

deeplearning.ai.

1. Neural Networks and Deep Learning
2. Improving deep neural networks: Hyperparameter tuning, regularization, optimization
3. Structuring your machine learning project
4. Convolutional Neural Networks
5. Natural Language Processing: Building sequence models

C1. Neural Networks & Deep Learning

W1. Introduction to DL

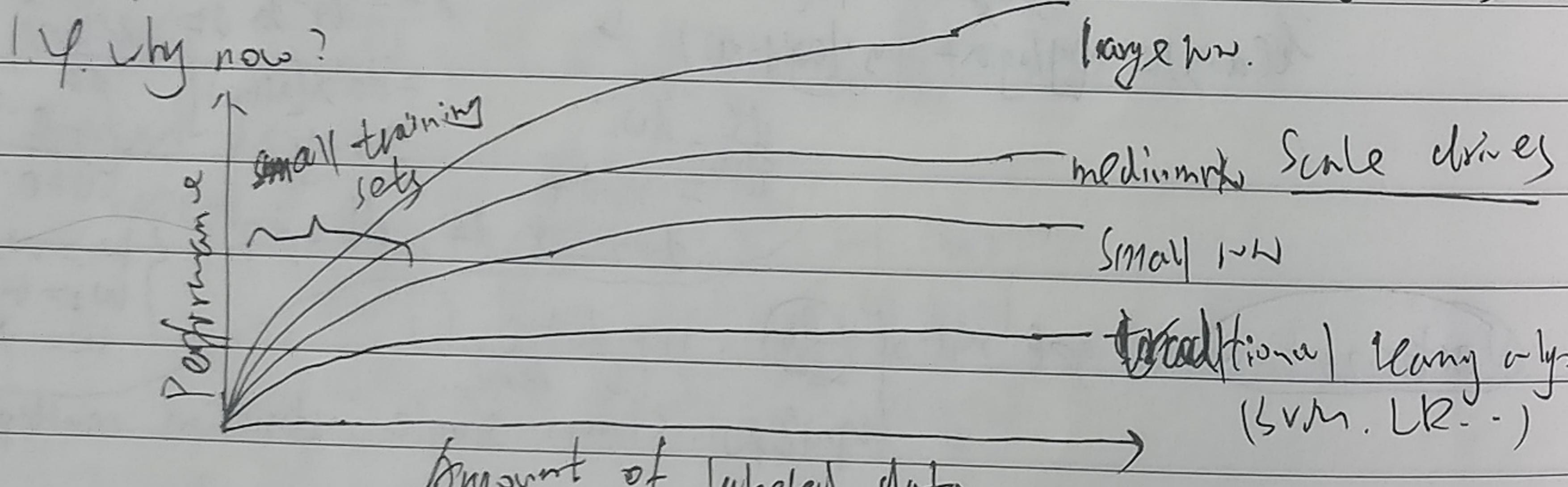
1.1 Supervised learning

Real estate, Online Advertising, Photo, Speech, translation, Auto-driving
Standard NN, CNN, RNN, Custom NN

structured data: list



unstructured data: Audio, Image, Text (human good at)



Data, computation, algorithms.
Signal → relu (faster gradient descent)
Idea
Experiment
Rule

W2. Basics of Neural Network Programming

(x, y) , $x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$
m training examples.

$$X = \begin{bmatrix} | & | & | \\ x_1^{(1)}, x_2^{(1)}, \dots, x_m^{(1)} & | \\ | & | & | \\ x_1^{(2)}, x_2^{(2)}, \dots, x_m^{(2)} & | \\ \vdots & \vdots & \vdots \\ | & | & | \\ x_1^{(m)}, x_2^{(m)}, \dots, x_m^{(m)} & | \end{bmatrix}_n^m$$

$X \in \mathbb{R}^{n_x \times m}$

$$Y = [y_1^{(1)}, y_2^{(1)}, \dots, y_m^{(1)}]$$

$Y \in \mathbb{R}^{1 \times m}$

logistic regression

$$\hat{y} = g(W^T x + b), W \in \mathbb{R}^{n_x \times n_y}$$

$$b \in \mathbb{R}^{n_y}$$

$$x_0 = 1, x \in \mathbb{R}^{n_x + 1}$$

$$\hat{y} = g(b^T x)$$

$$\theta = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}^T$$

Skeleton of deeplearning.ai

C1. Neural Networks and Deep Learning

W1. Introduction to DL

supervised learning, why now?

W2. Basics of neural network programming.

review activation, gradient descent for NN.

W3. DNN: Forward and Backward Propagation

gradient descent for NN, parameter & hyperparam

C2: Improving DNN

W1. Setting up your ML applications

Data, bias & variance, regularization,

Dropout, augmentation, early stopping,

normalizing training set, gradient vanish/explode,

weight initialize, grad check

W2. Optimization Algos

minibatch, moving average,

Momentum, RMSprop, Adam,

learning rate decay,

local optimal

W3. Hyperparameter tuning

tuning process, BN,

softmax, DL frameworks,

C3: Structuring ML projects

W1. intro to ML strategy

orthogonalization, evaluation metric,

optimizing and satisfying metrics

human performance, avoidable bias,

improving model performance.

W2. Error Analysis

categorize error analysis, incorrectly labeled data,

different distribution, bias & variance, addressing data mismatch,

transfer learning, multi-task learning, end-to-end DL

C4: Convolutional Neural Networks

W1. CNN

edge detect, padding, stride

Max Pooling, why convolutional?

W2. Case Studies

classic networks (LeNet, AlexNet, VGG, ResNet, Inception), transfer learning, Data augmentation, state-of-the-art

bounding box, Non-max suppression, Anchor Boxes, Yolo, region proposal

W3. Object detection

landmark detection, object detection, sliding window, Non-max suppression, NMS, Anchors, YOLO, region proposal

W4. Special applications of Convnet

face recognition, one-shot learning, Siamese network, triplet loss, neural style transfer

loss function (content, style)

C5: Sequence Models

W1. Recurrent Neural Networks

RNN, backpropagation, types of RNN,

language model, sequence generation, vanishing grads,

Gated RNN, LSTM, bidirectional networks

Deep RNN

W2. NLP and Word Embeddings

word representation/embeddings, embedding matrix, word2vec,

Negative sampling, glove word vectors,

Sentiment classification, Debias

W3. Sequence to Sequence Model

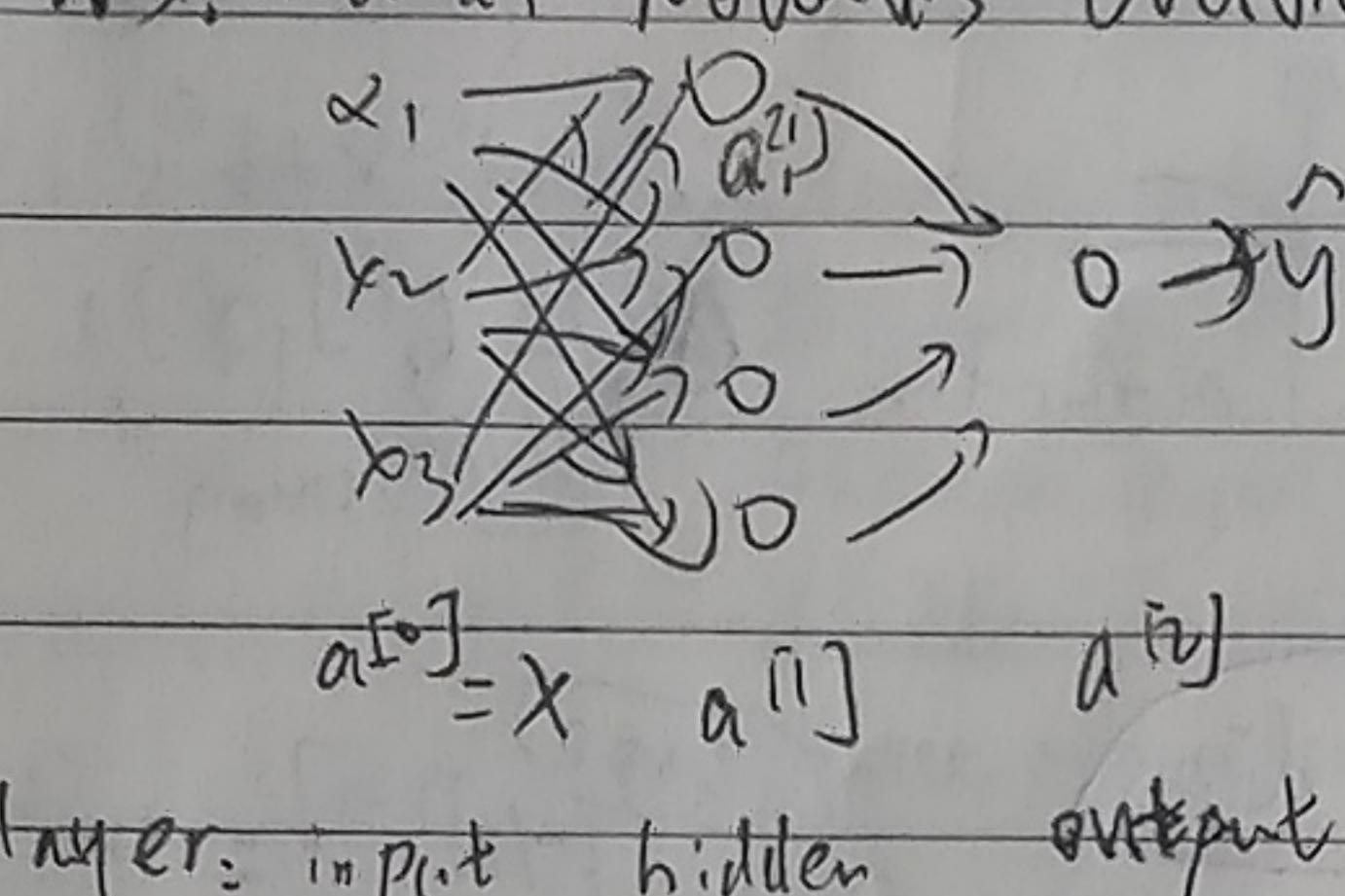
Basic model, Pick the most likely sentence,

Beam search, Attention model,

Speech recognition, Trigger word detection.

$$\begin{aligned} & \text{if } y=1: p(y|x) = \hat{y} \\ & \text{if } y=0: p(y|x) = 1-\hat{y} \end{aligned} \quad \left\{ \begin{array}{l} \text{lost on in samples} \\ \text{log P(labels on training set)} = \log \prod_{i=1}^m p(y^{(i)}|x^{(i)}) \end{array} \right.$$

W3: Neural Networks Overview



$$\begin{aligned} & \text{Input layer: } a^{(0)} = x, a^{(1)}, a^{(2)} \\ & \text{Hidden layer 1: } z^{(1)} = -w_1^{(1)T} - b_1^{(1)} \quad a^{(1)} = g(z^{(1)}) \\ & \quad z^{(1)} = W^{(1)T} x + b^{(1)} \\ & \quad g^{(1)} = \sigma(z^{(1)}) \\ & \text{Hidden layer 2: } z^{(2)} = -w_2^{(2)T} - b_2^{(2)} \quad a^{(2)} = g(z^{(2)}) \\ & \quad z^{(2)} = W^{(2)T} x + b^{(2)} \\ & \quad g^{(2)} = \sigma(z^{(2)}) \end{aligned}$$

activation function

$$\begin{aligned} & \text{sigmoid function: } \frac{1}{1+e^{-z}} \quad \text{Output layer: } \{0, 1\} \\ & \text{tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad \text{rescale.} \\ & \text{ReLU: } = \max(0, z) \quad \text{Leaky ReLU: } = \max(0, \alpha z + z) \end{aligned}$$

Gradient descent for NN

$$\text{Cost func.: } J(w^{(0)}, w^{(1)}, b^{(1)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)}, \hat{y}^{(i)})$$

Gradient descent:

$$\begin{aligned} & \text{repeat:} \quad \left\{ \begin{array}{l} \text{compute predictions } (\hat{y}^{(i)}, i=1 \dots m) \\ dJ/dw^{(1)} = \frac{\partial J}{\partial w^{(1)}} \cdot dJ/dw^{(1)}, dJ/db^{(1)}, dJ/db^{(2)} \\ w^{(1)} = w^{(1)} - \alpha dJ/dw^{(1)}, w^{(2)}, b^{(1)}, b^{(2)} \end{array} \right. \end{aligned}$$

Formulas for computing derivatives

Forward propagation:

$$\text{eg.: } \hat{y}^{(1)} = \hat{y}^{(1)T} X + b^{(1)}$$

$$d\hat{y}^{(1)} = d\hat{y}^{(1)T} \cdot I \quad (Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}])$$

