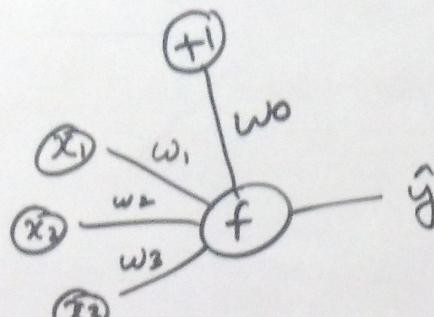
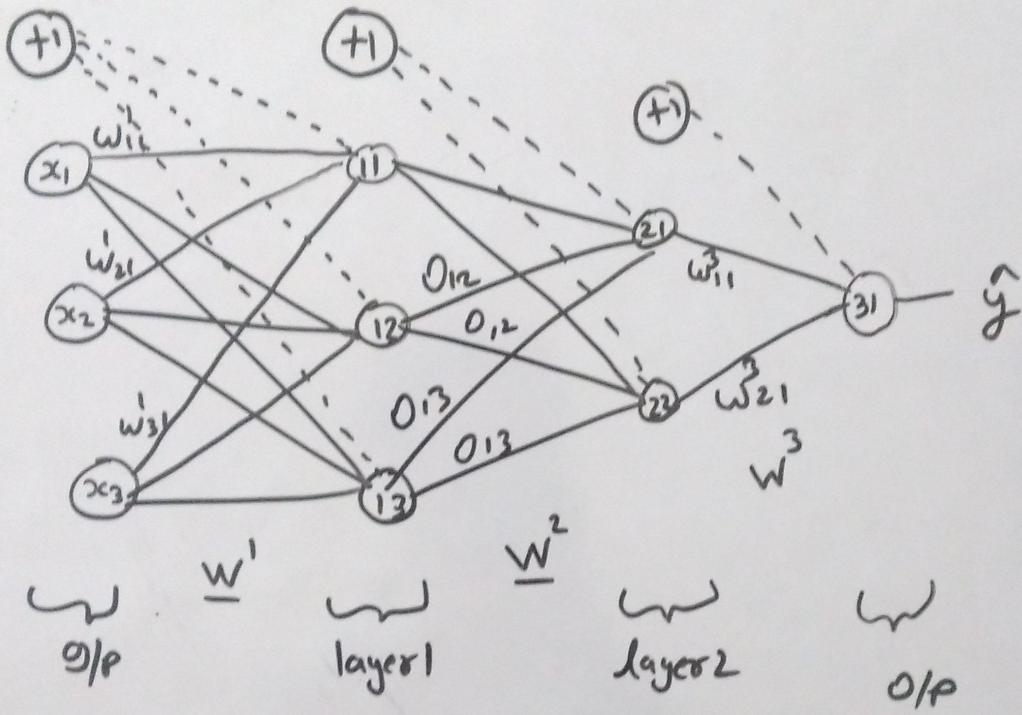


DL



Logistic regression unit
(single neuron)

$$\hat{y} = f(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3)$$



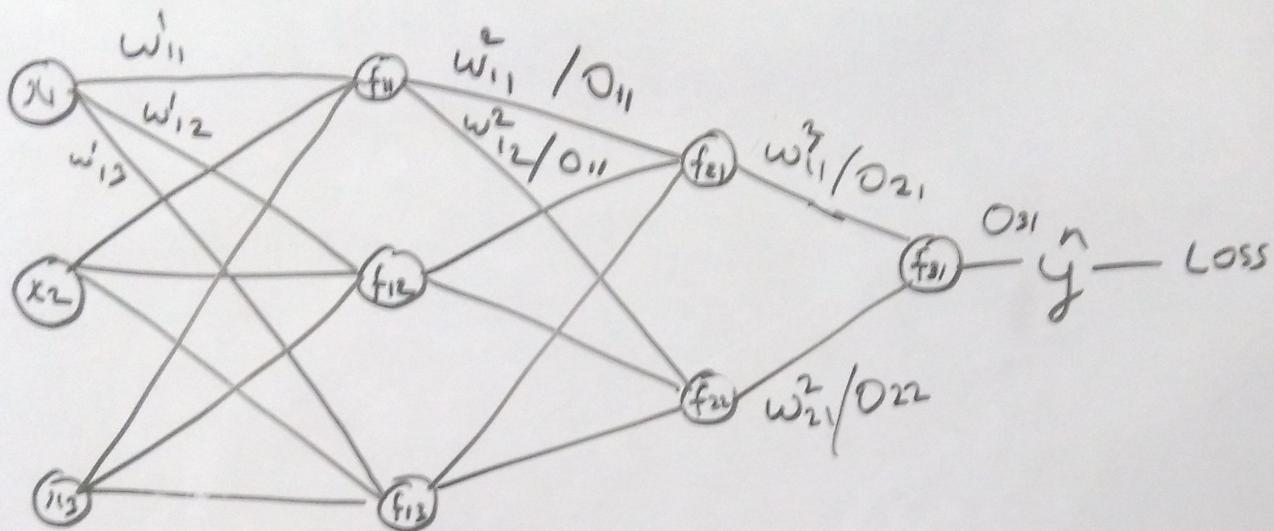
$$W^1 = \begin{bmatrix} w_{01} & w_{02} & w_{03} \\ w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \rightarrow \text{Bias} \quad (3+1) \times 3$$

$$W^2 = \begin{bmatrix} w_{01} & w_{02} \\ w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} \downarrow \begin{matrix} f_{10} \\ f_{11} \\ f_{12} \\ f_{13} \end{matrix} \quad \downarrow (1+3) \times 2$$

$$W^3 = \begin{bmatrix} w_{01} \\ w_{11} \\ w_{22} \end{bmatrix} \downarrow \begin{matrix} f_{20} \\ f_{21} \\ f_{22} \end{matrix} \quad (1+2) \times 1$$

