x_d ត្រូវបានបញ្ជូនមកកាន់ណឺរ៉ូន (រូបទី១) ដែលមានទម្ងន់ផ្ទាល់នៃធាតុចូលរៀងគ្នា w_1, w_2, \dots, w_d នោះ ផលបូកនៃធាតុចូលសរុបកំណត់ដោយ u និងលទ្ធផលបញ្ជូនបន្តនៃណឺរ៉ូននោះត្រូវបានកំណត់ដោយ z ដូចទម្រង់ខាងក្រោម។ នៅទីនេះ b ហៅថា bias ។

$$u = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

$$z = f(u) = f(w_1x_1 + w_2x_2 + \dots + w_dx_d + b)$$



រូបទី១ ធាតុចូលនិងលទ្ធផលនៃណឺរ៉ូនមួយ



រូបទី២ FNNទម្រង់២ថ្នាក់

ក្នុងករណីទម្រង់២ថ្នាក់ដូចក្នុងរូបទី២ សញ្ញាណត្រូវបានបញ្ជូនបន្តបន្ទាប់។ សន្មតថានៅថ្នាក់ទី១មានណឺរ៉ូនចំនួន d និង នៅថ្នាក់ទី២មានណឺរ៉ូនចំនួន k នោះលទ្ធផលបញ្ជូនបន្តនៃណឺរ៉ូននៅថ្នាក់លទ្ធផលទី២ត្រូវបានបង្ហាញក្នុងទម្រង់ខាងក្រោម ($i = 1, 2, \dots, k$) ។

$$u_i = \sum_{j=1}^d w_{ij}x_j + b_i$$

$$z_i = f(u_i)$$

យើងក៏អាចបង្ហាញជាទម្រង់វ៉ិចទ័រនិងម៉ាទ្រីសដូចខាងក្រោមផងដែរ។

$$\mathbf{u} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$\mathbf{z} = \mathbf{f}(\mathbf{u})$$

$$\mathbf{u} = \begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix}, \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}, \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_k \end{bmatrix}, \mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_k \end{bmatrix}, \mathbf{f}(\mathbf{u}) = \begin{bmatrix} f(u_1) \\ \vdots \\ f(u_k) \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & \cdots & w_{1d} \\ \vdots & \ddots & \vdots \\ w_{k1} & \cdots & w_{kd} \end{bmatrix}$$

2. អនុគមន៍សកម្ម Activation Function

អនុគមន៍សកម្ម (activation function) គឺជាអនុគមន៍ដែលគ្រប់គ្រងលើកម្រិតនៃការបញ្ជូនបន្តនូវលទ្ធផលរបស់ណឺរ៉ូននីមួយៗ។ ជាទូទៅអនុគមន៍សកម្មមានទម្រង់ជាអនុគមន៍មិនលីនេអ៊ែរ កើនជាប់ខាត។ មានទម្រង់ជាច្រើនត្រូវបានប្រើជាអនុគមន៍សកម្មដូចជា អនុគមន៍Sigmoid, អនុគមន៍Tanh, អនុគមន៍Softmax, អនុគមន៍ReLU (rectified linear function) ជាដើម។

អនុគមន៍Sigmoid $\sigma(\cdot)$ មានដែនកំណត់លើសំណុំចំនួនពិត $(-\infty, \infty)$ និងយកសំណុំរូបភាពលើចន្លោះបើក $(0, 1)$ ។ អនុគមន៍បែបនេះអាចផ្តល់នូវលទ្ធផលដែលយើងអាចបកស្រាយបានជាតម្លៃប្រូបាបនៃការបញ្ជូនបន្តឬមិនបញ្ជូនបន្ត $(1/0)$ ។

$$f(u) = \sigma(u) = \frac{1}{1 + e^{-u}}$$

អនុគមន៍Tanh $\tanh(\cdot)$ មានដែនកំណត់លើសំណុំចំនួនពិត $(-\infty, \infty)$ និងយកសំណុំរូបភាពលើចន្លោះបើក $(-1, 1)$ ។ អនុគមន៍បែបនេះមានលក្ខណៈស្រដៀងនឹងអនុគមន៍Sigmoidដែរ ដែលអាចផ្តល់នូវលទ្ធផលដែលមានបម្រែបម្រួលតិចតួចក្បែរតម្លៃថេរនៅពេលដែលតម្លៃនៃធាតុចូលធំ ខ្លាំងដល់កម្រិតណាមួយ និងប្រែប្រួលខ្លាំងនៅពេលដែលធាតុចូលមានតម្លៃក្បែរ 0។

$$f(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

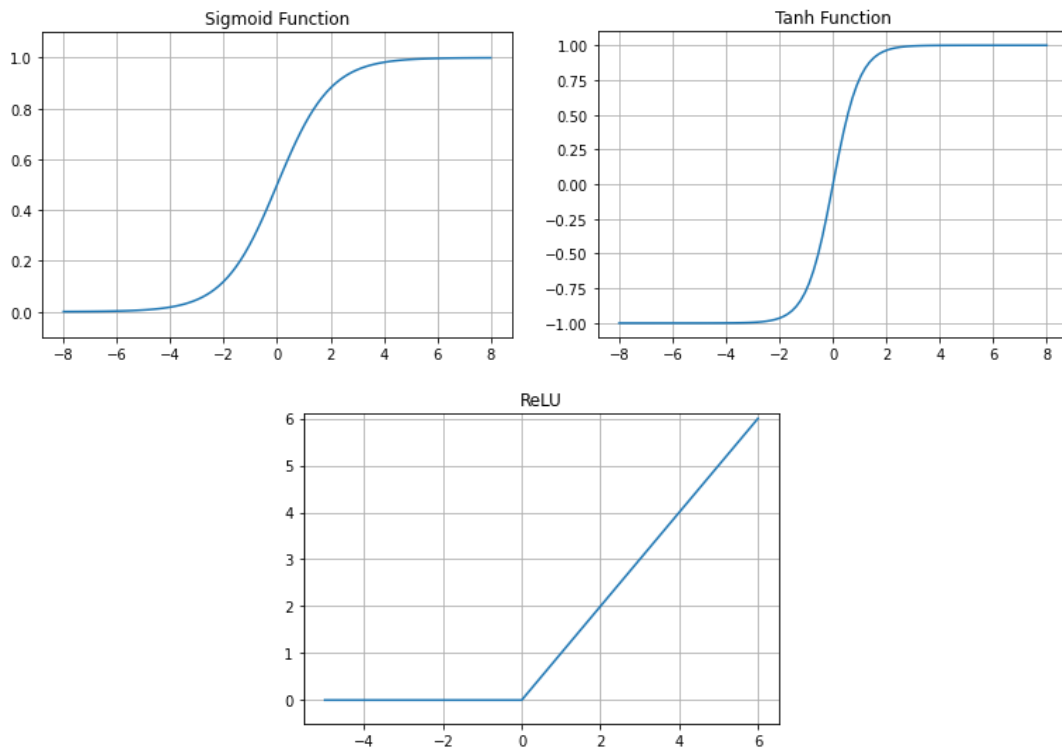
អនុគមន៍Softmax $\text{softmax}(\cdot)$ មានដែនកំណត់លើវ៉ិចទ័រចំនួនពិត \mathbb{R}^d និងយកសំណុំរូបភាពលើចន្លោះបើក $(0, 1)^d$ ដែលមានផលបូកគ្រប់កំប៉ូសង់ស្មើ 1។ អនុគមន៍បែបនេះអាចផ្តល់នូវលទ្ធផលដែលយើងអាចបកស្រាយបានជាតម្លៃប្រូបាបនៃលទ្ធផលដែលអាចចេញជា d ប្រភេទផ្សេងៗគ្នាបាន។ ក្នុងករណី $d = 2$ អនុគមន៍នេះសមមូលនឹងអនុគមន៍Sigmoid។ ចំពោះ $\mathbf{u} = (u_1 \cdots u_d)^\top$

$$\mathbf{f}(\mathbf{u}) = \text{Softmax}(\mathbf{u}) = (\text{Softmax}(u_1) \cdots \text{Softmax}(u_d))^\top$$

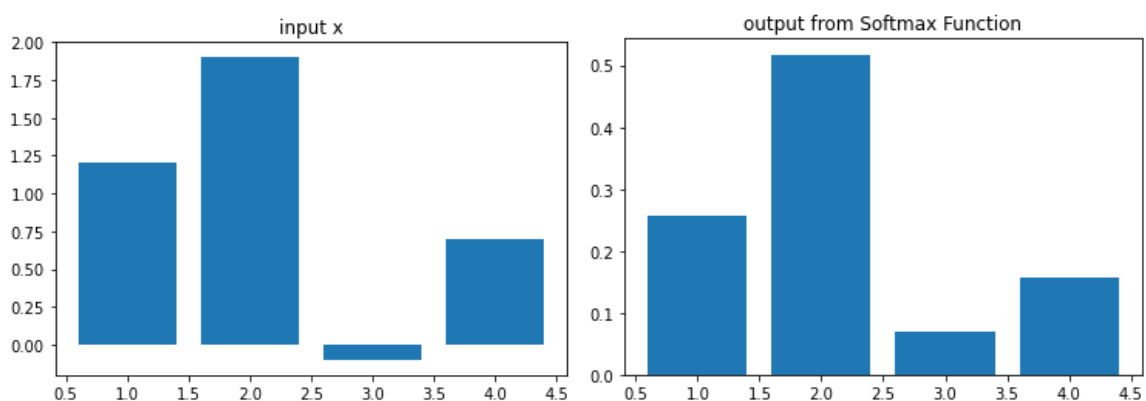
$$\text{Softmax}(u_i) = \frac{e^{u_i}}{\sum_{j=1}^d e^{u_j}}$$

អនុគមន៍ rectifier linear: $\text{ReLU}(\cdot)$ មានដែនកំណត់លើសំណុំចំនួនពិត $(-\infty, \infty)$ និងយកសំណុំរូបភាពលើចន្លោះបើក $(0, \infty)$ ។ នៅពេលដែលធាតុចូលមានតម្លៃតូចជាងឬស្មើសូន្យ លទ្ធផលបញ្ជូនបន្តត្រូវបានកំណត់ដោយ 0 និងបញ្ជូនបន្តនូវតម្លៃដូចធាតុចូលដែលបើធាតុចូលមានតម្លៃធំជាងឬស្មើសូន្យ ។ អនុគមន៍បែបនេះមានលក្ខណៈស្រដៀងនឹងអនុគមន៍ដឺក្រេទី១ (លីនេអ៊ែរ) ដែលអាចប្រើក្នុងករណីប៉ាន់ស្មានទាំងម៉ូឌែលលីនេអ៊ែរនិងមិនលីនេអ៊ែរបានល្អ ។

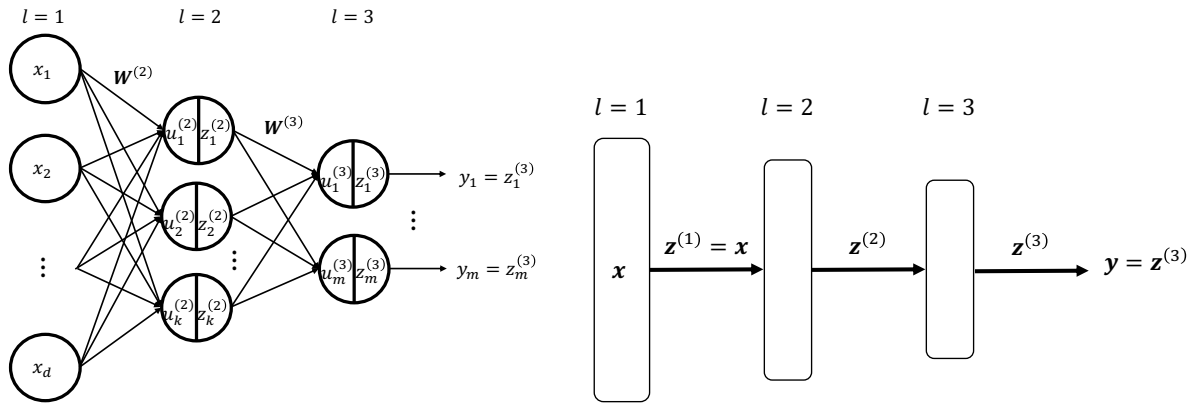
$$f(u) = \text{ReLU}(u) = \max(u, 0)$$



រូបទី៣ ក្រាបតាងអនុគមន៍សកម្ម sigmoid, tanh, ReLU



រូបទី៤ ទំនាក់ទំនងរវាងធាតុចូលនិងលទ្ធផលដោយ Softmax Function



រូបទី៥ Multilayer Network

3. បណ្តាញច្រើនថ្នាក់ (Multilayer Network)

នៅចំណុចនេះយើងពិនិត្យលើករណីមួយដែលមានច្រើនថ្នាក់ដូចក្នុងរូបទី៥។ ព័ត៌មាន (សញ្ញាណឬធាតុចូល) ត្រូវបានបញ្ជូនតាមលំដាប់លំដោយពីថ្នាក់នៅខាងឆ្វេងទៅស្តាំ។ នៅទីនេះ យើងកំណត់ហៅថ្នាក់នីមួយៗដោយ $l = 1, 2, 3, \dots$ ។ ក្នុងរូបខាងលើថ្នាក់ $l = 1$ ពោលគឺថ្នាក់នៅខាងឆ្វេងបំផុតហៅថាថ្នាក់នៃធាតុចូល (input layer), ថ្នាក់ $l = 2$ ហៅថាថ្នាក់នៃធាតុកណ្តាល (internal layer, hidden layer), ថ្នាក់ $l = 3$ ពោលគឺថ្នាក់នៅខាងស្តាំបំផុត ហៅថាថ្នាក់នៃលទ្ធផល (output layer)។ លទ្ធផលនៅថ្នាក់នីមួយៗអាចសរសេរជាទម្រង់ដូចខាងក្រោម។

$$\mathbf{u}^{(2)} = \mathbf{W}^{(2)}\mathbf{x} + \mathbf{b}^{(2)} \quad , \quad \mathbf{z}^{(2)} = \mathbf{f}(\mathbf{u}^{(2)})$$

$$\mathbf{u}^{(3)} = \mathbf{W}^{(3)}\mathbf{z}^{(2)} + \mathbf{b}^{(3)} \quad , \quad \mathbf{z}^{(3)} = \mathbf{f}(\mathbf{u}^{(3)})$$

ជាទូទៅ លទ្ធផលនៅថ្នាក់កណ្តាលកំណត់ដោយ

$$\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)}\mathbf{z}^{(l)} + \mathbf{b}^{(l+1)} \quad , \quad \mathbf{z}^{(l+1)} = \mathbf{f}(\mathbf{u}^{(l+1)})$$

និង លទ្ធផលនៅថ្នាក់លទ្ធផលចុងក្រោយកំណត់ដោយ

$$\mathbf{y} \equiv \mathbf{z}^{(l+1)}$$

ដូចដែលបានឃើញក្នុងទម្រង់គណនាខាងលើ ក្នុង FNN សញ្ញាណត្រូវបានបញ្ជូនបន្តបន្ទាប់ដោយការគណនាក្នុងរបៀបដូចគ្នាពីមួយថ្នាក់ទៅមួយថ្នាក់។ នេះគឺជាប្រភពដែលម៉ូដែលនេះត្រូវបានហៅថា feedforward ។ ទម្ងន់ផ្ទាល់នៃឈ្លីនចំពោះធាតុចូលតាមថ្នាក់នីមួយៗ $\mathbf{W}^{(l)}$ និង bias $\mathbf{b}^{(l)}$ ត្រូវបានហៅជារួមថាជា ប៉ារ៉ាម៉ែត្រនៃបណ្តាញ។ ក្នុងអត្ថបទនេះនិងបន្តបន្ទាប់យើងកំណត់ហៅដោយងាយនូវ បណ្តាញដែលមានប៉ារ៉ាម៉ែត្រ (ហៅជារួម) \mathbf{w} និងធាតុចូល \mathbf{x} ដោយ $\mathbf{y}(\mathbf{x}; \mathbf{w})$ ។

4. ការកំណត់ទម្រង់ណឺរ៉ូននៅថ្នាក់លទ្ធផលនិងអនុគមន៍លម្អៀង

4.1. ប្រភេទនៃបញ្ហានិងវិធីសាស្ត្ររៀន (Learning)

ដូចដែលបានរៀបរាប់ពីអត្ថបទមុននិងចំណុចខាងលើ បណ្តាញដែលបង្ហាញដោយទម្រង់នៃអនុគមន៍ច្រើនអថេរ $y(x; w)$ នឹងប្រែប្រួលនៅពេលដែលប៉ារ៉ាម៉ែត្ររបស់វាត្រូវបានផ្លាស់ប្តូរ។

ការជ្រើសរើសប៉ារ៉ាម៉ែត្របានល្អ នឹងធ្វើឱ្យបណ្តាញ (network) អាចបង្ហាញនូវអនុគមន៍ឬបញ្ហាដែលមានបានល្អប្រសើរ។

សន្មតថា អនុគមន៍ឬបញ្ហាជាគោលដៅដែលយើងចង់បង្ហាញដោយ Neural Network មិនប្រែប្រួលសណ្ឋានខាងក្នុងរបស់វាឡើយ ហើយទទួលបានធាតុចូល x និងបញ្ជូនចេញនូវលទ្ធផល t ។ គូនៃទិន្នន័យបែបនេះជាច្រើនត្រូវបានផ្តល់ឱ្យ $\{(x_1, t_1), \dots, (x_N, t_N)\}$ ។ នៅក្នុងអត្ថបទនេះ និងអត្ថបទបន្តបន្ទាប់ គូនីមួយៗហៅថាជាគម្រួសសម្រាប់រៀន (training sample) ហើយសំណុំទាំងមូលហៅថាសំណុំទិន្នន័យសម្រាប់រៀន (training data) ។

ដោយធ្វើការកែសម្រួលនិងកំណត់នូវតម្លៃប៉ារ៉ាម៉ែត្រ យើងអាចធ្វើការបង្ហាញទំនាក់ទំនងដែលមានក្នុងទិន្នន័យឡើងវិញបានដោយបណ្តាញ (network) របស់យើង។ ពេលគឺចំពោះគម្រួសសម្រាប់រៀន (x_n, t_n) នីមួយៗ យើងចង់កំណត់នូវប៉ារ៉ាម៉ែត្រណាដែលធ្វើឱ្យយើងអាចទទួលបាន $y(x_n; w)$ ដែលមានតម្លៃជិតបំផុតនឹង t_n ។ ដំណើរការកំណត់រកនូវប៉ារ៉ាម៉ែត្រដោយប្រើសំណុំទិន្នន័យសម្រាប់រៀនបែបនេះសន្មតហៅថាជាដំណើរការរៀន (learning process) ។