$$d\hat{y}^{(2)} = d\hat{y}^{(2)T} \cdot I$$

$$d\hat{y}^{(3)} = d\hat{y}^{(3)T} \cdot I$$

$$d\hat{y}^{(4)} = d\hat{y}^{(4)T} \cdot I$$

$$d\hat{y}^{(5)} = d\hat{y}^{(5)T} \cdot I$$

$$d\hat{y}^{(6)} = d\hat{y}^{(6)T} \cdot I$$

$$d\hat{y}^{(7)} = d\hat{y}^{(7)T} \cdot I$$

$$d\hat{y}^{(8)} = d\hat{y}^{(8)T} \cdot I$$

$$d\hat{y}^{(9)} = d\hat{y}^{(9)T} \cdot I$$

$$d\hat{y}^{(10)} = d\hat{y}^{(10)T} \cdot I$$

$$d\hat{y}^{(11)} = d\hat{y}^{(11)T} \cdot I$$

$$d\hat{y}^{(12)} = d\hat{y}^{(12)T} \cdot I$$

$$d\hat{y}^{(13)} = d\hat{y}^{(13)T} \cdot I$$

$$d\hat{y}^{(14)} = d\hat{y}^{(14)T} \cdot I$$

$$d\hat{y}^{(15)} = d\hat{y}^{(15)T} \cdot I$$

$$d\hat{y}^{(16)} = d\hat{y}^{(16)T} \cdot I$$

$$d\hat{y}^{(17)} = d\hat{y}^{(17)T} \cdot I$$

$$d\hat{y}^{(18)} = d\hat{y}^{(18)T} \cdot I$$

$$d\hat{y}^{(19)} = d\hat{y}^{(19)T} \cdot I$$

$$d\hat{y}^{(20)} = d\hat{y}^{(20)T} \cdot I$$

$$d\hat{y}^{(21)} = d\hat{y}^{(21)T} \cdot I$$

$$d\hat{y}^{(22)} = d\hat{y}^{(22)T} \cdot I$$

$$d\hat{y}^{(23)} = d\hat{y}^{(23)T} \cdot I$$

$$d\hat{y}^{(24)} = d\hat{y}^{(24)T} \cdot I$$

$$d\hat{y}^{(25)} = d\hat{y}^{(25)T} \cdot I$$

$$d\hat{y}^{(26)} = d\hat{y}^{(26)T} \cdot I$$

$$d\hat{y}^{(27)} = d\hat{y}^{(27)T} \cdot I$$

$$d\hat{y}^{(28)} = d\hat{y}^{(28)T} \cdot I$$

$$d\hat{y}^{(29)} = d\hat{y}^{(29)T} \cdot I$$

$$d\hat{y}^{(30)} = d\hat{y}^{(30)T} \cdot I$$

$$d\hat{y}^{(31)} = d\hat{y}^{(31)T} \cdot I$$

$$d\hat{y}^{(32)} = d\hat{y}^{(32)T} \cdot I$$

$$d\hat{y}^{(33)} = d\hat{y}^{(33)T} \cdot I$$

$$d\hat{y}^{(34)} = d\hat{y}^{(34)T} \cdot I$$

$$d\hat{y}^{(35)} = d\hat{y}^{(35)T} \cdot I$$

$$d\hat{y}^{(36)} = d\hat{y}^{(36)T} \cdot I$$

$$d\hat{y}^{(37)} = d\hat{y}^{(37)T} \cdot I$$

$$d\hat{y}^{(38)} = d\hat{y}^{(38)T} \cdot I$$

$$d\hat{y}^{(39)} = d\hat{y}^{(39)T} \cdot I$$

$$d\hat{y}^{(40)} = d\hat{y}^{(40)T} \cdot I$$

$$d\hat{y}^{(41)} = d\hat{y}^{(41)T} \cdot I$$

$$d\hat{y}^{(42)} = d\hat{y}^{(42)T} \cdot I$$

$$d\hat{y}^{(43)} = d\hat{y}^{(43)T} \cdot I$$

$$d\hat{y}^{(44)} = d\hat{y}^{(44)T} \cdot I$$

$$d\hat{y}^{(45)} = d\hat{y}^{(45)T} \cdot I$$

$$d\hat{y}^{(46)} = d\hat{y}^{(46)T} \cdot I$$

$$d\hat{y}^{(47)} = d\hat{y}^{(47)T} \cdot I$$

$$d\hat{y}^{(48)} = d\hat{y}^{(48)T} \cdot I$$

$$d\hat{y}^{(49)} = d\hat{y}^{(49)T} \cdot I$$

$$d\hat{y}^{(50)} = d\hat{y}^{(50)T} \cdot I$$

$$d\hat{y}^{(51)} = d\hat{y}^{(51)T} \cdot I$$

$$d\hat{y}^{(52)} = d\hat{y}^{(52)T} \cdot I$$

$$d\hat{y}^{(53)} = d\hat{y}^{(53)T} \cdot I$$

$$d\hat{y}^{(54)} = d\hat{y}^{(54)T} \cdot I$$

$$d\hat{y}^{(55)} = d\hat{y}^{(55)T} \cdot I$$

$$d\hat{y}^{(56)} = d\hat{y}^{(56)T} \cdot I$$

$$d\hat{y}^{(57)} = d\hat{y}^{(57)T} \cdot I$$

$$d\hat{y}^{(58)} = d\hat{y}^{(58)T} \cdot I$$

$$d\hat{y}^{(59)} = d\hat{y}^{(59)T} \cdot I$$

$$d\hat{y}^{(60)} = d\hat{y}^{(60)T} \cdot I$$

$$d\hat{y}^{(61)} = d\hat{y}^{(61)T} \cdot I$$

$$d\hat{y}^{(62)} = d\hat{y}^{(62)T} \cdot I$$

$$d\hat{y}^{(63)} = d\hat{y}^{(63)T} \cdot I$$

$$d\hat{y}^{(64)} = d\hat{y}^{(64)T} \cdot I$$

$$d\hat{y}^{(65)} = d\hat{y}^{(65)T} \cdot I$$

$$d\hat{y}^{(66)} = d\hat{y}^{(66)T} \cdot I$$

$$d\hat{y}^{(67)} = d\hat{y}^{(67)T} \cdot I$$

$$d\hat{y}^{(68)} = d\hat{y}^{(68)T} \cdot I$$

$$d\hat{y}^{(69)} = d\hat{y}^{(69)T} \cdot I$$

$$d\hat{y}^{(70)} = d\hat{y}^{(70)T} \cdot I$$

$$d\hat{y}^{(71)} = d\hat{y}^{(71)T} \cdot I$$

$$d\hat{y}^{(72)} = d\hat{y}^{(72)T} \cdot I$$

$$d\hat{y}^{(73)} = d\hat{y}^{(73)T} \cdot I$$

$$d\hat{y}^{(74)} = d\hat{y}^{(74)T} \cdot I$$

$$d\hat{y}^{(75)} = d\hat{y}^{(75)T} \cdot I$$

$$d\hat{y}^{(76)} = d\hat{y}^{(76)T} \cdot I$$

$$d\hat{y}^{(77)} = d\hat{y}^{(77)T} \cdot I$$

$$d\hat{y}^{(78)} = d\hat{y}^{(78)T} \cdot I$$

$$d\hat{y}^{(79)} = d\hat{y}^{(79)T} \cdot I$$

$$d\hat{y}^{(80)} = d\hat{y}^{(80)T} \cdot I$$

$$d\hat{y}^{(81)} = d\hat{y}^{(81)T} \cdot I$$

$$d\hat{y}^{(82)} = d\hat{y}^{(82)T} \cdot I$$

$$d\hat{y}^{(83)} = d\hat{y}^{(83)T} \cdot I$$

$$d\hat{y}^{(84)} = d\hat{y}^{(84)T} \cdot I$$

$$d\hat{y}^{(85)} = d\hat{y}^{(85)T} \cdot I$$

$$d\hat{y}^{(86)} = d\hat{y}^{(86)T} \cdot I$$

$$d\hat{y}^{(87)} = d\hat{y}^{(87)T} \cdot I$$

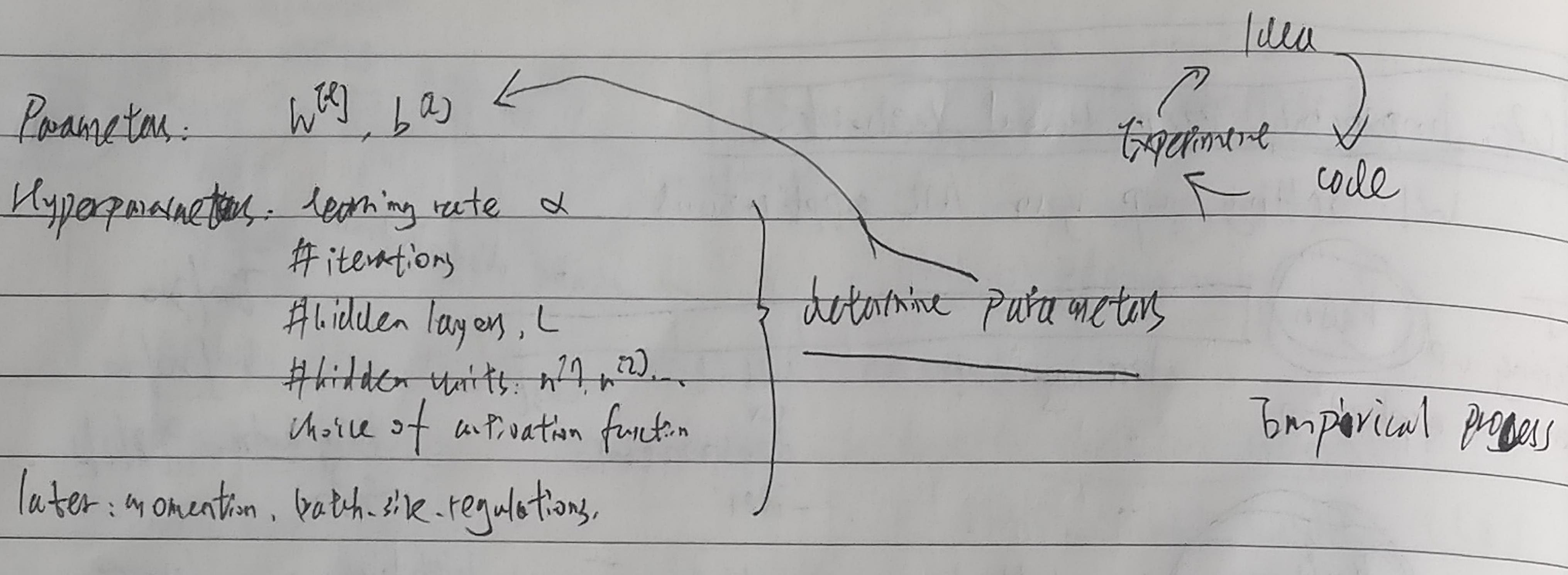
$$d\hat{y}^{(88)} = d\hat{y}^{(88)T} \cdot I$$

$$d\hat{y}^{(89)} = d\hat{y}^{(89)T} \cdot I$$

$$d\hat{y}^{(90)} = d\hat{y}^{(90)T} \cdot I$$

$$d\hat{y}^{(91)} = d\hat{y}^{(91)T} \cdot I$$

$$d\$$



3.1D. Computing gradients.

(Logistic regression)

$$x \rightarrow z = w^T x + b \rightarrow a = g(z) \rightarrow L(a, y)$$

$$\begin{aligned} dz &= a - y \\ da &= \frac{d}{da}(a, y) = -y \log a - (1-y) \log(1-a) \\ db &= da \cdot \frac{d}{db}(g(z)) = -\frac{y}{a} + \frac{1-y}{1-a} \end{aligned}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = da \cdot db = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot b(g(z)) \cdot (1-b(g(z)))$$

(Neural network gradients)

$$x \rightarrow z^{(1)} = w^{(1)} x + b^{(1)} \rightarrow a^{(1)} = g(z^{(1)}) \rightarrow L(a^{(1)}, y)$$

$$\begin{aligned} dz^{(1)} &= w^{(1)T} d\alpha^{(1)} \cdot g'(z^{(1)}) \\ dw^{(1)} &= d\alpha^{(1)} \cdot X^T \\ db^{(1)} &= d\alpha^{(1)} \end{aligned}$$

$$\begin{aligned} dz^{(2)} &= w^{(2)T} d\alpha^{(2)} \cdot g'(z^{(2)}) \\ dw^{(2)} &= d\alpha^{(2)} \cdot X^T \\ db^{(2)} &= d\alpha^{(2)} \end{aligned}$$

$$\vdots$$

$$w: [-] \quad \alpha: [-]$$

random initialize: $w^{(1)} = np.random.randn(m, n^{(1)}) * 2 \pi / 3$

$b^{(1)} = np.zeros((1, n^{(1)}))$

for sig. init. terms: $z^{(1)} = w^{(1)} x + b^{(1)}$

activation: $a^{(1)} = g^{(1)}(z^{(1)})$

slope: small

Why Deep L -layer neural network?

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

loop

$a^{(L)} = g^{(L)}(z^{(L)})$

check dims

$$w^{(L)} \in (n^{(L)}, n^{(L-1)})$$

$$z^{(L)} \in (n^{(L)}, 1)$$

vectorized

$$z^{(L)} = w^{(L)} + b^{(L)}$$

$$(n^{(L)}, m) \quad (n^{(L)}, n^{(L)}) \quad (n^{(L)}, m) \quad (n^{(L)}, m)$$

$$z^{(L)}, a^{(L)} \in (n^{(L)}, 1)$$

$$z^{(L)}, a^{(L)} \in (n^{(L)}, m)$$

$$dz^{(L)}, da^{(L)} = (n^{(L)}, m)$$

$$dz^{(L)}, da^{(L)} = (n^{(L)}, m)$$

human brain

backprop theory

$$y = x_1 \times x_2 \times x_3 \times \dots \times x_n$$

$$\begin{matrix} x_1 & \times & x_2 & \times & x_3 & \times & \dots & \times & x_n \\ \text{OP1} & & \text{OP2} & & \text{OP3} & & \dots & & \text{OPn} \end{matrix}$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

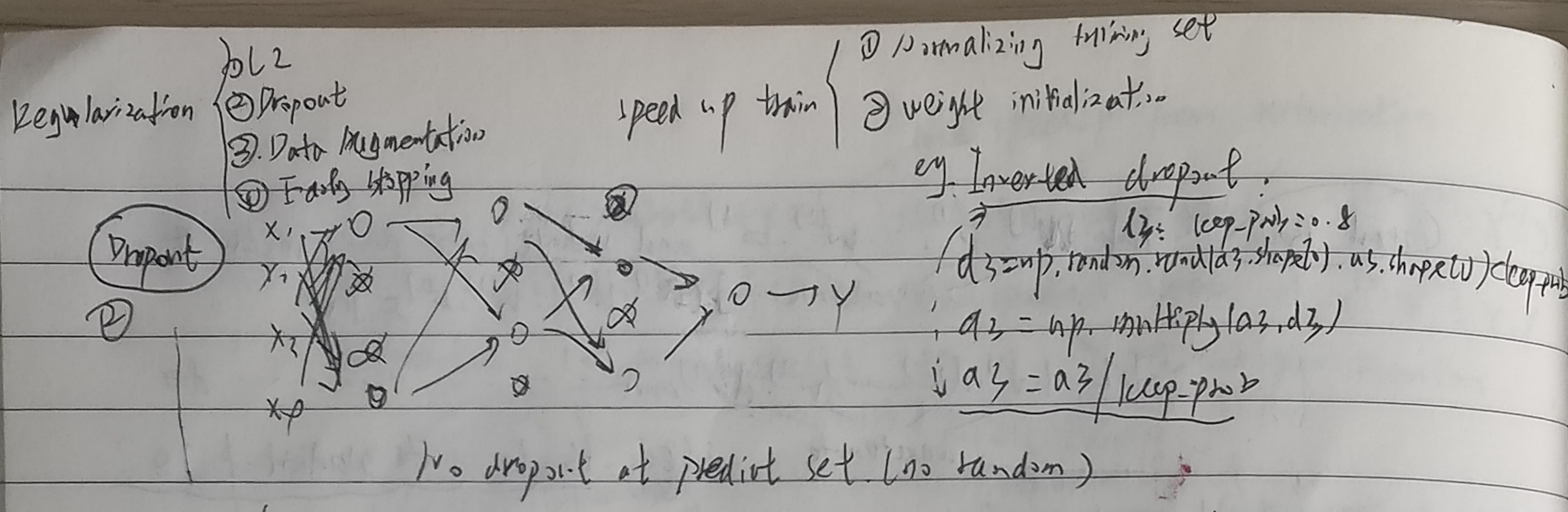
$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$

$$y = \text{OP}_1 \circ \text{OP}_2 \circ \text{OP}_3 \circ \dots \circ \text{OP}_n$$



downside: J is no longer well-defined. (trick: off dropout - J plot (than dropout))

Data augmentation

- ① flip, rotate, distort
- ② (CV)

Setting up your ML application

optimization problem

① subtract mean

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

2) Normalize variance

$$\sigma^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2$$

$$X = X - \bar{x}$$

$$X = X / \sigma$$

Same n, o for train/test

Weight initialization

L1

$$w_1 = 0$$

$$x_1 = 0$$

$$y_1 = 0$$

$$w_2 = 1$$

$$x_2 = 1$$

$$y_2 = 1$$

Gradient checking

$$f(\theta) = \theta^T \cdot \theta$$

$$f'(\theta) = 2\theta$$

$$\theta = 1, \epsilon = 0.01$$

$$\theta = 1, \epsilon = 0.001$$

$$\theta = 1, \epsilon = 0.0001$$

$$\theta = 1, \epsilon = 0.00001$$

$$\theta = 1, \epsilon = 0.000001$$

$$\theta = 1, \epsilon = 0.0000001$$

$$\theta = 1, \epsilon = 0.00000001$$

$$\theta = 1, \epsilon = 0.000000001$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta_i) - f(\theta_{i-1})$$

$$f'(\theta) = \frac{1}{2} \sum_{i=1}^n f'(\theta_i) - f'(\theta_{i-1})$$

$$\text{error: } D(f')$$

$$\text{error: } D(f)$$

$$f(\theta) = \frac{1}{2} \sum_{i=1}^n f(\theta$$

lose speed-up from vectorization

choosing mini-batch size

- if mini-batch size = 1: stochastic gradient descent $x^{(t)}, y^{(t)} = x^{(i)}, y^{(i)}$
- if mini-batch size = m: Batch gradient descent $x^{(t)}, y^{(t)} = X, Y$

in practice, $1 < \text{mini-batch size} < m$ → Too long per iteration

fastest learning:
vectorization
make progress without processing the entire set

$m=2000$: Batch gradient descent
 $m \gg m$: mini-batch

Exponentially weighted averages (moving average)

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

V_t : an average over $\frac{1}{1-\beta}$ days

$$\text{eg.: } V_{100} = 0.9 V_{99} + \theta_1 \theta_{100}$$

$$\begin{aligned} V_{100} &= 0.9 V_{99} + 0.1 \theta_{100} \\ &\quad \Rightarrow V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1 \times 0.9^2 \theta_{98} + 0.1 \times 0.9^3 \theta_{97} + \dots \\ V_{99} &= 0.9 V_{98} + 0.1 \theta_{99} \\ &\quad \Rightarrow V_{99} \approx 0.9 V_{98} + 0.1 \theta_{99} \end{aligned}$$

(statistics)

$$\begin{aligned} F_{t+1} &= (V_t + (1-\beta) F_t) \\ F_0 &= 0 \quad (\alpha = 1-\beta) \\ V_0 &= \theta_0 + (1-\beta)^0 t \quad (\text{computation, memory}) \end{aligned}$$

$$\begin{aligned} F_{t+1} &= \alpha F_t + \theta_{t+1} \quad (\text{bias correction}) \\ V_0 &= 0 \quad (\alpha = 1-\beta) \\ V_1 &= 0.98 V_0 + 0.02 \theta_1 \\ V_2 &= 0.98 V_1 + 0.02 \theta_2 \\ &= 0.96 \theta_0 + 0.02 \theta_2 \end{aligned}$$

Gradient descent with momentum

Value $\approx \theta_0 \approx \theta_t$ (initialization)

on iteration t : compute $d\theta$ on mini-batch.

