[모두를위한 딥러닝 시즌1. Sung Kim]

## Lecture 1 ML Basic.

- Imitation of explicit programming.
- $\Big[$ supervised learning : learning with labeled example
- $\Big[$ unsupervised learning : un-labeled data

## MLlab1 Tensorflow.

→ open source software library for numerical computation using data flow graphs.

- Graph ; (Node) & (edges).
  $$\Big( \text{mathematical operations} \quad \text{tensors communicated between them} \Big)$$

- TensorFlow mechanics
  ① Build graph using TF operations
    ↳ can use placeholders !! (& feeddict)
  ② feed data and run graph (operation)
    sess.run(op, feed_dict = {x: x_data}
  ③ update variables in the graph
    (& return values)

- Tensor rank / shape / Types

| | | |
|---|---|---|
| 0 | [] | Scalar |
| 1 | [D0] | 1 Dvector |
| 2 | [D0,D1] | matrix |
| 3 | [D0,D1,D2] | n tensor. |
| ⋮ | | |

## Lecture 2,3,4 Linear regression.

- $H(x) = Wx + b$ & cost function. $\frac{1}{m}\sum_{i}^{m}(H-y)^2$
  $$\boxed{\begin{array}{c}\text{minimize cost}(W,b)\\ W,b\end{array}}$$

- gradient descent — manual update
  for convex — Gradient Descent Optimizer
  cost fth.

---

DL ① (Sung Kim)

## Lecture 5 logistic classification?

~~Regression~~ Regression R → R
Classification R → B^n $(B = \{0,1\})$
  ↳ Facebook feed
    Credit Card Fraudulent, Spam
- (0/1) encoding Radiology, Finance,



- Logistic Hypothesis → $g(z) = \frac{1}{1+e^{-z}}$ ; sigmoid.
  $= \frac{1}{1+e^{W^T x}}$

- cost function? = cross-entropy
  for sigmoid function?
  $$C(H(x), y) = y\log(H(x)) - (1-y)\log(1-H(x))$$

- Grad. Desc. ⇒ $\Delta W = -\alpha \frac{\partial}{\partial W}\text{cost}(W)$

## (Lab5) logistic Classifier!

hypothesis = tf.sigmoid(tf.matmul(X,W)+b)
cost = tf.reduce_mean(Y * tf.log(hypoth))
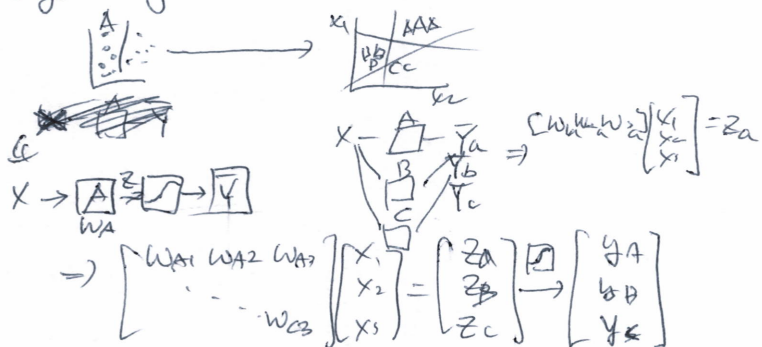                    + (1-Y) * tf.log(1-hyp))
train = tf.train.grad. optimizer(learning rate
predicted = tf.cast(hyp > 0.5)          = 0.01)
sess.run(tf.global_init ∨)  minimize(cost)
for step in range(~);
    cost_val, _ = sess.run[cost, train], feed_dict;

## Lecture 6 : Softmax classification (multinomial)

logistic regression   Multinomial classification



$\Rightarrow \begin{bmatrix} W_{A1} & W_{A2} & W_{A3} \\ & & \\ & & W_{C3} \end{bmatrix}\begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} Z_A \\ Z_B \\ Z_C \end{bmatrix} \boxed{S} \begin{bmatrix} Y_A \\ Y_B \\ Y_C \end{bmatrix}$

- softmax → $y \begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix} \rightarrow \left[ S(y_i) = \dfrac{e^{y_i}}{\sum_j e^{y_i}} \right] \rightarrow \begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$ (probabilities)

  Scores

  $\vdots$

  (one-hot encoding) $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

- Cost function = Cross entropy.

  $S(Y) \rightarrow \boxed{D(S, L) = -\sum_i L_i \log(S_i)}$

  $\begin{bmatrix} S(Y) & Y \\ \overset{=}{\hat{Y}} & \end{bmatrix}$ , data label or target , predicted probability logit

  (label) $Y = L = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow B$

  (case1) $\bar{Y} = \begin{matrix} A \\ B \end{matrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \rightarrow B$ ; $L$ (label) $\log$ (prediction) (cost) cost $\begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot \left( -\log \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \to 0$

  (case2) $\bar{Y} = \begin{matrix} A \\ B \end{matrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \rightarrow A$ ; $\begin{bmatrix} 0 \\ 1 \end{bmatrix} \odot \left( -\log \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0 \\ \infty \end{bmatrix} \to \infty$

  (same) $\begin{cases} \text{logistic cost} ; & C(H(x), y) = y\log(H(x)) - (1-y)\log(1-H) \\ \text{cross entropy} ; & D(S, L) = -\sum_i L_i \log(S_i) \end{cases}$

  (prediction) (label)

  ↳ total loss

  $\mathcal{L} = \dfrac{1}{N} \sum_i D_i(S, L) = \dfrac{1}{N} \sum_i \sum_j L_{ij} \log(S_{ij})$

  training data index.