DL

Training of MLP

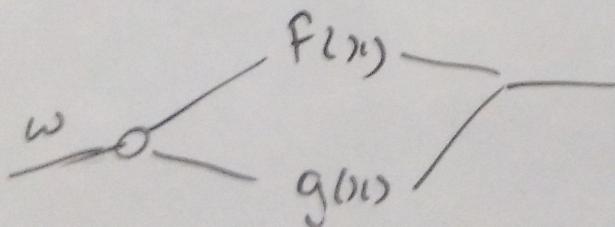


$$\text{here } O_{31} = \hat{y}$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{21}} \cdot \frac{\partial O_{21}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}} + \frac{\partial L}{\partial O_{31}} \cdot \frac{\partial O_{31}}{\partial O_{22}} \cdot \frac{\partial O_{22}}{\partial O_{11}} \cdot \frac{\partial O_{11}}{\partial w_{11}}$$

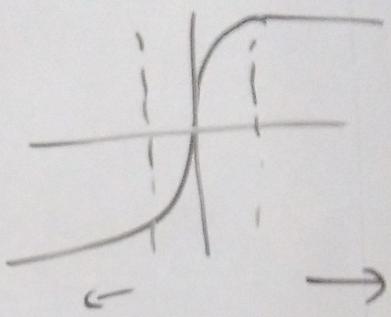


DL

Training of MLP

Problem of Sigmoid in Backpropagation

Max slope of Sigmoid = 0.5



$$\frac{\partial \sigma}{\partial w} = 0.5 \times 0.1 \times 0.2 \times \dots$$

- lead to very very small update to initial layer. It is known as vanishing gradient problem.

Exploding gradient problem:

Gradient become very-very large. hence lead to overshooting and failed to converge.

$$\frac{\partial L}{\partial w} = 10 \times 10 \times 10 \dots$$

These occur due to chain rule in Backpropagation and our activation functions structure.

30

Deep Learning

Activation Unit

$$\textcircled{1} \text{ Sigmoid: } \sigma(x) = \frac{1}{1+e^{-x}} ; \quad \sigma'(x) = \sigma(x) \cdot (1-\sigma(x))$$

[0,1]

$$\textcircled{2} \text{ tanh: } f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} ; \quad f'(x) = 1 - f(x)^2$$

$$\textcircled{3} \text{ Relu: } f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} ; \quad f'(x) = \begin{cases} 1 & ; x \geq 0 \\ 0 & ; x < 0 \end{cases}$$

$$\textcircled{4} \text{ Leaky Relu: } f(x) = \begin{cases} x & ; x \geq 0 \\ 0.01x & ; x < 0 \end{cases}$$

$$\textcircled{5} \text{ Softplus: } f(x) = \log(1+e^x) ; \quad f'(x) = \frac{1}{1+e^{-x}}$$

30

Deep Learning

① Uniform Initialization

$$w_{ij}^k \sim \text{Unif}\left(-\frac{1}{\sqrt{\text{fanin}}}, \frac{1}{\sqrt{\text{fanin}}}\right)$$

② Xavier (Glorot Initialization) (Sigmoid)

- Uniform $w_{ij}^k \sim \text{Unif}\left(-\frac{\sqrt{6}}{\sqrt{\text{fanin}+\text{fanout}}}, \frac{\sqrt{6}}{\sqrt{\text{fanin}+\text{fanout}}}\right)$

- Normal

$$w_{ij}^k \sim \text{Normal}(0, \alpha) ; \alpha = \frac{2}{\text{fanin} + \text{fanout}}$$

③ He Initialization (ReLU)

- Uniform $w_{ij}^k \sim \text{Unif}\left[-\sqrt{\frac{6}{\text{fanin}}}, \sqrt{\frac{6}{\text{fanin}}}\right]$

- Normal

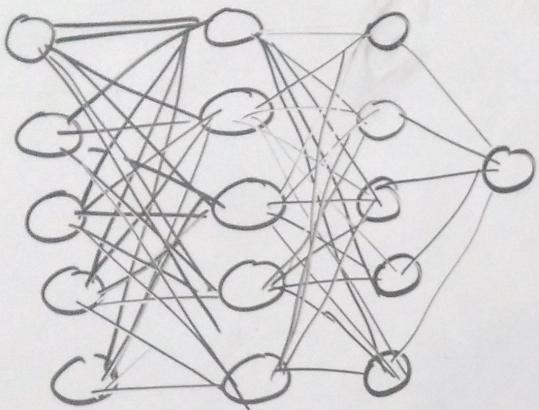
$$w_{ij}^k \sim N(0, \alpha) ; \alpha = \sqrt{\frac{2}{\text{fanin}}}$$