ហេតុនេះការប្រៀបធៀបរវាងតម្លៃលទ្ធផល $y(x_n; w)$ ដែលផ្តល់ដោយបណ្តាញនិងតម្លៃ t_n ត្រូវបានធ្វើឡើង។ ដើម្បីបង្ហាញពីកម្រិតជិតគ្នានៃតម្លៃទាំងពីរយើងកំណត់រង្វាស់សម្រាប់វាស់កម្រិតនេះ។ រង្វាស់នេះយើងសន្មតហៅថាជា អនុគមន៍កម្រិតលម្អៀង (error function, loss function) ។ ការកំណត់ប្រភេទនៃអនុគមន៍កម្រិតលម្អៀងខុសគ្នាទៅតាមប្រភេទចំណោទបញ្ហាដែលយើងចង់ដោះស្រាយ។

តារាងទី១ ប្រភេទនៃអនុគមន៍សកម្មនិងកម្រិតលម្អៀងតាមប្រភេទចំណោទ

ប្រភេទចំណោទបញ្ហា	អនុគមន៍សកម្មនៅថ្នាក់លទ្ធផល	អនុគមន៍កម្រិតលម្អៀង
តម្រូវតម្រង់ (regression)	Identity function $y(x) = x$	ផលបូកការេនៃលម្អៀង
ចំណាត់ថ្នាក់២ក្រុម	Sigmoid function	Cross entropy (2 classes)
ចំណាត់ថ្នាក់ច្រើនក្រុម	Softmax function	Cross entropy

4.2. ចំណោទតម្រូវតម្រង់

ចំណោទតម្រូវតម្រង់ (Regression) គឺជាប្រភេទចំណោទដែលធ្វើការកំណត់នូវអនុគមន៍ដើម្បីបង្ហាញនូវទំនាក់ទំនងរវាងធាតុចូលនិងលទ្ធផលដែលមានទម្រង់ជាអថេរជាប់។ ហេតុនេះក្នុង

ករណីនៃចំណោទតម្រិតម្រង់ យើងកំណត់យកអនុគមន៍សកម្ម (activation function) នៅថ្នាក់លទ្ធផលនៃបណ្តាញ (FNN) ដោយអនុគមន៍ដែលផ្តល់នូវសំណុំរូបភាពដូចដែននៃលទ្ធផលរបស់សំណុំទិន្នន័យសម្រាប់រៀន។ ឧទាហរណ៍ ក្នុងករណីដែលសំណុំទិន្នន័យសម្រាប់រៀនមានតម្លៃនៃលទ្ធផលលើចន្លោះ $[-1, 1]$ នោះយើងជ្រើសយកអនុគមន៍ \tanh សម្រាប់ជាអនុគមន៍សកម្ម។ ក្នុងករណីដែលសំណុំទិន្នន័យសម្រាប់រៀនមានតម្លៃនៃលទ្ធផលលើចន្លោះ $(-\infty, \infty)$ នោះយើងជ្រើសយកអនុគមន៍ដែលផ្តល់តម្លៃដូចជាតុចូល Identity function សម្រាប់ជាអនុគមន៍សកម្ម។

ក្នុងករណីចំណោទតម្រិតម្រង់នេះដើម្បីប្រៀបធៀបកម្រិតជិតគ្នារវាងលទ្ធផលនៃបណ្តាញនិងតម្លៃលទ្ធផលនៃគម្រូសម្រាប់រៀន យើងប្រើផលបូកការេនៃតម្លៃលម្អៀង (sum of squared residuals) រវាងតម្លៃ $y(x_n; \mathbf{w})$ និង t_n ចំពោះគ្រប់គម្រូសម្រាប់រៀនក្នុងសំណុំទិន្នន័យសម្រាប់រៀនទាំងអស់។ ពោលគឺអនុគមន៍កម្រិតលម្អៀងសម្រាប់សំណុំទិន្នន័យសម្រាប់រៀនកំណត់ដោយ ។

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|t - y(x_n; \mathbf{w})\|^2$$

នៅទីនេះការដាក់មេគុណ $\frac{1}{2}$ គឺដើម្បីសម្រួលដល់ការគណនាដល់ការធ្វើដេរីវេនាពេលខាងមុខ។

គោលដៅរបស់យើងគឺកំណត់ប៉ារ៉ាម៉ែត្រនៃបណ្តាញយ៉ាងណាដើម្បីឱ្យអនុគមន៍កម្រិតលម្អៀងខាងលើនេះមានតម្លៃអប្បបរមា។

4.3. ចំណោទធ្វើចំណាត់ថ្នាក់២ក្រុម

ក្នុងចំណោទធ្វើចំណាត់ថ្នាក់២ក្រុម ទិន្នន័យធាតុចូល x នឹងត្រូវបែងចែកទៅក្នុងក្រុមមួយក្នុងចំណោមពីរក្រុម។ ឧទាហរណ៍ករណីធាតុចូលជារូបថតមួយសន្លឹក។ នៅពេលនោះក្រោយពីទាញយកលក្ខណៈសម្គាល់របស់រូបថតនោះជាទម្រង់វ៉ិចទ័ររួច បណ្តាញនឹងទទួលយកវ៉ិចទ័រនោះជាធាតុចូលរួចធ្វើការបែងចែកថាជារូបមុខមនុស្សឬមិនមែន។ ក្នុងករណីនេះ លទ្ធផលនៃទិន្នន័យសម្រាប់រៀន t យកតម្លៃជាអថេរដាច់ $\{1(\text{មុខមនុស្ស}), 0(\text{មិនមែនមុខមនុស្ស})\}$ ។ បែបនេះចំណោទចំណាត់ថ្នាក់២ក្រុមក៏ជាចំណោទដែលទទួលធាតុចូល x និងប៉ាន់ស្មានលទ្ធផល t ដូចតម្រិតម្រង់ដែរគ្រាន់តែប្រភេទនៃតម្លៃលទ្ធផលជាអថេរដាច់។

ក្នុងករណីនេះ ដើម្បីប៉ាន់ស្មានតម្លៃលទ្ធផល យើងសិក្សាលើម៉ូដែលប្រូបាប ពោលគឺសិក្សាលើប្រូបាបដែលថាលទ្ធផល $t = 1$ នៅពេលធាតុចូល x ត្រូវបានទទួល $p(t = 1|x)$ ។ គំនិតនៅទីនេះគឺថា បើប្រូបាបនេះមានតម្លៃធំជាង 0.5 នោះយើងសន្និដ្ឋានថាលទ្ធផលគឺ $t = 1$ និង សន្និដ្ឋានថាលទ្ធផលគឺ $t = 0$ ក្នុងករណីផ្ទុយពីនេះ។

ដោយពិនិត្យលើគំនិតបែបនេះ អ្នកអាចនឹកឃើញដល់លក្ខណៈនៃអនុគមន៍ Sigmoid ដែលបានបង្ហាញខាងលើ។ បើយើងប្រើអនុគមន៍ Sigmoid ជាអនុគមន៍សកម្មសម្រាប់បណ្តាញ FNN $y(x; \mathbf{w})$ នោះយើងអាចបង្ហាញម៉ូដែលប្រូបាបខាងលើដោយ FNN បាន។

$$p(t = 1|x) \approx y(x; \mathbf{w})$$

ហេតុនេះដើម្បីកំណត់ប៉ារ៉ាម៉ែត្រនៃ FNN យើងអាចសិក្សាពីសំណុំអថេរសម្រាប់រៀន (training data) $\{(x_n, t_n)\}_{n=1}^N$ បានដោយកំណត់យកប៉ារ៉ាម៉ែត្រដែលធ្វើឱ្យបំណែងចែកប្រូបាប $p(t|x; \mathbf{w})$ មានភាពប្រហាក់ប្រហែលគ្នាបំផុតជាមួយនឹងរបាយនៃសំណុំអថេរសម្រាប់រៀន។ ក្នុងការកំណត់ប៉ារ៉ាម៉ែត្រនៃម៉ូដែលប្រូបាបបែបនេះ យើងហៅថា ការប៉ាន់ស្មានកម្រិតសាកសមបំផុតនៃទិន្នន័យ (Maximum Likelihood Estimation, MLE) ។

ដោយប្រើតម្លៃនៃ $p(t = 0|x; \mathbf{w})$, $p(t = 1|x; \mathbf{w})$ យើងអាចបង្ហាញ $p(t|x; \mathbf{w})$ ជាមួយតាមទម្រង់ខាងក្រោម។

$$p(t|x; \mathbf{w}) = p(t = 1|x; \mathbf{w})^t p(t = 0|x; \mathbf{w})^{1-t}$$

ដោយការសន្មតខាងលើ $p(t = 1|x) = y(x; \mathbf{w})$ នោះ $p(t = 1|x) = 1 - y(x; \mathbf{w})$ ។ ក្រោមការសន្មតនៃម៉ូដែលបែបនេះ ការប៉ាន់ស្មានកម្រិតសាកសមបំផុតនៃទិន្នន័យ MLE គឺជាការកំណត់នូវកម្រិតសាកសមនៃទិន្នន័យ (likelihood) សម្រាប់រៀនរបស់ប៉ារ៉ាម៉ែត្រ \mathbf{w} និងជ្រើសយកតម្លៃនៃ \mathbf{w} ណាដែលធ្វើឱ្យកម្រិតសាកសមនោះមានតម្លៃអតិបរមា។ កម្រិតសាកសមនៃទិន្នន័យសម្រាប់រៀន (training data) របស់ប៉ារ៉ាម៉ែត្រ \mathbf{w} ត្រូវបានកំណត់ដូចទម្រង់ខាងក្រោម។

$$L(\mathbf{w}) = \prod_{n=1}^N p(t_n|x; \mathbf{w}) = \prod_{n=1}^N \{y(x_n; \mathbf{w})\}^{t_n} \{1 - y(x_n; \mathbf{w})\}^{1-t_n}$$

ដើម្បីសម្រួលដល់ការធ្វើបរិមាណ យើងអនុវត្តអនុគមន៍លោការីតលើកន្សោមខាងលើ។ ការធ្វើបែបនេះមិនប៉ះពាល់ដល់អថេរភាព (ភាពកើនចុះ) នៃអនុគមន៍ឡើយ។ ក្នុងករណីនេះ អនុគមន៍កម្រិតលម្អៀងនៃចំណោទចំណាត់ថ្នាក់ ២ ក្រុមកំណត់ដោយ $E(\mathbf{w})$ ដូចខាងក្រោម។

$$E(\mathbf{w}) = -\log L(\mathbf{w}) = -\sum_{n=1}^N \{t_n \log y(x_n; \mathbf{w}) + (1 - t_n) \log(1 - y(x_n; \mathbf{w}))\}$$

ដូចដែលបានបង្ហាញខាងលើ អនុគមន៍ Sigmoid ត្រូវបានប្រើសម្រាប់ជាអនុគមន៍សកម្មក្នុងថ្នាក់លទ្ធផលនៃ FNN ។ ចំណុចនេះអាចបកស្រាយដូចខាងក្រោម។

ប្រូបាប $p(t = 1|x)$ អាចសរសេរជាទម្រង់ប្រូបាបមានលក្ខខណ្ឌដូចខាងក្រោម ។

$$p(t = 1|x) = \frac{p(x; t = 1)}{p(x; t = 0) + p(x; t = 1)}$$

ដោយយក

$$u \equiv \log \frac{p(x; t = 1)}{p(x; t = 0)}$$

នោះយើងបាន

$$p(t = 1|x) = \frac{1}{1 + \exp(-u)} = \sigma(u)$$

ពោលគឺ ម៉ូដែលប្រូបាប $p(t = 1|x)$ ដែលសិក្សាខាងលើសមមូលទៅនឹងអនុគមន៍ Sigmoid ។

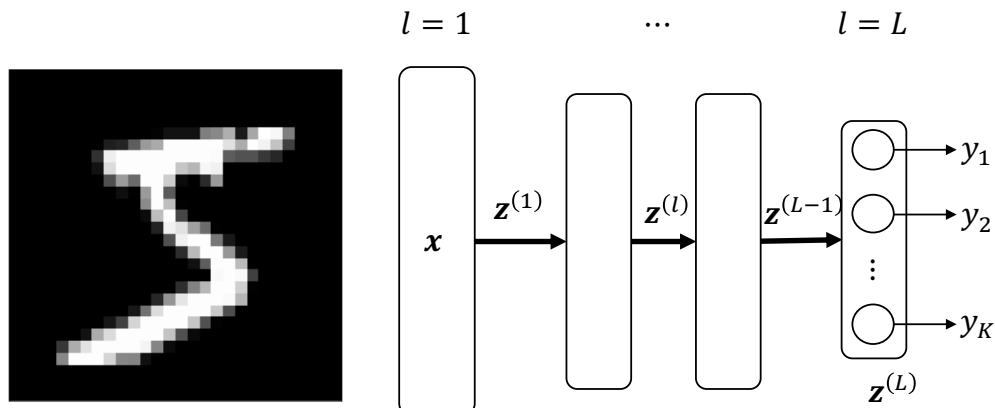
4.4. ចំណោទធ្វើចំណាត់ថ្នាក់ច្រើនក្រុម

ក្នុងចំណោទធ្វើចំណាត់ថ្នាក់ច្រើនក្រុម ទិន្នន័យធាតុចូល x នឹងត្រូវបែងចែកទៅក្នុងក្រុមមួយ ក្នុងចំណោមក្រុមមានកំណត់ច្រើន ។ ឧទាហរណ៍ករណីធាតុចូលជារូបថតនៃលេខសរសេរដោយដៃ មួយសន្លឹក ។ នៅពេលនោះក្រោយពីទាញយកលក្ខណៈសម្គាល់របស់រូបថតនោះជាទម្រង់វ៉ិចទ័ររួច បណ្តាញនឹងទទួលយកវ៉ិចទ័រនោះជាធាតុចូលរួចធ្វើការបែងចែកថាជាលេខណាមួយក្នុងចំណោម ០ ដល់ ៩ ។ ក្នុងករណីនេះ លទ្ធផលនៃទិន្នន័យសម្រាប់រៀន t យកតម្លៃជាអថេរដាច់ $\{0, 1, \dots, 9\}$ ។

ក្នុងករណីនេះ ការប្រើបណ្តាញ FNN សម្រាប់ចំណាត់ថ្នាក់ច្រើនក្រុម អាចធ្វើបានដោយ កំណត់យកថ្នាក់លទ្ធផលមានចំនួនណឺរ៉ូនស្មើនឹងចំនួននៃក្រុមដែលត្រូវបែងចែក ។ សន្មតថាចំនួន ក្រុមដែលត្រូវបែងចែកក្នុងចំណោទមាន K ។ នៅទីនេះ យើងប្រើបណ្តាញ FNN ដែលមាន L ថ្នាក់និង ចំនួនណឺរ៉ូននៅថ្នាក់លទ្ធផលមានចំនួន K ។ លទ្ធផលដែលផ្តល់ដោយណឺរ៉ូននីមួយៗក្នុងថ្នាក់លទ្ធផល អាចសរសេរដោយទម្រង់ខាងក្រោម ។

$$y_k \equiv z_k^{(L)} = \frac{\exp(u_k^{(L)})}{\sum_{i=1}^K \exp(u_i^{(L)})}$$

ដូចដែលបានបកស្រាយក្នុងចំណុចអនុគមន៍សកម្មខាងលើយើងអាចប្រើអនុគមន៍ Softmax ដើម្បីបម្លែងលទ្ធផលជាប្រូបាបនៃករណីបែងចែកចូលក្នុងក្រុមនីមួយៗ ។ យើងនឹងពិនិត្យលើភាពត្រឹម ត្រូវនៃការសន្មតនេះដោយប្រើម៉ូដែលប្រូបាបដូចករណី២ក្រុមដែរ ។



រូបទី៦ FNN ក្នុងចំណោមចំណាត់ថ្នាក់ច្រើនក្រុម (K-ក្រុម)

សន្មតថាថ្នាក់នីមួយៗនៃចំណោទខាងលើគឺ c_1, c_2, \dots, c_K លទ្ធផលនៃណឺរ៉ូន k នៅថ្នាក់ លទ្ធផលចុងក្រោយនៃបណ្តាញ FNN កំណត់ដោយ $y_k (= z_k^{(L)})$ ជាប្រូបាបនៃព្រឹត្តិការណ៍ដែលធាតុ ចូល x ត្រូវកំណត់ថានៅក្នុងក្រុម c_k ។

$$p(c_k|x) = y_k = z_k^{(L)}$$

ជាលទ្ធផល ធាតុចូល x ត្រូវកំណត់ថានៅក្នុងក្រុម c_k បើតម្លៃប្រូបាប $p(c_k|x)$ មានតម្លៃធំជាងគេក្នុង ចំណោមក្រុមទាំងអស់។

ក្នុងករណីចំណាត់ថ្នាក់ច្រើនក្រុមនេះ យើងកំណត់សរសេរលទ្ធផលពិត t_n នៃអថេរ x_n ដោយ ទម្រង់វ៉ិចទ័រ (one-hot vector) $t_n = [t_{n1} \ \dots \ t_{nk}]^T$ ។

ក្នុងករណីចំណាត់ថ្នាក់រូបថតលេខសរសេរដោយដៃខាងលើនោះ $K = 10$ ។
បើ x_n ជាលេខ 0 ដែលស្ថិតនៅក្នុងក្រុម c_1 នោះ $t_n = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]^T$ និង
បើ x_n ជាលេខ 5 ដែលស្ថិតនៅក្នុងក្រុម c_6 នោះ $t_n = [0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]^T$ ។