$$V_t d\theta = \beta V_{t-1} d\theta + (1-\beta) d\theta$$

$$\begin{aligned} \theta' &= \theta - \frac{1}{1-\beta} V_t d\theta \\ &= \theta - \alpha V_t d\theta \quad (\text{friction}) \end{aligned}$$

$$\alpha' = \alpha \cdot \frac{1}{1-\beta} \quad (\text{hyperparameters: } \alpha, \beta)$$

average on last T_0 gradient

$$(\beta = 0.9)$$

RMSprop: root mean square prop

In iteration t :

compute $d\theta$, $d\theta$ on current mini-batch,

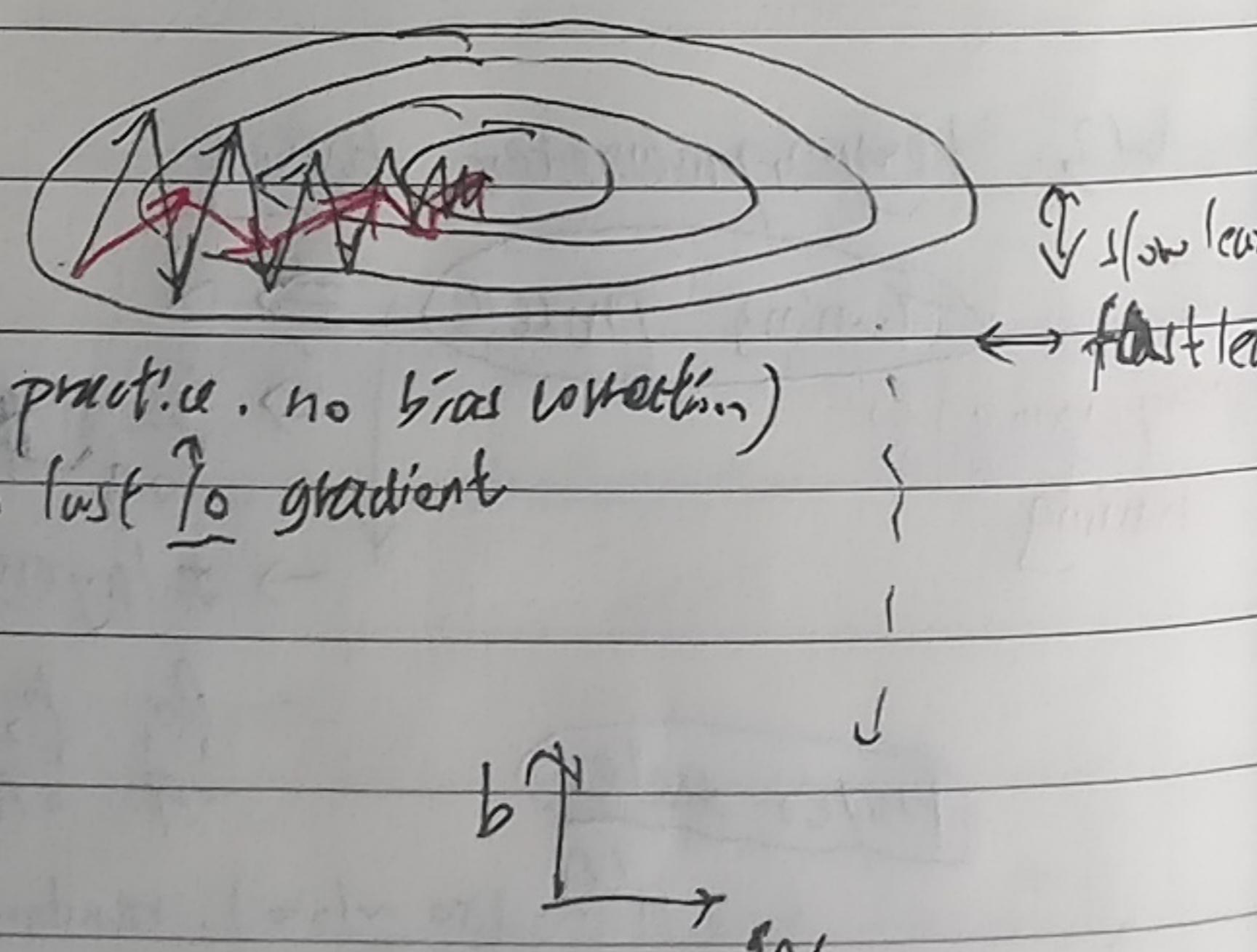
$$Sd\theta = \beta_2 Sd\theta + (1-\beta_2)(d\theta)^2 \quad (\text{element-wise})$$

$$Sd\theta = \beta_2 Sd\theta + (1-\beta_2)(d\theta)^2 \quad (\text{large})$$

$$W = W - \alpha \frac{d\theta}{\sqrt{Sd\theta + \epsilon}} \quad (\text{faster})$$

$$b = b - \alpha \frac{d\theta}{\sqrt{Sd\theta + \epsilon}} \quad (\text{slow})$$

$$10^{-8}, \text{ stability} \quad (\text{stability})$$



could range α

Grad check

Take $(W^{(1)}, b^{(1)}, \dots, W^{(L)}, b^{(L)})$

concatenate $[W^{(1)}, b^{(1)}, \dots, W^{(L)}, b^{(L)}]^\top$

Take $dW^{(1)}, db^{(1)}, \dots, dW^{(L)}, db^{(L)}$ and concatenate

check: Is $d\theta$ is the gradient of $J\theta$?

for each i :

$$\rightarrow d\theta_{\text{approx}}^{(i)} = \frac{J(\theta, \theta_1, \dots, \theta_i + \epsilon, \dots, \theta_L) - J(\theta, \theta_1, \dots, \theta_i - \epsilon, \dots, \theta_L)}{2\epsilon}$$

$$\approx d\theta^{(i)} = \frac{\partial J}{\partial \theta_i}$$

$$\text{check: } \|d\theta_{\text{approx}} - d\theta\|_2 \approx 10^{-7} \text{ - great!}$$

$$\|d\theta_{\text{approx}}\|_2, \|d\theta\|_2 \quad (\epsilon = 10^{-7}) \quad 10^{-5} \text{ - wrong.}$$

Notes:

- Don't use in training. - only for debugging.

a) If fails, check component.

b) Remember regularization: $J\theta = \frac{1}{m} \sum_{i=1}^m (y_i^{(i)} - \theta^\top x_i^{(i)})^2 + \frac{\lambda}{2m} \|\theta\|^2$

c) doesn't work with dropout.

d) Run at random initialization; perhaps check again after some training. $w, b \neq 0$

W2: Optimization Algorithm. (learning factor)

mini-batch gradient descent

$$x^{(t)}, y^{(t)} \quad \text{batch size: for example } 1000$$

$$x = (n^*, m), y = (1, m)$$

repeat

for i in m/batch-size:

Forward prop on $x^{(t)}$:

$$z^{(i)} = W^{(1)} x^{(i)} + b^{(1)}$$

$$\delta^{(i)} = g^{(1)}(z^{(i)})$$

backward pass

$$\delta^{(1)} = g^{(1)}(z^{(1)})$$

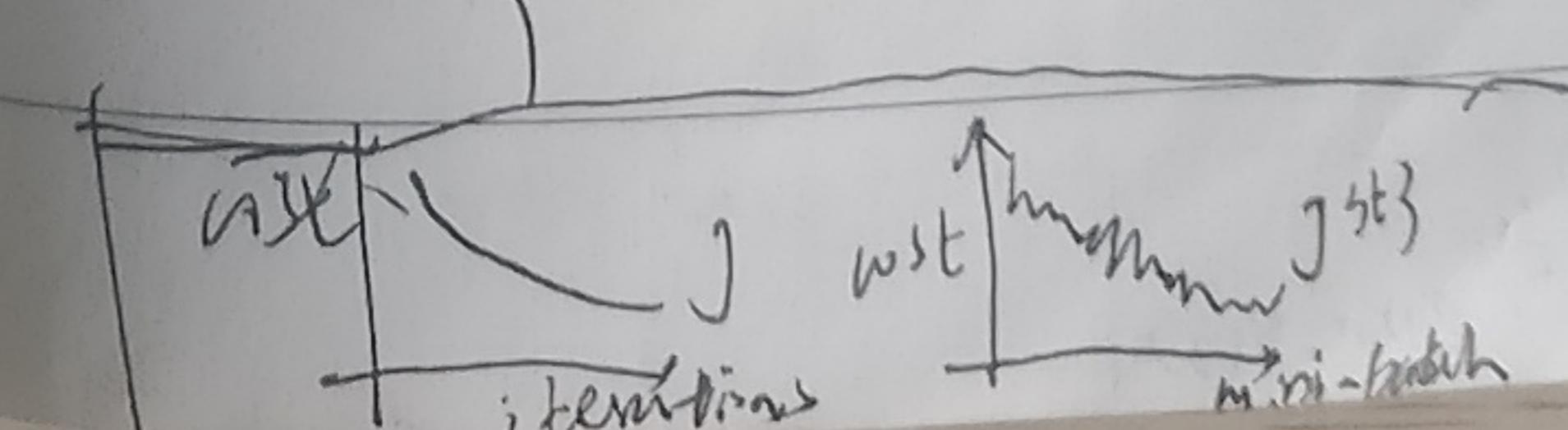
compute cost $J^{(t)}$ for $x^{(t)}, y^{(t)}$

$$J^{(t)} = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, z^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|W^{(l)}\|_F^2$$

Back to compute gradients $J^{(t)}$ (using $x^{(t)}, y^{(t)}$)

$$W^{(t)} = W^{(t)} - \alpha dW^{(t)}, b^{(t)} = b^{(t)} - \alpha db^{(t)}$$

epoch



BN Batch normalization BN

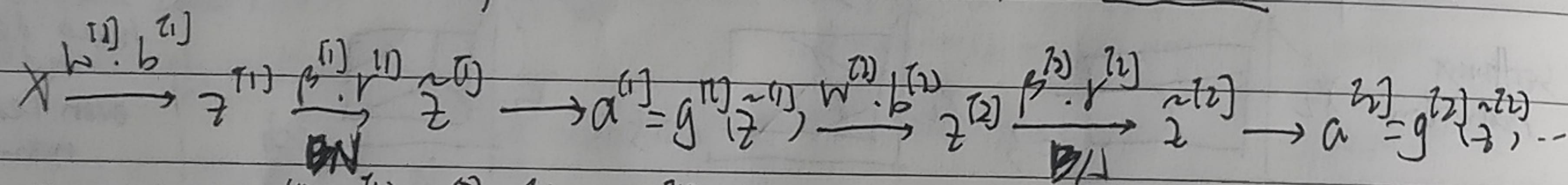
Gives some intermediate values in NN: $\tilde{z}^{(1)}, \dots, \tilde{z}^{(L)}$

$$- M = \frac{1}{m} \sum_{i=1}^m \tilde{z}^{(1)i}, \quad \hat{b} = \frac{1}{m} \sum_{i=1}^m (\tilde{z}^{(1)i} - M)^2$$

$$\tilde{z}^{(1)i} = \frac{\tilde{z}^{(1)i} - M}{\sqrt{\hat{b}} + \epsilon}$$

$$\tilde{z}^{(1)i} = r \tilde{z}^{(1)i}_{norm} + \beta \quad \text{learnable parameters (control range)}$$

Use $\tilde{z}^{(1)i}$, instead of $\tilde{z}^{(1)i}$.



Parameters: $(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)})$

$$\begin{aligned} & \text{minibatch } \\ & X^{(1)} \xrightarrow{W^{(1)}, b^{(1)}} \tilde{z}^{(1)} \xrightarrow{g^{(1)}, \gamma^{(1)}} z^{(1)} \xrightarrow{W^{(2)}, b^{(2)}} \tilde{z}^{(2)} \xrightarrow{g^{(2)}, \gamma^{(2)}} z^{(2)} \xrightarrow{\dots} \\ & X^{(2)} \xrightarrow{W^{(2)}, b^{(2)}} \tilde{z}^{(2)} \xrightarrow{\dots} \dots \end{aligned}$$

Parameters: $w^{(1)}, \beta^{(1)}, \gamma^{(1)}$

$$\begin{matrix} \tilde{z}^{(1)} \\ \downarrow \\ z^{(1)} \\ \downarrow \\ z^{(1)norm} \\ \downarrow \\ z^{(1)} \end{matrix}$$

$$\begin{aligned} \tilde{z}^{(1)} &= W^{(1)} a^{(1)} + \beta^{(1)} \\ z^{(1)} &= (\tilde{z}^{(1)} - M^{(1)}) / \sqrt{\hat{b}^{(1)}} \\ z^{(1)norm} &= \gamma^{(1)} z^{(1)norm} + \beta^{(1)} \end{aligned}$$

Why BN works?

Learning on shifting input distribution
disjoint adjacent layers

BN as regularization. (Scaled by mean/variance) inst on mini-batch. (batch, regularization)
noise similar to dropout.

$\beta \neq 0$ at test time

no mini-batch.

$$M, \hat{b} \text{ (moving average). } \tilde{z}^{(1)}_i, \tilde{z}^{(2)}_i, \dots, \tilde{z}^{(L)}_i$$

$$\tilde{z}^{(1)}_i = \frac{\tilde{z}^{(1)}_i - M}{\sqrt{\hat{b}} + \epsilon}$$

$$\tilde{z}^{(1)}_i, \tilde{z}^{(2)}_i, \dots \rightarrow M \text{ same for } \hat{b}$$

"covariate shift": black cat \rightarrow colored cat

noise similar to dropout.

Adam Adaptive moment estimation

$$V_{dw} = 0, V_{db} = 0, S_{dw} = 0, S_{db} = 0$$

On iteration t:

compute d_w, d_b using current minibatch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) d_w, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) d_b \quad \leftarrow \text{"momentum": } \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) d_w^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) d_b^2 \quad \leftarrow \text{"rmsprop": } \beta_2$$

$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{corrected} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{corrected} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{corrected} = S_{db} / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$$

Learning rate decay

$$\text{① } \alpha = \frac{\text{fix decay rate} * \text{epoch num}}{\text{epoch num}} \quad ; \quad \text{② } \alpha = 0.95^\text{epoch num} \cdot \alpha_0, \text{ exponentially decay}$$

$$\text{③ } \alpha = \frac{k}{\text{epoch num}} \cdot \alpha_0, \text{ or } \alpha = \frac{k}{nt} \cdot \alpha_0 \quad \text{④ } \alpha \quad \begin{matrix} \uparrow \\ \text{--- discrete.} \end{matrix}$$

⑤ by hand

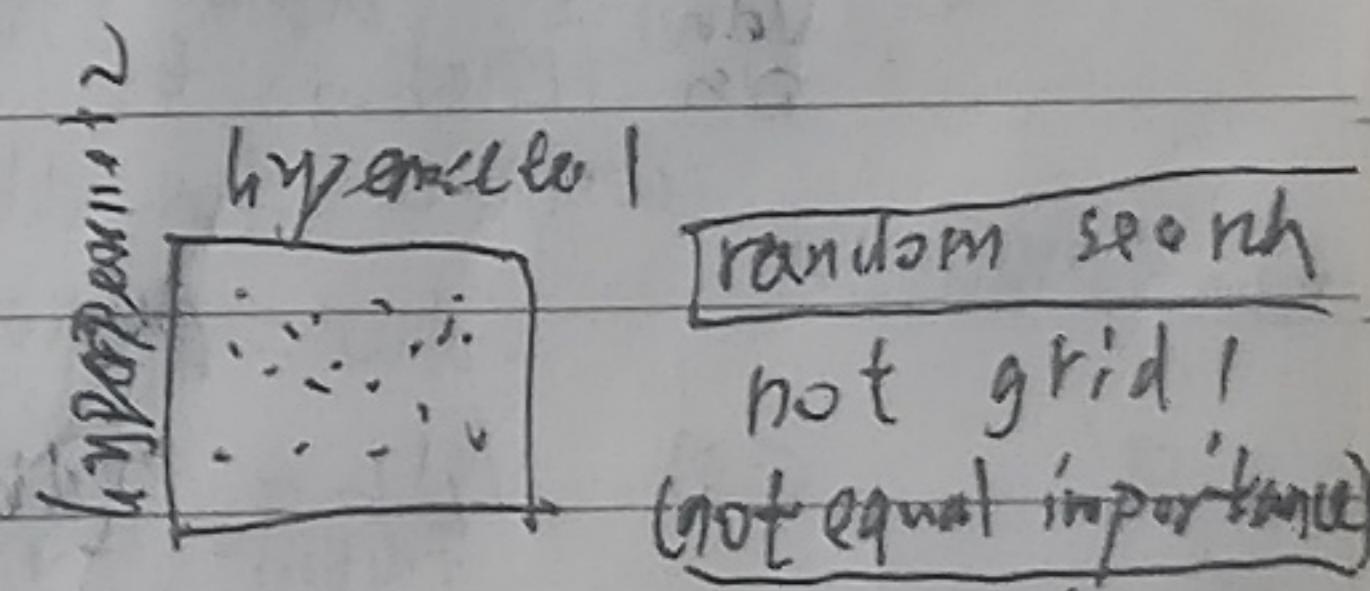
Total optimal

low-dimension \rightarrow high dimension

→ unlikely to get stuck in a bad local optima → more saddle, not local optimal

2): Plateaus can make learning slow. (rmsprop, adam)

all features are $\sqrt{2^{-20000}}$



W3: Hyperparameter tuning

Tuning process

Hyperparameters
Tuning

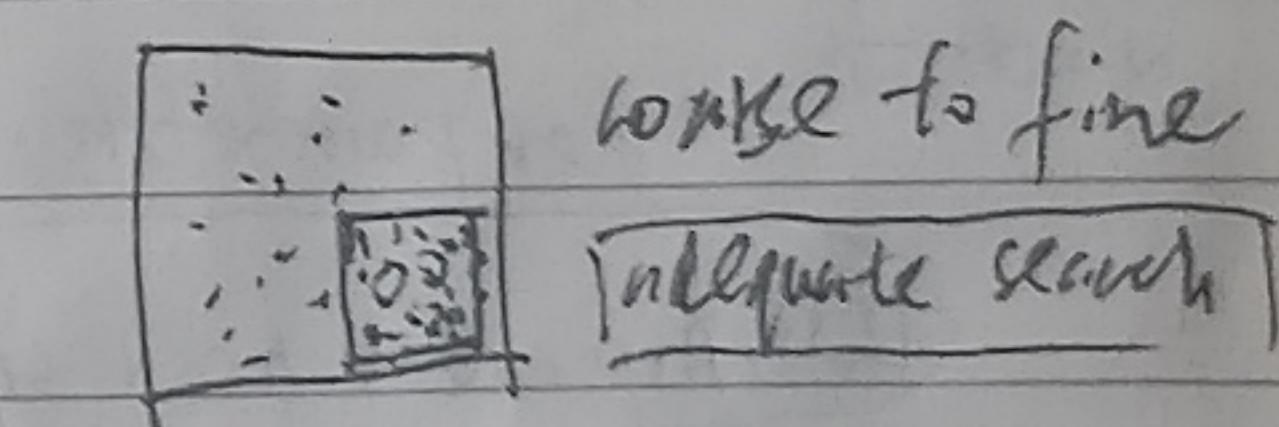
$\Rightarrow \alpha$

$\Rightarrow \beta$: #hidden units, mini-batch size

$\Rightarrow \# \text{ layers, learning rate decay}$

$\dots, \beta_1, \beta_2, \dots$

$\dots, 0.9, 0.99, 10^{-8}$



Proper scale

$\alpha^{(1)}: (0.01 \sim 1)$: random search

#layers: $(2 \sim 4)$: grid search

$\alpha: (0.001 \sim 1)$: log scale

$(r = np.random.rand() * 10^{-4})$

$\alpha = 10^r$

$\beta = (0.9 \sim 0.99)$

$1 - \beta: (0.1 \sim 0.001)$

$10^{-1} \sim 10^{-3}$

$r = np.random.rand()$

$\beta = 1/2 * (1-r)$

building one model

Training many models in parallel

resources

Punjas

Carina

10

C3. Structuring your ML strategy

W1: Introduction to ML strategy

(why ML strategy?)