Deep Learning

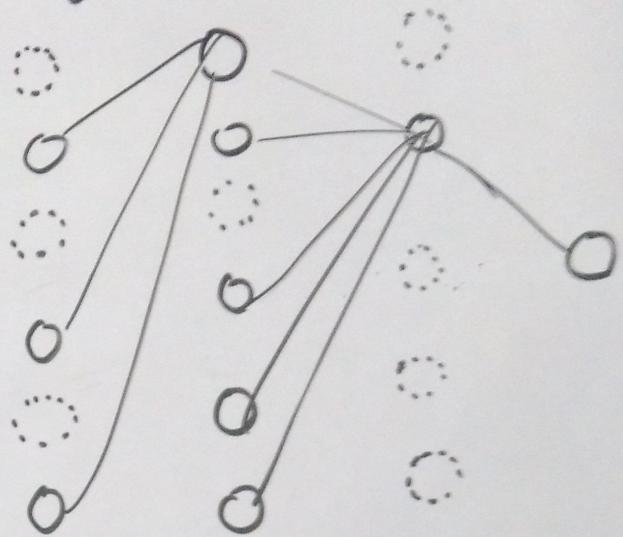
Dropout

(Randomization for Regularization)

Dropout
probability $P=0$

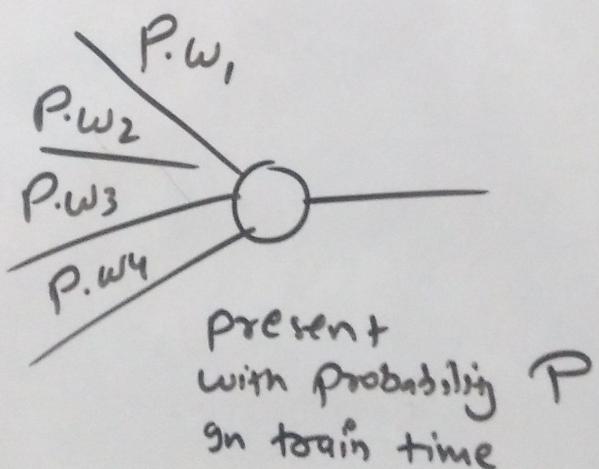


$$=0 =0 =0.$$



$$0.5 \quad 0.2 \quad 0.8$$

Test
Time



present
with probability P
in train time

Deep Learning

Batch Normalization

- Protect from Bad Initialization
- Solve Problem of Internal Covariance shift.

B₁:

B₂: $(x) \rightarrow L_1 \rightarrow L_2 \rightarrow L_3 \rightarrow L_4 \rightarrow L_5 \rightarrow (O/r)$

B₃:

for each Batch (Parameter to learnt):

$$\mu_B = \frac{1}{m} \sum_{i=1}^B x_i$$

$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^B (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \therefore \epsilon: 0.0001 \\ (\text{Protect from div by zero})$$

$\text{BN}_{\gamma, \beta}(x_i): \boxed{\tilde{x}_i = \gamma \cdot \hat{x}_i + \beta}$ Batch Normalized data.

Test time

μ_B : running avg of μ_B at train time.

σ_B^2 : Running avg. of σ_B^2 at train time

β, γ : learn by Optimization.

Deep Learning - Optimizers

- Simple Model (SUM, Log. Reg.) : Convex loss
- Deep Learning Model : Non-convex losses / Saddle points

SGD with Momentum:

Idea: exponential weighted Average.

$$v_t = \alpha_t + \gamma \cdot \alpha_{t-1} + \gamma^2 \cdot \alpha_{t-2} + \gamma^3 \cdot \alpha_{t-3}$$

\rightarrow Reduced impaced.
if $\gamma < 1$

$$\begin{aligned} v_1 &= \eta \cdot g_t \quad // \text{initial} \\ v_t &= \gamma \cdot v_{t-1} + \eta \cdot g_t \\ w_t &= w_{t-1} - v_t \end{aligned}$$

