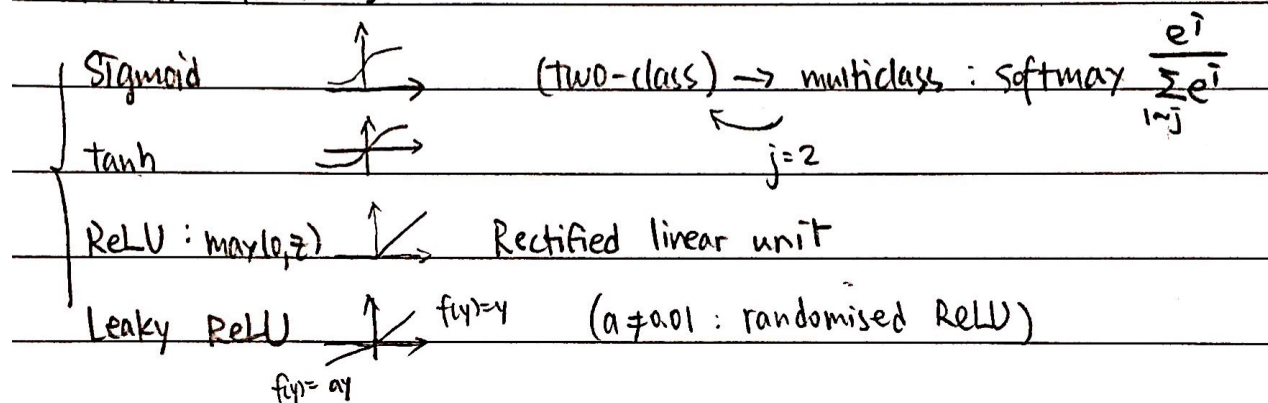


Linear regression \rightarrow Logistic regression

$$\sigma(z) = \frac{1}{1+e^{-z}} \quad \text{logistic activation (sigmoid)}$$

Activation functions



Gradient descent

Loss function $\mathcal{L}(\hat{y}^{(i)}, y^{(i)})$ e.g. logistic loss $= -[\hat{y} \log \hat{y} + (1-\hat{y}) \log (1-\hat{y})]$

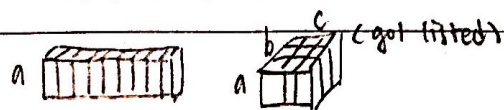
Cost function $J(w, b) = \frac{1}{m} \sum \mathcal{L}$

Python

rank 1 array : column/row vector

flatten - $x = x.reshape(x.shape[0], -1).T$

$a \times b \times c \leftarrow a \times b \times c \times d$



Max likelihood estimation

	High bias (underfit)	→ (✓) →	High variance (overfit)
model complexity	low		high
training error	high	reasonably low	low
validation error	(high)	reasonably low	high
solutions	bigger network longer training		regularization more data

Gradient checking

Regularization ① $J(w, b) = \frac{1}{m} \sum \mathcal{L} + \lambda R(w_i)$

l_1 regularization $\frac{\lambda}{2m} \|w\|_1$ $\|w\|_1 = \sum w_j$

l_2 regularization $\frac{\lambda}{2m} \|w\|_2^2$ $\|w\|_2^2 = \sum w_j^2 = w^T w$

→ Direct effect: smaller/less $w_j \rightarrow$ less complex ← thus

ie, as $\lambda \uparrow$ $w \downarrow$ $z = w \cdot a + b \downarrow$ with activation function, every layer \approx linear

② Dropout keep-prob (= 1 is no dropout)

inverted dropout { train time : reassign probability }
test time : no change

mask matrix $=(np.random.rand(P, shape) < p)/p$

③ Early stopping

Weight Initialization of W (partially helps vanishing gradient problem)

$W^{[L]} = np.random.rand(shape) * \sqrt{\frac{2}{n^{[L-1]} + n^{[L]}}}$ (Xavier init)

when give $var(w_i) = \frac{2}{n}$

Optimization

Mini-batch Gradient Descent

$$X = \underbrace{[x^{(1)} x^{(2)} \dots x^{(n)}]}_{x^{(1)}} \dots \underbrace{[x^{(1)} x^{(2)} \dots x^{(n)}]}_{x^{(n)}}$$

$$\text{total} = a \text{ (epoch)} \times n \text{ (mini-batch)}$$

$$\hookrightarrow m \leq 2000 : \text{batch GD} \quad \hookrightarrow \text{usually } 2^k$$

size $a = \text{all}$

$\alpha = 1$: stochastic (one by one)

$$\text{Exponentially weighted average } \bar{x}_n = \beta \bar{x}_{n-1} + (1-\beta) \theta_n$$