  ### ▢ Lab6 Softmax classifier.

  - hypothesis $= tf.matmul(X, W) + b$

    $\downarrow$ (linear) logit

    $tf.nn.softmax(tf.matmul(X.W) + b)$

    logit.

  - loss function

    cost = tf.reduce_mean(-tf.reduce_sum

    (Y * tf.log(hyp) (axis=1))

    opt = tf.train.Grad.optim(learnrate=0.1)

    - minimize (cost)

    [axis=0 ⇒ data index]

- Test & one-hot encoding

  hyp = tf.nn.softmax(tf.matmul(X.W) + b)

  a = sess.run(hyp, feed-dict X )

  print (a, sess.run(tf.arg-max(a, 1)))

- Fancy Softmax Classifier

  $\begin{cases} logits = tf.matmul(X.W) + b, \quad hyp = softmax \\ cost = tf.reduce\_mean(-tf.reduce\_sum( \\ \qquad Y * tf.log(hyp), axis=1) \end{cases}$

  cost_i = tf.nn.softmax_cross_entropy_with_logits

  (logits = logits, labels = Y_one_hot)

  cost = tf.reduce_mean(cost_i)

- Discrete label into onehot vector!

  ① Y = tf.placeholder(tf.int32, [None, 1]) ← shape=(?, 1)

  ② Y_one_hot = tf.one_hot(Y, nb_classes)

  ③ Y_one_hot = ~~(Y_one_hot,~~ ⇒ shape=(?, 1, 7) 9=7 (label index size)

     tf.reshape(Y_one_hot,

  ① [1, 3, 0, 6] [-1, nb_classes])

  ② [[0100000] [0001000] [1000000] [0000001]] ⇒ Shape (?, 7)

  ③ [0100000, 0001000, 1000001, 0000001]

  (? = # of training examples)

  ## Lecture 7. Learning Rate // Data Preprocessing // Overfitting

  (LR)

- Gradient Descent

  ↳ Overshooting problem : large LR.

  Small learning rate ; ↓ takes too long. stops at local minimum

  try several LR.

- Data (X) Preprocessing for grad des.

  

  → $X_1 \& X_2$ scales too different

  ⭐ Standardization $x'_j = \dfrac{x_j - \mu_j}{\sigma_j}$

  X_std[:, 0] = (X[:, 0] - X[:, 0].mean(1)) / X[:, 0].std()
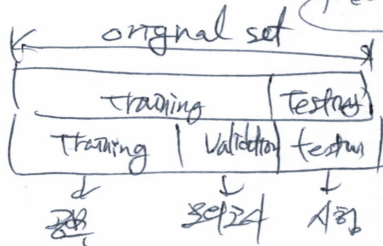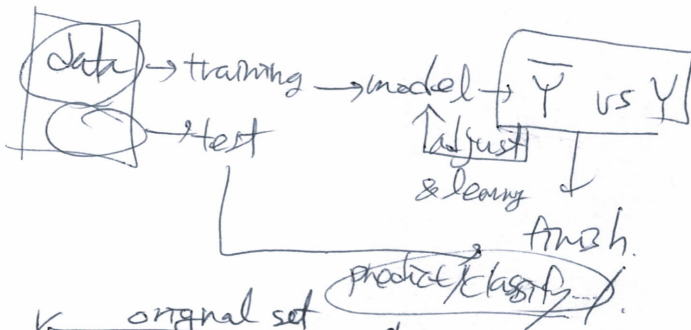
• **Overfitting Problem**

solution ↳/ More training data
  Reduce the # of parameters
  Regularization.

  ↳ # Let's not have too big numbers in the weight

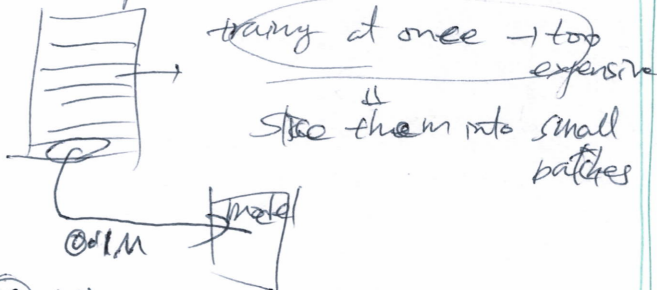  $$L = \frac{1}{N}\sum_i D(S(Wx_i+b), L_i) + \lambda \sum W^2$$

  ↑ # regularization strength.

  l2reg = 0.001 X
  tf.reduce_sum(tf.square(W))

• **Performance Evaluation ?**

Data → training → model → Ŷ vs Y
  → test        ↑ Adjust
              & learn ↓
                  Finish.
              product/classify...

  original set
  | Training | Testing |
  | Training | Validation | testing |
  학습      검증      시험

• **Online learning ?**

100M examples

  train at once → too expensive
         "
  Slice them into small batches

@01M  model

ex MNIST → Accuracy ?!