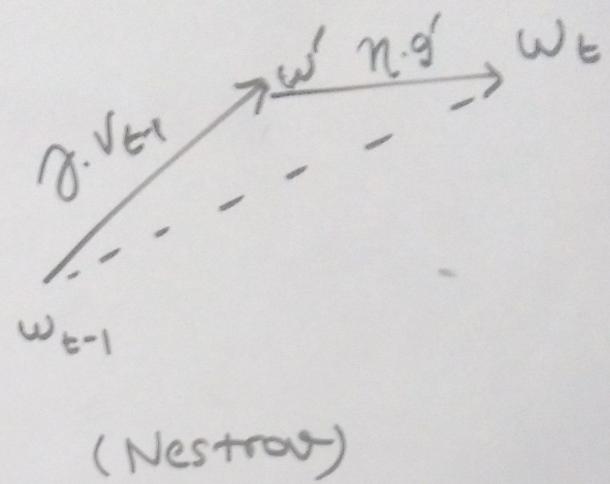
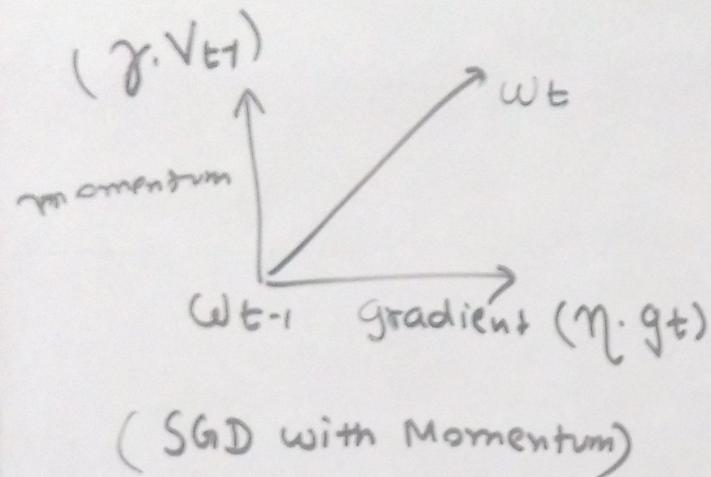
here $\eta = \text{learning rate}$
 $g_t = \text{gradient update}$
to make w_t .

$$\gamma = 0.9 \quad // \text{generally.}$$

$$w_t = w_{t-1} - (\eta \cdot g_t + \gamma \eta \cdot g_{t-1} + \gamma^2 \eta \cdot g_{t-2} \dots)$$

Deep Learning - Optimizers

Nestrov Accelerated Gradient (NAG)



- first go with Momentum then calculate gradient then get w_t .

$$\omega' = w_{t-1} - \gamma \cdot v_{t-1}$$

$$g' = \left(\frac{\partial L}{\partial w} \right)_{w_{t-1} - \gamma \cdot v_{t-1}}$$

$$w_t = \omega' - \eta \cdot g'$$

or

$$w_t = w_{t-1} - (\gamma \cdot v_{t-1} + \eta \cdot g')$$

Deep Learning - Optimizers

Adagrad (adaptive gradient)

each weight need to have its own adaptive learning rate.

$$w_t = w_{t-1} - \eta \cdot g_t$$

$$\eta = \frac{c}{\sqrt{\alpha_{t-1} + \epsilon}} \quad \therefore \text{here } c = 0.01$$

$$\alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2$$

as iteration increases learning rate automatically decrease.

Adadelta (Similar RMSprop)

for long iteration η is very low hard to converge

all same as Adagrad except α_{t-1}

$$\alpha_{t-1} = \sum_{i=1}^{t-1} g_i^2 \rightarrow \boxed{(1-\gamma) \cdot \hat{g}_{t-1}^2 + \gamma \cdot \alpha_{t-2}}$$

here $\gamma = 0.95$ // generally.

Deep Learning - Optimizers

Adam (Adaptive Moment Estimation)

mean - 1st Moment

Variance - 2nd Moment

$$\| m_0 \leftarrow 0$$

$$\| v_0 \leftarrow 0$$

$$m_t = \beta_1 \cdot m_{t-1} + (1-\beta_1) g_t$$

$$v_t = \beta_2 \cdot v_{t-1} + (1-\beta_2) g_t^2$$

// generally $\beta_1 = 0.9$; $\beta_2 = 0.99$

$$\hat{m}_t = \frac{m_t}{1-(\beta_1)^t} ; \hat{v}_t = \frac{v_t}{1-(\beta_2)^t}$$

$$w_t = w_{t-1} - \eta \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

Deep Learning - Others

Gradient Clipping

$$W = \boxed{2000 \ 1 \ 1 \ 1 \ 1}$$

(mainly monitor initial layer for gradient vanishing and exploding problem.)

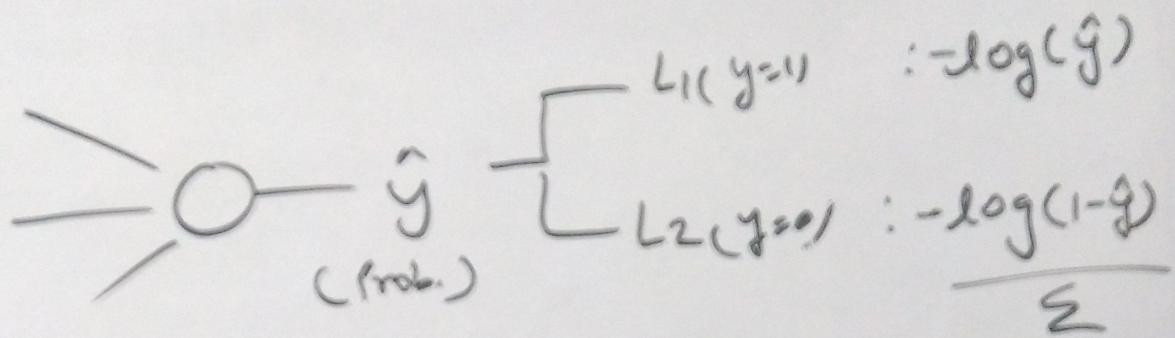
$$G_{\text{new}} = \frac{G}{\|G\|} \cdot Z$$

$$\therefore \|G\| = \sqrt{G_1^2 + G_2^2}$$

- It makes value b/w 0 and 1. So resolve problem of vanishing gradient.