ដោយសន្មតសរសេរបែបនេះ យើងអាចបង្ហាញ $p(t|x)$ ជាអនុគមន៍ទម្រង់ខាងក្រោម។

$$p(t|x) = \prod_{k=1}^K p(c_k|x)^{t_k}$$

ហេតុនេះចំពោះសំណុំទិន្នន័យសម្រាប់រៀន $\{(x_n, t_n)\}_{n=1}^N$ កម្រិតសាកសមនៃទិន្នន័យ សម្រាប់រៀន (training data) របស់ប៉ារ៉ាម៉ែត្រ w ត្រូវបានកំណត់ដូចទម្រង់ខាងក្រោម។

$$L(w) = \prod_{n=1}^N p(t_n|x_n; w) = \prod_{n=1}^N \prod_{k=1}^K p(c_k|x)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K (y_k(x_n; w))^{t_{nk}}$$

ដើម្បីសម្រួលដល់ការធ្វើបរមាកម្ម យើងអនុវត្តអនុគមន៍លោការីតលើកន្សោមខាងលើ។
ក្នុងករណីនេះ អនុគមន៍កម្រិតលម្អៀងនៃចំណោទចំណាត់ថ្នាក់ច្រើនក្រុមកំណត់ដោយ $E(\mathbf{w})$
ដូចខាងក្រោម។ អនុគមន៍កម្រិតលម្អៀងបែបនេះហៅថា cross entropy ។

$$E(\mathbf{w}) = -\log L(\mathbf{w}) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_k(\mathbf{x}_n; \mathbf{w})$$

ការប្រើអនុគមន៍Softmax សម្រាប់ជាម៉ូដែលនៃបំណែងចែកច្រើនថ្នាក់នេះអាចបកស្រាយ
ដូចខាងក្រោម។

ប្រូបាបដែលធាតុចូល \mathbf{x} ត្រូវកំណត់ថានៅក្នុងក្រុម C_k អាចគណនាដោយ

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)}{\sum_{i=1}^K p(\mathbf{x}|C_i)}$$

នៅទីនេះដោយតាង $u_k = \log(p(\mathbf{x}, C_k))$ នោះ $p(\mathbf{x}, C_k) = \exp(u_k)$ ហេតុនេះ

$$p(C_k|\mathbf{x}) = \frac{\exp(u_k)}{\sum_{i=1}^K \exp(u_i)}$$

កន្សោម $p(C_k|\mathbf{x})$ ដែលទាញបានខាងលើនេះដូចគ្នាទៅនឹងអនុគមន៍Softmaxដែរ។ ដូច្នេះភាពត្រឹម
ត្រូវនៃការប្រើប្រាស់អនុគមន៍Softmax ជាអនុគមន៍សកម្មសម្រាប់ចំណោទចំណាត់ថ្នាក់ច្រើនក្រុម
ដោយFNNត្រូវបានផ្ទៀងផ្ទាត់។

Learning Process in Neural Network

ក្នុងអត្ថបទមុន យើងបានដឹងរួចមកហើយថាដោយការផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្រនៃFNN (Feedforward Neural Network) អាចឱ្យយើងបង្ហាញពីទំនាក់ទំនងផ្សេងៗនៃទិន្នន័យដែលយើងមាន។ ចំណុចសំខាន់ក្នុងការកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រគឺការស្វែងរកតម្លៃណាដែលសាកសមបំផុតក្នុងការបង្កើតបានជាបណ្តាញFNNដែលអាចពណ៌នាទំនាក់ទំនងក្នុងសំណុំទិន្នន័យសម្រាប់រៀន។ ការកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រដោយផ្អែកលើសំណុំទិន្នន័យសម្រាប់រៀន ហៅថា ដំណើរការរៀន (Learning Process) ។ មានវិធីសាស្ត្រច្រើនដែលត្រូវបានប្រើក្នុងដំណើរការរៀននៃFNN។ ក្នុងអត្ថបទនេះ យើងនឹងណែនាំវិធីសាស្ត្រកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រដោយវិធីគណនាច្រំដែលលើតម្លៃលេខតាមប្រមាណវិធីងាយៗគឺ Stochastic Gradient Descent (SGD) ។ ដើម្បីងាយស្រួលស្វែងយល់អំពីSGD ជាដំបូងយើងនឹងណែនាំអំពីគំនិត និងការគណនាក្នុងវិធីសាស្ត្រ Gradient Descent ជាមុន។

1. វិធីសាស្ត្រ Gradient Descent

ដូចដែលបានបង្ហាញក្នុងអត្ថបទមុន ប៉ារ៉ាម៉ែត្រដែលសាកសមបំផុតក្នុងការបង្កើតបានជាបណ្តាញFNNដែលអាចពណ៌នាទំនាក់ទំនងក្នុងសំណុំទិន្នន័យសម្រាប់រៀន គឺជាតម្លៃណាដែលធ្វើឱ្យអនុគមន៍កម្រិតលម្អៀងរវាងលទ្ធផលពីFNNនិងទិន្នន័យក្នុងសំណុំសម្រាប់រៀនតូចបំផុត។ ហេតុនេះគោលដៅរបស់យើងក្នុងដំណាក់កាលរៀននៃNeural Network គឺចង់កំណត់ប៉ារ៉ាម៉ែត្រនៃFNNដែលធ្វើឱ្យអនុគមន៍កម្រិតលម្អៀងមានតម្លៃតូចបំផុត ពោលគឺតម្លៃផលបូកការនែលម្អៀងក្នុងករណីនៃចំណោទតម្រៃតម្រង់ ឬ Cross Entropyក្នុងករណីចំណោទចំណាត់ថ្នាក់ក្រុម មានតម្លៃតូចបំផុត។

$$\text{Regression : } E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \| \mathbf{t} - \mathbf{y}(x_n; \mathbf{w}) \|^2$$

$$\text{Classification : } E(\mathbf{w}) = - \sum_{n=1}^N \{ y(x_n; \mathbf{w})^{t_n} + (1 - y(x_n; \mathbf{w}))^{1-t_n} \}$$

គោលគំនិតក្នុងGradient Descent គឺផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្របន្តិចម្តងៗទៅតាមទិសដៅដែលធ្វើឱ្យតម្លៃនៃអនុគមន៍កម្រិតលម្អៀងមានការថយចុះ។ អ្នកអាចធ្វើការប្រដូចវិធីនេះទៅនឹងការចុះជំរាលឬចុះពីទីភ្នំដោយរំកិលខ្លួនអ្នកបន្តិចម្តងៗទៅកាន់ទីដែលទាបជាងកន្លែងដែលអ្នកនៅ។ ពេលដែលអ្នករំកិលខ្លួនដល់ទីដែលលែងមានបម្រែបម្រួលនៃរយៈកម្ពស់ អ្នកអាចសន្និដ្ឋានបានថាអ្នកដល់ទីដែលទាបបំផុតហើយ។ ដូចគ្នានេះដែរ នៅក្នុងវិធីសាស្ត្រGradient Descent តាមលក្ខណៈគណិតវិទ្យានៃ gradient (តម្លៃដេរីវេនៃអនុគមន៍ត្រង់ចំណុចណាមួយ) តម្លៃgradientត្រង់ចំណុចណាមួយគឺជាតម្លៃមេគុណប្រាប់ទិសនៃខ្សែកោងត្រង់ចំណុចនោះហើយក៏ជាតម្លៃជំហាននៃបម្រែបម្រួលតម្លៃអនុគមន៍ពេលអ្នកធ្វើបម្រែបម្រួលលើអថេរមិនអាស្រ័យ។



រូបទី១ គំនិតក្នុង Gradient Descent

ពេលនេះ យើងពិនិត្យលើការគណនាក្នុងវិធីសាស្ត្រ Gradient Descent ។ យើងសិក្សា លើ ករណីសំណុំទិន្នន័យសម្រាប់រៀន $\mathcal{D} = \{(x_1, t_1), \dots, (x_N, t_N)\}$ និងអនុគមន៍ កម្រិត លម្អៀង $E(\mathbf{w})$ ពេលគឺ ក្នុងករណីនៃចំណោទតម្រៃតម្រង់

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|t - y(x_n; \mathbf{w})\|^2$$

និងក្នុងករណីនៃចំណោទចំណាត់ថ្នាក់ច្រើនក្រុម

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_k(x_n; \mathbf{w})$$

។ គោលដៅរបស់យើងគឺកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រ \mathbf{w} ដែលធ្វើអប្បបរមាកម្មលើ $E(\mathbf{w})$ ។

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w})$$

សន្មតថាអនុគមន៍នេះយកតម្លៃអប្បបរមាត្រង់ចំណុច $\mathbf{w}^* \in \mathbb{R}^M$ ។ វិធីសាស្ត្រ Gradient Descent អាចឱ្យយើងគណនាតម្លៃ (ប្រហែល) នៃ \mathbf{w}^* បានដោយចាប់ផ្តើមពីតម្លៃ $\mathbf{w}^{(0)}$

ណាមួយ រួចធ្វើការផ្លាស់ប្តូរតម្លៃនេះតាមការគណនាដូចខាងក្រោម។

$$\nabla E(\mathbf{w}) \equiv \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \left[\frac{\partial E}{\partial w_1} \quad \dots \quad \frac{\partial E}{\partial w_M} \right]^T$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \nabla E(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

នៅទីនេះ $t = 0, 1, \dots$ គឺជាលេខរៀងនៃការផ្លាស់ប្តូរតម្លៃអថេរ x ។ $\nabla E(\mathbf{w})$ គឺជាដេរីវេដោយផ្នែកនៃ អនុគមន៍ E ធៀបនឹងអថេរ \mathbf{w} ឬហៅថា gradient ។ η_t គឺជាកម្រិតនៃការផ្លាស់ប្តូរតម្លៃអថេរដោយ គ្រប់គ្រងលើឥទ្ធិពលនៃតម្លៃ gradient ។ η_t ត្រូវបានហៅថាជា អត្រារៀនឬ learning rate ។ ជាទូទៅ តម្លៃនៃ η_t ត្រូវបានកំណត់យកចន្លោះ ០ និង ១ ដោយតម្លៃយ៉ាងតូច ។

យើងអាចកំណត់លក្ខខណ្ឌសម្រាប់បញ្ចប់ការផ្លាស់ប្តូរតម្លៃនៃអថេរបាន ដោយយកពេលវេលា លក្ខណៈជាចំនួននៃ gradient យកតម្លៃសូន្យឬក្បែរសូន្យ ។

ពិនិត្យលើករណីគម្រោងមួយ $f(x) = x^2 - 2x - 3$ ។ ករណីនេះយើងដឹងច្បាស់ថាតម្លៃអប្បបរមានៃអនុគមន៍គឺ -4 នៅពេលដែល $x^* = 1$ ។ យើងនឹងធ្វើផ្ទៀងផ្ទាត់ជាមួយតម្លៃដែលគណនាតាមរយៈ Gradient Descent ។

ដំបូងយើងគណនាអនុគមន៍ដេរីវេ $\frac{df(x)}{dx} = 2x - 2$ និង កំណត់យកអត្រា $\eta = 0.1$ បើ ។ យើងចាប់ផ្តើមពីចំណុច $x^{(0)} = 0$, $f(x^{(0)}) = -3$ ។ ដោយផ្លាស់ប្តូរតម្លៃអប្បបរមាតាមរយៈ Gradient Descent ខាងលើយើងបានបម្រែបម្រួលនៃតម្លៃអប្បបរមានិងតម្លៃអនុគមន៍ដូចតារាងខាងក្រោម ។

តារាងទី១ បម្រែបម្រួលនៃតម្លៃអប្បបរមានិងអនុគមន៍តាម Gradient Descent

t	$x^{(t)}$	$\frac{df(x)}{dx}$	$f(x)$
0	0.00	-2.00	-3.00
1	0.20	-1.60	-3.36
2	0.36	-1.28	-3.59
⋮	⋮	⋮	⋮
44	0.999946	-0.000109	-4.00
45	0.999956	-0.000087	-4.00

2. វិធីសាស្ត្រ Stochastic Gradient Descent (SGD)

ការធ្វើបរមាកម្មលើតម្លៃអនុគមន៍ដោយប្រើ Gradient Descent តម្លៃអនុគមន៍កម្រិតលម្អៀងនៃគ្រប់ទិន្នន័យទាំងអស់ $E(\mathbf{w})$ ក្នុងសំណុំទិន្នន័យសម្រាប់រៀន (training data) ត្រូវបានធ្វើអប្បបរមាកម្ម ។ ទាំងក្នុងករណីចំណោទតម្រិតម្រង់ (Regression) និងចំណោទចំណាត់ថ្នាក់ក្រុមទិន្នន័យ (Classification) អនុគមន៍កម្រិតលម្អៀងនៃគ្រប់ទិន្នន័យអាចសរសេរបានជាផលបូកនៃគ្រប់តម្លៃកម្រិតលម្អៀងក្នុងករណីគម្រោងសម្រាប់រៀននីមួយៗ $E_n(\mathbf{w})$ ។

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

ការផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្រដូចបានបង្ហាញក្នុង Gradient Descent ដោយប្រើអនុគមន៍កម្រិតលម្អៀងនៃគ្រប់ទិន្នន័យទាំងអស់ $E(\mathbf{w})$ ហៅថា ការរៀនជាក្រុម/ជាបាច់ (batch learning) ។ ផ្ទុយពីនេះ វិធីសាស្ត្រនៃការធ្វើអប្បបរមាអនុគមន៍កម្រិតលម្អៀងដោយធ្វើការផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្រដោយប្រើគម្រោងសម្រាប់រៀនម្តងមួយៗនិងតម្លៃអនុគមន៍កម្រិតលម្អៀងលើគម្រោងនោះ $E_n(\mathbf{w})$ ហៅថា stochastic gradient descent (SGD) ។ ក្នុងវិធី SGD គម្រោងទិន្នន័យសម្រាប់រៀន (training sample) ម្តងមួយៗ ត្រូវបានជ្រើសយកដោយចៃដន្យដើម្បីគណនា gradient នៃអនុគមន៍ $E_n(\mathbf{w})$ រួចធ្វើការផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រតែម្តង ដោយមិនចាំបាច់ធ្វើការបូកសរុបគ្រប់ទិន្នន័យដែលមាននោះឡើយ ។

$$\nabla E_n(\mathbf{w}) \equiv \frac{\partial E_n(\mathbf{w})}{\partial \mathbf{w}} = \left[\frac{\partial E_n}{\partial w_1} \quad \dots \quad \frac{\partial E_n}{\partial w_M} \right]^T$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \nabla E_n(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

ដូចដែលអ្នកអាចធ្វើការកត់សម្គាល់បាន ដោយប្រើgradient descent ពេលខ្លះយើងអាចនឹងទទួលបានតម្លៃប៉ារ៉ាម៉ែត្រដែលធ្វើឱ្យតម្លៃអនុគមន៍កម្រិតលម្អៀងធ្លាក់ចុះទៅក្នុងទីតាំងដែលជាបរមាធៀបតែមិនមែនជាកន្លែងអប្បបរមាពិតប្រាកដប្រសិនបើទីតាំងនៃការចាប់ផ្តើមរបស់អ្នកមិនប្រសើរ។ ប៉ុន្តែជាមួយ SGD ដោយសាររាល់ការផ្លាស់ប្តូរទិន្នន័យ គម្រិតទិន្នន័យនីមួយៗត្រូវបានជ្រើសរើសដោយចៃដន្យ ហេតុនេះភាពប្រថុយប្រថាននៃការធ្លាក់ចូលទៅក្នុងអប្បបរមាធៀបអាចត្រូវបានដោះស្រាយមួយកម្រិត។

ផ្ទុយពី ការរៀនជាក្រុម/ជាបាច់ (batch learning) ការរៀនដោយប្រើវិធីសាស្ត្រSGD ត្រូវបានហៅថា ការរៀនអនឡាញ (online learning) ។ ក្រៅពីនេះ ការរៀនដោយប្រើវិធីសាស្ត្រផ្លាស់ប្តូរប៉ារ៉ាម៉ែត្រដូចSGD ប៉ុន្តែជំនួសការរើសយកគម្រិតទិន្នន័យមួយមួយៗ ការរើសយកទិន្នន័យចំនួន ($< N$) សម្រាប់ប្រើគណនា gradient ក៏ត្រូវបានអនុវត្តដែរ។ ក្នុងករណីនេះគេហៅថា ការរៀនជាក្រុមតូច/ជាបាច់តូច (minibatch learning) ។ ក្នុងការរៀនជាក្រុមតូច ការផ្លាស់ប្តូរប៉ារ៉ាម៉ែត្រត្រូវធ្វើឡើងដូចខាងក្រោម។

$$E_t(\mathbf{w}) = \frac{1}{|\mathcal{D}_t|} \sum_{n \in \mathcal{D}_t} E_n(\mathbf{w}) , \mathcal{D}_t \subset \mathcal{D}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \nabla E_t(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^{(t)}}$$

3. សមត្ថភាពទូទៅ (Generalization Performance) និងស្ថានភាពរៀនហួសកំរិត (Overfitting)

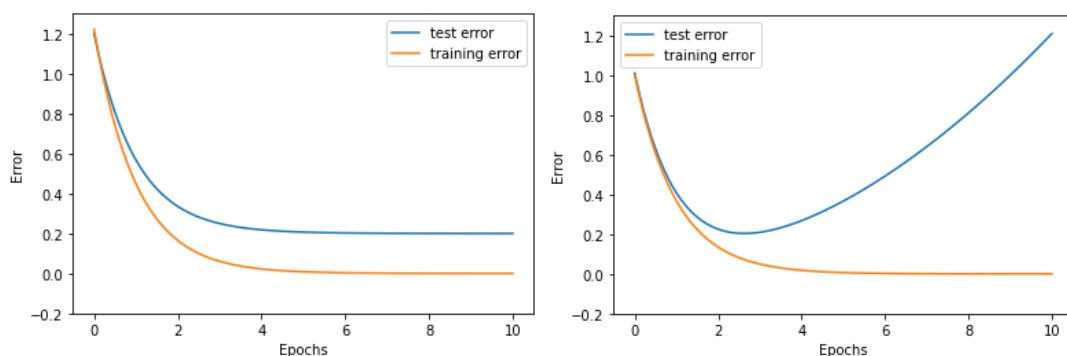
តាមការរៀបរាប់ពីដំណើរការរៀនក្នុងចំណុចខាងលើ យើងបានសិក្សាពីការកំណត់ប៉ារ៉ាម៉ែត្រណាដែលធ្វើឱ្យតម្លៃនៃអនុគមន៍កម្រិតលម្អៀងចំពោះសំណុំទិន្នន័យសម្រាប់រៀនមានតម្លៃតូចបំផុត។ ប៉ុន្តែតាមពិតទៅក្រៅពីសំណុំទិន្នន័យសម្រាប់រៀន យើងក៏ចង់បានបណ្តាញFNNដែលធ្វើឱ្យកម្រិតលម្អៀងក្នុងករណីទិន្នន័យផ្សេងក្រៅពីទិន្នន័យសម្រាប់រៀនមានតម្លៃតូចផងដែរ។