**Lab7 Learning rate & Evaluation**

• ↳ X-data = [[...],[...],[...]...[...]]
  y-data = [[∼],[∼],[∼]...[∼]] (one hot)

② X_test = ∼
  y_test = ∼

• learning rate → in  tf.train.GradDes(learrate)
                      = xx.min(cost)

---

• non-normalized input ?    M③

  xy = MinMaxScaler (xy)

• **MNIST.**
  [28×28×1] Image → X = tf.placeholder(tf.float, 32,
                  [None, 784])
                  Y = tf.placeholder(tf.float,
                  [None, nb_classes])
                              = 10

from tensorflow.examples.tutorials.mnist import
                                  input_data
mnist = input_data.read_data_sets('MNIST_data/,
                          onehot=True,
batch_xs, batch_ys = mnist.train.next_batch (100
  ⋮
print("Accuracy #:", accuracy.eval(session=sess,
            feed_dict = {X: mnist.test.images,
                        Y: mnist.test.labels}

testing
is_correct = tf.equal(tf.arg_max(hyp,1),
                  tf.arg_max(Y,1))

• accuracy = tf.reduce_mean(tf.cast  T/F 1 ↦0.0
                          (is_correct, tf.float))

training epochs = 15
batch_size = 100  → Full data set execution.
                    training
                  → Training unit for
with tf.Session() as sess:    a single weight
  sess.run(tf ∼ initialize)      update.
  for epoch in range(training epochs):
    # avg_cost = 0
    total_batch = int(mnist.train.num_examples/
    for i in range(total_batch): batch size
      batch_xs, batch_ys = mnist.train.next_batch
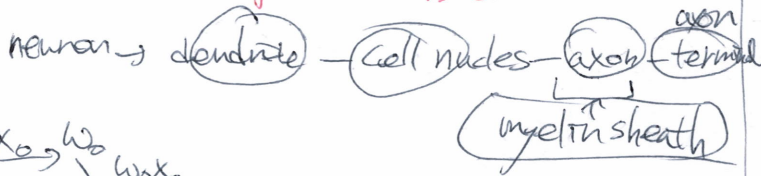                      (batch_size)
      c, _ = sess.run[cost,opt] feed_dict
      avgcost+=            [X=xs, y_batch
            c / total_batch         batch  ys]

sess.run (node) = node.eval()

mnist visualization ⇒ plt.imshow (
        mnist.test.images [r : r+].
        reshape(28,28), cmap='Greys',
        interpolation = 'nearest')
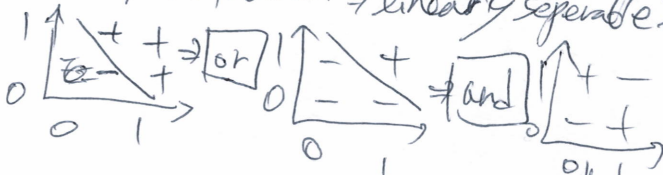        plt.show ()

---

# Lecture 8 Deep Neural Network

neuron → dendrite — Cell nudes — axon — terminal
                               myelin sheath



$X_0$ → $W_0$ $w_0 x_0$
$\sum W_i x_i + b$ → $f$ → output axon
               activation function

Hardware?
  └ Frank Rosenblatt ~1957 : Perceptron
  ( Widrow & Hoff. ~1960 : Adaline/Madaline

(False Promises) of Rosenblatt.
  └ July 08, 1958...

(Sample) AND/OR problem ⇒ linearly seperable?



• Perceptrons (1969)
  └ Marvin Minsky.
    XOR → x
              └ Need we MLP (multi layer peceptron)
              └ Not possible to train it...

(XOR) not linearly seperable.

"No one on earth had found a viable way to train" 1969

• Backpropagation
(1974, 1982 by Paul Werbos, 1986 Hinton)


      forward, "dog"
      backward = ? ← human face — label.
training example   error

• Convolutional Neural Network (Hubel & Wiesel 1959)
⇒ local connection (1980. LeCun)
inspired by

---

• Nav Lab 1984 1994 : 자율주행..
• A Big Problem ☆    Vanishing Gradient problem!
  └ Back Propagation → (Not work for many layers)
  &) Other existing algorithm!
  ⇒ SVM, Random Forrest.

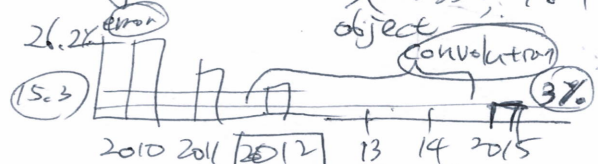(1995) by LeCun → New approach! ⟍ agree

★ CIFAR → HINTON!

| 2006 | Hinton "A fast learning algorithm for deep belief nets" |
| 2007 | Y.Bengio |

        "Greedy layer-wise traing of deep network"

• If the weights are (initialized) in a (clever way) rather than randomly, deep neural network can be trained well

• Imagenet → 1000 class, 1.4 M images
                 object
26.2% error         (convolution)



(SoA)                         3%
2010 2011 2012 13 14 2015
         └ Alexnet.

• Deep API learning ( Dr. Sung Kim's Lab)..
  alphago , Algae

• G. Hinton's summary of finding up to today
  ① Our labeled datasets were thousands of times too small.
  ② Our computers were millions of times too slow
  ③ We initialized the weights in a stupid way
  ④ We used the wrong type of non-linearity.

  (from "a brief history of neural-nets and deep learning part 4")

Why DNN?  youtube (translation), Facebook (user training), Google (web search), netflix (recommendation), amazon
⋮

**Lab8. Tensor manipulation**

① t = np.array [0, 1, 2, ... 6])
(t.ndim
(t.shape        t[:2], t[2:]
t[0], t[6] ...        ↑ slicing

② t = np.array ( [[1,2,3],
[4,5,6],
[ : ],
[ ]])  → t.ndim = 2  →rank
t.shape = (4,3)

**Tensor**

4 • t = tf.costant ([1,2,3,4])
tf.shape(t).eval() → array([4], dtype=int)

• t = tf.costant ([[1,2][3,4]])
number of [] = rank   ← tf.shape(t).eval()
⇒ [2,2]
                                x  y

AXIS? ⇒ rank = 4? ⇒ 4 axes!!