- Ideas: collect more data
- Collect more diverse training set
- Train algo longer with gradient descent
- Try adam instead of gradient descent
- Try bigger network

(Orthogonalization)

- each knob kind of does only one thing
- Fit training set well on cost function
↓
(≈ human level performance)
- Fit dev set well on cost func.
- Fit test set well on cost func.
- Performs well in real world

(Using a single number evaluation metric)

Setting up your goal: F1 instead of Precision or Recall. $F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}}$
sometimes average.

Dev set + single real number evaluation metric. speed up iterating.

(Optimizing and satisfying metrics)

eg. car classification example:		optimizing metric	satisfying metric
classifier	accuracy	running time	maximize accuracy
1	9%	8 hours	subject to running time < 10 hrs
2	92%	9 hours	
3	95%	15 hours	

N metrics { 1 optimizing
2 satisfying

(train/dev/test distribution)

Guideline: choose a dev and test set to reflect data you expect to get in the future

↓
same distribution and consider important to do well on.
Randomly shuffle to dev/test

set size: old way: 70% / 30% or 60% / 20% / 20%

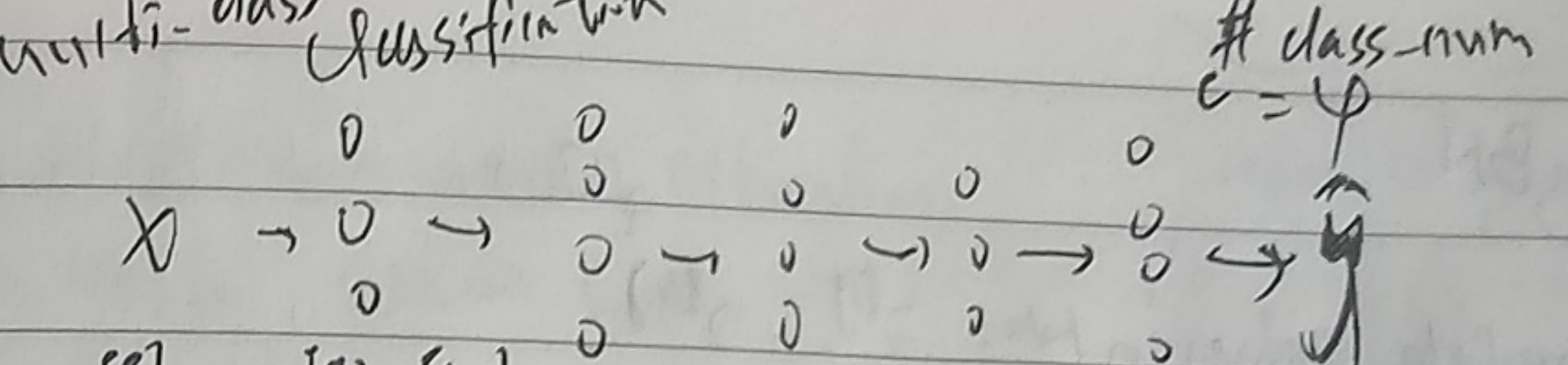
big data norm: 98% / 1% / 1%

guidelines:

- set test set big enough to give high confidence in the overall performance in your system.
- Test set may < 30%
- The dev set has to be big enough to evaluate different ideas

multi-class classification

Multi-class classification
S2 framework



hard max:

$$\text{eg: } z = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{3} \\ \vdots \\ \frac{1}{c} \end{bmatrix}$$

$$z(x) = w^{(1)}(x) + b^{(1)}$$

activation function:

$$\rightarrow t = e^{z(x)}$$

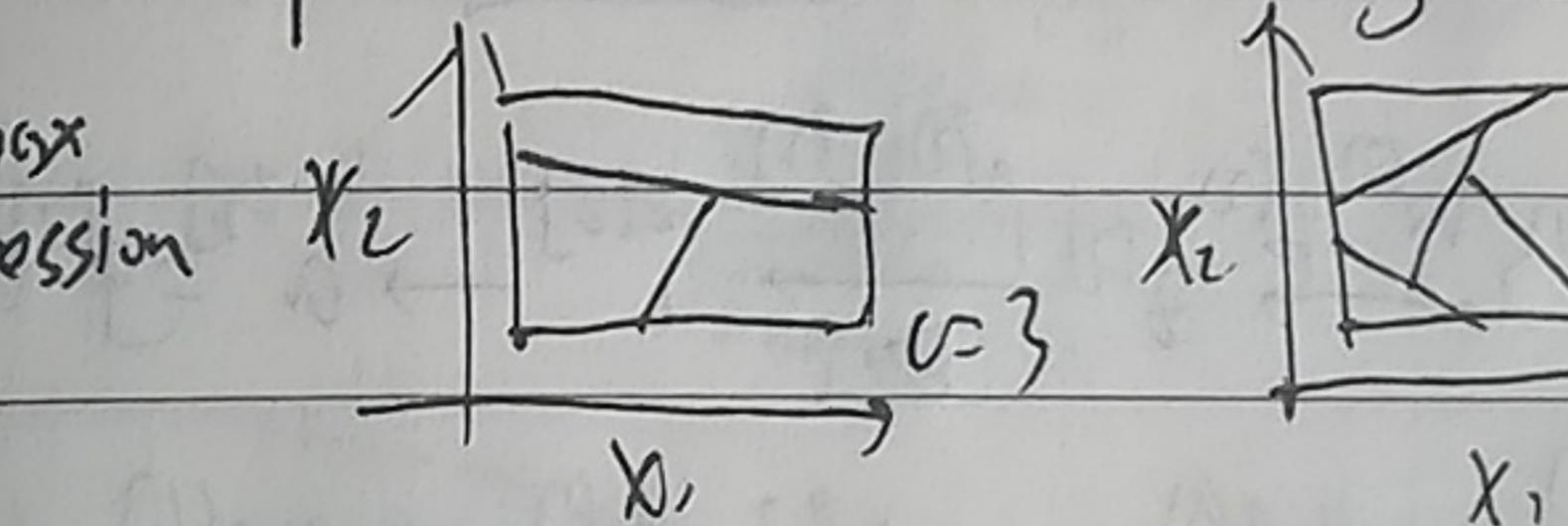
$$\rightarrow a^{(1)} = \frac{e^{z(x)}}{\sum_i e^{z_i}}, a_i^{(1)} = \frac{e^{z_i}}{\sum_i e^{z_i}}$$

$$t = \begin{bmatrix} e^{\frac{1}{2}} \\ e^{\frac{1}{3}} \\ \vdots \\ e^{\frac{1}{c}} \end{bmatrix}, \sum_i t_i = 176.3$$

$$a^{(1)} = \frac{t}{176.3} \rightarrow \begin{bmatrix} 0.842 \\ 0.046 \\ \dots \\ 0.099 \end{bmatrix}$$

examples: decision boundary.

softmax regression



loss function

$$\text{eg: } y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \hat{y}^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.3 \\ 0.1 \end{bmatrix}, c=3.$$

$$J(w^{(1)}, b^{(1)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)})$$

$$\ell(y^{(i)}, \hat{y}^{(i)}) = -\sum_{j=1}^c y_j^{(i)} \log \hat{y}_j^{(i)} = -y_2^{(i)} \log \hat{y}_2^{(i)} = \log \hat{y}_2^{(i)}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}], \hat{Y} = [\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)}]$$

$$= \begin{bmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \in \mathbb{R}^{4 \times m}, \quad = \begin{bmatrix} 0.3 & 0.2 & \dots \\ 0.2 & 0.3 & \dots \end{bmatrix} \in \mathbb{R}^{4 \times m}$$

$$\text{Back Prop. } d\ell/dw^{(1)} = \hat{y} - y = \frac{d\ell}{d\hat{y}^{(1)}}$$

Intro to Deep learning framework
programming framework

Caffe/Caffe2

TensorFlow

Keras

Lasagne

mxnet

PaddlePaddle

TensorFlow

Theano

PyTorch

Criteria:

- Ease of programming (development and deployment)
- Running speed
- Truly open (good governance)

W2. Error Analysis

J Confident error analysis

Ideas for cat detection:

- 1. Fix pictures of dogs being recognized as cats
- 2. Fix great cats (tigers, panthers, etc...) being misrecognized.
- 3. Improve performance on blurry images.

Error Analysis

Index	dog	great cat	blurred	Unknown
1	✓			
2		✓		
3		✓	✓	
4				✓
5				
Total	8%	43%	61%	6%

Error analysis:
 ① Get ~100 images misrecognized on dev set
 ② Count up how many are others.
 → choose direction to improve.

Clean up incorrectly labeled data

DL is quite robust to random errors in the training set.
 not systematic error,

Gold Incorrectly labeled rate, if overall dev/test err 10% | 2%
 Errors due to incorrect labels 0.6% | 0.6%

Errors due to other causes 9.4% | 1.4%

not important right now more worthwhile to fix

Count of dev set:
 help you select between classifier A & B

Guides: • Apply same process to dev/test sets to make sure they continue to come from same distribution

• Consider examining examples your algo got right as well as ones it got wrong.

• Train and dev/test set may now come from slightly different distributions.

Build your first system quickly, then iterate

Speech algorithm example.

Background noise:

Cafe noise

Car noise

Accented speech

Traffic from microphone

Young children's speech

Setup dev/test set and metric

Build initial system quickly.

• Use bias/variance analysis & Error analysis

to prioritize next steps.

unless familiar with this domain or more performance, the more (PAPER...)

J Mismatched training and testing on different distributions

Cat app example

Data from web: 200k ; Data from app: 10k

shuffle

X option1: train: 200k 2.5k 2.5k

✓ option2: train: web: 200k 2.5k 2.5k

app: 5k 2.5k 2.5k

speech recognition example

Training

Published

Smart speaker control

Voice keyboard

500k

Dev/Test

Speech acted review filter

10k ✓ 5k ✓ 5k

When to change dev/test set and metrics

↳ doing well on your metrics & dev/test set,
 but not doing well on application

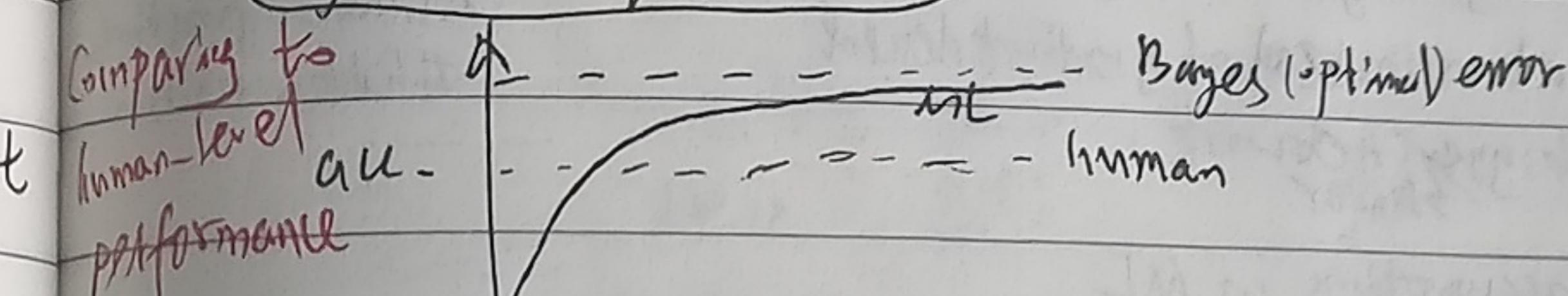
→ how to define a metric
 orthogonalization

→ 2. how to do well on this metric: How to look

$$J = \frac{1}{m} \sum_{i=1}^m \ell(y_i^{(i)}, \hat{y}_i^{(i)}) \Rightarrow J = \frac{1}{2m} \sum_{i=1}^m \text{w}(y_i^{(i)}, \hat{y}_i^{(i)})$$

push note

J Why human performance?



- human are quite good at some tasks.
- as long as m is worse than human, you can get labeled data from humans
- gain insight from manual error analysis
- using human get this right??
- Better analysis of bias/variance

Unavoidable bias

Human (x Bayes) 1% ; 7.5%

Training error 8% ; 8%

Dev error 1% ; 1%

↓ avoidable bias

↓ Variance

Focus on bias Focus on variance

When approach human-level, not know how far away you are from Bayes error,
 therefore how much you should be trying to reduce avoidable bias

harder to tease out the bias and variance effects,

Surpassing human-level performance

Team of human	0.5%	↑	0.1% ?	0.1% ?
Human	1%	↑	0.2% ?	0.2% ?
Train error	0.6%	↓	0.3%	0.3%
Dev error	0.8%	↓	0.4%	0.4%

ML significantly surpassing human
 Online advertising : Speech recg.
 Product recommendations : Same w/
 Logistics (Predict transit time) : Media:
 Loan approval : Bots
 ↓
 Structured data.
 not natural perception

Improving your model performance

two fundamental assumptions of supervised learning

1. fit train set well

~ avoidable bias

2. train set performs well on dev/test set.

~ variance

Human-level
Avoidable error

Training error
Variance

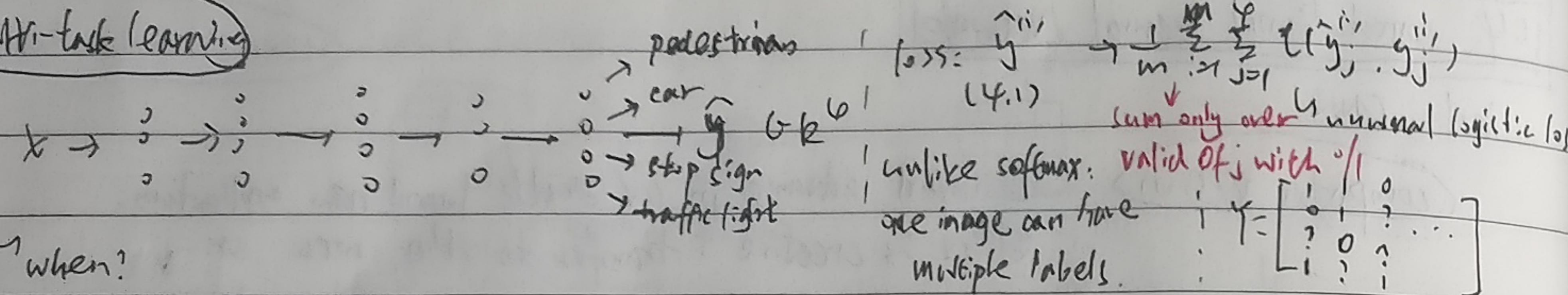
Dev error

Train bigger model
 Train longer/faster optimization
 - Momentum, backprop, adam
 NN architecture
 Hyperparameter search

More data
 Regularization
 - L2, No point, Data augmentation

NN architecture
 Hyperparameter search

Multi-task learning

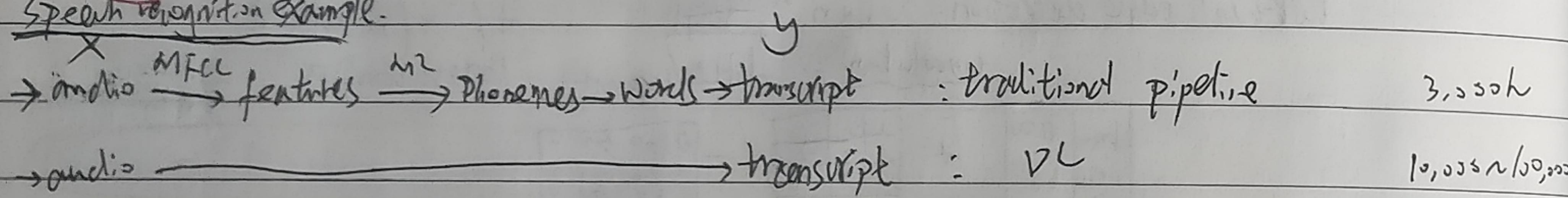


when?