កម្រិតលម្អៀងចំពោះសំណុំទិន្នន័យសម្រាប់រៀនហៅថា កម្រិតលម្អៀងពេលរៀន (training error) និងតម្លៃសង្ឃឹមគណិតនៃកម្រិតលម្អៀងចំពោះសំណុំទិន្នន័យទាំងអស់ (population) ទាំងទិន្នន័យដែលបានរៀននិងទាំងទិន្នន័យដែលមិនមានក្នុងសំណុំសម្រាប់រៀនហៅថាកម្រិតលម្អៀងទូទៅ (generalization error) ។ គោលដៅយើងចង់បានបណ្តាញដែលផ្តល់ឱ្យនូវកម្រិតលម្អៀងទូទៅតូច

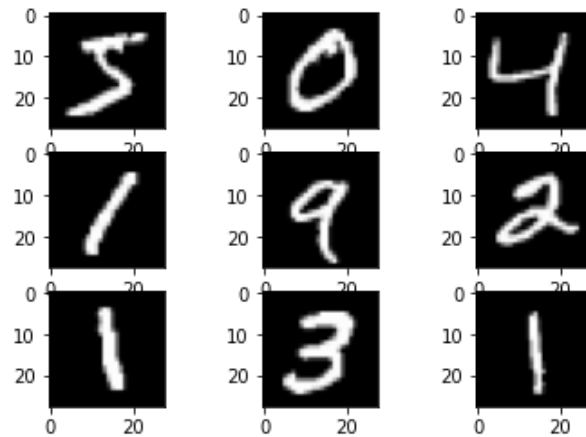
បំផុត ប៉ុន្តែដោយសារកម្រិតលម្អៀងទូទៅជាតម្លៃសង្ខេបគណិតលើសំណុំទិន្នន័យទាំងអស់ នោះការគណនាដោយផ្ទាល់មិនអាចធ្វើបានឡើយ ។ ហេតុនេះ ជំនួសឱ្យការសិក្សាលើសំណុំទិន្នន័យទាំងអស់ គេសិក្សាលើសំណុំតម្រូវទិន្នន័យមួយផ្នែកដែលមិនមានក្នុងសំណុំសម្រាប់រៀន ។ សំណុំនេះហៅថា សំណុំទិន្នន័យសម្រាប់វាយតម្លៃ (test data set) ហើយកម្រិតលម្អៀងនៃសំណុំទិន្នន័យសម្រាប់វាយតម្លៃនេះហៅថា កម្រិតលម្អៀងពេលវាយតម្លៃ (test error) ។ ពេលគឺជំនួសឱ្យគោលដៅនៃអប្បបរមាកម្មលើកម្រិតលម្អៀងទូទៅយើងសិក្សាលើអប្បបរមាកម្មលើកម្រិតលម្អៀងពេលវាយតម្លៃវិញ ។

រូបទី២ខាងក្រោមបង្ហាញពីបម្រែបម្រួលនៃអនុគមន៍កម្រិតលម្អៀងនៅពេលប៉ារ៉ាម៉ែត្រត្រូវបានផ្លាស់ប្តូរក្នុងដំណើរការរៀន ។ ក្រាបប្រភេទនេះហៅថាខ្សែកោងបង្ហាញដំណើរការរៀន (learning curve) ។ ជាទូទៅអនុគមន៍កម្រិតលម្អៀងពេលរៀនថយនៅពេលដំណើរការរៀនត្រូវបានអនុវត្ត ។ ផ្ទុយទៅវិញ អនុគមន៍កម្រិតលម្អៀងពេលវាយតម្លៃថយចុះដូចអនុគមន៍កម្រិតលម្អៀងពេលរៀនដែរ នៅដំណាក់កាលដំបូង ប៉ុន្តែនៅត្រង់ចំណុចណាមួយ វាក៏ចាប់ផ្តើមមានតម្លាតរវាងករណីរៀននិងការវាយតម្លៃ ។ ក្នុងករណីដែលមិនល្អ នៅពេលដែលដំណើរការរៀនបន្តទៅមុខ អនុគមន៍កម្រិតលម្អៀងពេលរៀនថយ ឯកម្រិតលម្អៀងវាយតម្លៃចាប់ផ្តើមកើនឡើង ។ ករណីបែបនេះបង្ហាញក្នុងរូបទី២ ផ្នែកខាងស្តាំ ។ ស្ថានភាពដែលកម្រិតលម្អៀងពេលរៀនថយចុះទាបខ្លាំងខណៈកម្រិតលម្អៀងវាយតម្លៃកើនឡើងហៅថា ស្ថានភាពរៀនហួសកំរិត ។ ស្ថានភាពបែបនេះអាចនិយាយបែបងាយបានថា បណ្តាញអាចផ្តល់លទ្ធផលបានល្អចំពោះតែទិន្នន័យណាដែលខ្លួនធ្លាប់បានរៀន រីឯទិន្នន័យដែលមិនមានក្នុងសំណុំសម្រាប់រៀន បណ្តាញមិនអាចផ្តល់លទ្ធផលបានល្អឡើយ ។

ក្នុងការដោះស្រាយបញ្ហាបែបនេះមានតិចតិចជាច្រើនត្រូវបានស្រាវជ្រាវនិងអនុវត្ត ។ យើងនឹងណែនាំលម្អិតក្នុងអត្ថបទខាងក្រោយ ប៉ុន្តែដំណោះស្រាយងាយគឺការបញ្ឈប់ដំណើរការរៀនមុនកំណត់ (early stopping) ដោយបញ្ឈប់ដំណើរការរៀនមុនពេលកម្រិតលម្អៀងវាយតម្លៃកើនឡើង ។



រូបទី២ អនុគមន៍កម្រិតលម្អៀងពេលរៀននិងពេលវាយតម្លៃក្នុងករណីធម្មតានិងករណីOverfitting



រូបទី៣ គម្រូទិន្នន័យក្នុងDataset (28 x 28)

4. សិក្សាឧទាហរណ៍៖ ការសម្គាល់លេខសរសេរដោយដៃដោយប្រើFNN

4.1. អំពីសំណុំទិន្នន័យនិង លក្ខណៈសម្គាល់នៃទិន្នន័យ

យើងប្រើសំណុំទិន្នន័យ MNIST¹ សម្រាប់ចំណោទធ្វើចំណាត់ថ្នាក់ក្រុមលេខសរសេរដោយដៃទៅតាមក្រុម(០ដល់៩)ដោយប្រើFNN។ គម្រូទិន្នន័យមានបង្ហាញក្នុងរូបទី៣។ ទិន្នន័យជារូបនៃលេខសរសេរដោយដៃទាំងនេះមានទំហំ28x28។ នៅទីនេះយើងបំប្លែងវាជាទម្រង់វ៉ិចទ័រដែលមានទំហំ784(=28x28) និង បំប្លែងលេខសំគាល់ថ្នាក់របស់វាជាទម្រង់one-hot vector។

ចំពោះដំណើរការរៀននិងវាយតម្លៃ យើងបែងចែក60000គម្រូទិន្នន័យសម្រាប់ដំណើរការរៀន (training data)និង10000គម្រូទិន្នន័យសម្រាប់ការវាយតម្លៃ(test data) ។

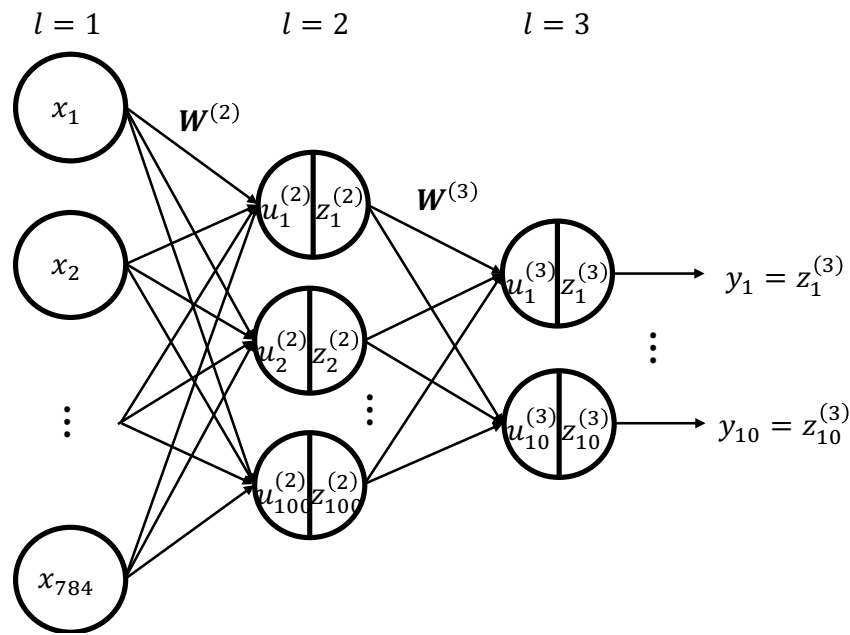
4.2. អំពីបណ្តាញFNN

យើងកំណត់យកFNN៣ថ្នាក់សម្រាប់ប្រើនៅទីនេះ ពោលគឺ input layer- hidden layer- output layer ។ ចំនួនណឺរ៉ូននៅinput layer គឺស្មើនឹងទំហំវ៉ិចទ័រដែលជាធាតុចូល(784)ឯចំនួនណឺរ៉ូននៅoutput layerគឺស្មើនឹងចំនួនក្រុម(10) ។ ចំពោះចំនួនណឺរ៉ូននៅhidden layerយើងកំណត់យកស្មើនឹង 100 ។

ចំពោះអនុគមន៍សកម្ម(activation function) យកកំណត់ប្រើ sigmoid function ចំពោះលទ្ធផលពីinput layerនិងsoftmax functionសម្រាប់លទ្ធផលចុងក្រោយ។

ចំពោះដំណើរការរៀនយើងនឹងប្រើ SGD (minibatch) ។

¹ <http://yann.lecun.com/exdb/mnist/> , <https://keras.io/api/datasets/>



រូបទី៤ ទម្រង់ FNN ដែលប្រើក្នុងឧទាហរណ៍

4.3. អំពីដំណើរការរៀនជាក្រុមតូច (minibatch learning)

Step-1 : ជ្រើសរើស minibatch

ជ្រើសយកគម្រូទិន្នន័យពីសំណុំទិន្នន័យសម្រាប់រៀនមួយផ្នែក (minibatch) ដោយចៃដន្យ។ គោលដៅយើងគឺធ្វើអប្បបរមាតម្លៃនៃអនុគមន៍កម្រិតលម្អៀងចំពោះទិន្នន័យទាំងនេះ។

Step-2: គណនា gradient

ដើម្បីអនុវត្ត Gradient Descent (minibatch) យើងគណនាតម្លៃ gradient ធៀបនឹងប៉ារ៉ាម៉ែត្រនីមួយៗ។

Step-3: ផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រ

ផ្លាស់ប្តូរតម្លៃនៃប៉ារ៉ាម៉ែត្រតាមទំនាក់ទំនងដែលបានបង្ហាញនៅចំណុច 2. ខាងលើ។

Step-4: អនុវត្ត Step-1, Step-2, Step-3 ម្តងហើយម្តងទៀតរហូតដល់ចំនួនដងដែលបានកំណត់ជាមុន (Iteration number)

នៅទីនេះយើងកំណត់យក អត្រារៀន (learning rate) $\eta = 0.1$, ទំហំក្រុមតូច (minibatch) $|\mathcal{D}_t| = 100$ និងចំនួនដងក្នុងដំណើរការរៀន Iteration number = 10000 ។

4.4. ការបង្កើតបណ្តាញFNNដោយប្រើPython

```
import numpy as np
```

```
class SimpleFNN:
    def __init__(self, input_size, hidden_size, output_size):
        self.params = {}
        self.params['W1'] = np.random.randn(input_size, hidden_size)
        self.params['b1'] = np.random.randn(hidden_size)
        self.params['W2'] = np.random.randn(hidden_size, output_size)
        self.params['b2'] = np.random.randn(output_size)

    def predict(self, x):
        W1, W2 = self.params['W1'], self.params['W2']
        b1, b2 = self.params['b1'], self.params['b2']
        u1 = np.dot(x, W1) + b1
        z1 = sigmoid(u1)
        u2 = np.dot(z1, W2) + b2
        y = softmax(u2)
        return y

    def loss(self, x, t):
        y = self.predict(x)
        return cross_entropy_error(y, t)

    def accuracy(self, x, t):
        y = self.predict(x)
        y = np.argmax(y, axis=1)
        t = np.argmax(t, axis=1)
        accuracy = np.sum(y == t) / float(x.shape[0])
        return accuracy

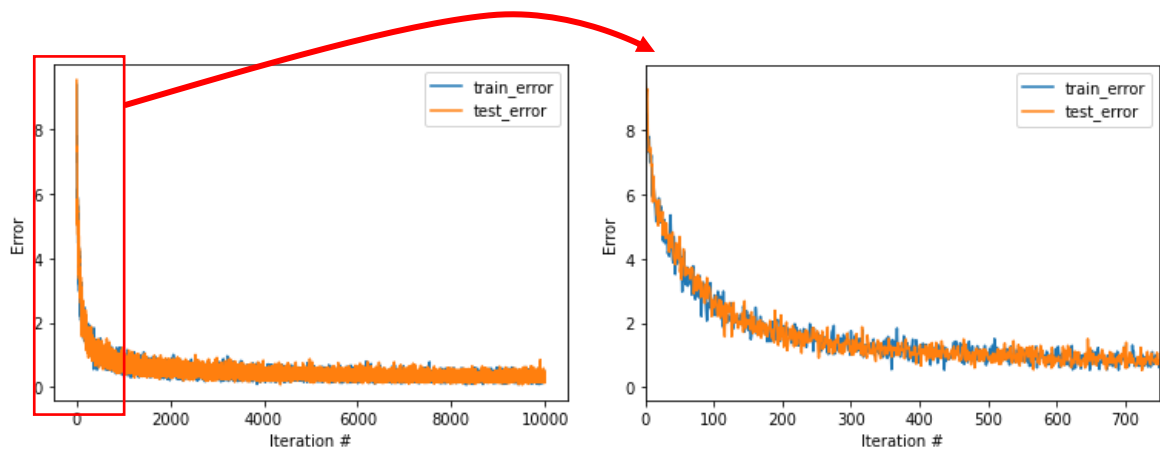
    def gradient(self, x, t):
        W1, W2 = self.params['W1'], self.params['W2']
        b1, b2 = self.params['b1'], self.params['b2']
        grads = {}
        batch_num = x.shape[0]
        a1 = np.dot(x, W1) + b1
        z1 = sigmoid(a1)
        a2 = np.dot(z1, W2) + b2
        y = softmax(a2)
        dy = (y - t) / batch_num
        grads['W2'] = np.dot(z1.T, dy)
        grads['b2'] = np.sum(dy, axis=0)
        dz1 = np.dot(dy, W2.T)
        dal = sigmoid_grad(a1) * dz1
        grads['W1'] = np.dot(x.T, dal)
        grads['b1'] = np.sum(dal, axis=0)
        return grads
```

```
def sigmoid(x):
    return 1 / (1 + np.exp(-x))

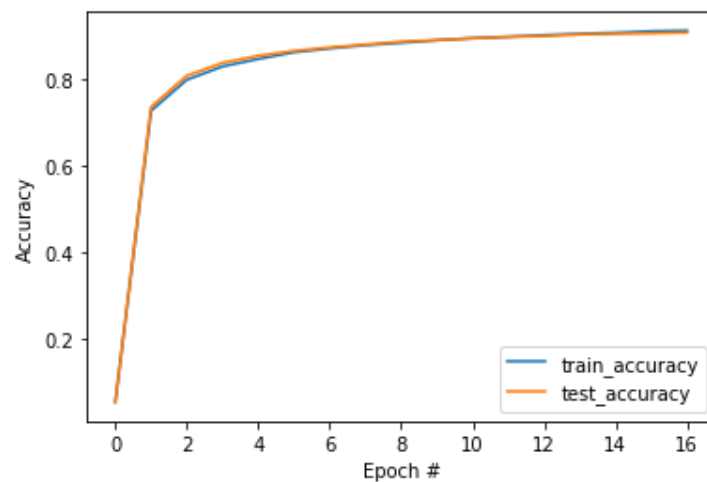
def sigmoid_grad(x):
    return (1.0 - sigmoid(x)) * sigmoid(x)

def softmax(x):
    x = x - np.max(x, axis=-1, keepdims=True) # To deal with overflow
    return np.exp(x) / np.sum(np.exp(x), axis=-1, keepdims=True)

def cross_entropy_error(y, t):
    batch_size = y.shape[0]
    t = t.argmax(axis=1)
    return -np.sum(np.log(y[np.arange(batch_size), t] + 1e-7)) / batch_size
```



រូបទី៥ កម្រិតលម្អៀង (Cross-entropy) ក្នុងដំណើរការរៀន



រូបទី៦ អត្រាត្រឹមត្រូវនៃការប៉ាន់ស្មានក្រុមចំពោះទិន្នន័យសម្រាប់រៀននិងទិន្នន័យវាយតម្លៃ

4.5. លទ្ធផល

រូបទី៥បង្ហាញពីបម្រែបម្រួលនៃអនុគមន៍កម្រិតលម្អៀងក្នុងដំណើរការរៀនចំពោះទិន្នន័យសម្រាប់រៀននិងទិន្នន័យវាយតម្លៃ។ រូបទី៦បង្ហាញពីអត្រាត្រឹមត្រូវនៃការធ្វើចំណាត់ថ្នាក់ក្រុមលើលេខសរសេរដៃដោយប្រើFNN ចំពោះទិន្នន័យសម្រាប់រៀននិងទិន្នន័យសម្រាប់វាយតម្លៃ។

ក្នុងអត្ថបទក្រោយយើងនឹងណែនាំអំពីវិធីសាស្ត្រកំណត់ប៉ារ៉ាម៉ែត្រក្នុងដំណើរការរៀនផ្សេងទៀតក្រៅពីSGD ដែលមានប្រសិទ្ធភាពក្នុងការគណនាចំពោះទម្រង់នៃNeural Network ។