• inner-most axis = 3 ⇒ -1 → [[[[0000],
• outer-most axis = 0 → axis0    [0000]]
AXIS 3                            [[0000],
                                 [0000]]
                                 [[....]]]

[
                              ~
                            ]]

• matmul vs multiply
⇒ A: shape (2,2)
B: shape (2,1) ⇒ tf.matmul(A,B).eval()
(2,2)×(2,1) ⇒ shape (2,1)
A * B ⇒ element-wise multiplication
But... if shape is not same
→ broadcasting → Be careful!

• Reduce_mean.
⇒ tf.reduce_mean([1,2], axis=0).eval()
                     ↑
                    int? → output = int

• stack

• tf.ones_like(x). & tf.zeros_like(x)

---

NLS

[[1, 2],     → reduce_mean(x, axis=0)
 [3, 4]]        ½ [[1,2] + [3,4]]

reduce_mean(x, axis=-1)
axis0   axis-1    (½[1+2],
                   ½[3+4])

• Argmax [[0,1,2],
          [2,1,0]]
— tf.argmax(x, axis=0).eval()
      ↳ [1,0,0]
— tf.argmax(x, axis=-1).eval()
      ↳ [2,0], argmax ≡ "position of maxima along the axis"

• Reshape  ↝ [[[0,1,2],
              [3,4,5]],
             [[6,7,8],
              [9,10,11]]]
tf.reshape(t, shape[-1,3]).eval()
      ↳ [[ 0 1 2]
         [ 3 4 5]
         [ 6 7 8]
         [ 9 10 11]]
tf.reshape(t, shape[-1,1,3]).eval()
      "remain the last two axes"

• tf.squeeze, expand
      ↳ = flatten    ↳ tf.expand_dims
• zip (python) numpy      ([0,1,2], 1).eval()
• one hot
tf.one_hot([[0][1][2][0]], depth=3).eval()
      ↳ [[[1 0 0]
         [0 1 0]
         [0 0 1]
         [1 0 0]]], dtype=float32)

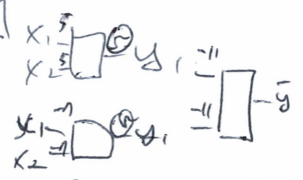& tf.reshape(t.shape=[-1,3]).eval()
(∵ redundant rank is generated after tf.onehot)

• Cast
⇒ tf.cast(X, int or float)

**Lecture 9. Deep learning 정리!**

✱ XOR?



neural net

✱ Backpropagation (Chain rule)



$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g}\frac{\partial g}{\partial w} = 1 \times 5 = 5$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g}\frac{\partial g}{\partial w}$$

$$f = wx + b$$
$$g = wx$$

∴ graph 구성이 BP도 바로 적용가능 !!

**LAB9**

```
w₁ = tf.va~
b₁ = tf.va~
layer1 = tf.sigmoid(tf.matmul(X,W₁)+b₁)
w₂ = tf.v—
b₂ = tf.va—
layer2 = tf.sigmoid(tf.matmul(layer1,W₂)+b₂)
       hypothesis
```

**TensorBoard** : TF logging/debugging tool

① From TF graph, decide which tensors you want to log.
  · w2_hist = tf.summary.histogram("weights2", W2)
  · cost_summ = tf.summary.scalar("cost", cost)

② Merge all summaries
  · summary = tf.summary.merge_all()

③ Create writer and add graph
  · writer = tf.summary.FileWriter('./logs')
  · writer.add_graph(sess.graph)

④ Run summary merge and add summary (매번)
  · s, _ = sess.run([summary, optimizer], feed_dict=feed_dict)
  · writer.add_summary(s, global_step=global step)

⑤ Launch TensorBoard
  · tensorboard --logdir=./logs
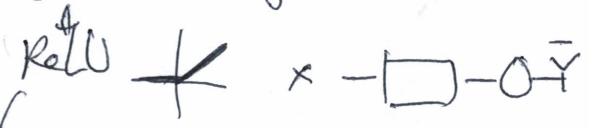
· what tensor?
  — Scalar tensor ⇒ tf.summary.scalar()
  — Histogram for multiple-valued tensor

· Add scope for better graph hierarchy
  with tf.name_scope("layer1") as scope:
  {
  with tf.name_scope("layer2") as scope:

**Lecture 10. ReLU**
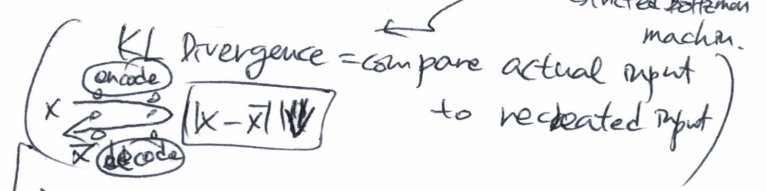
Sigmoid ; Vanishing Gradient. problem!



$$L1 = tf.nn.relu(tf.matmul(X,W1)+b1)$$

Leaky ReLU & Maxout & ELU

**Initialization (smart way)**

· Not all 0's
· challenging issue
· Hinton et al. (2006) DBN with (RBM)
  restricted Boltzman machine



KL Divergence = compare actual input to recreated input

RBM RBM RBM RBM → supervised learning! (Backprop)

· No need to use RBM !
· simple method ⓞ → Xavier initialization (2010) (Glorot)
              → He's initialization (2015)

$$\left[ \begin{array}{l} \text{Xavier: } W = np.random.randn(fan\_in, fan\_out) \\ \qquad\qquad // np.sqrt(fan\_in) \\ He : W = np.ra... \qquad\qquad ] \\ \qquad\qquad // np.sqrt(fan\_in/2) \end{array} \right.$$

○ prettytensor implementation.