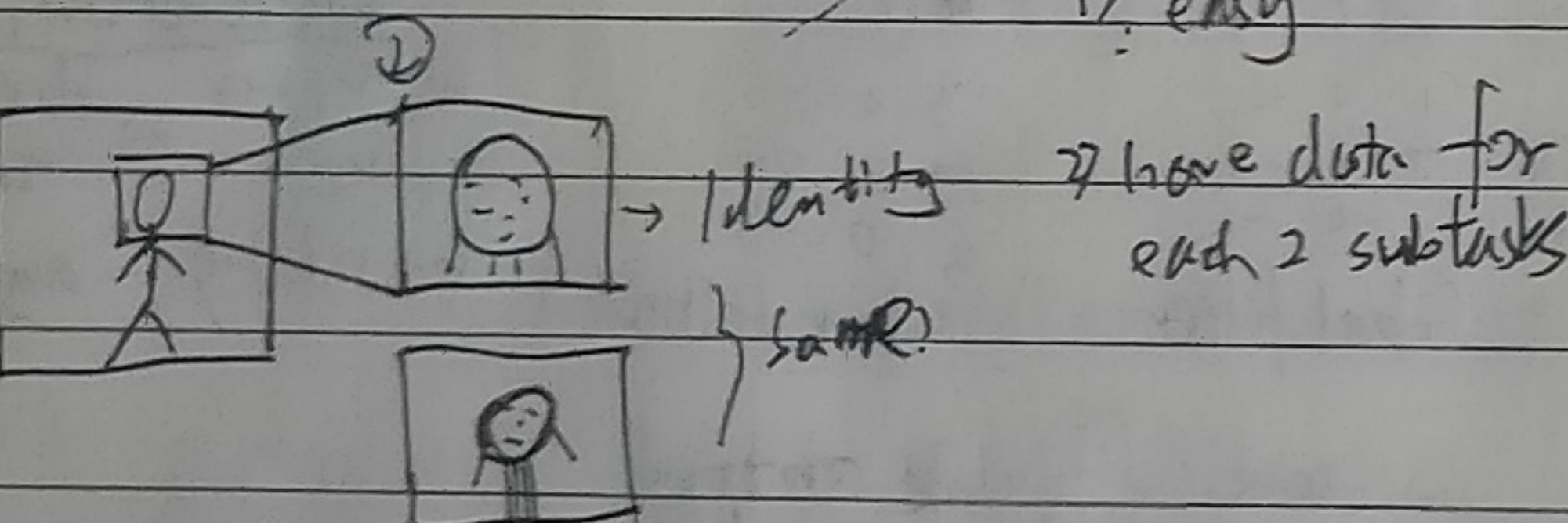
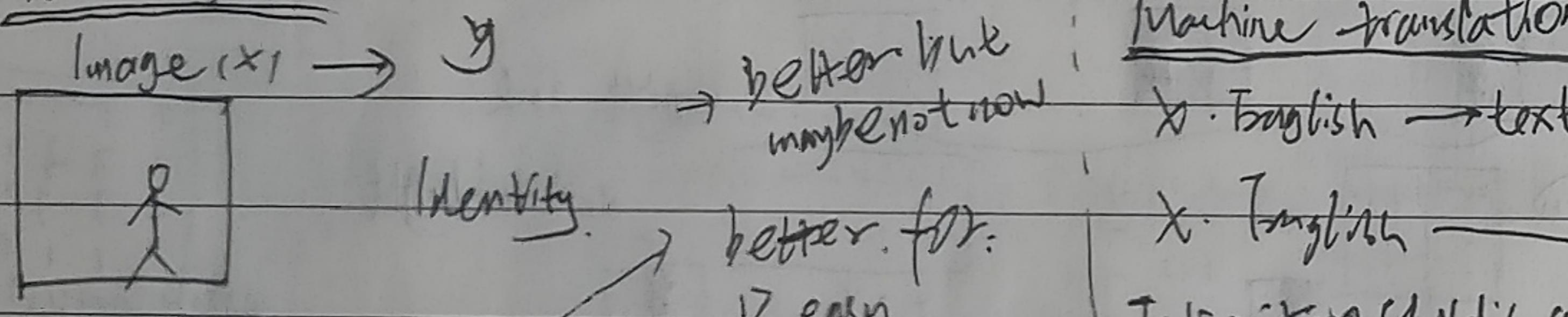
- Train on a set of tasks that share low-level features.
- Usually, amount of data you have for each task is quite similar.
- Can train a big enough neural network to do well on all the tasks.

End-to-end DL

Speech recognition example



Fake recognition

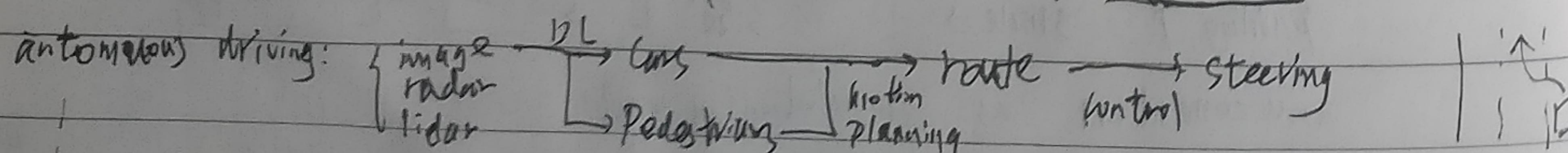


End-to-end DL

- Pros: ① let the data speak
 - ② Less hand-designing of components needed.
- Cons:
- ① May need large amount of data
 - ② Excludes potentially useful hand-designed components.

When to apply end-to-end DL?

key question: do you have sufficient data for a function of the complexity needed to map x to y ?



- Use DL to learn individual components.
- Carefully choose $x \rightarrow y$ depending on what task you get data for

End-to-End is not mature now.

Bias and Variance with mismatched data distribution

Cat classifier example

"Human level" $\approx 0\%$

Train error 1%

Train-dev error 1%

Dev error 10%

Test error 10%

Same distribution

Train error 1%

Variance 1%

Train-dev error 0%

Dev error 10%

Data mismatch 10%

Test error 10%

Human level \downarrow avoidable bias 0%

Train set error \downarrow variance 1%

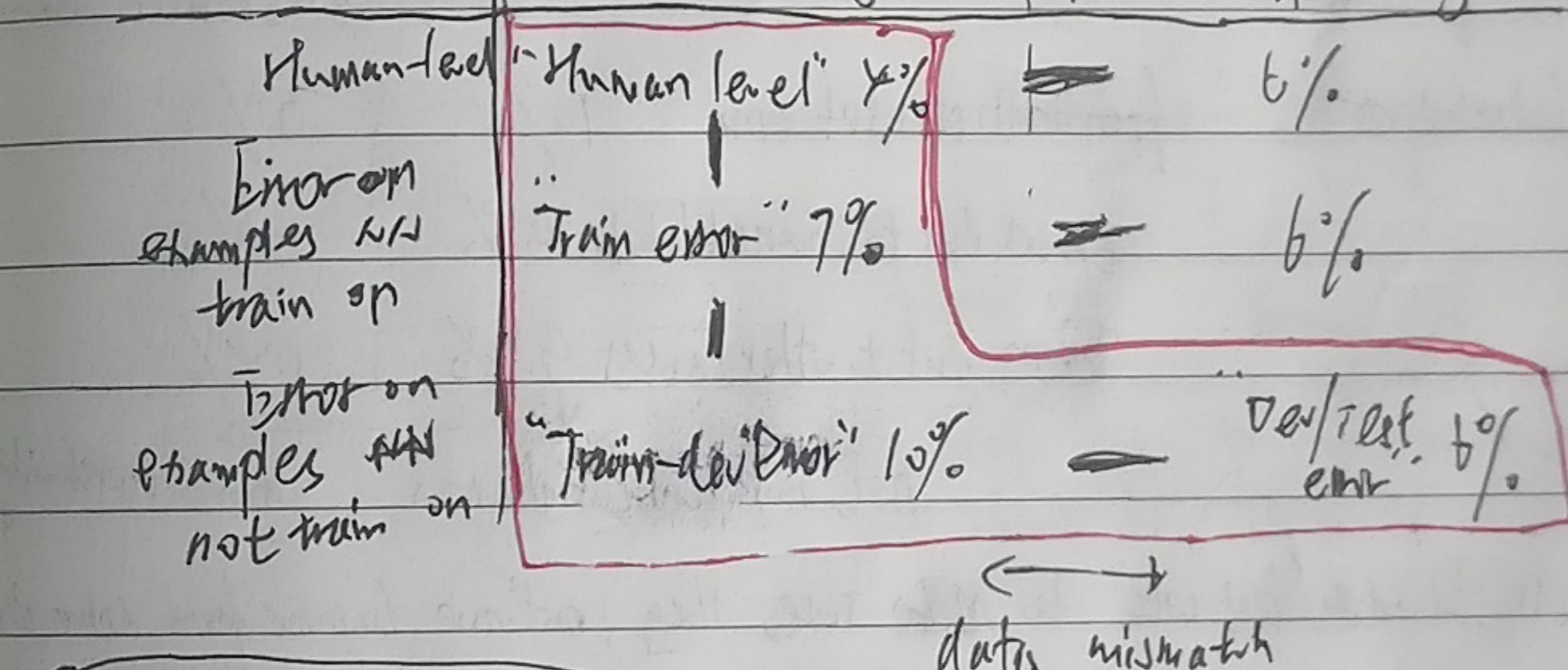
Train-dev set error \downarrow Data mismatch 1%

Dev set error \downarrow degree of overfitting to dev set 10%

Test error \downarrow degree of overfitting to dev set 10%

more general formulation (rearview mirror example)

general speech recognition | Rearview speech recognit...



Addressing data mismatch

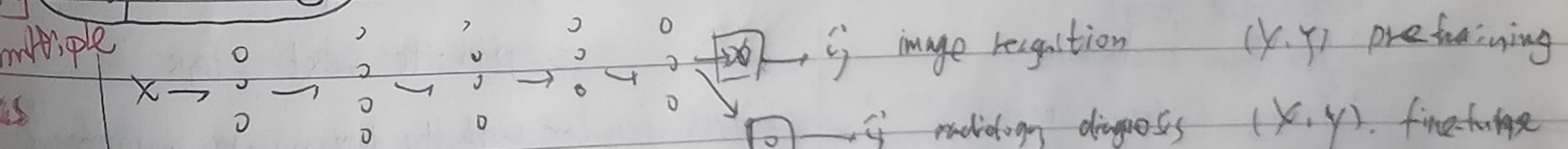
- Carry out manual error analysis to try to understand difference between training and dev/test sets.
e.g. noisy car noise, street numbers.
- Make training data more similar; or collect more data similar to dev/test set.
e.g. simulate noisy in-car data, clean speech + car noise = synthesized in-car audio

Artificial data synthesis

clean speech + car noise = synthesized in-car audio

1h. synthesis of all audio in car may overfitting

Transfer learning



When? $A \rightarrow B$

- Task A, B have same input X .
- But more data for A than B.
- Low level features from A would be helpful for B.

W2 : case studies

Why look at case studies?
case studies

classic networks : LeNet-5

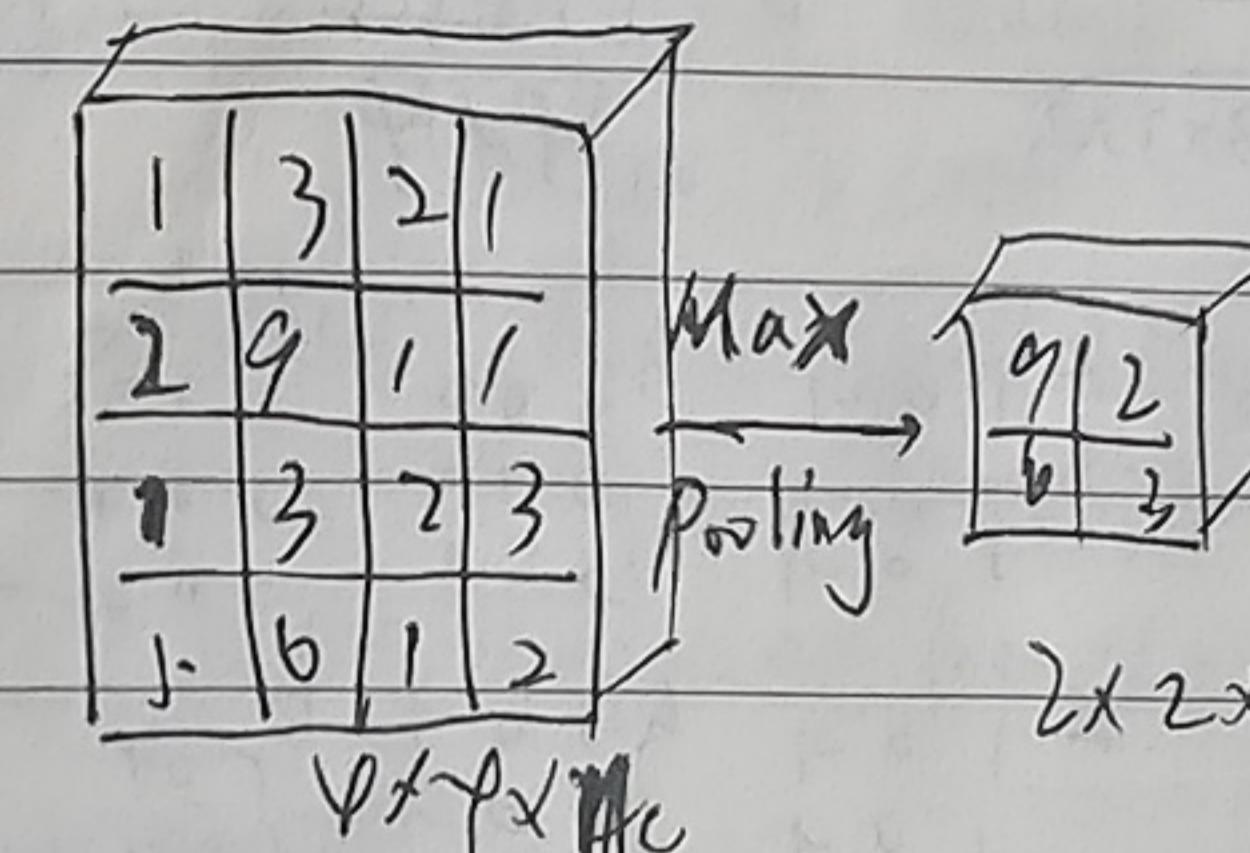
AlexNet
VGG

ResNet
Inception

Types of layers in a convolutional network

Convolution (CONV)
Pooling (POOL)
Fully-connected (FC)

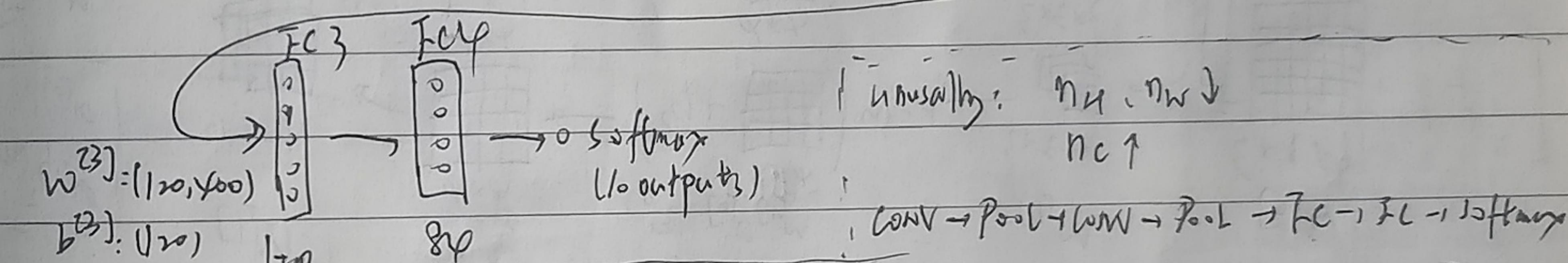
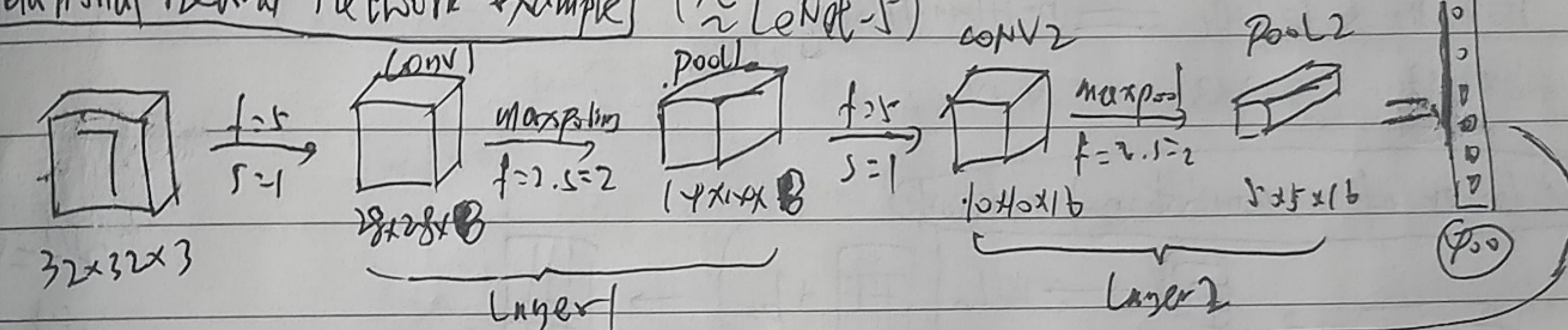
Pooling layers