Backpropagation (BP)

ក្នុងអត្ថបទមុនយើងបានសិក្សាអំពីដំណើរការរៀនដើម្បីកំណត់តម្លៃប៉ារ៉ាម៉ែត្រនៃ FNN តាមរយៈវិធីសាស្ត្រ stochastic gradient descend (SGD) ។ ដូចដែលអ្នកអាចចាប់អារម្មណ៍បានក្នុង SGD ការគណនា gradient ឬដេរីវេនៃអនុគមន៍កម្រិតលម្អៀងត្រូវបានធ្វើឡើង ។ ជាទូទៅការគណនានេះអាចធ្វើបានតាមរយៈការគណនាតម្លៃប្រហែលនៃដេរីវេដោយប្រើតម្លៃអនុគមន៍ផ្ទាល់ ។ ប៉ុន្តែការគណនាបែបនេះចំណាយពេលច្រើន ។ ដើម្បីគណនា gradient ឬដេរីវេនៃអនុគមន៍ប្រកបដោយប្រសិទ្ធភាព វិធីសាស្ត្រគណនា gradient ដោយប្រើ backpropagation ត្រូវបានប្រើប្រាស់ជាទូទៅ ។

1. ភាពលំបាកក្នុងការគណនា gradient

នៅក្នុង stochastic gradient descend ការគណនាតម្លៃ gradient នៃអនុគមន៍កម្រិតលម្អៀង ($\nabla E(\mathbf{w}) = \partial E(\mathbf{w}) / \partial \mathbf{w}$) គឺជាដំណាក់កាលដ៏សំខាន់ ។ ចំពោះ FNN ច្រើនថ្នាក់ ការគណនា gradient ចំពោះប៉ារ៉ាម៉ែត្រមានភាពស្មុគស្មាញខ្លាំង ។

ជាឧទាហរណ៍តម្លៃកម្រិតលម្អៀងចំពោះគម្រិតទី n នៃចំណោទតម្រៃតម្រង់ (regression) កំណត់ដោយ $E_n = \frac{1}{2} \|\mathbf{y}(x_n) - \mathbf{t}_n\|^2$ ។ យើងសាកល្បងគណនាដេរីវេធៀបនឹងប៉ារ៉ាម៉ែត្រទំនងផ្ទាល់ $w_{ji}^{(l)}$ នៃថ្នាក់ទី l ។

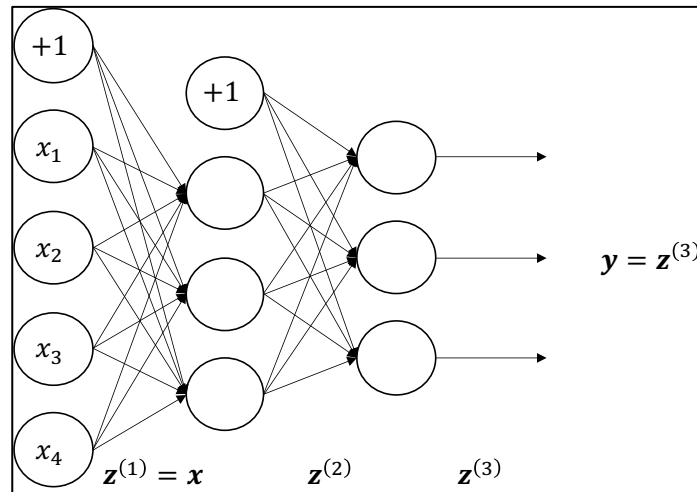
ដំបូងយើងពិនិត្យឃើញថា

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = (\mathbf{y}(x_n) - \mathbf{t}_n)^T \frac{\partial \mathbf{y}}{\partial w_{ji}^{(l)}}$$

បន្ទាប់មកទៀតយើងត្រូវគណនាដេរីវេ $\frac{\partial \mathbf{y}}{\partial w_{ji}^{(l)}}$ ។ ដោយលទ្ធផលនៃ FNN $\mathbf{y}(x)$ កំណត់ដោយទម្រង់ខាងក្រោម នោះយើងអាចមើលឃើញបានថាការគណនាដេរីវេតាមវិធីបែបនេះមិនមានប្រសិទ្ធភាពឡើយ ពោលគឺត្រូវចំណាយពេលច្រើនក្នុងការគណនាដោយប្រើ Programming ។

$$\begin{aligned} \mathbf{y}(x) &= \mathbf{f}(\mathbf{u}^{(L)}) \\ &= \mathbf{f}(\mathbf{W}^{(L)} \mathbf{z}^{(L-1)} + \mathbf{b}^{(L)}) \\ &= \mathbf{f}(\mathbf{W}^{(L)} \mathbf{f}(\mathbf{W}^{(L-1)} \mathbf{z}^{(L-2)} + \mathbf{b}^{(L-1)}) + \mathbf{b}^{(L)}) \\ &= \mathbf{f}(\mathbf{W}^{(L)} \mathbf{f}(\mathbf{W}^{(L)} \mathbf{f}(\dots \mathbf{f}(\mathbf{W}^{(l)} \mathbf{z}^{(l-1)} + \mathbf{b}^{(l)}) \dots)) + \mathbf{b}^{(L)}) \end{aligned}$$

វិធីសាស្ត្រ backpropagation អាចជួយដោះស្រាយបញ្ហាគណនាដេរីវេនៃអនុគមន៍បណ្តាក់ច្រើនជាន់បែបនេះបាន ។ នៅក្នុងការបង្ហាញខាងក្រោម ដើម្បីសម្រួលដល់ការសរសេរ យើងកំណត់សរសេរ bias ជាផ្នែកមួយនៃប៉ារ៉ាម៉ែត្រទំនងផ្ទាល់នៃណឺរ៉ូនដែរ ពោលគឺ $w_{0j}^{(l)} = b_j^{(l)}$ ។ ហេតុនេះដោយកំណត់ណឺរ៉ូនទី 0 នៃថ្នាក់ $(l-1)$ ឱ្យបញ្ចេញនូវលទ្ធផល $z_0^{(l-1)} = 1$ ជានិច្ចនោះយើងអាចសរសេរលទ្ធផលនៃណឺរ៉ូនដោយទម្រង់ខាងក្រោម ។



រូបទី១ FNNមានប៉ារ៉ាម៉ែត្រពីរថ្នាក់ (ណឺរ៉ូន៣ថ្នាក់)

$$u_j^{(l)} = \sum_{i=1}^k w_{ji}^{(l)} z_i^{(l-1)} + b_j = \sum_{i=0}^k w_{ji}^{(l)} z_i^{(l-1)}$$

2. ការគណនាតាម backpropagation ករណីFNNមានប៉ារ៉ាម៉ែត្រពីរថ្នាក់ (ណឺរ៉ូន៣ថ្នាក់)

រូបទី១បង្ហាញទម្រង់នៃFNNមានប៉ារ៉ាម៉ែត្រពីរថ្នាក់ (ណឺរ៉ូន៣ថ្នាក់) នៃចំណោទតម្រេតម្រង់ (regression) ។ អនុគមន៍សកម្ម (activation function) នៃថ្នាក់លទ្ធផលចុងក្រោយកំណត់ដោយអនុគមន៍identity ($f(x) = x$) ។

សន្មតធាតុចូលនៃបណ្តាញនេះដោយ $x = [x_1 \ x_2 \ x_3 \ x_4]^T$ ។ លទ្ធផលនៃថ្នាក់ណឺរ៉ូនទីមួយពេលថ្នាក់ធាតុចូលគឺ $z_i^{(1)} = x_i$ និងលទ្ធផលនៃថ្នាក់កណ្តាល $z_j^{(2)}$ ព្រមទាំងលទ្ធផលនៃថ្នាក់លទ្ធផលចុងក្រោយ $y_j(x) = z_j^{(3)}$ កំណត់ដោយទម្រង់ខាងក្រោម ។

$$z_j^{(2)} = f(u_j^{(2)}) = f\left(\sum_i w_{ji}^{(2)} z_i^{(1)}\right)$$

$$y_j(x) = z_j^{(3)} = u_j^{(3)} = \sum_i w_{ji}^{(3)} z_i^{(2)}$$

សន្មតយកផលបូកការនៃលម្អៀងជាអនុគមន៍កម្រិតលម្អៀងនៃបណ្តាញនេះ ។

$$E_n = \frac{1}{2} \|y(x_n) - t_n\|^2 = \frac{1}{2} \sum_j (y_j(x_n) - t_n)^2$$

ពេលនេះយើងពិនិត្យលើការគណនាដេរីវេនៃអនុគមន៍នេះធៀបនឹងប៉ារ៉ាម៉ែត្ររបស់វា ។ ដំបូងយើងគណនាដេរីវេធៀបប៉ារ៉ាម៉ែត្រនៃផ្នែកថ្នាក់លទ្ធផលចុងក្រោយនៃបណ្តាញ $\frac{\partial E_n}{\partial w_{ji}^{(3)}}$ ។

$$\frac{\partial E_n}{\partial w_{ji}^{(3)}} = \frac{\partial}{\partial w_{ji}^{(3)}} \left\{ \frac{1}{2} \|\mathbf{y}(\mathbf{x}_n) - \mathbf{t}_n\|^2 \right\} = (\mathbf{y}(\mathbf{x}_n) - \mathbf{t}_n)^\top \frac{\partial \mathbf{y}}{\partial w_{ji}^{(3)}}$$

ដោយ

$$y_j(\mathbf{x}) = z_j^{(3)} = u_j^{(3)} = \sum_i w_{ji}^{(3)} z_i^{(2)}$$

នោះ

$$\frac{\partial \mathbf{y}}{\partial w_{ji}^{(3)}} = [0 \dots 0 z_i^{(2)} 0 \dots 0]^\top$$

ដូច្នេះ

$$\frac{\partial E_n}{\partial w_{ji}^{(3)}} = (\mathbf{y}(\mathbf{x}_n) - \mathbf{t}_n)^\top \frac{\partial \mathbf{y}}{\partial w_{ji}^{(3)}} = (y_j(\mathbf{x}_n) - t_{nj}) z_i^{(2)}$$

បន្ទាប់ពីនេះ យើងពិនិត្យលើថ្នាក់កណ្តាល $\frac{\partial E_n}{\partial w_{ji}^{(2)}}$ ។

$$\frac{\partial E_n}{\partial w_{ji}^{(2)}} = \frac{\partial E_n}{\partial u_j^{(2)}} \frac{\partial u_j^{(2)}}{\partial w_{ji}^{(2)}}$$

ដោយ

$$u_j^{(2)} = \sum_i w_{ji}^{(2)} z_i^{(1)}$$

នោះយើងបាន

$$\frac{\partial u_j^{(2)}}{\partial w_{ji}^{(2)}} = z_i^{(1)}$$

ម្យ៉ាងទៀត បើយក k ជាចំនួនណាមួយនៅ ថ្នាក់លទ្ធផលចុងក្រោយនោះ E_n មាន $u_1^{(3)}, \dots, u_k^{(3)}$ ជាអថេរដែលយើងអាចសរសេរអនុគមន៍ដេរីវេដូចខាងក្រោម។

$$\begin{aligned} \frac{\partial E_n}{\partial u_j^{(2)}} &= \sum_k \frac{\partial E_n}{\partial u_k^{(3)}} \frac{\partial u_k^{(3)}}{\partial u_j^{(2)}} \\ \frac{\partial E_n}{\partial u_k^{(3)}} &= \frac{\partial}{\partial u_k^{(3)}} \left\{ \frac{1}{2} \sum_j (y_j(\mathbf{x}_n) - \mathbf{t}_n)^2 \right\} \\ &= \frac{\partial}{\partial u_k^{(3)}} \left\{ \frac{1}{2} \sum_j (u_j^{(3)} - \mathbf{t}_n)^2 \right\} = u_k^{(3)} - t_{nk} \end{aligned}$$

ដោយ

$$\frac{\partial u_k^{(3)}}{\partial u_j^{(2)}} = \frac{\partial}{\partial u_j^{(2)}} \left\{ \sum_i w_{ki}^{(3)} z_i^{(2)} \right\}$$

$$\begin{aligned}
&= \frac{\partial}{\partial u_j^{(2)}} \left\{ \sum_i w_{ki}^{(3)} f(u_j^{(2)}) \right\} \\
&= w_{kj}^{(3)} f'(u_j^{(2)})
\end{aligned}$$

ហេតុនេះ យើងបាន

$$\frac{\partial E_n}{\partial w_{ji}^{(2)}} = \left(f'(u_j^{(2)}) \sum_k w_{kj}^{(3)} (u_k^{(3)} - t_{nk}) \right) z_i^{(1)}$$

3. ករណីទូទៅ៖ FNNច្រើនថ្នាក់

ជំបូង យើងពិនិត្យលើប៉ារ៉ាម៉ែត្រ $w_{ji}^{(l)}$ នៃថ្នាក់ទី l ។

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \frac{\partial E_n}{\partial u_j^{(l)}} \frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}}$$

ដោយពិនិត្យរូបទី២ (ឆ្វេង) បម្រែបម្រួលនៃ $u_j^{(l)}$ ជះឥទ្ធិពលលើតម្លៃនៃ E_n តាមរយៈតម្លៃនៃ $z_j^{(l)}$ និង តម្លៃលទ្ធផលនៃណឺរ៉ូននៅថ្នាក់ទី $(l+1)$ ។ ហេតុនេះយើងបាន

$$\frac{\partial E_n}{\partial u_j^{(l)}} = \sum_k \frac{\partial E_n}{\partial u_k^{(l+1)}} \frac{\partial u_k^{(l+1)}}{\partial u_j^{(l)}}$$

។ ដោយពិនិត្យលើកន្សោមខាងលើ យើងឃើញថា $\partial E_n / \partial u_j^{(l)}$ បង្ហាញនៅអង្គទាំងសង្ខេប ។ នៅទីនេះ យើងសន្មតតាង

$$\delta_j^{(l)} \equiv \frac{\partial E_n}{\partial u_j^{(l)}}$$

។ ដោយប្រើទំនាក់ទំនង $u_k^{(l+1)} = \sum_j w_{kj}^{(l+1)} z_j^{(l)} = \sum_j w_{kj}^{(l+1)} f(u_j^{(l)})$ យើងបាន $\partial u_k^{(l+1)} / \partial u_j^{(l)} = w_{kj}^{(l+1)} f'(u_j^{(l)})$ ។ ហេតុនេះ

$$\delta_j^{(l)} = \frac{\partial E_n}{\partial u_j^{(l)}} = \sum_k \frac{\partial E_n}{\partial u_k^{(l+1)}} \frac{\partial u_k^{(l+1)}}{\partial u_j^{(l)}} = \sum_k \delta_k^{(l+1)} \left(w_{kj}^{(l+1)} f'(u_j^{(l)}) \right)$$

តាមទំនាក់ទំនងនេះបង្ហាញថាយើងអាចគណនា $\delta_j^{(l)}$ បានដោយប្រើតម្លៃនៃ $\delta_k^{(l+1)}$ ($k = 1, 2, \dots$) ។ មានន័យថា បើយើងដឹងតម្លៃ δ របស់ថ្នាក់ខាងលើពេលគឺ $(l+1)$ នោះយើងអាចគណនាតម្លៃ δ នៅថ្នាក់ក្រោមបន្តបន្ទាប់តាមទំនាក់ទំនងនេះ ។ រូបទី២ (ស្តាំ) បង្ហាញពីដំណើរការនៃការគណនា

δ ពីថ្នាក់លើឆ្ពោះទៅថ្នាក់ក្រោមបែបនេះ។ ទំនាក់ទំនងខាងលើនេះពិតជានិច្ចចំពោះគ្រប់ថ្នាក់ទាំងអស់នៃបណ្តាញ។

ដូច្នេះបើ δ នៅថ្នាក់លទ្ធផលចុងក្រោយត្រូវបានគណនា នោះ យើងអាចគណនា δ ពេលគឺដេរីវេនៃប៉ារ៉ាម៉ែត្រនៅថ្នាក់ក្រោមជាបន្តបន្ទាប់បានដោយអនុវត្តតាមទំនាក់ទំនងដោយខាងលើ។ ដោយសារតែលំដាប់នៃការគណនា δ នៅទីនេះមានទិសដៅផ្ទុយពីទិសដៅបញ្ជូនសញ្ញាណក្នុងការប៉ាន់ស្មានលទ្ធផលរបស់បណ្តាញ ដូចនេះទើបគេហៅឈ្មោះវិធីសាស្ត្រនេះថាជា backpropagation ។

ត្រលប់មកផ្នែកនៅសល់នៃ $\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \frac{\partial E_n}{\partial u_j^{(l)}} \frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}}$ ។ ដោយ $\frac{\partial E_n}{\partial u_j^{(l)}}$ អាចគណនាបានដោយគណនា δ ដូច្នេះយើងពិនិត្យលើ $\frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}}$ ។ តាមទំនាក់ទំនង $u_j^{(l)} = \sum_{i=0}^k w_{ji}^{(l)} z_i^{(l-1)}$ នោះ $\frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}} = z_i^{(l-1)}$ ហេតុនេះយើងបាន