✳ Batch normalization ✳
　　○ Layer sequential uniform variance.

---

( Dropout & Model Ensemble )



＋＋＋＋＋ overfitting.

test
train
epoch

→ Regularization.
　　　　↳① L2 Regularization
　　　　↳② Dropout (2014)

○ Tensorflow implementation.

dropout_rate = tf.placeholder("float")

$L_1$ = tf.nn.relu(tf.add(tf.matmul(X,W1),B1)

$L_1$ = tf.nn.dropout(_$L_1$, dropout_rate)

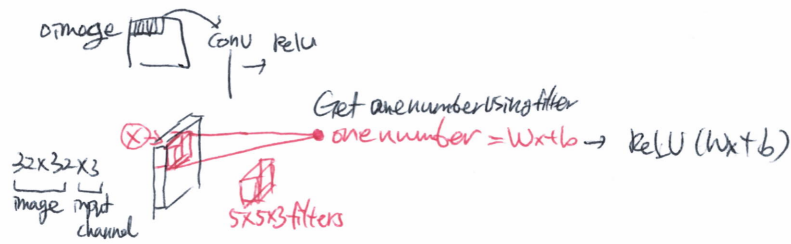　✳② evaluation 할때 dropout rate = 1.0

○ Ensemble ?

Data set

Trai... Set 1 　set 2 　　set 3
(Net) 　(Net) 　　(Net)

Combine
#
ensemble
prediction

•Build DNN !

－Fast forward － □→□→□→□→□→□...
　(He's ResNet)

－Split & Merge x →□<□□□>□→?

－Recurrent →□→□→□→□

OPTIMIZERS
　↳ Gradiendescent, AdadeltaOpt
Stoch. Adagrad Opt. Momentum Opt, (Adam Opt
✳
　　FtrlOpt, ....

(CNN) ✗✗

↳ Next page.

# ☆ CNN

- Hubel & Wiesel 1959 Receptive fields of single neurons in the cat's striate cortex.

○image ⟶ conv relu →

Get one number using filter
• one number = Wx+b → ReLU (Wx+b)

22×32×3 image input channel

5×5×3 filters

→ How many numbers we get?

(ex)

η×η input (spatially)
assume 3×3 Filter → (1 stride) → 5×5 output
$$(= η - [3-1])$$
→ (2 stride) → 3×3 output

Output size $\left[\dfrac{N-F}{stride}\right] + 1$

$$\left(\dfrac{η-3}{2} + 1\right)$$

In practice: Common to zero pad the board.

eg: input 7×7, 3×3 filter, stride 1 ⇒ pad with 1 pixel boarder
⇒ what is the output?
⇒ 7×7 output.

common — stride 1, filter F×F, zero padding $\dfrac{(F-1)}{2}$

convolution layer
(5×5) 6 filters ×3
32×32×3
activation maps
28×28×6

- Convolution layers.

(conv ReLU) 6 of (5×5×3) filters
32×32×3
(conv ReLU) 10 of (5×5×6) filters 10

Q. How many would variables?
$(5×5×3)×6$
$+ (5×5×6)×10$
$(F_1 \cdot F_1 \cdot 10)^2 × N_3$

- pooling (= sampling)
resize (sampling) → (How?)
output
max pool with 2×2 filter, and stride 2.

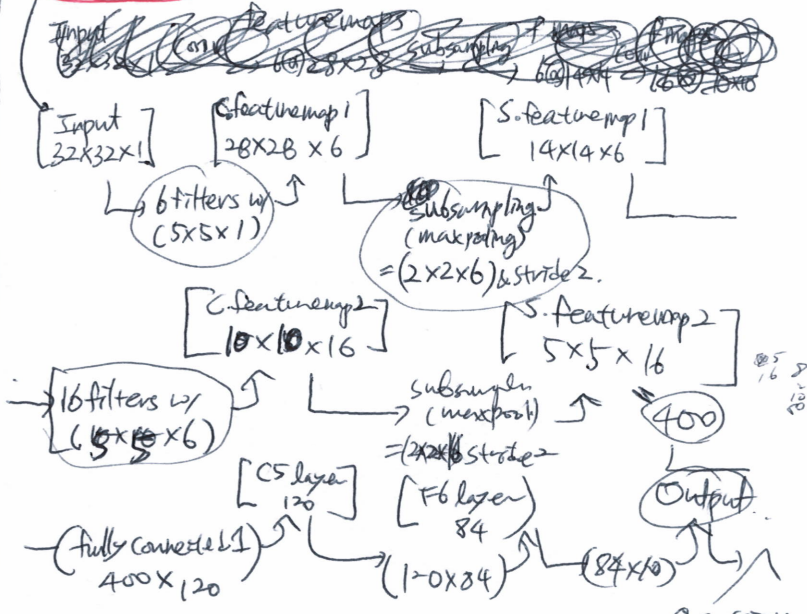| 1 | 2 | 4 |
| 6 | 7 | 8 |
| 3 | 2 | 1 |
| 2 | 4 |

2×2  | 6 | 8 |
     | 3 | 4 |

- Fully connected layer (FC layer)