\downarrow old average \downarrow new input

for approx over $\frac{1}{1-\beta}$ day

\rightarrow bias correction, for warm-up

① GD with momentum

$$\begin{cases} v_{dw} = \beta_1 v_{dw} + (1-\beta_1) dW \\ v_{db} = \beta_1 v_{db} + (1-\beta_1) db \end{cases} \rightarrow \text{usually } 0.9$$

intuitive:
controlling learning rate
 \downarrow vertically
 \uparrow horizontally

$$\text{then update } \begin{cases} w = w - \alpha v_{dw} \\ b = b - \alpha v_{db} \end{cases} \rightarrow \text{to be tuned}$$

② RMSprop (root mean square)

$$\begin{cases} s_{dw} = \beta_2 s_{dw} + (1-\beta_2) (dW)^2 \\ s_{db} = \beta_2 s_{db} + (1-\beta_2) (db)^2 \end{cases} \rightarrow \text{usually } 0.999$$

$$\text{then update } \begin{cases} w = w - \alpha \frac{dw}{\sqrt{s_{dw}} + \epsilon} \\ b = b - \alpha \frac{db}{\sqrt{s_{db}} + \epsilon} \end{cases}$$

ADAM optimization (adaptive moment estimation)

$$\textcircled{1} v_{dw} = s_{dw} = v_{db} = s_{db} = 0 \text{ (init)}$$

$$\textcircled{2} \text{ mini-batch iteration: } \begin{aligned} v_{dw}^{\text{corrected}} &= \frac{v_{dw}}{1-\beta_1^t} & v_{db}^{\text{corrected}} &= \frac{v_{db}}{1-\beta_1^t} \\ s_{dw}^{\text{corrected}} &= \frac{s_{dw}}{1-\beta_2^t} & s_{db}^{\text{corrected}} &= \frac{s_{db}}{1-\beta_2^t} \end{aligned}$$

$$w = w - \alpha \frac{v_{dw}^{\text{corrected}}}{\sqrt{s_{dw}^{\text{corrected}} + \epsilon}} \quad b = b - \alpha \frac{v_{db}^{\text{corrected}}}{\sqrt{s_{db}^{\text{corrected}} + \epsilon}}$$

usually 10^{-8}
(to avoid 0 denominator)

Learning rate decay

$$\alpha = \begin{cases} 1 + \text{decayrate} \times \text{epoch num} & \alpha_0 \\ 0.95^{\text{epoch num}} & \alpha_0 \quad (\text{exponentially}) \\ \frac{k}{\sqrt{\text{e.n.}}} & \alpha_0 \\ \frac{k}{\sqrt{t}} & \alpha_0 \\ \text{stair case} & \uparrow \end{cases}$$

or manually

Saddle point, plateau (rather than local optima)

Hyperparameter tuning

Panda vs Caviar (many in parallel) model

Batch normalization $\tilde{z}_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$ for non-input layer(s) $\left\{ \begin{array}{l} \text{pro: covariate shift} \\ \text{(not quite matching train} \leftrightarrow \text{test)} \end{array} \right.$

* but don't want norm distribution at all time $\tilde{z}^{(i)} = \gamma \tilde{z}_{\text{norm}}^{(i)} + \beta$ $\left\{ \begin{array}{l} \text{con: like dropout, will} \\ \text{introduce noise} \end{array} \right.$

and at test time:

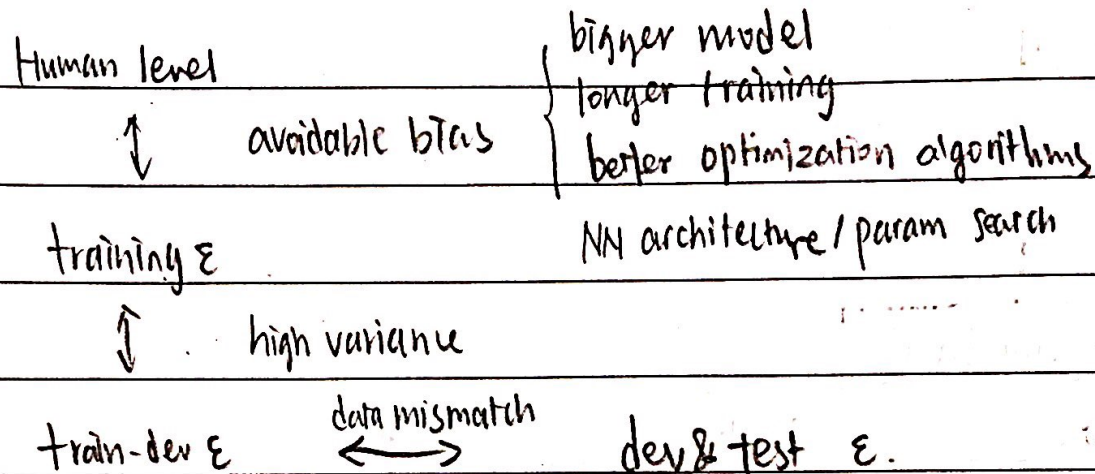
use μ, σ^2 from exponentially weighted average
(across mini-batches), i.e. "running average"