Deep Learning - Others

Logistic Regression [Sigmoid + log loss]



log loss:

$$\frac{1}{m} \sum [-y \log(\hat{y}) - (1-y) \log(1-\hat{y})]$$

Multiclass

(2)	\circ	\hat{c}_1	$L: -\log(\hat{c}_1)$
(2)	\circ	\hat{c}_2	$L: -\log(\hat{c}_2)$
(2)	\circ	\hat{c}_3	$L: -\log(\hat{c}_3)$
		<u>①</u>	<u>Σ</u>

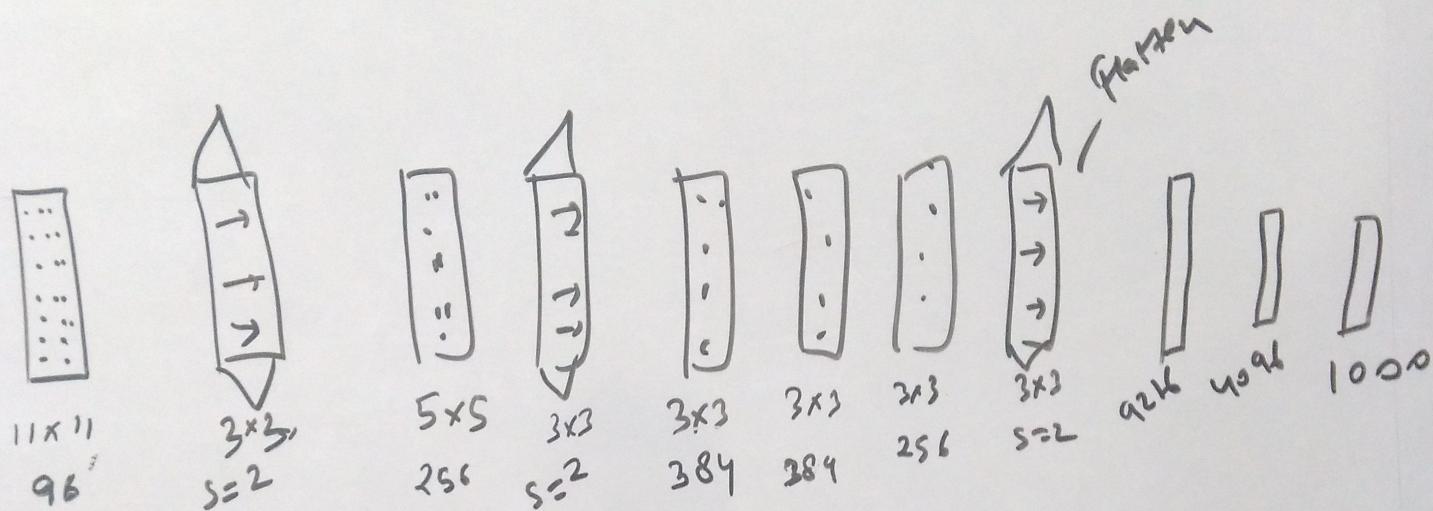
Softmax
(Normalized exponent)