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)}$$

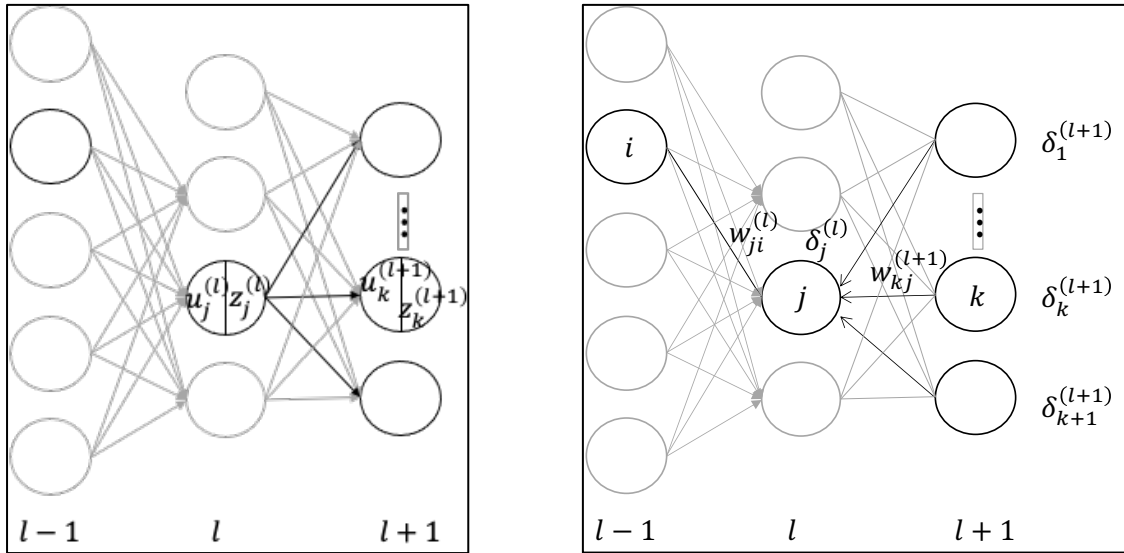
ដូចបង្ហាញក្នុងទំនាក់ទំនងដែលទាញបាននេះ ការគណនាដេរីវេដោយផ្នែកនៃអនុគមន៍កម្រិតលម្អៀងធៀបប៉ារ៉ាម៉ែត្រ $w_{ji}^{(l)}$ ដែលភ្ជាប់ណឺរ៉ូនទី i នៃថ្នាក់ទី $(l-1)$ និងណឺរ៉ូនទី j នៃថ្នាក់ទី (l) អាចគណនាបានយ៉ាងងាយដោយប្រើ $\delta_j^{(l)}$ នៃណឺរ៉ូនទី j នៃថ្នាក់ទី (l) និងលទ្ធផល $z_i^{(l-1)}$ នៃណឺរ៉ូនទី i នៃថ្នាក់ទី $(l-1)$ ។ ដូចបានបញ្ជាក់ខាងលើ $\delta_j^{(l)}$ អាចគណនាបានដោយគណនាជាបន្តបន្ទាប់ពីថ្នាក់ខាងលើតាមទំនាក់ទំនងដែលបានទាញពីខាងដើម។ ក្នុងករណីនេះ $\delta_j^{(L)}$ នៃថ្នាក់ខាងលើបំផុត (ថ្នាក់លទ្ធផលចុងក្រោយ) អាចកំណត់បានដោយ

$$\delta_j^{(L)} = \frac{\partial E_n}{\partial u_j^{(L)}}$$

ដែលការគណនាជាក់ស្តែងប្រែប្រួលទៅតាមប្រភេទនៃអនុគមន៍កម្រិតលម្អៀង (ទៅតាមចំណោទ) ។

ដោយបូកសរុបការគណនាខាងលើ នៅពេលដែលគម្រោងទិន្នន័យសម្រាប់រៀន (x_n, t_n) ត្រូវបានផ្តល់ gradient ឬដេរីវេនៃអនុគមន៍កម្រិតលម្អៀងអាចគណនាបានតាមលំដាប់ដោយខាងក្រោម។ ក្នុងករណីនៃការរៀនជាក្រុមតូច (minibatch) ផលបូកនៃ gradient បានមកពីការគណនាចំពោះគម្រោងទិន្នន័យនីមួយៗត្រូវបានយកជា gradient នៃក្រុមតូចនោះ

$$\frac{\partial E}{\partial w_{ji}^{(l)}} = \sum_n \frac{\partial E_n}{\partial w_{ji}^{(l)}}$$



រូបទី២ ករណីទូទៅនៃ FNN ច្រើនថ្នាក់

Input : គម្រូទិន្នន័យសម្រាប់រៀន (x_n, t_n)

Output : gradient ឬដេរីវេនៃអនុគមន៍កម្រិតលម្អៀង $\frac{\partial E_n}{\partial w_{ji}^{(l)}} \quad (l = 2, \dots, L)$

- (1) ដោយយក $z^{(1)} = x$ គណនាតម្លៃនៃ $u^{(l)}, z^{(l)} \quad (l = 2, \dots, L)$ តាមលំដាប់
- (2) គណនា $\delta_j^{(L)}$ (តាមធម្មតាវាត្រូវបានកំណត់ដោយ $\delta_j^{(L)} = z_j - t_j$)
- (3) ចំពោះថ្នាក់កណ្តាល $(l = L - 1, L - 2, \dots, 2)$ គណនាតម្លៃ $\delta_j^{(l)}$ តាមលំដាប់ដោយ

$$\delta_j^{(l)} = \sum_k \delta_k^{(l+1)} \left(w_{kj}^{(l+1)} f'(u_j^{(l)}) \right)$$

- (4) ចំពោះថ្នាក់ទី $l \quad (l = 2, \dots, L)$ គណនាតម្លៃ $\frac{\partial E_n}{\partial w_{ji}^{(l)}}$ ដោយ

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)}$$

4. ការគណនា $\delta_j^{(L)}$ នៃថ្នាក់លទ្ធផលចុងក្រោយ

ដូចបានបង្ហាញខាងលើ ការគណនា $\delta_j^{(L)}$ នៃថ្នាក់លទ្ធផលចុងក្រោយអាស្រ័យនឹងប្រភេទអនុគមន៍កម្រិតលម្អៀងដែលប្រើ ពោលគឺអាស្រ័យនឹងប្រភេទនៃចំណោទ ។

4.1. ករណីចំណោទតម្រេតម្រង់ Regression

ក្នុងករណីចំណោទតម្រេតម្រង់អនុគមន៍កម្រិតលម្អៀងគឺជាផលបូកការេនៃលម្អៀងលើគម្រូទិន្នន័យនីមួយៗ ។

$$E_n = \frac{1}{2} \|y(x_n) - t_n\|^2 = \frac{1}{2} \sum_j (y_j - t_j)^2$$

ក្នុងករណីនេះ អនុគមន៍សកម្ម (activation) នៃថ្នាក់លទ្ធផលចុងក្រោយនៃ FNN គឺជាអនុគមន៍ identity: $y_j = z_j^{(L)} = u_j^{(L)}$ ។ ហេតុនេះ $\delta_j^{(L)}$ នៃថ្នាក់លទ្ធផលចុងក្រោយគឺ

$$\delta_j^{(L)} = u_j^{(L)} - t_{nj} = z_j^{(L)} - t_j = y_j - t_j$$

ពោលគឺគម្លាតរវាងតម្លៃលទ្ធផលនៃបណ្តាញ (ណឺរ៉ូន) និងលទ្ធផលក្នុងទិន្នន័យសម្រាប់រៀន ។

4.2. ករណីចំណោទចំណាត់ថ្នាក់ក្រុម Classification

ក្នុងចំណោទចំណាត់ថ្នាក់ក្រុម អនុគមន៍កម្រិតលម្អៀងត្រូវបានកំណត់ដោយ

$$E_n = -t \log y - (1 - t) \log(1 - y)$$

ក្នុងករណីនេះ អនុគមន៍សកម្មនៃថ្នាក់លទ្ធផលចុងក្រោយនៃ FNN គឺជាអនុគមន៍ sigmoid ហេតុនេះ

$$y = \frac{1}{1 + \exp(-u)}, \quad \frac{dy}{du} = -\frac{\exp(-u)}{(1 + \exp(-u))^2} = y(1 - y)$$

ដូច្នេះយើងបាន

$$\begin{aligned} \delta_j^{(L)} &= -\frac{t}{y} \cdot \frac{dy}{du} - \frac{1-t}{1-y} \left(-\frac{dy}{du} \right) \\ &= -t(1-y) - (1-t)y = y - t \end{aligned}$$

ក្នុងចំណោទចំណាត់ថ្នាក់ច្រើនក្រុម អនុគមន៍កម្រិតលម្អៀងត្រូវបានកំណត់ដោយ

$$E_n = -\sum_k t_k \log y_k = -\sum_k t_k \log \left(\frac{\exp(u_k^{(L)})}{\sum_i \exp(u_i^{(L)})} \right)$$

ក្នុងករណីនេះ អនុគមន៍សកម្មនៃថ្នាក់លទ្ធផលចុងក្រោយនៃ FNN គឺជាអនុគមន៍ softmax ហេតុនេះ

$$y = \frac{\exp(u_k^{(L)})}{\sum_i \exp(u_i^{(L)})}$$

ដូច្នេះយើងបាន

$$\begin{aligned} \delta_j^{(L)} &= -\sum_k t_k \frac{1}{y_k} \cdot \frac{\partial y_k}{\partial u_j^{(L)}} \\ &= -t_j(1 - y_j) - \sum_{k \neq j} t_k(-y_j) \\ &= (y_j - t_j) \sum_k t_k \\ &= y_j - t_j \left(\because \sum_k t_k = 1 \text{ (one-hot vector)} \right) \end{aligned}$$

តាមលទ្ធផលខាងលើ ទាំងករណីចំណោទតម្រៃតម្រង់ ទាំងករណីចំណោទចំណាត់ថ្នាក់
ក្រុម $\delta_j^{(L)}$ នៃថ្នាក់លទ្ធផលចុងក្រោយគឺ

$$\delta_j^{(L)} = y_j - t_j$$

ពោលគឺគម្លាតរវាងតម្លៃលទ្ធផលនៃបណ្តាញ (ណឺរ៉ូន) និងលទ្ធផលក្នុងទិន្នន័យសម្រាប់រៀន ។

Technique in Learning Process

យើងបានសិក្សាអំពីការកំណត់តម្លៃប៉ារ៉ាម៉ែត្រនៃ Neural Network នៅក្នុងដំណើរការរៀន ដូចជាវិធី stochastic gradient descent និង វិធីក្នុងការគណនាដេរីវេដោយផ្នែកប្រកបដោយប្រសិទ្ធភាពដូចជា backpropagation ។ ក្នុងអត្ថបទនេះយើងនឹងលើកយកល្បិចនិងវិធីសាស្ត្រមួយចំនួនដែលត្រូវបានប្រើដើម្បីបង្កើនប្រសិទ្ធភាពនៃការរៀនក្នុង Neural Network ។ ប្រសិទ្ធភាពនៃការរៀននៅទីនេះសំដៅដល់ភាពល្អនៃអត្រាត្រឹមត្រូវក្នុងការប៉ាន់ស្មានលទ្ធផលរបស់បណ្តាញស្មុគស្មាញដូចជា FNN ច្រើនថ្នាក់ដែលហៅថា Deep Learning ព្រមទាំងល្បឿននៃការរៀនរបស់វា ។

1. ការកាត់បន្ថយស្ថានភាពរៀនហួសកម្រិត (Overfitting)

1.1. ការដាក់កម្រិតលើទម្ងន់ផ្ទាល់នៃណឺរ៉ូន (weight decay)

បញ្ហារៀនហួសកម្រិតគឺជាស្ថានភាពដែលបណ្តាញអាចផ្តល់លទ្ធផលល្អតែចំពោះទិន្នន័យសម្រាប់រៀនប៉ុណ្ណោះ មិនអាចផ្តល់ចម្លើយល្អប្រសើរចំពោះទិន្នន័យផ្សេងពីទិន្នន័យសម្រាប់រៀនឡើយ ។ ស្ថានភាពនេះអាចបកស្រាយបានថាជាករណីដែលការធ្វើបម្រាមមូលដ្ឋានលើអនុគមន៍កម្រិតលម្អៀងបានធ្លាក់ចូលទៅចំណុចនៃអប្បបរមាជ្រៀបដែលមិនមែនជាចំណុចអប្បបរមានៃដែនកំណត់ទាំងមូល ។ បញ្ហានេះកាន់តែងាយកើតឡើងនៅពេលដែលកម្រិតសេរីភាពនៃបណ្តាញពិសេសចំនួននៃប៉ារ៉ាម៉ែត្រមានកាន់តែច្រើន ។ ដោយសារតែកម្រិតសេរីភាពនៃបណ្តាញមួយបង្ហាញពីកម្រិតនៃសមត្ថភាពរបស់បណ្តាញ ដូច្នេះការកាត់បន្ថយកម្រិតសេរីភាពដើម្បីដោះស្រាយស្ថានភាពរៀនហួសកម្រិតមិនមែនជារឿងល្អឡើយបើពិនិត្យលើសមត្ថភាពទាំងមូលរបស់បណ្តាញ ។ ហេតុនេះដើម្បីកាត់បន្ថយការកើតឡើងនៃស្ថានភាពរៀនហួសកម្រិតនិងរក្សានូវសមត្ថភាពខ្ពស់របស់បណ្តាញផង ការគ្រប់គ្រងទៅលើកម្រិតសេរីនៃទម្ងន់ផ្ទាល់របស់ណឺរ៉ូនត្រូវបានសិក្សាព្រមគ្នានៅពេលដំណើរការការរៀន ។ វិធីនេះហៅថា Regularization ។ វិធីសាស្ត្រងាយមួយនៃ Regularization ប្រើចំពោះ Neural Network គឺការដាក់កម្រិតលើទម្ងន់ផ្ទាល់នៃណឺរ៉ូន (weight decay) ។

ការដាក់កម្រិតលើទម្ងន់ផ្ទាល់នៃណឺរ៉ូនអាចធ្វើបានដូចខាងក្រោម ។

$$E_t(\mathbf{w}) = \frac{1}{|\mathcal{D}_t|} \sum_{n \in \mathcal{D}_t} E_n(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

នៅទីនេះ λ ជាប៉ារ៉ាម៉ែត្រដើម្បីគ្រប់គ្រងកម្រិតនៃការដាក់កម្រិតលក្ខខណ្ឌលើទំហំនៃប៉ារ៉ាម៉ែត្រទម្ងន់ផ្ទាល់នៃណឺរ៉ូន \mathbf{w} ។ ជាទូទៅតម្លៃនេះត្រូវបានកំណត់ដោយចំនួនពិតតូចខ្លាំង $\lambda = 0.1 \sim 0.00001$ ។ ក្នុងករណីនេះការផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រតាម SGD ត្រូវបានធ្វើឡើងដូចខាងក្រោម ។

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \left(\frac{1}{|\mathcal{D}_t|} \sum_{n \in \mathcal{D}_t} \nabla E_n + \lambda \mathbf{w}^{(t)} \right)$$

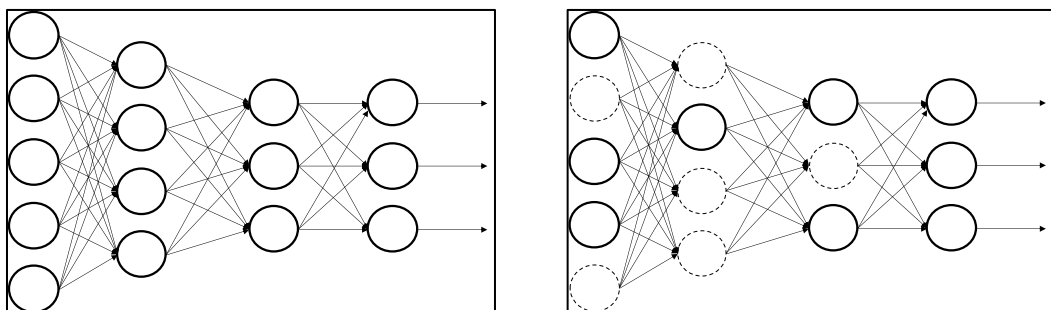
តាមកន្សោមខាងលើនេះ យើងអាចមើលឃើញថាតម្លៃនៃ ប៉ារ៉ាម៉ែត្រនឹងត្រូវបានបន្ថយជានិច្ចរាល់ពេល
នៃការផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រនៃដំណើរការរៀនត្រូវបានធ្វើឡើង ។ ហេតុនេះហើយទើបវិធីនេះហៅថា
weight decay ។ ក្រៅពីវិធីខាងលើនេះក៏មានវិធីផ្សេងដែរដូចជាការកំណត់លក្ខខណ្ឌអតិបរមានៃ
ទំហំប៉ារ៉ាម៉ែត្រទៅក្នុងចំណោមអប្បបរមានៃអនុគមន៍កម្រិតលម្អៀងដោយទម្រង់ខាងក្រោម ។

$$\|w\|^2 = \sum_i w_{ji}^2 < c$$

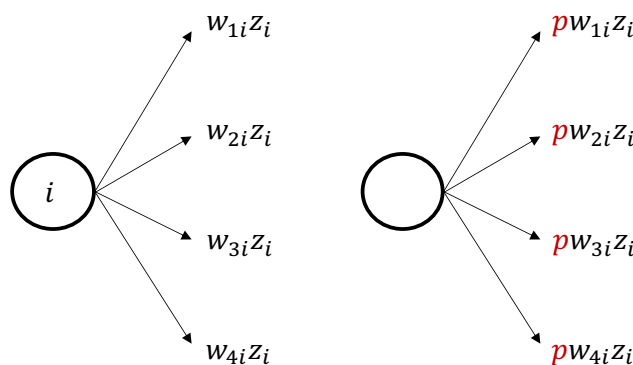
1.2. Dropout

Dropout គឺជាវិធីសាស្ត្ររៀនដែលណឺរ៉ូនត្រូវប្រើត្រូវបានជ្រើសដោយចៃដន្យជាមួយតម្លៃប្រូបាបដែលបានកំណត់ ។ រូបទី១បង្ហាញពីទម្រង់នៃFNNនៅពេលដែលណឺរ៉ូនមួយចំនួនត្រូវបានជ្រើស
ដើម្បីរៀន ។ នៅពេលរៀន (learning process) ណឺរ៉ូននៃថ្នាក់ធាតុចូលនិងថ្នាក់កណ្តាលត្រូវបាន
ជ្រើសរើសដោយប្រូបាប p ដោយចៃដន្យ ណឺរ៉ូនក្រៅពីនេះត្រូវបានចាត់ទុកថាគ្មាន រួចដំណើរការ
រៀនដូចធម្មតា ។ ប្រូបាប p អាចកំណត់ដោយតម្លៃផ្សេងគ្នានៅតាមថ្នាក់នីមួយៗបាន ។

ក្រោយពេលដំណើរការរៀនបានបញ្ចប់ ការប៉ាន់ស្មានលទ្ធផលត្រូវបានធ្វើជាមួយនឹងណឺរ៉ូន
ទាំងអស់ដោយគ្មានការជ្រើសរើសតែមួយផ្នែកដូចកាលដំណាក់កាលរៀនឡើយ ។ ប៉ុន្តែនៅករណីនេះ
តម្លៃនៃលទ្ធផលនៅថ្នាក់នីមួយៗត្រូវបានគុណនឹងតម្លៃប្រូបាបដែលបានកំណត់កាលពីពេលរៀន ។