Hyperparameters: $f=2$, $s=2$; max/average
(no need to learn) pooling rate

$$\frac{n+2p-f}{s} + 1, n_h \times n_w \times n_c \Rightarrow \left[\frac{n_h - f + 1}{s} \right] \times \left[\frac{n_w - f + 1}{s} \right]$$

Convolutional neural network example (\approx LeNet-5)

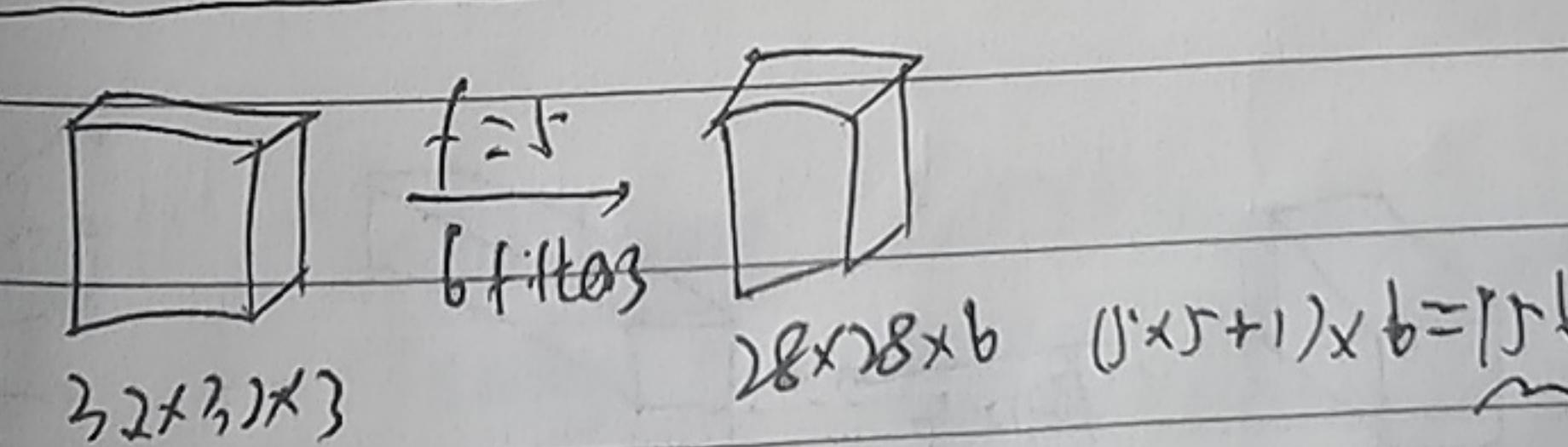


unusually: $n_h, n_w \downarrow$
 $n_c \uparrow$

CONV → Pool → CONV → Pool → FC → FC → Softmax

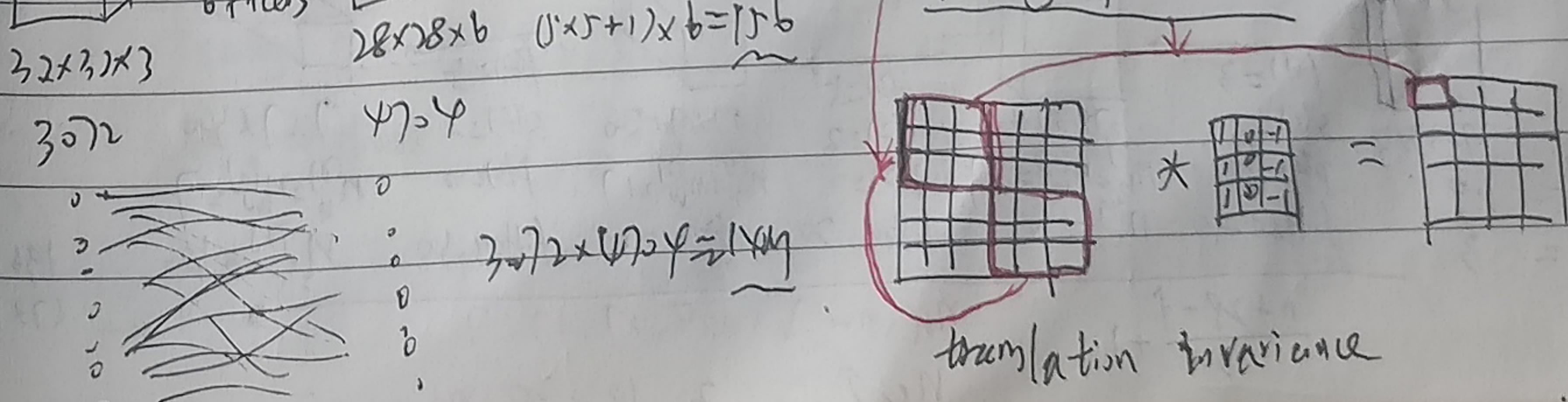
	Activation shape	Activation size	# Parameters
Input	(32, 32, 3)	3072	0
CONV (f=5, s=1)	(28, 28, 8)	6272	$208(5 \times 5 \times 3) \times 8$
POOL1	(14, 14, 8)	1568	0
CONV (f=5, s=1)	(10, 10, 16)	1600	$416(5 \times 5 \times 8) \times 16$
POOL2	(5, 5, 16)	400	0
FC3	(120, 1)	120	$4800(400 \times 16 + 1)$
FC4	(84, 1)	84	$108(120 \times 16 + 1)$
Softmax	(10, 1)	10	84(84 \times 1 + 1)

Why convolution?

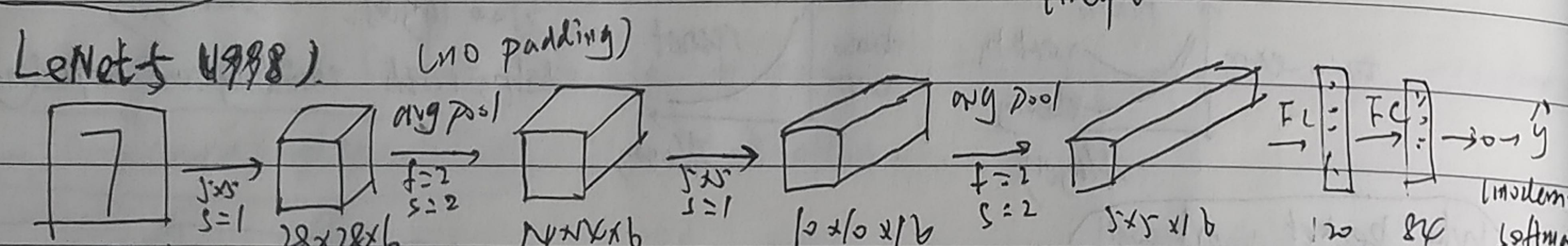


Parameter sharing

Sparsity of connections

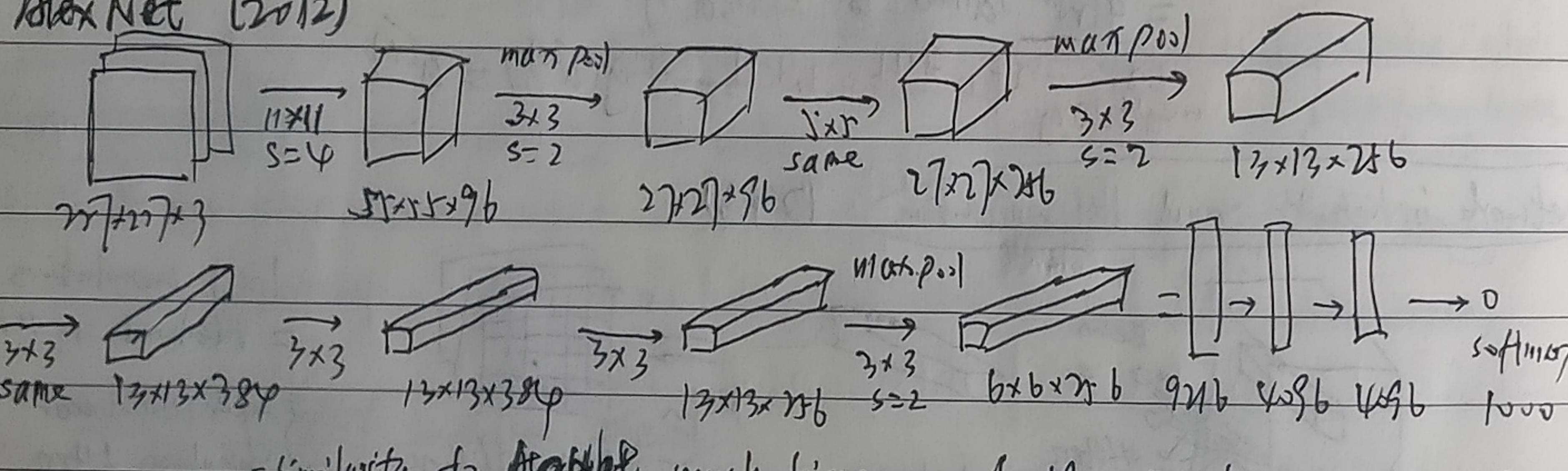


Classic networks



- 60k parameters, $n_h, n_w \downarrow$, $n_c \uparrow$
- conv - pool - conv - pool - fc - fc - output
- ① Sigmoid/Tanh, not ReLU.
- ② save computation, different filter paper.
- ③ non-linearity after pooling different (ReLU) section II

AlexNet (2012)

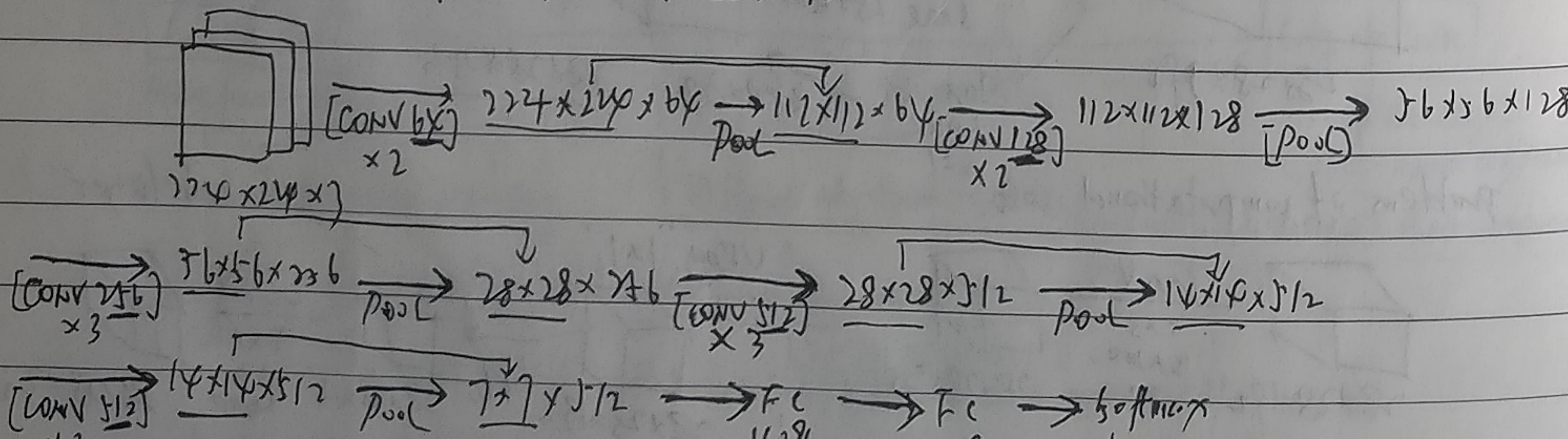


- similarity to LeNet, much bigger number of parameters.
- ReLU
- multiple GPUs
- local response normalization (LRN) (doesn't help so much)

(start convise)
in W

VGG16 (2014)

CONV = 3x3 filter, $s=1$, same; MAX-POOL = 2×2 , $s=2$

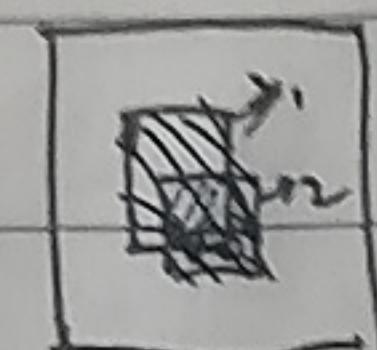


→ 16: layers with weights

- simple principle, uniform layers

~ 138M

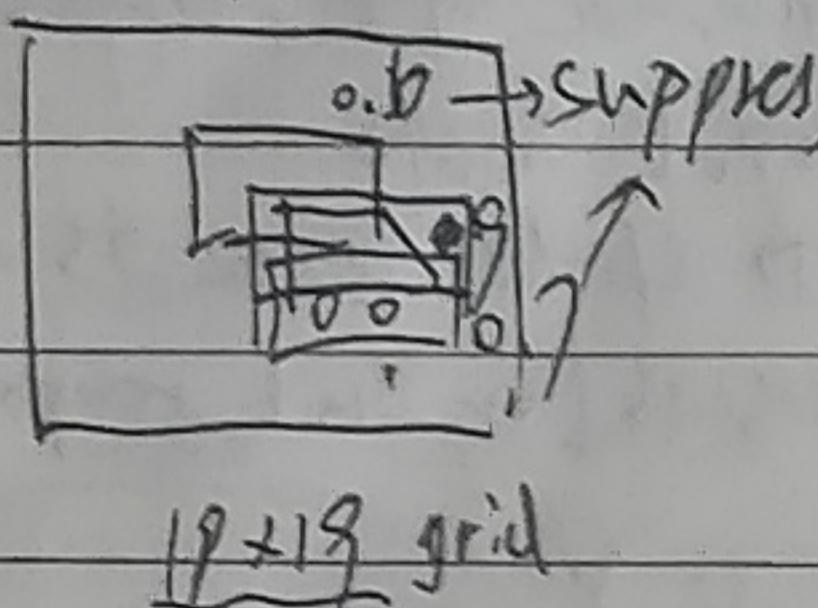
Intersection over Union (IoU)



$$IoU = \frac{\text{size of intersection}}{\text{size of union}}$$

"Correct" if $IoU > 0.5$

Non-max suppression



Each output prediction is $\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$. Output maximal probabilities classifications, but suppress it close by ones that are non-maximum.

Discard all boxes with $P_c \leq 0.6$.

while remaining boxes:

- Pick the box with the largest P_c output that has a prediction.
 - Discard any remaining boxes with $IoU \geq 0.5$ with previous box.
- (multiple object, independently implement non-max.)

Anchor Boxes

overlapping objects

Previously:

output: $3 \times 3 \times 8$

with anchor box:
anchor box 1 anchor box 2

output: $3 \times 3 \times 1$ to $(3 \times 3 \times 2 \times 8)$. (grid cell, anchor box)

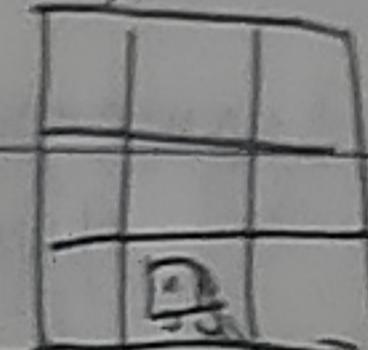
two anchor but 3 object, two object same grid \Rightarrow 19x19 grid take for the grid cell with highest IoU.

anchor box: by hand, by algo (k-means)

YOLO Algorithms

Train

1. pedestrian
2. car
3. motorcycle

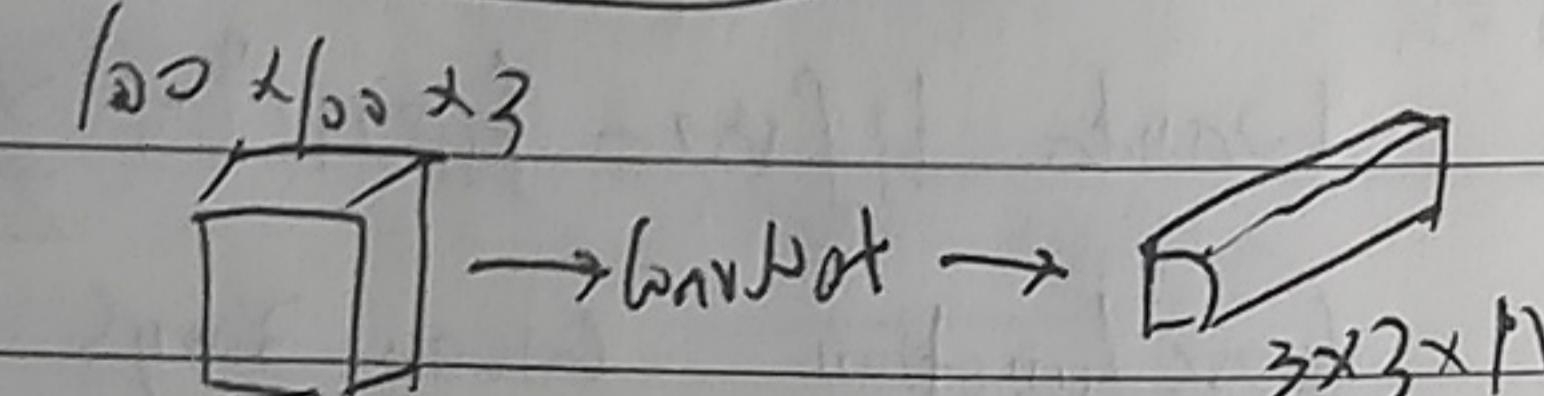


$$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \\ P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

grid: $3 \times 3 \times 2 \times 8$
grid anchors (real size: $19 \times 19 \times 5 \times 8$)

shape predictions

outputting the non-max suppressed output



• For each grid cell, get 2 predicted bounding boxes.