$$= \frac{e^{z_1}}{\sum_{i=1}^k e^{z_i}}$$

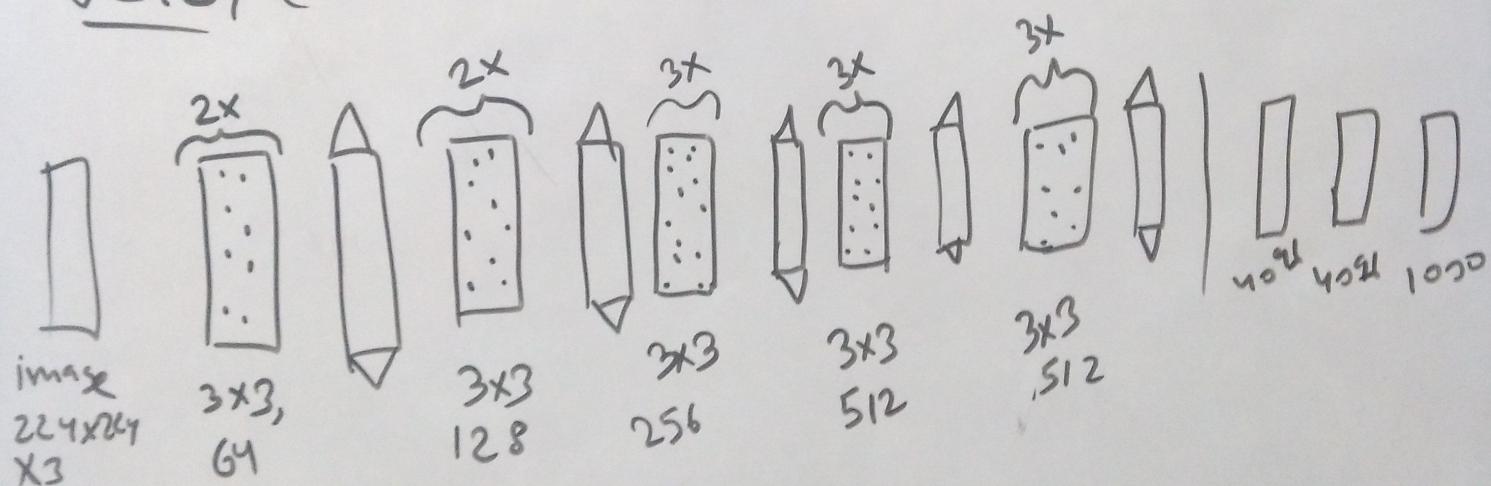
use
cross
entropy
loss

Deep Learning - Others

AlexNet (2012)

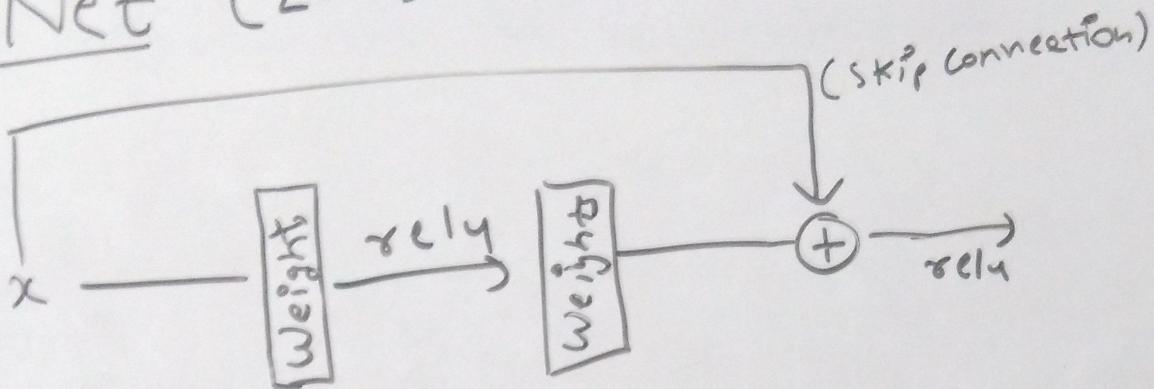


VGG (2014/2015)



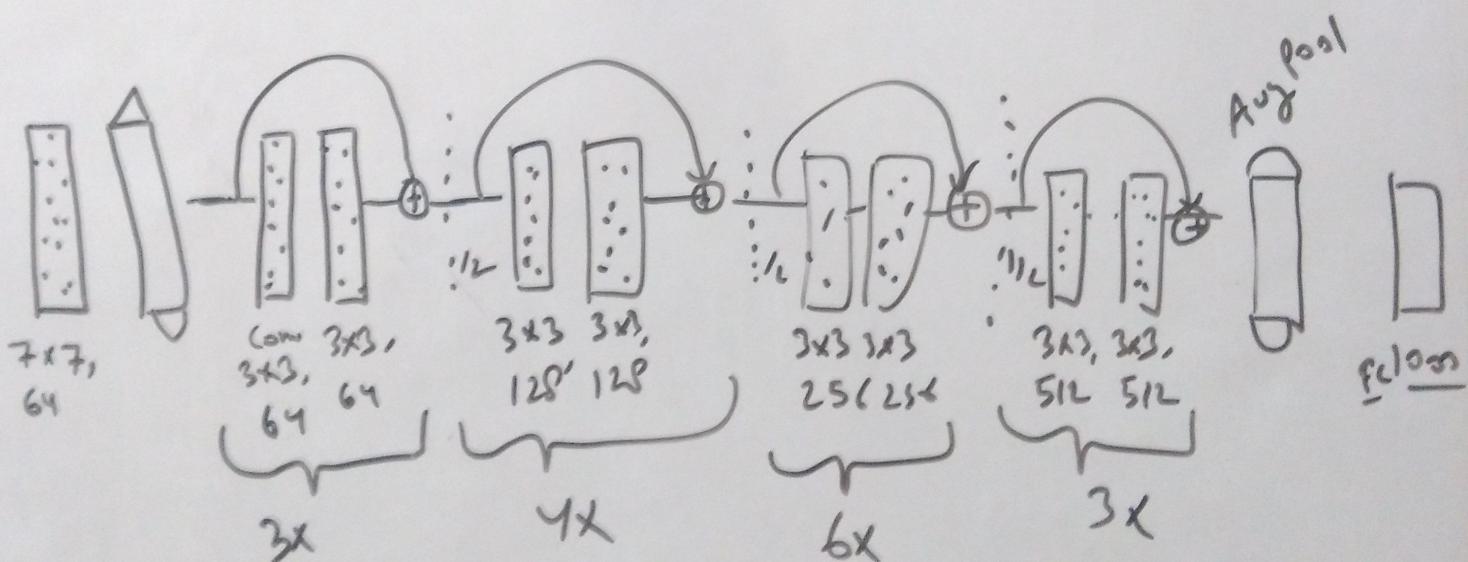
Deep Learning - Others

ResNet (2015)