- ConvNet JS demo : training on CIFAR-10

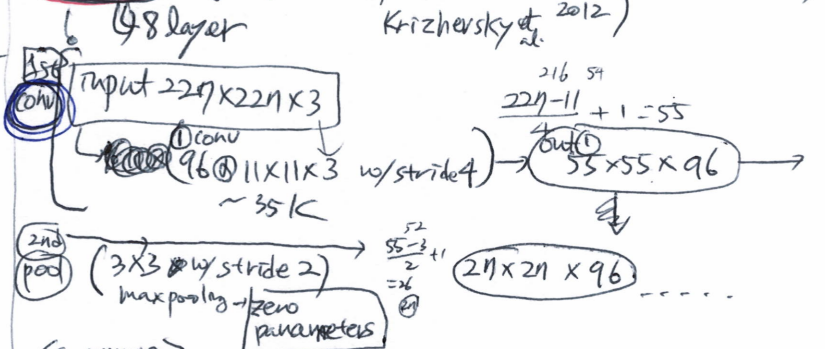( cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html )

---

[case study]

LeNet-5 (LeCun, 1998)

[Input 32×32×1] [C. feature map 1 28×28 ×6] [S. feature map 1 14×14×6]

6 filters w/ (5×5×1)
subsampling (maxpooling) = (2×2×6) & stride 2.

[C. feature map2 10×10×16] [S. feature map 2 5×5×16]

16 filters w/ (5×5×6)
subsampling (maxpool) = (2×2×16) stride 2

400

[C5 layer 120] [F6 layer 84] Output

(fully connected1) 400×120
(120×84) (84×10)

gaussian connection. (or Softmax)

---

AlexNet (ImageNet competition 1st place, Krizhevsky et al. 2012)
8 layer

conv | Input 227×227×3

① conv 96@ 11×11×3 w/stride4 → out① 55×55×96
~35K

$\dfrac{227-11}{4} + 1 = 55$

2nd pool (3×3 w/ stride 2) max pooling → zero parameters

$\dfrac{55-3}{2}+1 = 6$  27×27×96

<summary>

[227 × 227 × 3] Input
[55 × 55 × 96] conv : 96@ 11×11(×3) filter w/s=4 & pad=0.
[27 × 27 × 96] Maxpool: 3×3 filter w/s=2
[27 × 27 × 96] norm1 : Normalization layer
[27 × 27 × 256] conv2 : 256@ 5×5(×96) filter s=1 pad=2
[13 × 13 × 256] maxpool : 3×3 filter w/ s=2.
[13 × 13 × 256] Norm2
[13 × 13 × 384] conv3 : 384@ 3×3(×256) s=1, pad1.
[13 × 13 × 384] conv4 : 384@ 3×3(384) s=1, p=1
[13 × 13 × 256] conv5 : 256@ 3×3(×384) s=1, p=1
[6 × 6 × 256] Maxpool3 : 3×3 filter, s=2.
[4096] FC6
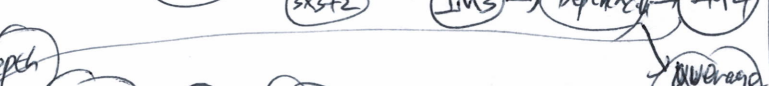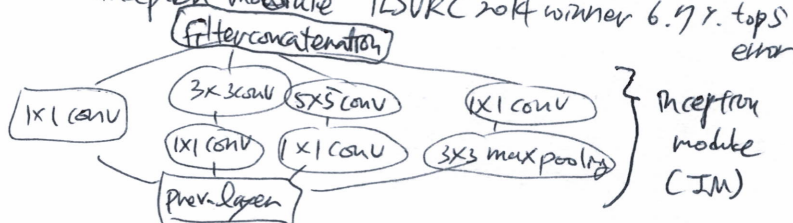[4096] FC7
[1000] FC8 (class scores)

AlexNet continued....

⇒ 1st use of (ReLU).

- Norm layers ( not common anymore)
- heavy data augmentation.
- dropout = 0.5
- batchsize = 128
- SGD Momentum = 0.9
- learning rate = 1e-2. reduced by 10 manually when val accuracy plateaus.
- L2 weight decay = 5e-4
- 7 CNN ensemble : 18.2% → 15.4%

(GoogLeNet) Szegody et. al. 2014

↳ Inception module  ILSVRC 2014 winner 6.7% top5 error



[Filter concatenation]
1x1 conv | 3x3 conv | 5x5 conv | 1x1 conv
1x1 conv | 1x1 conv | 3x3 maxpooling
Prev layer

} Inception module (IM)

Input → Conv 7x7 + 2S → MaxPool 3x3 + 2S → Local Res Norm
→ Conv (1x1 (v)) → Conv 3x3 + 1S → Local Res Norm → IM1 → depth concat
→ IM2 → depth concat → maxP 3x3+2 → IM3 → Depth c → IM4
↓ MP 3x3+2
Depth c → IM5 → DC → IM6 → DC

→ average pool → Conv 1x1 → FC → FC → Softmax

→ average p73 → Conv 1x1 → FC → FC → Softmax

(ResNet) (He et al, 2015)

ILSVRC 2015 winner (3.6% top 5 error)

1st place in 5 main tracks.
→ Imagenet classification (152 layers)
       detection
       localization
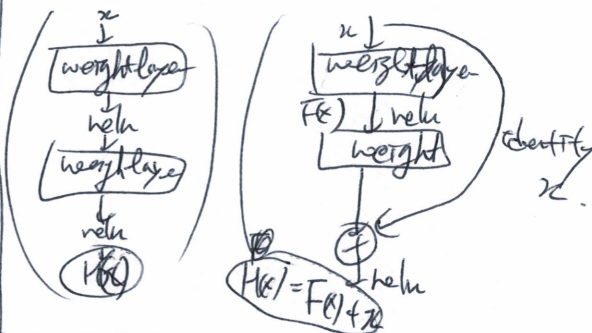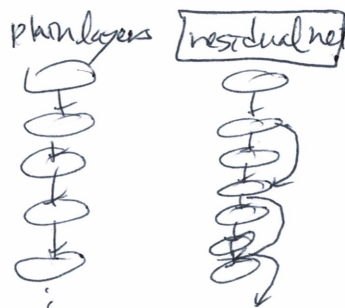  COCO  detection
  COCO  segmentation