• Get rid of probability predictions.

• For each class, use non-max-suppression to generate final predictions

Region proposal

→ R-CNN: Propose regions, classify proposed regions one at a time. output: label + bounding box
(slow)

→ Fast R-CNN: Propose regions. Use convolutional implementation of sliding window to classify all regions.

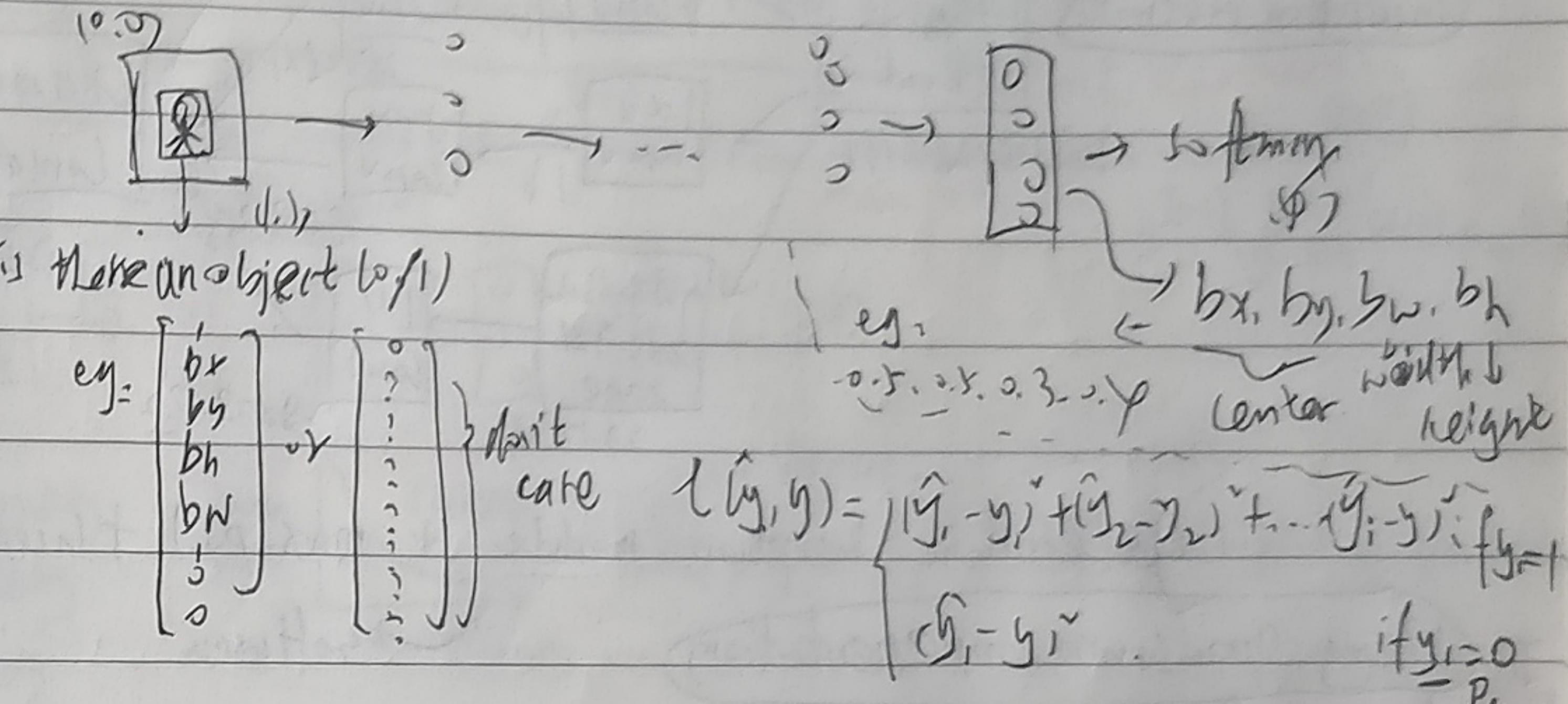
→ Faster R-CNN: Use convolutional networks to propose regions

faster than YOLO

instead of traditional segmentation algs.

W3. object detection

Object localization

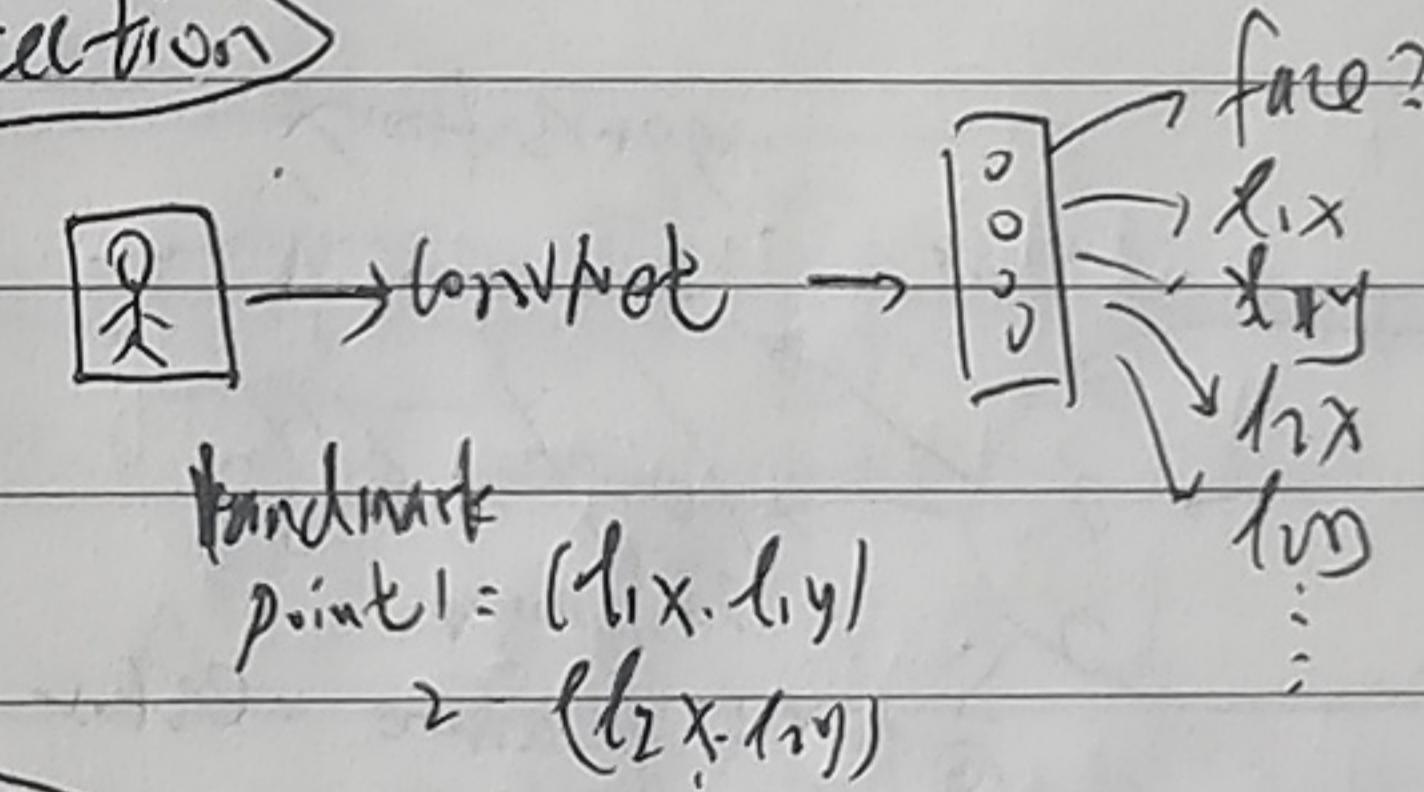


$$\text{target } y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \text{is there an object } (0/1)$$

$$y_g = \begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

$$\text{cate } \ell(y_g, y) = \ell(\hat{y}_1 - y_1 + \hat{y}_2 - y_2 + \dots + \hat{y}_n - y_n) ; \hat{y}_i = \begin{cases} \hat{y}_i & \text{if } y_i = 0 \\ P_c & \text{otherwise} \end{cases}$$

Landmark detection



AP hat
emotion, face recognition,
posture...

Object detection

Sliding windows detection

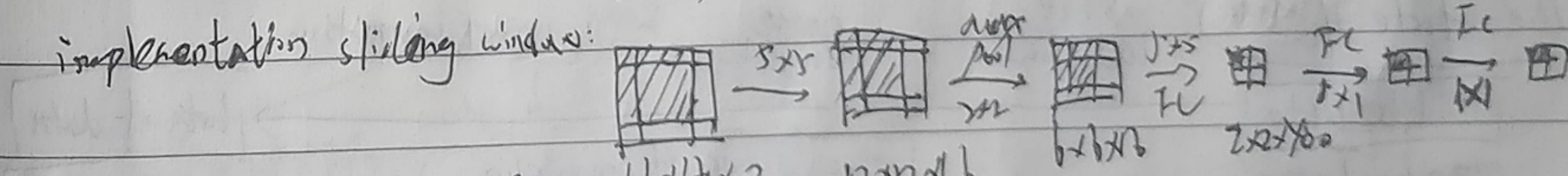
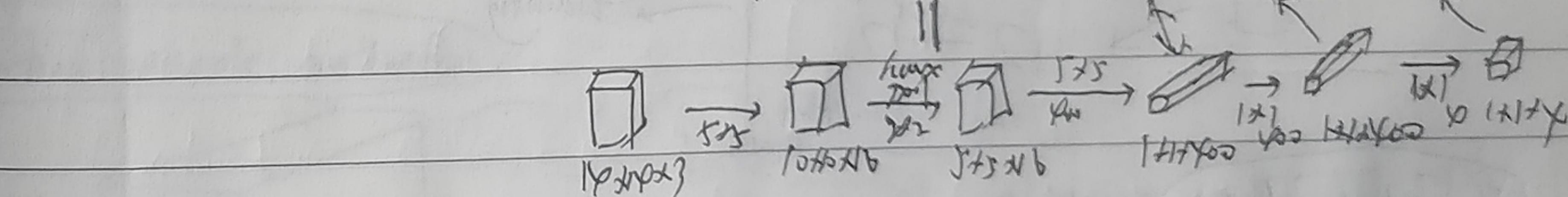
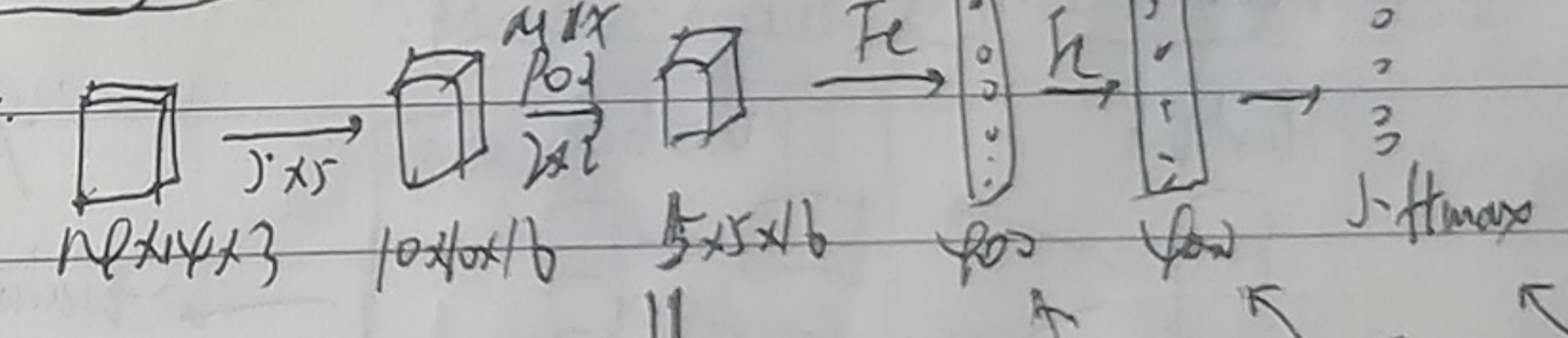
first: properly cut image (full object)

$(x, y) \rightarrow$ convert. then \rightarrow sliding windows detection.

cons:
① computation
② stride

Convolutional implementation of sliding window

Turning FC to conv layer:

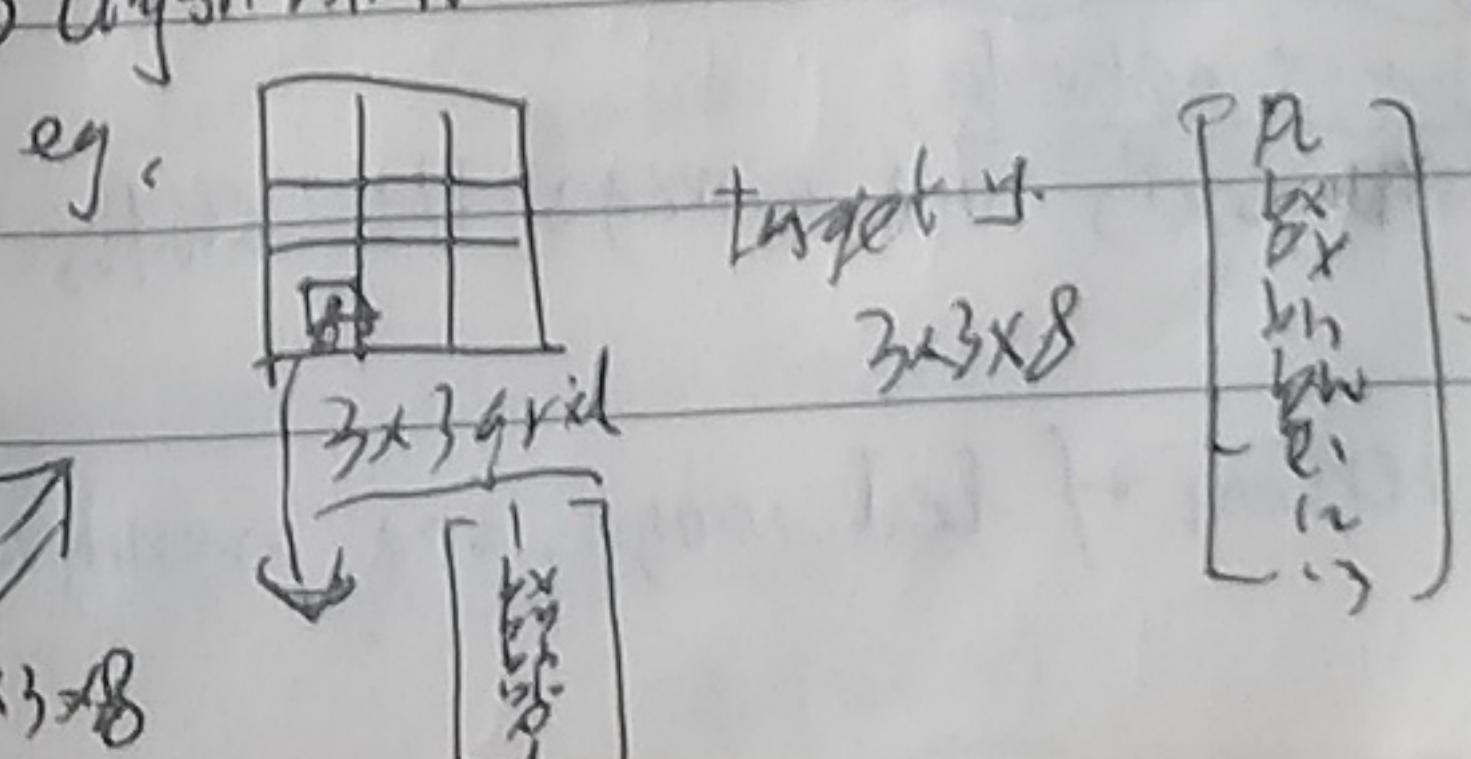


implementation sliding window:

eg: $16 \times 16 \times 3 \rightarrow 16 \times 16 \rightarrow 8 \times 8 \times 400 \rightarrow 8 \times 8 \times 400 \rightarrow 8 \times 8 \times 400$

Bounding box prediction

YOLO algorithm



Specify the bounding box

$$\begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \rightarrow y = \begin{bmatrix} b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \begin{cases} 0.8 & y \in (0, 1) \\ 0.3 & \\ 0.9 & \\ 0.5 & \} \text{ valid } > 1 \end{cases}$$

Bleu score: how well of Machine translation

(Papineni et al., 2002) - Bleu: A method for automatic

French

Error analysis on beam search.

$$\text{Case 1. } P(y^*|x) > P(\hat{y}|x)$$

Beam search is at fault.

$$\text{Case 2. } P(y^*|x) \leq P(\hat{y}|x)$$

Beam is at fault. (\hat{y} is better than y^* actually)

do list: Human Algo $P(y^*|x)$ $P(\hat{y}|x)$ at fault?

$$\begin{matrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{matrix} \quad \begin{matrix} a^{(0)} \\ a^{(1)} \\ \vdots \\ a^{(T_x)} \end{matrix} \quad \begin{matrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(T_y)} \end{matrix}$$

B: R: regularization, more data, network...