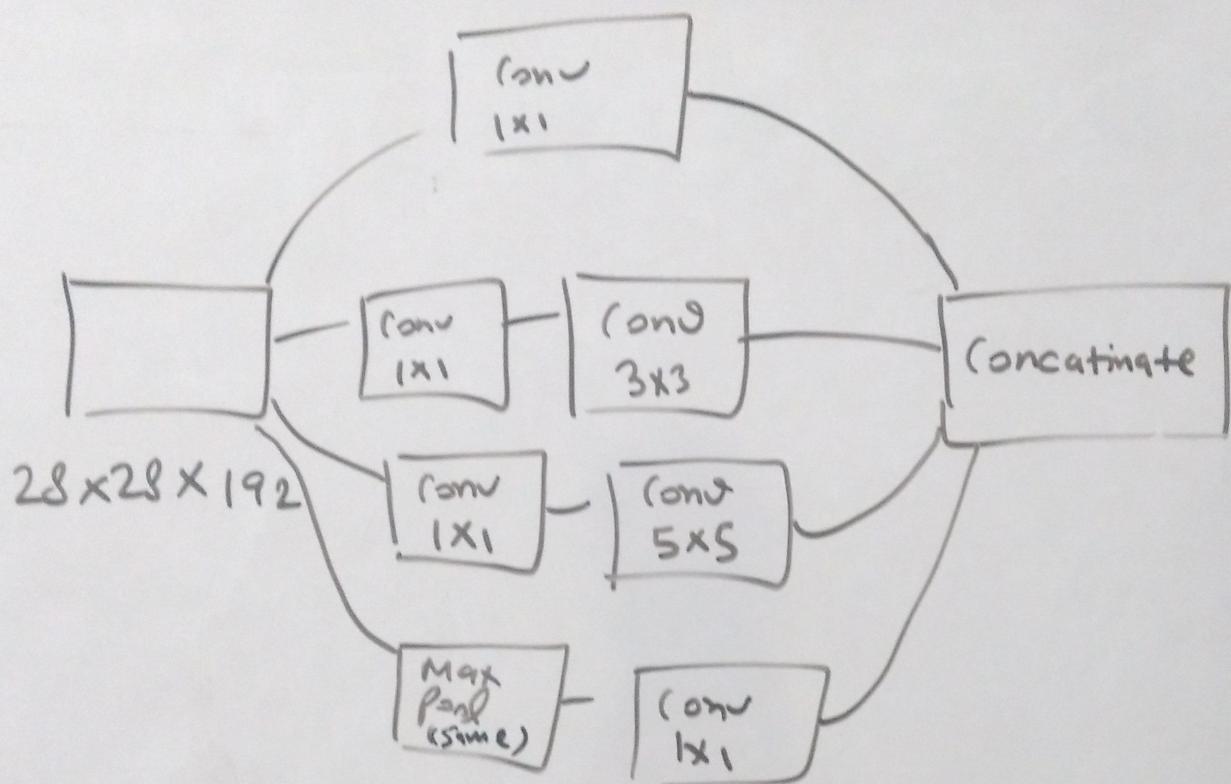
- Deeper Network lead to Reduce performance Solved this.

34 layer Residual



Deep Learning - Others

Inception (GoogLeNet) (2015)

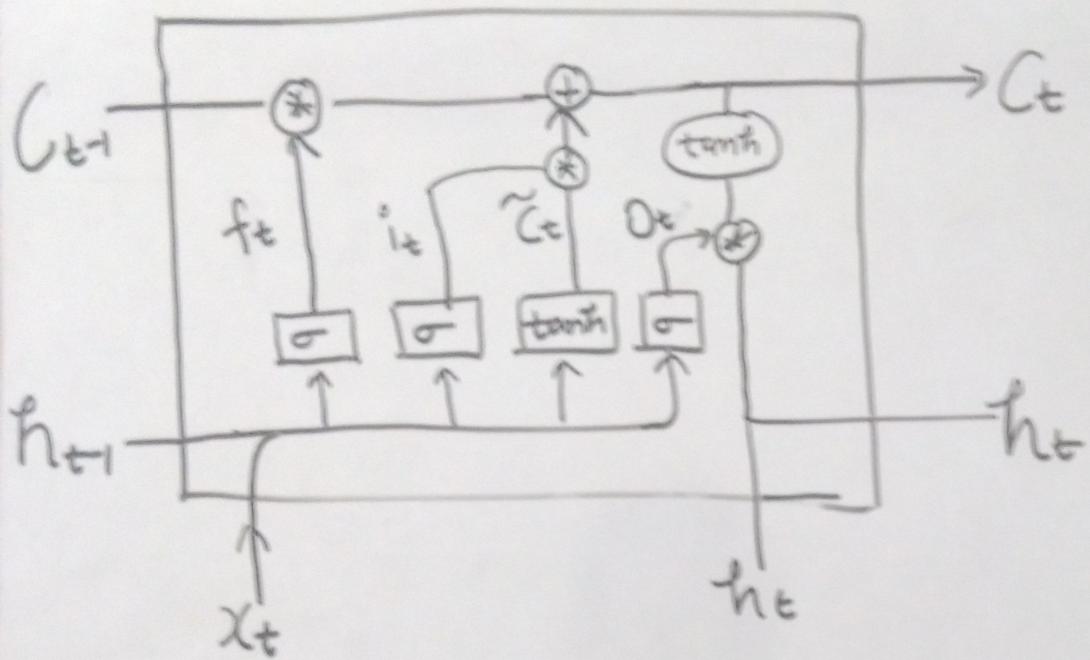


Inception v2

(googlenet + ResNet)