2017 = 2.3%

---

• Revolution of depth.  (ML9)

AlexNet 8 layers → VGG 19 layer → ResNet 152 layer
  2012            19  2014          2015

(2~3 week training GPU)
                     ↓ 8
• runtime = faster than VGG Net

plain layers | residual net



weight layer
  ↓ relu
weight layer
  ↓ relu
H(x)

x
weight layer
  F(x) ↓ relu
  weight          identity x
H(x) = F(x) + x ← relu

---

• CNN for Sentence classification (Yoon Kim) 2014
• Alphago

---

< Tensorflow CNN >

input → convolution (subsampling) → FC → output
           ↓                         ↓
      [= feature extraction] [classification]

image vector → filter w/stride → subsampling (pooling)

• simple convolution layer.
  3x3x1 image
  3x2x1 filter(W)   →  2x2x1 output.

< Tensor shape >
• Image  "-1" is good.
  y (#data index (??), x, y, color)
• filter → (x_f, y_f, color, #filters)
• stride → (S_x, S_y)
• pading → "VALID"

→ tf.nn.conv2d (image, weight, strides=[1,1,1,1], pading = VALID)

→ pading = SAME → output map size = input map size

• pooling
↳ tf.nn.max_pool (image ksize = [1,2,2,1], stride = [1,1,1,1], pading = SAME)

# ① MNIST Conv. layer

```
sess = tf. InteractiveSession()
img = img. reshape (-1, 28, 28, 1)
W1 = tf. Variable( tf.random_normal([3,3,1,5],     #filters
                              stddev=0.01))          ↓

conv2d = tf.nn.conv2d (img, W1, strides=[1,2,2,1],
                              padding = 'SAME')

sess.run (global_initd ~ ())
conv2d_img = conv2d.eval()
   conv2d_img = np. swapaxes(conv2d_img, 0, 3)
for i, one_img in enumerate (conv2d_img):
    plt. subplot (1,5, i+1), plt.imshow(one_img.
                              reshape(14,14)
                              cmap='grey')
```

# ① MNIST Maxpooling.

```
pool = tf.nn. maxpool (conv2d, ksize=[1,2,2,1],
                              strides=[1,2,2,1] ,padding
                                        ='same')
print (pool)

sess.run ( tf. globa~)
pool_img = pool.eval()
pool_img = np.swapaxes (pool_img, 0, 3)
for ~.
```

RNN ~ next page.

# ⟨RNN in Tensorflow⟩



① $^{\uparrow}_{\downarrow}$38

Cell = tf. contrib. rnn. BasicRNNCell

(num_units = hidden_size)

outputs, _states = tf.nn. dynamic_rnn (

cell, x_data, dtype=tf.float 32 )

cell= tf. contrib. rnn .BasicLSTMCell(~)

(ex) One node : 4 (input_dim) in 2 (hidden-size)

$$h = [1 \ 0 \ 0 \ 0]$$
$$e = [0 \ 1 \ 0 \ 0]$$
$$l = [0 \ 0 \ 1 \ 0]$$
$$o = [0 \ 0 \ 0 \ 1]$$

(X) → [[[1,0,0,0]]]

shape = (1,1,4)

(h)

hidden size = 2.

shape (output) = (1,1,2)

[[[ x, x]]]

(code)

hidden Size = 2

cell = tf. contrib. rnn. Basic LSTMCell (

num_units =hidden_size)

x-data = np. array ([[[1,0,0,0]]], dtype=np. float32)

outputs, -stats = tf.nn. dynamic _rnn ( cell,

x-data, dtype = tf. float 32)

sess. run ( tf. goal _variables _ initializer())

pp. print ( outputs .eval())

array ([[[ -0.424~ , 0.64 ~ ]]])

output shape(1,5,2) [[XX], [X,X] [X,X] [X,X] [XX]]



# of seq. input dim

(unfold)

Input

shape=(1,5,4) :[[[1000] [0100] [0010] [0010] [0001]]]

# of sequence input dim

Coding?

---

\# One cell input_dim(4) → output_dim(2). seq =5
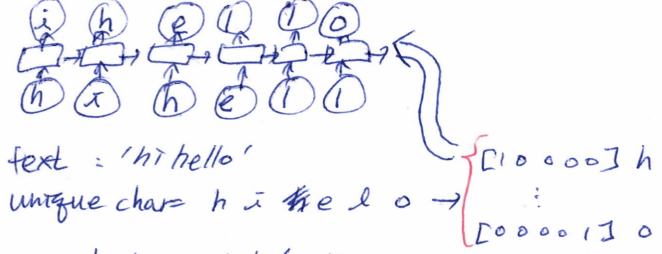
RNN

hiddensize = 2.     h=[1,0,0,0]

cell ~

x-data = np. array([[ h e l l o ]], dtype. float 32)

outputs. & states = tf. nn. dynamic _rnn (cell, x_data,

dtype=?)

sess. run ( tf. glob ~ )

pp. print ( outputs .eval())

• Batching Input.

batch size =3

shape = (3, 5, 2)  ————— (input dim)

batch-size   sequence (or hiddensize )

---

✱ Teach RNN "hihello"



text : 'hi hello'

unique chars h i e l o → { [1 0 0 0 0] h
                              ⋮
                            [0 0 0 0 1] o }

∴ input shape = [[1,6,5]]

output shape = [[1,6,5]]

single sentence .