រូបទី១ ដំណើរការនៃរៀនដោយប្រើDropout ។ (ឆ្វេង) បណ្តាញដើម (ស្តាំ) បណ្តាញពេលរៀន



រូបទី២ ដំណើរនៅពេលប៉ាន់ស្មានលទ្ធផលក្រោយពេលរៀនដោយប្រើDropout

បំណងនៃការអនុវត្ត Dropout គឺកាត់បន្ថយកម្រិតសេរីភាពនៃបណ្តាញនៅពេលដំណើរការ រៀនឱ្យតូចដើម្បីកាត់បន្ថយឱកាសនៃការកើតឡើងនូវស្ថានភាពរៀនហួសកម្រិត (overfitting) ។ ការធ្វើបែបនេះសមមូលនឹងការរៀនដោយប្រើបណ្តាញច្រើនដែលឯករាជ្យគ្នា និងដែលត្រូវបានកាត់បន្ថយ ចំនួនណឺរ៉ូន រួចធ្វើមធ្យមភាគលើលទ្ធផលនៅដំណាក់កាលប៉ាន់ស្មានលទ្ធផល។

2. វិធីសាស្ត្រជំនួយក្នុងការរៀន

2.1. ការផ្តល់ទម្រង់ស្តង់ដារលើទិន្នន័យ (standardization of data)

ក្នុងករណីដែលទិន្នន័យសម្រាប់រៀន (training data) មានស្ថានភាពលំអៀងទៅកាន់ក្រុម តម្លៃណាមួយ ការរៀននិងប៉ាន់ស្មានលទ្ធផលនឹងបានល្អប្រសើរឡើយ។ ដើម្បីដោះស្រាយបញ្ហាបែប នេះការធ្វើដំណើរការដំណាក់កាលដំបូង (preprocessing) លើសំណុំទិន្នន័យមុនពេលអនុវត្តដំណើរ ការរៀនត្រូវបានធ្វើឡើង។ ដំណើរការដំណាក់កាលដំបូងសាមញ្ញមួយគឺការធ្វើឱ្យទិន្នន័យទៅជាទម្រង់ ស្តង់ដារដែលហៅថា standardization។ ដំណើរការនេះត្រូវធ្វើឡើងដោយការបម្លែងទិន្នន័យនីមួយៗ តាមទម្រង់ខាងក្រោម។ នៅទីនេះ x_{ni} គឺជាកំប៉ូសង់ទី i នៃទិន្នន័យ x_n និង $\bar{x}_i = \frac{1}{N} \sum_{n=1}^N x_{ni}$ ។

$$\sigma_i = \sqrt{\frac{1}{N} \sum_{n=1}^N (x_{ni} - \bar{x}_i)^2}$$

$$x_{ni} \leftarrow \frac{x_{ni} - \bar{x}_i}{\sigma_i}$$

ការអនុវត្តបែបនេះជាការបន្លែងទូទៅនៅក្នុងវិធីសាស្ត្រស្ថិតិដែលក្រោយបម្លែងរួចទិន្នន័យនឹងមានតម្លៃ មធ្យមនៃកំប៉ូសង់នីមួយៗស្មើ 0 និងរ៉ាង្ស្យង់ស្មើ 1 ។

2.2. របៀបកំណត់កម្រិតរៀន (learning rate)

នៅក្នុងវិធីសាស្ត្ររៀនដូចជា gradient descend ឬ stochastic gradient descend (SGD) ការផ្លាស់ប្តូរនិងកំណត់តម្លៃនៃប៉ារ៉ាម៉ែត្រប្រែប្រួលទៅតាមតម្លៃកម្រិតរៀន (learning rate) ដែលបានកំណត់។ នៅក្នុងដំណើរការរៀនមានវិធីសាស្ត្រចំបង២ដែលត្រូវបានប្រើដើម្បីកំណត់តម្លៃ កម្រិតរៀន។

វិធីសាស្ត្រទី១ គឺការកំណត់តម្លៃកម្រិតរៀនឱ្យធំនៅពេលចាប់ផ្តើមរៀន រួចបន្ថយតម្លៃនេះ បន្តិចម្តងៗនៅពេលដំណើរការរៀនបន្ត។ ជាឧទាហរណ៍ការកំណត់កម្រិតរៀននៅដំណាក់កាលរៀនទី t ដោយ $\eta_t = \eta_0 - \alpha t$ ដែល η_0 ជាតម្លៃកម្រិតរៀននៅពេលចាប់ផ្តើមនិង α កម្រិតនៃការបន្ថយ។

វិធីសាស្ត្រទី២ គឺការកំណត់តម្លៃកម្រិតរៀនខុសៗគ្នានៅតាមថ្នាក់ផ្សេងៗគ្នានៃបណ្តាញ ។ របៀបនៃការកំណត់អាចប្រែប្រួលទៅតាមទម្រង់នៃបណ្តាញ ។ ជាឧទាហរណ៍ការកំណត់កម្រិតរៀន ឱ្យធំនៅថ្នាក់ដែលកែវនឹងថ្នាក់ធាតុចូល (input layer) និងកំណត់តម្លៃកាន់តែតូចនៅថ្នាក់ដែលនៅ កែវថ្នាក់លទ្ធផលចុងក្រោយ (output layer) ។

វិធីសាស្ត្រនៃការកំណត់តម្លៃកម្រិតរៀនដោយស្វ័យប្រវត្តិដូចជា AdaGrad ក៏ត្រូវបានប្រើជា ញឹកញាប់ផងដែរ ។ ឧបមាថា gradient នៃអនុគមន៍កម្រិតលម្អៀងកំណត់ដោយ $g_t = \nabla E_n$ និង កំប៉ូសង់ទី i កំណត់ដោយ $g_{t,i}$ ។ បើតាម SGD ប៉ារ៉ាម៉ែត្រនឹងត្រូវផ្លាស់ប្តូរតម្លៃដោយ $-\eta g_{t,i}$ ប៉ុន្តែ AdaGrad វាផ្លាស់ប្តូរដោយតម្លៃខាងក្រោម ។ ការផ្លាស់ប្តូរបែបនេះអាចបកស្រាយដោយងាយថាជាការ ផ្ដោតការយកចិត្តទុកដាក់លើតម្លៃកំប៉ូសង់ណាដែលការផ្លាស់ប្តូរមិនសូវបានអនុវត្តជាញឹកញាប់ ។

$$-\frac{\eta}{\sqrt{\sum_{k=1}^t g_{k,i}^2}} g_{t,i}$$

2.3. វិធីសាស្ត្រ Momentum

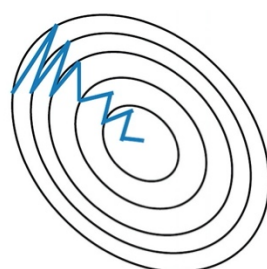
វិធីសាស្ត្រនៃការកំណត់តម្លៃកម្រិតរៀនដោយស្វ័យប្រវត្តិផ្សេងទៀតដែលមានប្រសិទ្ធភាពគឺ វិធីសាស្ត្រ Momentum ។ វិធីនេះត្រូវបានដឹងថាជួយឱ្យដំណើរការផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រក្នុង gradient descend ឆាប់ចប់ ។ ក្នុងវិធីសាស្ត្រនេះប៉ារ៉ាម៉ែត្រត្រូវបានផ្លាស់ប្តូរតាមទម្រង់ខាងក្រោម ។

$$\begin{aligned} \Delta w^{(t-1)} &= \Delta w^{(t)} - \Delta w^{(t-1)} \\ w^{(t+1)} &= w^{(t)} - \eta \nabla E_t + \mu \Delta w^{(t-1)} \end{aligned}$$

នៅទីនេះ μ ច្រើនត្រូវបានកំណត់ដោយ $\mu = 0.5 \sim 0.9$ ។ វិធីនេះត្រូវបានបកស្រាយថាជាវិធីដែល ជួយគ្រប់គ្រងដល់ការប្រែប្រួលនៃការផ្លាស់ប្តូរតម្លៃប៉ារ៉ាម៉ែត្រកុំឱ្យមានការប្រែប្រួលខ្លាំងដែលប៉ះពាល់ ដល់ការរួចតូចនៃតម្លៃកម្រិតលម្អៀងនៅដំណាក់កាលរៀន ។



Stochastic Gradient Descent **without** Momentum



Stochastic Gradient Descent **with** Momentum

រូបទី៣ ភាពខុសគ្នានៃបម្រែបម្រួលប៉ារ៉ាម៉ែត្រដោយប្រើSGDដែលមានឬគ្មាន momentum

(Ref: Stochastic Gradient Descent on your microcontroller by SIMONE

<https://eloquentarduino.github.io/2020/04/stochastic-gradient-descent-on-your-microcontroller/>)

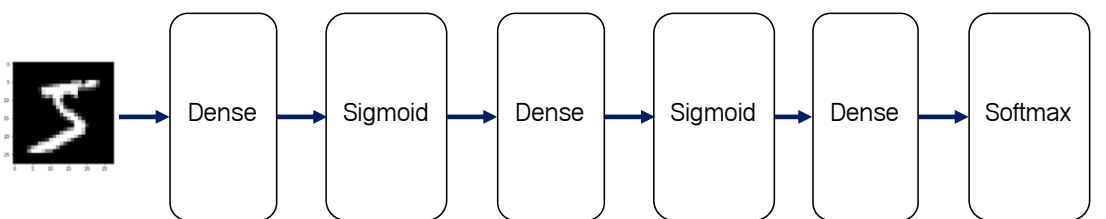
Convolutional Neural Network (CNN)

ក្នុងអត្ថបទនេះ យើងនឹងសិក្សាអំពីប្រភេទមួយផ្សេងទៀតនៃ Neural Network ដែលហៅថា Convolutional Neural Network (CNN) ។ CNN ត្រូវបានប្រើប្រាស់យ៉ាងទូលំទូលាយក្នុងការអនុវត្តលើបច្ចេកវិទ្យាកំណត់សម្គាល់រូបភាពឬកំណត់សម្គាល់សម្លេង និងប្រកបដោយប្រសិទ្ធភាពខ្ពស់បើធៀបនឹងបណ្តាញ FNN ធម្មតា ។ ក្នុងបច្ចេកវិទ្យាកំណត់សម្គាល់រូបភាពឬ បច្ចេកវិទ្យានានាទាក់ទងនឹងការប្រើប្រាស់រូបភាពក្នុងបញ្ហាសិប្បនិម្មិតត្រូវបានមើលឃើញថាមានការប្រើប្រាស់ CNN ជាមូលដ្ឋាន ។ ថ្មីៗនេះទៀត សូម្បីក្នុងបច្ចេកវិទ្យាភាសា (Natural Language Processing) បណ្តាញ CNN ក៏ត្រូវបានប្រើប្រាស់ផងដែរ ។

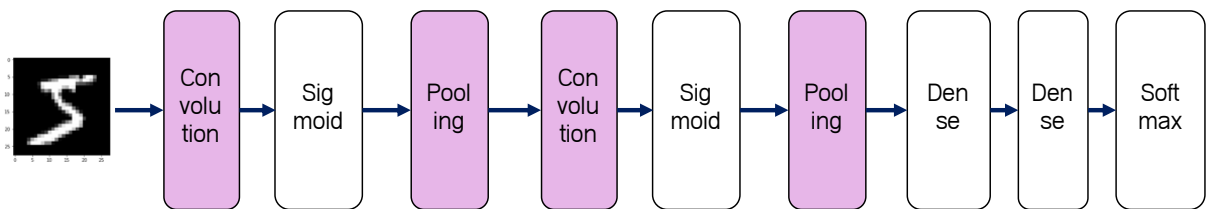
1. ទម្រង់និងដំណើរការទូទៅរបស់ CNN

ដើម្បីស្វែងយល់អំពី CNN ជាដំបូងយើងនឹងពិនិត្យលើទម្រង់ជាទូទៅរបស់វាជាមុន ។ CNN ក៏ដូចជាបណ្តាញ Neural Network ដទៃដែរគឺកើតឡើងពីការផ្គុំបញ្ចូលគ្នានូវថ្នាក់ណឺរ៉ូនជាច្រើនថ្នាក់ ។ ប៉ុន្តែចំណុចពិសេសគឺនៅក្នុង CNN មានថ្នាក់ពិសេសពីដែលមិនមាននៅក្នុង Feedforward Neural Network (FNN) ដែលយើងសិក្សាកន្លងមក គឺ ថ្នាក់ Convolution និង ថ្នាក់ Pooling ។ យើងនឹងសិក្សាលម្អិតអំពីថ្នាក់ទាំងពីរនេះនៅចំណុចបន្ទាប់ ។

ក្នុងបណ្តាញ FNN ណឺរ៉ូននៃថ្នាក់ដែលនៅជាប់គ្នាត្រូវបានភ្ជាប់គ្នាទាំងអស់ដោយទម្ងន់ផ្ទាល់របស់វា ។ ទម្រង់បែបនេះគេហៅថា តំណភ្ជាប់ពេញទី (fully-connected) ។ ក្នុងរូបនិងការបកស្រាយខាងក្រោម យើងនឹងបង្ហាញតំណភ្ជាប់ពេញទីបែបនោះដោយ Dense រីឯអនុគមន៍សកម្មនឹងបង្ហាញដោយឈ្មោះនៃអនុគមន៍នោះ ។ ក្នុងករណីនេះ ទម្រង់នៃបណ្តាញ FNN អាចបង្ហាញបានដូចរូបទី១ ។ បណ្តាញ FNN អាចត្រូវបាននិយាយថាជាបណ្តុំនៃតំណភ្ជាប់ពេញទីនិងអនុគមន៍សកម្មច្រើនថ្នាក់បន្តគ្នា ។ ផ្ទុយពីនេះ ទម្រង់នៃបណ្តាញ CNN អាចបង្ហាញបានដូចរូបទី២ ។ ដូចដែលអ្នកបានឃើញថ្នាក់ថ្មី (Convolution និង Pooling) ត្រូវបានប្រើជាបន្តបន្ទាប់ មុនពេលដែលបណ្តុំនៃតំណភ្ជាប់ពេញទីនិងអនុគមន៍សកម្មត្រូវបានអនុវត្ត ។ ទម្រង់ដែលថ្នាក់ Convolution និង Pooling ត្រូវបានប្រើមុនត្រូវបានមើលឃើញជាទូទៅក្នុងបណ្តាញ CNN ។



រូបទី១ ទម្រង់នៃបណ្តាញ FNN ដោយប្រើតំណភ្ជាប់ពេញទី (fully-connected layer)



រូបទី២ ទម្រង់នៃCNN

2. អំពីថ្នាក់ Convolution

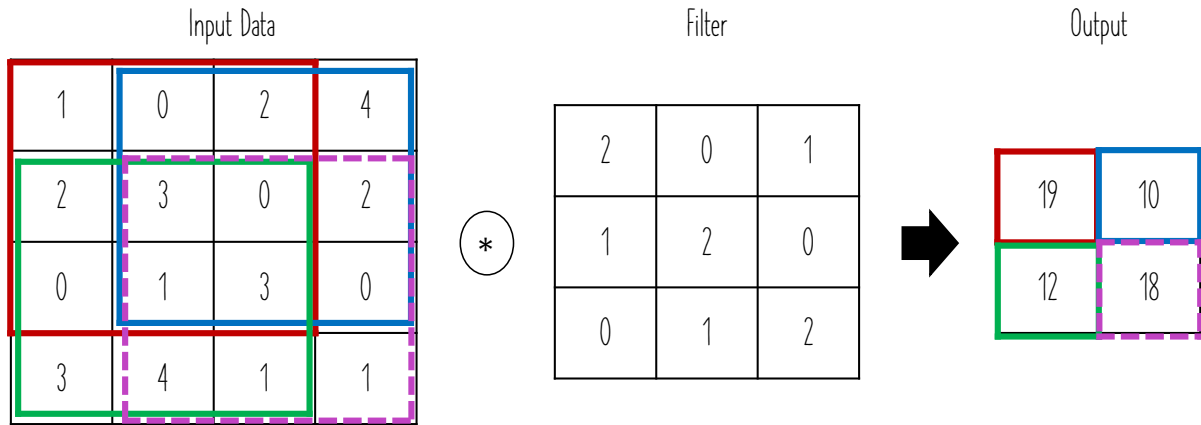
ជាដំបូងយើងនឹងពិនិត្យលើភាពខុសគ្នារវាងថ្នាក់តំណភ្ជាប់ពេញទី (fully-connected) ដែលត្រូវបានប្រើក្នុងFNNនិងថ្នាក់Convolutionនេះ ។

2.1. បញ្ហានៃថ្នាក់តំណភ្ជាប់ពេញទី

នៅក្នុងបណ្តាញFNNណឺរ៉ូននៃថ្នាក់តំណភ្ជាប់ពេញទីត្រូវបានភ្ជាប់គ្នាគ្រប់ណឺរ៉ូនទាំងអស់ដែលមាននៅថ្នាក់ជាប់គ្នា រីឯចំនួននៃណឺរ៉ូននៅថ្នាក់លទ្ធផលអាចត្រូវបានកំណត់ដោយសេរី។

បញ្ហានៃថ្នាក់តំណភ្ជាប់ពេញទីគឺ ទម្រង់នៃទិន្នន័យមិនត្រូវបានប្រើប្រាស់ឱ្យសមស្របនោះឡើយ ។ ឧទាហរណ៍បើធាតុចូលជារូបថតដូចរូបទី១និងទី២ នោះធាតុចូលស្ថិតនៅក្នុងទម្រង់វិមាត្រ៣ពេល គឺ បណ្តោយ ទទឹង និង channel (RGB) ។ បើធាតុចូលនេះ ត្រូវបានបញ្ជូនទៅក្នុងFNNដែលមានថ្នាក់តំណភ្ជាប់ពេញទី នោះដូចដែលបានបង្ហាញនៅឧទាហរណ៍នៃមេរៀនអំពីFNN ទិន្នន័យនេះត្រូវបានពន្លាតជា១វិមាត្រដោយគម្រៀបគ្រប់តម្លៃនៃPixelទាំងអស់ជាកំប៉ូសង់នៃវ៉ិចទ័រ១វិមាត្រនោះ។ ក្នុងឧទាហរណ៍ដែលយើងបានសិក្សាកន្លងមក រូបថតដែលមានទម្រង់ (28,28,1) ត្រូវបានប្រើជាវ៉ិចទ័រ១វិមាត្រដែលមានកំប៉ូសង់ចំនួន784 (=28x28x1) ។ ហេតុនេះការប្រើប្រាស់ថ្នាក់តំណភ្ជាប់ពេញទីធ្វើឱ្យខាតបង់ព័ត៌មានដែលយើងអាចទាញយកពីទម្រង់ដើមនៃទិន្នន័យដូចជាទិន្នន័យអំពីចម្ងាយនៃរបស់ក្នុងរូបដែលបង្ហាញដោយទម្រង់ RGB ច្រើនchannelជាដើម ។