Deep Learning - Others

LSTM (Long short term memory)



$$f_t = \sigma(W_f[x_t, h_{t-1}] + b_f)$$

$$i_t = \sigma(W_i[x_t, h_{t-1}] + b_i)$$

$$\tilde{C}_t = \tanh(W_c[x_t, h_{t-1}] + b_c)$$

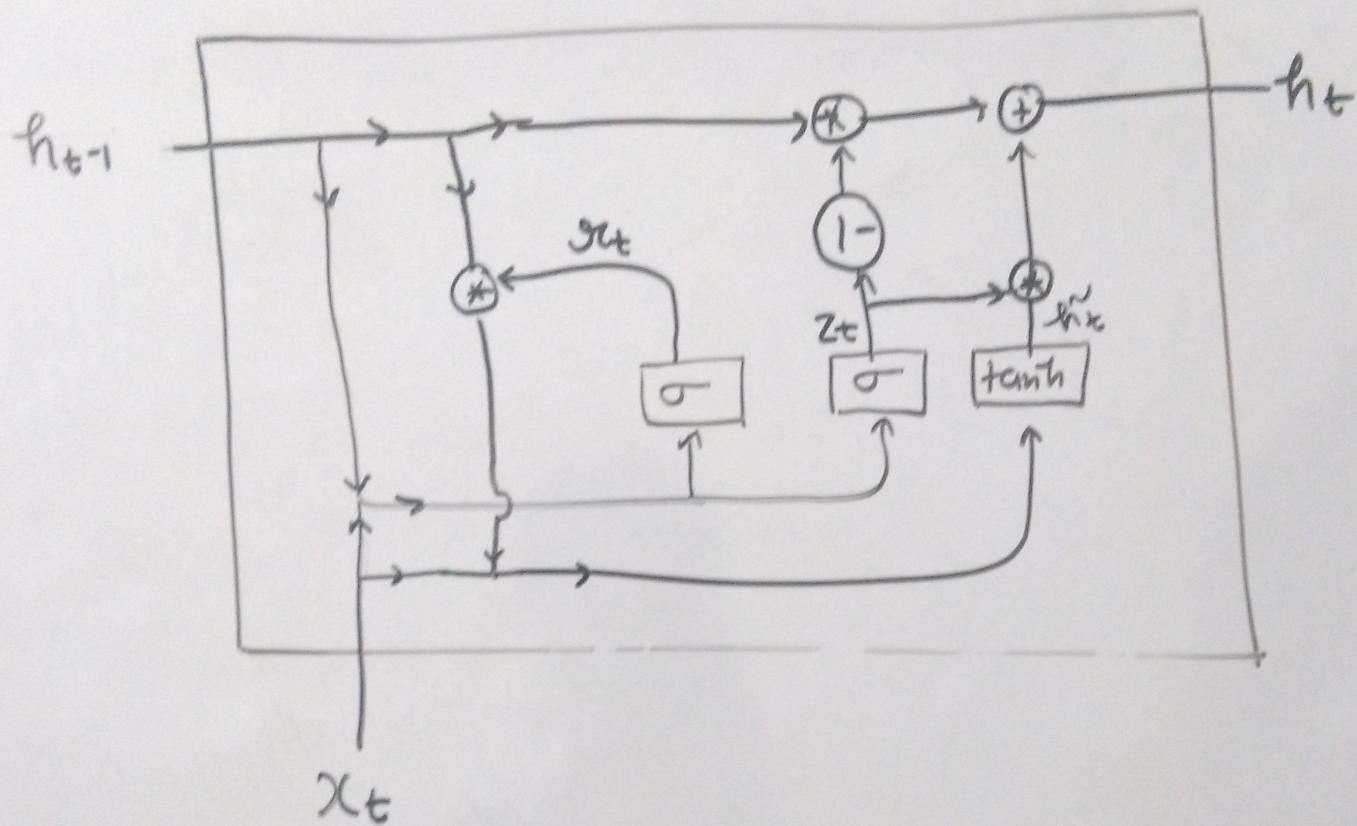
$$O_t = \sigma(W_o[x_t, h_{t-1}] + b_o)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$h_t = O_t * \tanh(C_t)$$

Deep Learning - Others

GRU (Gated Recurrent Unit)



$$g_t = \sigma(W_g \cdot [h_{t-1}, x_t] + b_g)$$

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

$$\tilde{h}_t = \tanh(W_h \cdot [g_t * h_{t-1}, x_t] + b_h)$$

$$h_t = h_{t-1} * (1 - z_t) + \tilde{h}_t * z_t$$