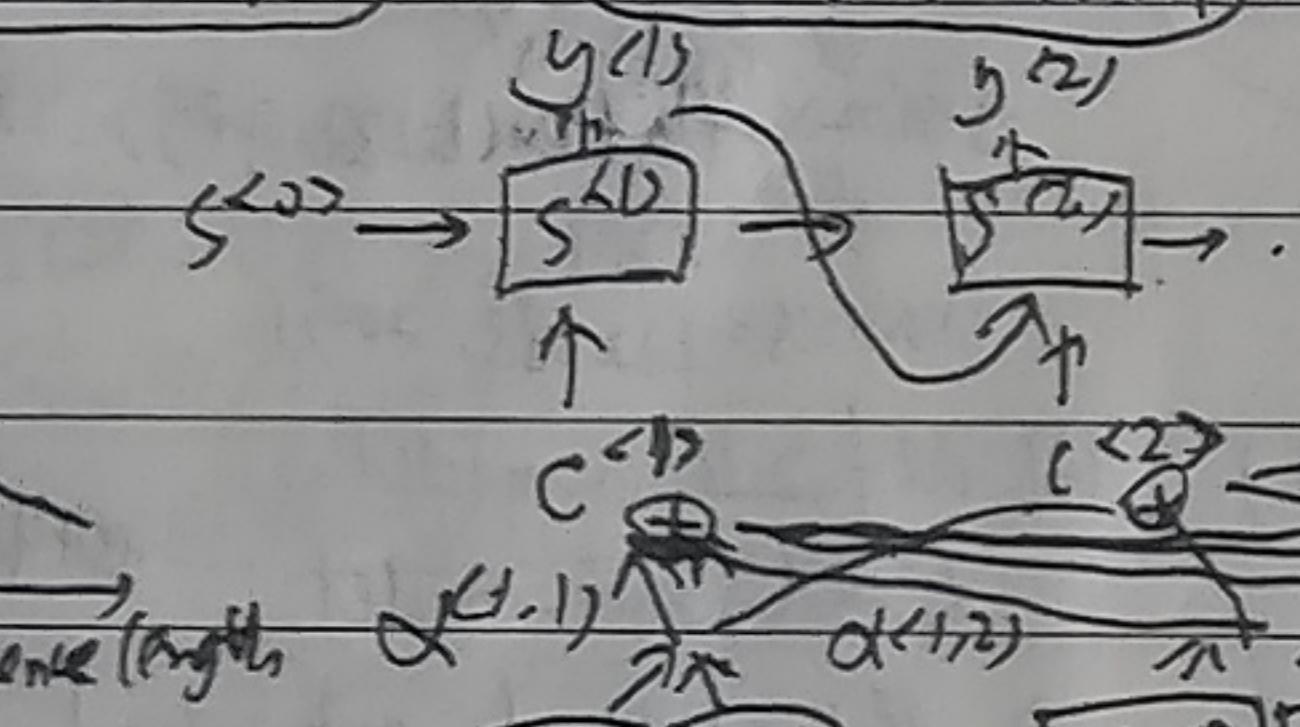
Attention model intuition

$$\alpha^{(t,t')} = \text{amount of attention } y^{(t')} \text{ should pay to } a^{(t')}$$
$$\sum_{t'} \alpha^{(t,t')} = 1, \alpha^{(t,t')} = \frac{\exp(a^{(t,t')})}{\sum_{t'} \exp(a^{(t,t')})}$$

Badrinarayanan et al., 2014. Neural

Machine translation by jointly

Xu et al., 2015 Show, attend and tell: visual image



drawback: computation

$$\alpha^{(t,t')} = \frac{\exp(a^{(t,t')})}{\sum_{t'} \exp(a^{(t,t')})}$$

Speech recognition

Graves et al., 2013 audio clip y transcript

LeCun et al., 1998 connectionist air pressure

Temporal fastfilter

tf.Variable(tf.get_variable(),
trainable=True, name='conv1')
自行标注
自行标注
自行标注

fig, axarr = plt.subplots(1, 2)
axarr[0].set_title('...')
axarr[0].imshow(x[0, :, :, 0])
axarr[1] = ...

[A2: optimization method]

① gradient descent: $W^{(t)} = W^{(t-1)} - \alpha \cdot dW^{(t)}$
 $b^{(t)} = b^{(t-1)} - \alpha \cdot db^{(t)}$

Sgd, mini-batch, ② D. 73.7%

② a. shuffle(x, y); b. Partition last batch!

③ momentum: $Vdw^{(t)} = \beta Vdw^{(t-1)} + (1-\beta) dW^{(t)}$

79.7% $W^{(t)} = W^{(t-1)} - \alpha \cdot Vdw^{(t)}$

$Vdb^{(t)} = \beta Vdb^{(t-1)} + (1-\beta) db^{(t)}$

$b^{(t)} = b^{(t-1)} - \alpha \cdot Vdb^{(t)}$

initialize $Vdw^{(0)} = np.zeros(W^{(0)})$

$Vdb^{(0)} = np.zeros(b^{(0)})$, shape = $b^{(0)}$

④ reformat initialize with 0. several iterations "building",

Beta, smoother update. $\alpha = 0.8 = 0.999, \beta = 0.9$

⑤ Adam: combines RMSprop & momentum.

94% $Vdw^{(t)} = \beta_1 Vdw^{(t-1)} + (1-\beta_1) dW^{(t)}$

t. steps of adam $Vdw^{(t)} = Vdw^{(t)} / (1-\beta_1^t)$

$Sdw^{(t)} = \beta_2 Sdw^{(t-1)} + (1-\beta_2) dW^{(t)2}$

$Sdw^{(t)} = Sdw^{(t)} / (1-\beta_2^t)$

$W^{(t)} = W^{(t-1)} - \alpha \cdot Vdw^{(t)} / \sqrt{Sdw^{(t)}} + \epsilon$

initialize $Vdw^{(0)}, Sdw^{(0)}, np.zeros, shape = W^{(0)}$

$Vdb^{(0)}, Sdb^{(0)}, np.zeros, shape = b^{(0)}$

⑥ momentum usually help but small α and simple dataset.

adam converge faster.

advantage: low memory requirements ($>$ Adam, momentum)

usually works well with little tuning of hyperparameters

(except α)

[C4: CNN]

[A1-1: CNN step by step]

① zero-pad: $np.pad(X, ((0, 0), (Pd, Pd), (Pd, Pd), (0, 0)), 'constant')$

② conv-forward: a. $Z = np.zeros(m, nH, nW, nC)$, edge...

b. $A_{prev} = zero_pad(A_{prev})$,

c. for i in range(nH):

for h in range(nH):

for w in range(nW):

for c in range(nC):

vert_start = stride * h

vert_end = vert_start + f

hor_start = stride * w

hor_end = hor_start + f

a slice prev = A_prev[:, :, :, c]

$Z[i, h, w, c] = a \cdot slice \cdot prev * W[:, :, c] + b[:, :, c]$

③ forward pass: $A = np.concatenate([Z], axis=3)$

④ backpropagation: $dW[:, :, :, c] = a \cdot slice \cdot prev * dz[:, :, c]$

$db[:, :, :, c] = dz[:, :, c] * dt[:, :, h, w, c]$

⑤ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑥ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑦ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑧ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑨ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑩ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑪ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑫ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑬ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑭ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑮ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑯ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑰ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑱ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑲ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

⑳ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉑ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉒ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉓ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉔ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉕ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉖ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉗ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉘ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉙ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉚ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉛ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉜ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

㉝ backpropagation-pool: $mask = (x == np.max(x))$

3) $da_prev_pad[h, :, :, c] += a \cdot slice \cdot dz[:, :, c] * mask$

$db[:, :, :, c] += dt[:, :, h, w, c] * mask$

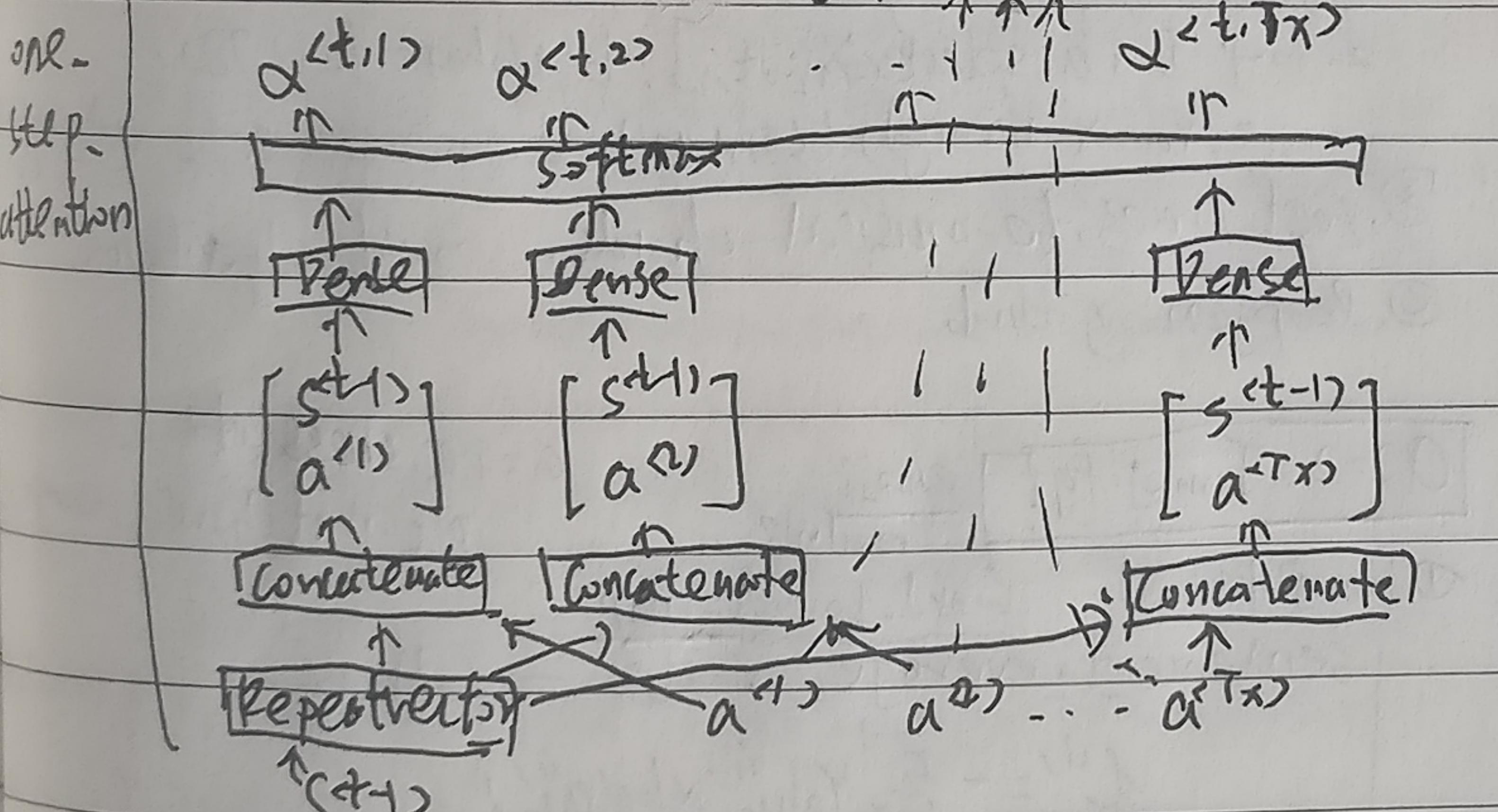
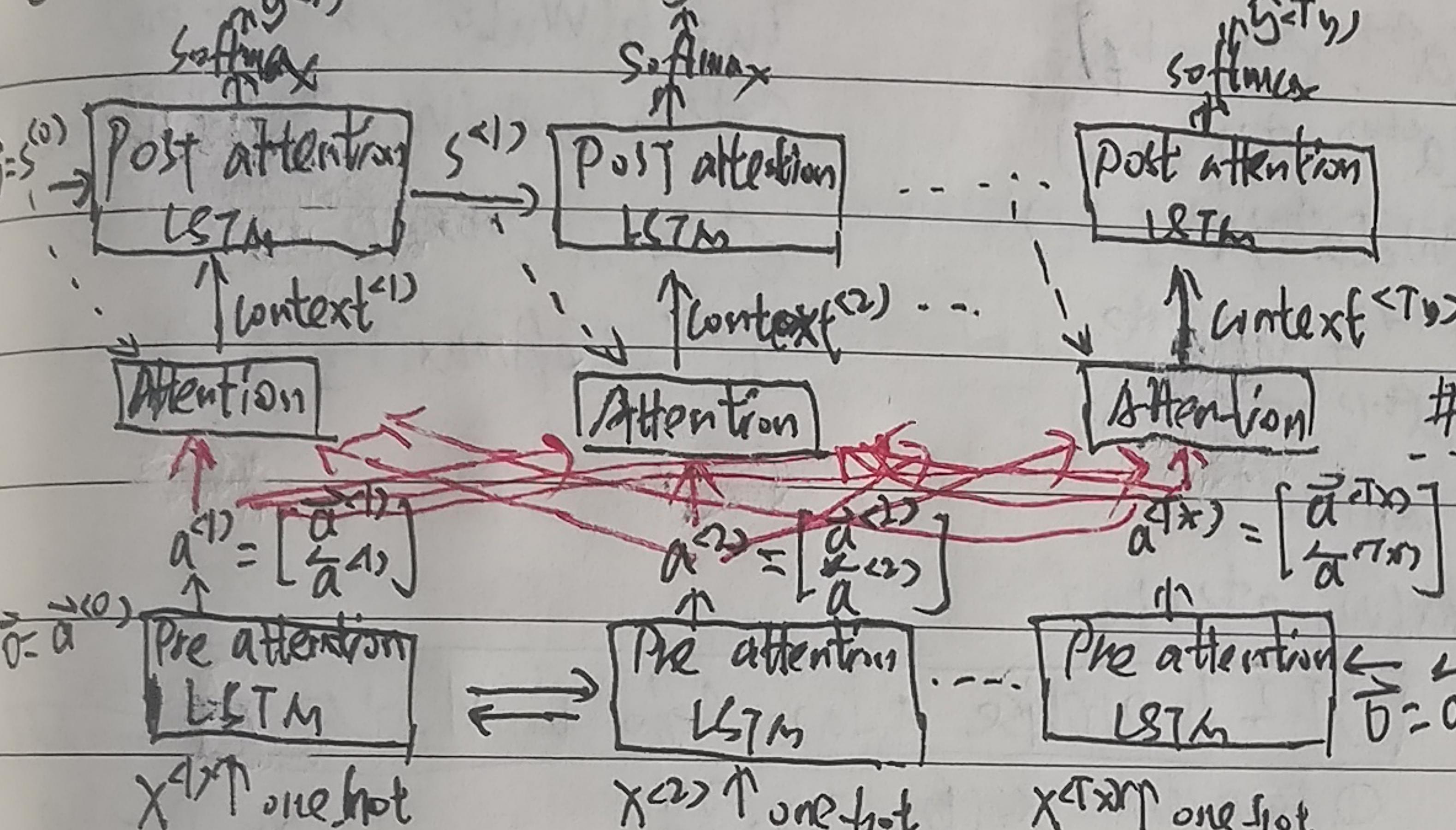
㉞ backpropagation-pool: $mask = (x == np.max(x))$

A3-1: Neural Machine Translation (NMT) date translation

① language translation require massive datasets and usually takes days of training on GPUs.

② preprocess, $X: \text{shape}((5000, 30))$; $Y_{\text{oh}}: \text{shape}((5000, 30, 37))$
 $m: \text{len}(\text{human_vocab})$
 $T_x: \text{len}(\text{human_vocab})$
 $T_y: \text{len}(\text{machine_vocab})$
 $\text{dict}: \{1, 2, 3, \dots, n\}$

③ attention mechanism.



④ model:

1. Propagate input to a Bidirectional LSTM.
2. Loop $t = 0 \dots T_y - 1$:
 - a. call one_step_attention on a and $s^{(t-1)}$, get $\text{context}^{(t)}$
 - b. give $\text{context}^{(t)}$ to post LSTM, with s, c .
 - c. apply softmax to $s^{(t)}$, get $\hat{s}^{(t)}$
 - d. keras: three inputs $[a, s, \hat{s}^{(t)}]$.

remember:

1. attention. $T_x \neq T_y$ can be different.
2. visualized $a^{(t,t)}$. see attention.

pydub
module

A3-2: Trigger word detection

① data synthesis. audio recording: long list of number measuring air pressure change

Spectrogram. $f \uparrow$ time $105 \rightarrow T_x = 5511$

10000: used by pydub to synthesize audio. $T_y = 1357$

$10/10000 = 1\text{ms}$

background. (10s)

negative

positive.

→ advantages: ①, easy to record lots of negative, positive/boundary

② convenient to generate labels.

③ Model. $\hat{y}^{(t)}$

sigmoid

Dense

↓

dropout 0.8

BN

dropout 0.8

GRU(128) →

↑

dropout 0.8

GRU(128) →

↑