ផ្ទុយពីនេះ ថ្នាក់Convolution ទទួលយកធាតុចូលដោយរក្សាទម្រង់ដើមនៃទិន្នន័យ។ ក្នុងករណីរូបថតទម្រង់៣វិមាត្រដូចពណ៌នាខាងលើ ថ្នាក់Convolutionនឹងទទួលយករូបថតវិមាត្រ៣នោះជាធាតុចូលនិងបញ្ចេញរូបថតវិមាត្រ៣ជាលទ្ធផល។ ហេតុនេះ CNN អាចមានលទ្ធភាពក្នុងការសិក្សាបានស៊ីជម្រៅលើលក្ខណៈពិសេសរបស់ទិន្នន័យបាន។ នៅក្នុងCNN ធាតុចូលនៃថ្នាក់Convolutionហៅថា ផែនទីលក្ខណៈចូល (input feature map) និងធាតុចេញ (លទ្ធផល) របស់វាហៅថា ផែនទីលក្ខណៈចេញ (output feature map) ។ ក្នុងករណីខ្លះយើងហៅវាជារួមថា feature map ។



រូបទី៣ ឧទាហរណ៍នៃប្រមាណវិធីConvolution(\otimes)

2.2. ប្រមាណវិធី Convolution

ក្នុងថ្នាក់Convolution ប្រមាណវិធីConvolutionត្រូវបានអនុវត្ត។ ប្រមាណវិធីនេះអាចប្រៀបបានជាដំណើរការនៃFilterនៅក្នុងបច្ចេកវិទ្យារូបភាព(Image Processing) ។ ការធ្វើប្រមាណវិធី Convolution លើរូបភាពដែលមានទំហំ $W \times W$ និងFilter ដែលមានទំហំ $H \times H$ គឺកំណត់ដោយទម្រង់ខាងក្រោម។ u_{ij} ជាតម្លៃលទ្ធផលនៅទីតាំង (i, j) និង x_{st} ជាតម្លៃនៅPixel (s, t) និង h_{pq} ជាតម្លៃFilterនៃទីតាំង (p, q) ។

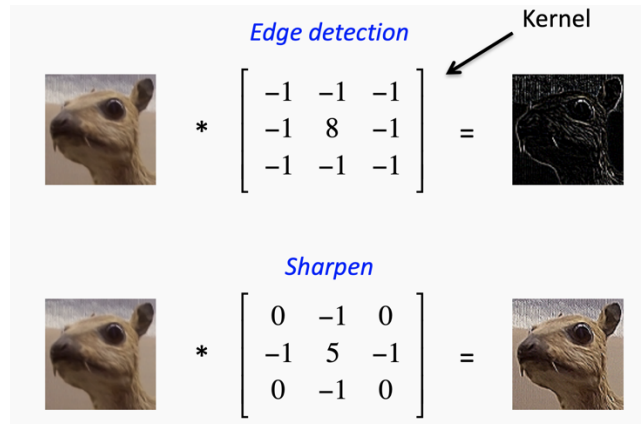
$$u_{ij} = \sum_{p=0}^{H-1} \sum_{q=0}^{H-1} x_{i+p, j+q} h_{pq}$$

រូបទី៣បង្ហាញពីឧទាហរណ៍នៃប្រមាណវិធីConvolution។ ឧបមាថាទិន្នន័យជាធាតុចូលមានទម្រង់ ២វិមាត្រ(height, width) ដែលមានទំហំ(4,4) និង Filterមានទំហំ(3,3)នោះលទ្ធផលធាតុចេញមានទំហំ(2,2)។ របៀបនៃការគណនាគឺ ចន្លោះនៃទិន្នន័យដែលមានទំហំដូចFilterត្រូវបានគណនាដោយធ្វើផលគុណគ្រប់ធាតុដែលមានទីតាំងត្រូវគ្នាបន្តបន្ទាប់ប្រមាណវិធីបូកបញ្ចូលគ្នា។ ឧទាហរណ៍នៅក្នុងចន្លោះពណ៌ក្រហម និងពណ៌បៃតងការគណនាត្រូវបានអនុវត្តដូចខាងក្រោម។

$$1 \times 2 + 0 \times 0 + 2 \times 1 + 2 \times 1 + 3 \times 2 + 0 \times 0 + 0 \times 0 + 1 \times 1 + 3 \times 2 = 19$$

$$2 \times 2 + 3 \times 0 + 0 \times 1 + 0 \times 1 + 1 \times 2 + 3 \times 0 + 3 \times 0 + 4 \times 1 + 1 \times 2 = 12$$

ដំណើរការConvolution នេះអាចទាញយកនូវលក្ខណៈឬទម្រង់ជាក់លាក់ដែលមាននៅទីតាំងផ្សេងៗនៃរូបភាពបានតាមការកំណត់Filterសមស្រប។ រូបទី៤បង្ហាញពីឧទាហរណ៍នៃលក្ខណៈនោះ។



រូបទី៤ ឧទាហរណ៍នៃដំណើរការទាញយកFeature MapដោយConvolution

(Image from: Simple Introduction to Convolutional Neural Networks

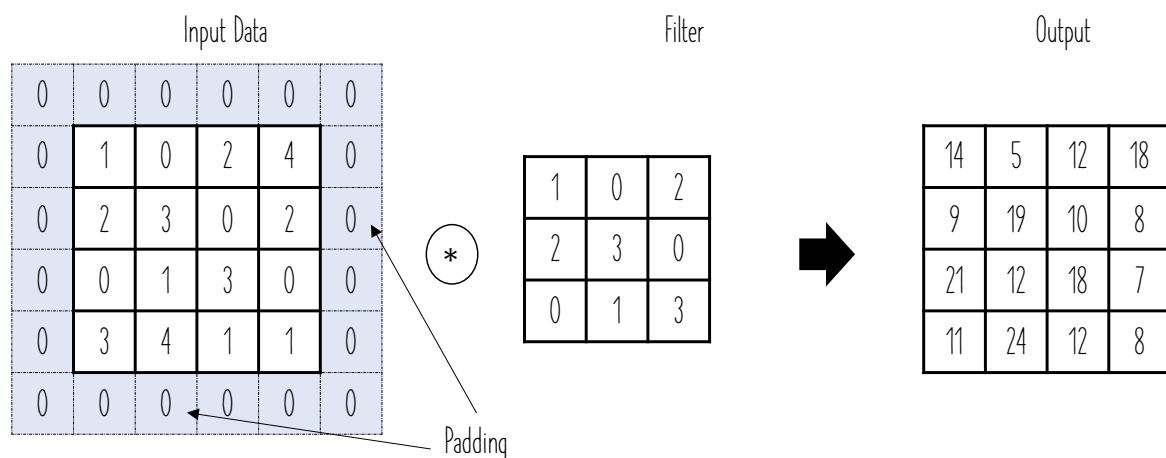
<https://towardsdatascience.com/simple-introduction-to-convolutional-neural-networks-cdf8d3077bac>)

2.3. ការតម្រឹម Padding

ដូចដែលបានបង្ហាញខាងលើ តាមរយៈប្រមាណវិធីConvolution ទំហំនៃលទ្ធផលដែលទទួលបានមានទំហំតូចជាធាតុចូលដើមជានិច្ច។ ទំហំនៃលទ្ធផលនោះកំណត់ដោយទម្រង់ខាងក្រោម។ នៅទីនេះ $\lfloor x \rfloor$ ជាតម្លៃចំនួនគត់ធំបំផុតដែលតូចជាង x ។

$$\left(w - 2 \left\lfloor \frac{H}{2} \right\rfloor\right) \times \left(w - 2 \left\lfloor \frac{H}{2} \right\rfloor\right)$$

ប៉ុន្តែក្នុងដំណើរការវិភាគទិន្នន័យ នៅដំណាក់កាលខ្លះការរក្សាទំហំនៃធាតុចូលនិងលទ្ធផលឱ្យដូចគ្នាជារឿងចាំបាច់។ ក្នុងករណីនេះដំណើរការនៃតម្រឹម (Padding) ត្រូវបានប្រើ។ ដំណើរការនេះគឺបន្ថែមតម្លៃនៅផ្នែកជ្រុងដោយជុំវិញនៃរូបភាពដោយតម្លៃសូន្យ។ រូបទី៥បង្ហាញពីគម្រោងដំណើរការនេះ។



រូបទី៥ ឧទាហរណ៍នៃដំណើរការ Padding

2.4. ប្រមាណវិធីរំកិល Sride

ក្នុងការគណនាខាងលើ តំបន់ដែលប្រើដើម្បីអនុវត្តFilterត្រូវបានរំកិលម្តងៗពីឆ្វេងទៅស្តាំ និងពីលើចុះក្រោម។ ចន្លោះឬទំហំដែលត្រូវប្រើក្នុងដំណើរការរំកិលតំបន់ដើម្បីអនុវត្តFilter នេះហៅថា Sride។ ក្នុងករណីសាមញ្ញខាងលើ តម្លៃនៃ Sride ត្រូវបានកំណត់ដោយ 1 ។ ជាទូទៅក្នុងករណីដែល តម្លៃ Sride គឺ s លទ្ធផលនៃ Convolution ត្រូវបានកំណត់ដោយទម្រង់ខាងក្រោម។

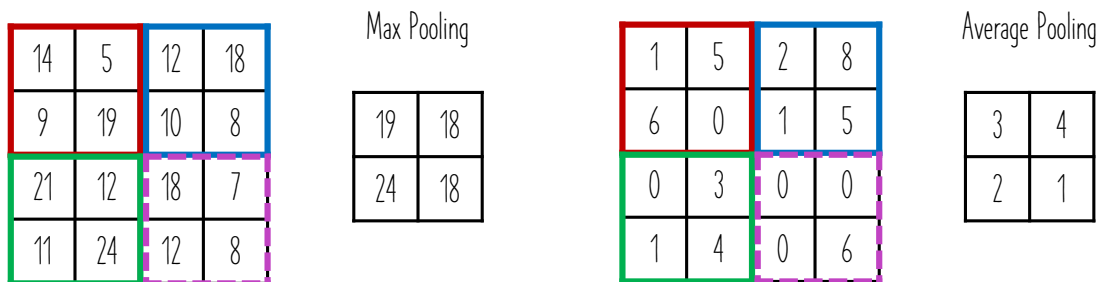
$$u_{ij} = \sum_{p=0}^{H-1} \sum_{q=0}^{H-1} x_{si+p, sj+q} h_{pq}$$

ក្នុងករណីដែលប្រើ Padding ជាមួយផង លទ្ធផលនឹងផ្តល់ឱ្យដោយទំហំខាងក្រោម។

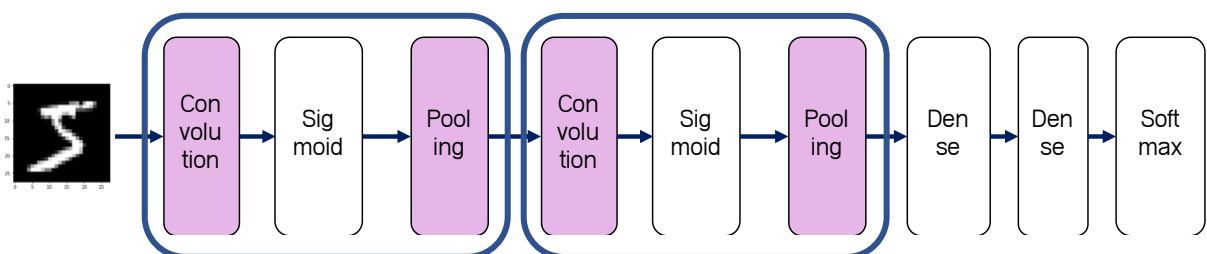
$$\left(\left\lfloor \frac{W-1}{s} \right\rfloor + 1\right) \times \left(\left\lfloor \frac{W-1}{s} \right\rfloor + 1\right)$$

3. ថ្នាក់ Pooling

ក្នុងថ្នាក់ Pooling ប្រមាណវិធីដែលអនុវត្តគឺការបង្រួមទំហំនៃទិន្នន័យឱ្យរួមតូចជាងមុនដោយ កំណត់យកតម្លៃតំណាងនៅក្នុងតំបន់។ ជាទូទៅថ្នាក់ Pooling មានពីរផ្នែកលើប្រភេទនៃតម្លៃតំណាង ដែលកំណត់យក ពោលគឺ តម្លៃអតិបរមាក្នុងតំបន់ និងតម្លៃមធ្យមក្នុងតំបន់។ ក្នុងករណីតម្លៃអតិបរមា ត្រូវបានប្រើប្រាស់ត្រូវបានហៅថា Max Pooling ឯករណីតម្លៃមធ្យមត្រូវបានប្រើហៅថា Average Pooling។ លក្ខណៈពិសេសនៅដំណាក់កាលនេះគឺ គ្មានប៉ារ៉ាម៉ែត្រដែលត្រូវរៀន។



រូបទី៦ ឧទាហរណ៍នៃដំណើរការ Pooling



រូបទី៧ ទម្រង់នៃ CNN

4. ទម្រង់នៃបណ្តាញCNNនិងការរៀន (Learning)

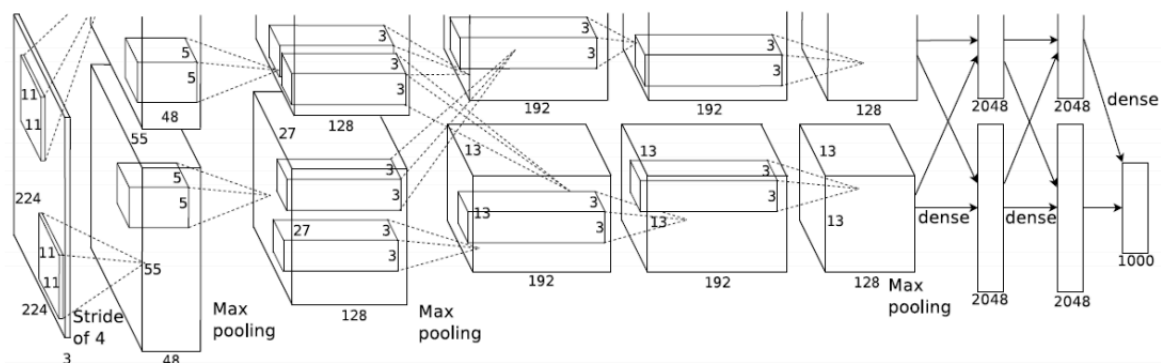
ដូចបង្ហាញខាងដើម CNN ត្រូវបានប្រើប្រើនក្នុងបច្ចេកវិទ្យារូបភាព ។ ទម្រង់នៃបណ្តាញនេះ ច្រើនកំណត់ដោយទម្រង់ដូចរូបទី៧ ដោយប្រើថ្នាក់ដែលកើតពីConvolutionនិងPoolingច្រើនដង និងប្រើថ្នាក់ដែលកើតពីតំណភ្ជាប់ពេញទីចុងក្រោយ ។ ចំពោះចំណោទធ្វើចំណាត់ថ្នាក់រូបភាពច្រើន ក្រុម យើងអាចបកស្រាយបានថា ថ្នាក់ដែលប្រើConvolutionនិងPooling ជាថ្នាក់ដែលទាញយក លក្ខណៈសម្គាល់ពិសេស (feature map) ពីរូបភាព រីឯថ្នាក់កើតពីតំណភ្ជាប់ (Dense) ត្រូវបានប្រើ ដើម្បីធ្វើចំណាត់ថ្នាក់ក្រុមដូចក្នុងFNNដែរ ។

ដូចដែលអ្នកអាចធ្វើការកត់សម្គាល់បាន ថ្នាក់នីមួយៗនៃCNNត្រូវបានតភ្ជាប់គ្នាជាបន្តបន្ទាប់ ប្រៀបដូចជាថ្នាក់នៃFNNដែរ ។ ហេតុនេះការកំណត់តម្លៃប៉ារ៉ាម៉ែត្រនៃថ្នាក់នីមួយៗ (ដំណើរការរៀន) អាចធ្វើបានដោយប្រើ Backpropagation Algorithmដូចក្នុងករណីFNNផងដែរ ។ អ្នកអាចប្រើប្រាស់ FrameWork ដែលផ្តល់ឱ្យប្រើជាសេរីជាច្រើនដើម្បីបង្កើតម៉ូដែលCNNបានដូចជា Keras, PyTorch, Tensorflow ជាដើម ។

5. អនុវត្តន៍នៃបណ្តាញCNN

CNNត្រូវបានអភិវឌ្ឍជាបណ្តាញជំនិងស្មុគស្មាញដើម្បីដោះស្រាយបញ្ហានានាពិសេសក្នុង បច្ចេកវិទ្យាធ្វើកំណត់សម្គាល់វត្ថុក្នុងរូបភាព ។ រូបទី៨បង្ហាញពីគម្រោងម៉ូដែលដែលប្រើCNNដើម្បីធ្វើ ចំណាត់ថ្នាក់1000ក្រុមនៃវត្ថុក្នុងរូបភាពដែលត្រូវបានស្គាល់ថាផ្តល់នូវលទ្ធផលល្អប្រសើរខ្លាំង ។ សម្រាប់ចំណេះដឹងបន្ថែមទាក់ទងនឹងការប្រកួតប្រជែងបង្កើតម៉ូដែលសម្រាប់ធ្វើចំណាត់ថ្នាក់ក្រុមវត្ថុ ក្នុងរូបភាព មិត្តអ្នកអានអាចចូលទៅកាន់តំណភ្ជាប់ខាងក្រោម ។

<http://www.image-net.org/challenges/LSVRC/>



Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. 2012. ImageNet classification with deep convolutional neural networks. In Proc. of NIPS, pp. 1097-1105.

រូបទី៨ គម្រោងម៉ូដែលដែលប្រើCNN