• Creating rnn cell

↳ rnn-cell= rnn-cell. Basic RNN Cell (rnn size)

Basic LSTMCell

GRU Cell

outputs, _states = tf. nn. dynamic _ rnn ( rnn-cell, X,

initial state = initial_state,
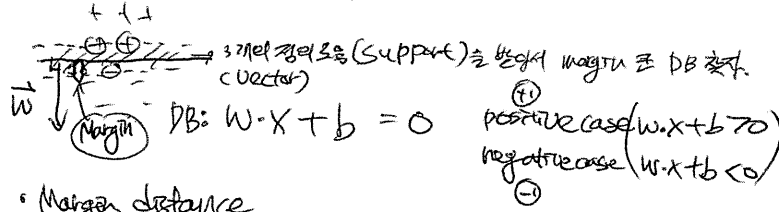
dtype = tf. float 32 )

• Loss !

sequence _loss =

# ⟨SUPPORT Vector Machine⟩ aailab kaist channel

- **Decision Boundary w/o Prob.**



여러개의 DB 가능함 → DB가 data point 에 너무가까우면 data의 noise 때문에 영향받음!!

- **Decision Boundary with Margin** ⟵ [ + & — region 에서 가능한 영향을 줄여야함. ]



3개의 정의하는 점들(support) 을 통해서 margin 또 DB 결정. (vector)

$DB: W \cdot X + b = 0$

positive case $(w \cdot x + b > 0)$ $(+)$
negative case $(w \cdot x + b < 0)$ $(-)$

- **Margin distance**

$f(x) = w \cdot x + b$ ⟹ $0 →$ on decision Boundary
$> 0 →$ A positive point $X$



- $X = x_p + r \frac{W}{\|W\|}$ , $f(x_p) = 0 \Rightarrow f(x) = w \cdot x + b$
  $= w(x_p + r\frac{w}{\|w\|}) + b$
  $= \overline{w x_p + b} + r \frac{w \cdot w}{\|w\|}$

∴ distance $r = \frac{f(x)}{\|W\|}$

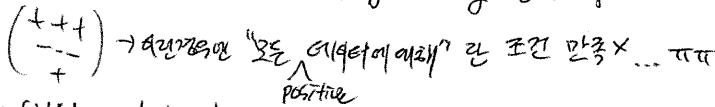- **Maximizing the Margin.**
  ↳ Good DB = maximum margin = $\max(r)$ !
  - Optimization problem $\{ \max_{w,b} 2r = \frac{2a}{\|w\|}$
    $s.t (w x_j + b) y_j \geq a, \forall_j \}$

  - $a =$ arbitrary number and can be normalized

  $\boxed{\min_{w,b} \|w\| \quad s.t \; (w x_j + b) y_j \geq 1, \forall_j}$

  ⟹ quadratic optimization problem.
  ↳ quadratic programming 이용해서 optimization.

  $\begin{pmatrix} + & + & + \\ - & - & - \\ & + \end{pmatrix}$ → 우리경우에 "모든 데이터에 의해서" 란 조건 만족X ... ㅠㅠ
  ∧ positive

- **SVM with hard margin**

  solution { — Hard margin : No errors case are allowed
  — Soft margin : 에러 허용.
  — Kernal trick : $(+ + + +)$ non-linear boundary 사용!!

- **"Error" cases in SVM**

  Error handling ! ⟹ (Option 1) $\min_{w,b} \|w\| + C \times$ #error
  s.t. $(w x_j + b) y_j \geq 1, \forall_j$


0-1 loss
(DB)

(Option 2) Slack variable $\xi_j \geq 1$ when misclassified.
error.
$\min_{w,b} \|w\| + C \sum \xi_j$
s.t $(w x_j + b) y_j \geq 1 - \xi_j, \forall_j (\xi_j \geq 0, \forall_j)$
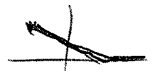
---

- **Soft Margin SVM.**
  { $\min_{w,b} \|w\| + C \sum_j \xi_j \; s.t. \; (w x_j + b) y_j \geq 1 - \xi_j, \forall_j$
  $\xi_j \geq 0, \forall_j$

- **Comparison to Logistic Regression**
  — Loss function $\xi_j = loss(f(x_j), y_j)$
  — SVM loss function : Hinge loss
  - $\xi_j = (1 - (w x_j + b) y_j)_+$
  — Logistic Regression loss function : log loss.
  - $\hat{\theta} = \arg\max_\theta \sum_{1 \leq i \leq N} \log(P(y_i | x_i ; \theta))$ ....
  — which loss fn is preferable?
    - Around the decision boundary?
    - overall place?

  $\xi_j = -\log(P(y_j | x_j, w, b))$

- **Strength of the Loss Function** $= \log(1 + e^{(w x_j + b) y_j})$

- **non-linear decision boundary.**

  $\boxed{\text{Feature Mapping to Expand Dim.}}$

  - $\min_{w,b,\xi_j} \|w\| + C \sum_j \xi_j \; s.t. \; (w \varphi(x_j) + b) y_j \geq 1 - \xi_j, \forall_j$
    $\xi_j \geq 0, \forall_j$
  - $\varphi(\langle x_1, x_2 \rangle) =$
    $\langle x_1, x_2, x_1^2, x_2^2, x_1 x_2, x_1^3 \cdots \rangle$