→ Week 3 intro to TensorFlow

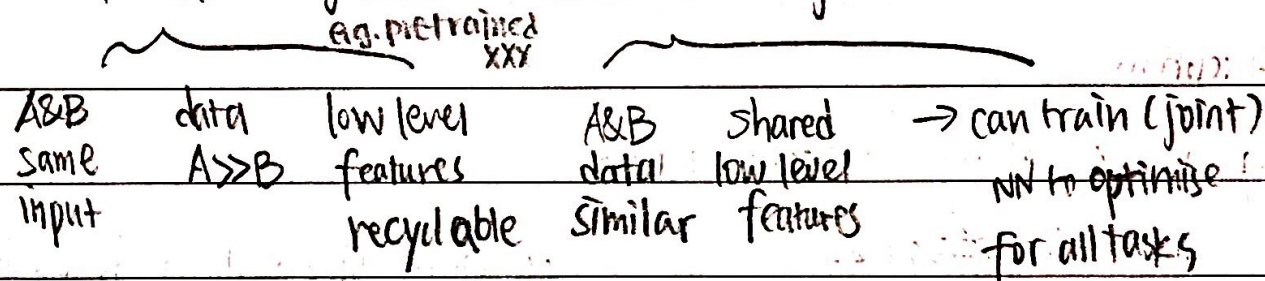
Orthogonalization : one variable tuned at a time

Evaluation metric

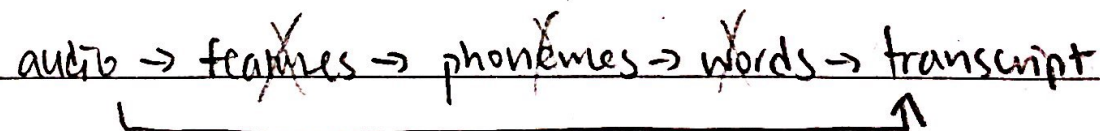
Performance



Transfer learning vs. multi-task learning



End-to-end: no intermediate nodes / constraints



CNN

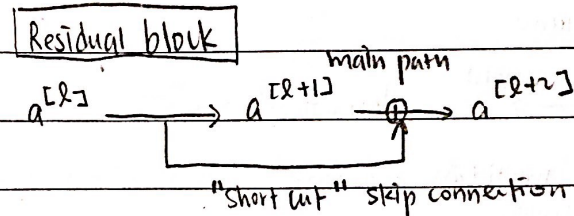
Case studies

LeNet-5

AlexNet

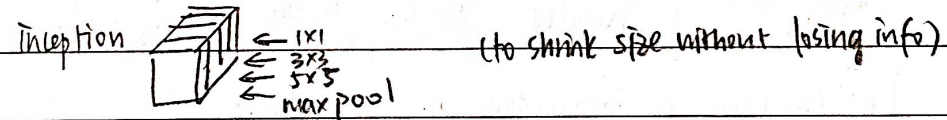
VGG

ResNet



Inception

1×1 Convolution : $n \times n$ flatten



Data Augmentation

- Mirroring
- Random cropping
- Rotation
- Shearing \rightarrow
- Local warping

Colour shifting e.g. PCA (principles component analysis)

(Computer vision, neural style transfer)

Object Localization

Landmark Detection

Sliding window

Bounding box - YOLO = you only look once

$$\text{intersection over union} = \frac{\text{intersection}}{\text{union}} \geq 0.5 - \text{correct}$$

non-max suppression

anchor boxes for each cell

region proposals : segmentation algorithm

- only run sliding windows on interesting blobs/regions

Face verification vs. recognition

1:1

1:K

much more difficult

One shot learning : input = 1 image

Siamese network

Triplet loss



Anchor - Positive

→ small

Negative

→ large

$$d(A, p) - d(A, N) + \alpha \leq 0$$

$$\|f(A) - f(p)\|^2 - \|f(A) - f(N)\|^2 + \text{margin param}$$

$$L = [\dots] \quad 1, 0, J_{\max}$$

Neural Style Transfer

content (C) + style (S) → generated picture (G)

cost function

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

content cost function

style cost function

difference in

k

~~activation~~ function

channel-channel similarity

$$\|a_i - a_j\|^2$$

$n_c \times n_c$