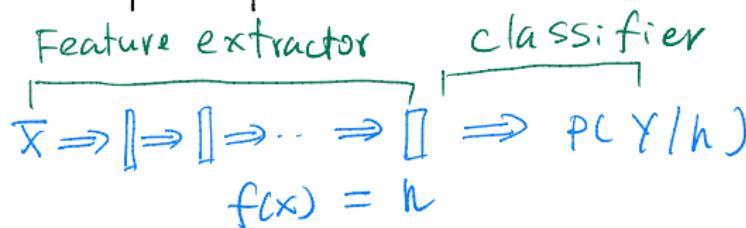


PART 2

PART 1: classic supervised learning

$$\{\bar{x}_i, y_i\}_{i=1}^N \Rightarrow p(y|\bar{x})$$

PART 2: Deep supervised learning



L06: Feedforward NN (FNN)

L07: Convolutional NN (CNN)

penultimate feature layer \downarrow logits

$$\underbrace{\left[\begin{array}{c} \vdots \\ \bar{x} \end{array} \right] \Rightarrow \cdots \Rightarrow \left[\begin{array}{c} \bar{z} \\ \bar{z}^T \bar{w} + b \end{array} \right]}_{p(\bar{y}|\bar{x}, \theta)} \rightarrow \bar{z} \quad p(\bar{y}|\bar{z})$$

$$J(\theta) = -\frac{1}{N} \sum_{i=1}^N p(y_i|\bar{x}_i; \theta)$$

- ① \bar{y} : real-valued vector \checkmark identity matrix
- $$p(\bar{y}|\bar{x}) = N(\bar{y}|\bar{z}, I)$$

per-sample cross entropy loss

$$L = -\log p(\bar{y}|\bar{x}; \theta) \propto \frac{1}{2} \|\bar{y} - \bar{z}\|_2^2$$
$$= \frac{1}{2} [(y_1 - z_1)^2 + (y_2 - z_2)^2 + \dots]$$

$$\frac{\partial L}{\partial z_k} = z_k - y_k \quad \underline{\text{Error}}$$

predicted observed

$$\frac{\partial L}{\partial z_2}$$

$$\textcircled{2} \quad Y \in \{0, 1\} \quad P(Y | \bar{x}) = \text{Ber}(Y | \sigma(z))$$

$$P(Y=1 | \bar{x}) = \sigma(z)$$

$$L = -\log P(Y | \bar{x}, \theta) = -[Y \log \sigma(z) + (1-Y) \log(1-\sigma(z))]$$

predicted *observed*

$$\frac{\partial L}{\partial z} = \sigma(z) - Y \quad \text{Error}$$

Sigmoid ok as output unit

$$\sigma(z) - Y \approx 0 \quad \left\{ \begin{array}{l} Y=1, \quad \sigma(z) \approx 1 \quad z \gg 0 \\ \qquad \qquad \qquad \text{Model correct} \\ Y=0, \quad \sigma(z) \approx 0, \quad z \ll 0 \\ \qquad \qquad \qquad \text{Model correct} \end{array} \right.$$

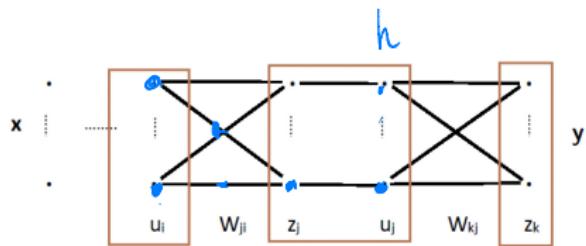
$$\textcircled{3} \quad Y \in \{1, 2, \dots, C\} \quad \bar{z} = (z_1, z_2, \dots, z_C)^T$$

$$P(Y=c | \bar{x}) = \frac{e^{z_c}}{\sum_c e^{z_c}}$$

$$L = -\log P(Y | \bar{x}, \theta)$$

$$\frac{\partial L}{\partial z_c} = P(Y=c | \bar{x}) - \mathbb{1}(Y=c) \quad \underline{\text{Error}}$$

↑ ↑
predicted observed



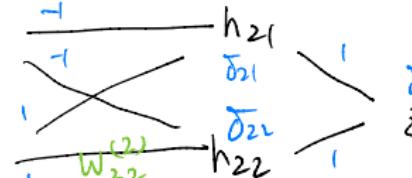
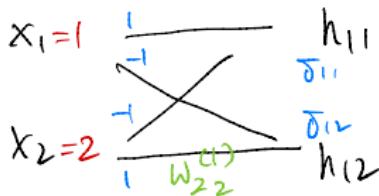
$$z_j = \sum_i u_i w_{ji} + \dots \quad u_j = g(z_j) \quad z_k = \sum_j u_j w_{kj} + \dots$$

$$\begin{array}{c} x_1 \quad x_2 \quad y \\ \hline 1 \quad 2 \quad 0 \end{array}$$

$$z_{11} = -1, u_{11} = 0$$

$$z_{21} = 1, u_{21} = 1$$

h : ReLU



$p(y|z)$:

Sigmoid

$$z = 2 \quad p(y=1|x_1=1, x_2=2) = \sigma(2)$$

Forward

$$z_{12} = 1, u_{12} = 1$$

$$z_{22} = 1, u_{22} = 1$$

$$\delta_{11} = \frac{\partial u_{11}}{\partial z_{11}} (\delta_{21} \times (-1) + \delta_{22} \times (-1))$$

$$= 0$$

$$\delta_{21} = \frac{\partial u_{21}}{\partial z_{21}} \delta_z \cdot 1$$

$$= 0.88$$

$$\delta_z = \sigma(2) - y$$

$$\approx 0.88$$

$$\delta_{12} = \frac{\partial u_{12}}{\partial z_{12}} (\delta_{21} \cdot 1 + \delta_{22} \cdot 1)$$

$$= 1 (0.88 \times 1 + 0.88 \times 1)$$

$$= 1.76$$

Backward

$$\delta_{22} = 0.88$$

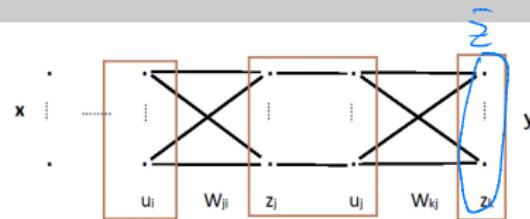
$$\frac{\partial L}{\partial w_{22}^{(1)}} = \text{input} \downarrow \frac{\partial L}{\partial w_{22}^{(2)}} = \text{error}$$

2022-02-04

Derivation of Back prop

$$L = -\log p(y|\bar{x}, \theta)$$

$$= -\log p(y|\bar{z}, \theta)$$



$$\underline{z}_k = \sum_j u_j W_{kj}, \quad u_j = g(z_j), \quad z_j = \sum_i u_i W_{ji}$$

$$\frac{\partial L}{\partial w_{kj}} = \frac{\partial L}{\partial z_k} \frac{\partial z_k}{\partial w_{kj}}$$

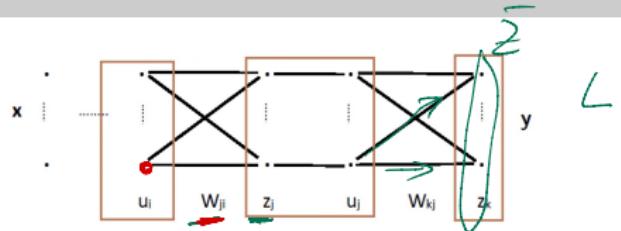
 $w_{kj} \rightarrow z_k \rightarrow L$

$$= \delta_k \quad u_j$$

$$= u_j \quad \delta_k$$

input error

$$w_{ji} \rightarrow z_j \rightarrow u_j \rightarrow \{z_k\} \rightarrow L$$



$$z_k = \sum_j u_j W_{kj}, \quad u_j = g(z_j), \quad z_j = \sum_i u_i W_{ji}$$

$$\begin{aligned} \frac{\partial L}{\partial w_{ji}} &= \sum_k \frac{\partial L}{\partial z_k} \frac{\partial z_k}{\partial w_{ji}} \\ &= \sum_k \delta_k \frac{\partial z_k}{\partial u_j} \frac{\partial u_j}{\partial z_j} \frac{\partial z_j}{\partial w_{ji}} \\ &= \sum_k \delta_k \frac{w_{kj}}{\frac{\partial u_j}{\partial z_j}} \frac{u_i}{\frac{\partial z_j}{\partial w_{ji}}} \\ &= \frac{u_i}{\text{input}} \frac{\frac{\partial u_j}{\partial z_j}}{\sum_k w_{kj} \delta_k} \delta_j \end{aligned}$$

$\begin{aligned} z_k &= u_j_1 w_{kj_1} \\ &\quad + u_j_2 w_{kj_2} \end{aligned}$

$\frac{\partial z_k}{\partial u_j_1} = w_{kj_1}$

SGD: $\bar{\theta} \leftarrow \bar{\theta} - \varepsilon \bar{g}$: push $\bar{\theta}$ in direction $-\bar{g}$
for ε distance, stop

Momentum

- * $\bar{\theta}$ has speed \bar{v}
- * At each iteration $0 < \alpha < 1$
- * current speed reduce $\bar{v} \leftarrow \alpha \bar{v}$
- * Get acceleration $-\varepsilon \bar{g}$
- * New speed: $\bar{v} \leftarrow \alpha \bar{v} - \varepsilon \bar{g}$
- * $\bar{\theta} \leftarrow \bar{\theta} + \bar{v}$

$$\begin{array}{c}
 t \quad V \\
 \hline
 0 \quad \bar{V}_0 = 0 \\
 1 \quad \bar{V}_1 = \alpha \bar{V}_0 - \varepsilon \bar{g}_1 = -\varepsilon \bar{g}_1 \\
 2 \quad \bar{V}_2 = \alpha \bar{V}_1 - \varepsilon \bar{g}_2 = -\varepsilon (\alpha \bar{g}_1 + \bar{g}_2) \\
 3 \quad \bar{V}_3 = \alpha \bar{V}_2 - \varepsilon \bar{g}_3 = -\varepsilon (\alpha^2 \bar{g}_1 + \alpha \bar{g}_2 + \bar{g}_3)
 \end{array}$$

$0 < \alpha < 1$

Exponentially decaying
 average of past
 gradients

AdaGrad (Adaptive Gradient)

$$\bar{\theta} = (\theta_1, \theta_2, \dots, \theta_{1,000,000})$$

$$\bar{g} = (g_1, g_2, \dots, g_{1,000,000})$$

Past change (gradient)

$$\bar{r} = (r_1, r_2, \dots, r_{1,000,000})$$

$$\bar{r} = 0$$

$$\bar{r} \leftarrow \bar{r} + \bar{g} \odot \bar{g}$$

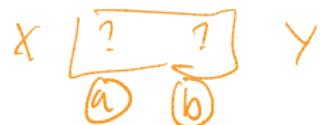
\uparrow \bar{g}^2

Componentwise
product

After t iterations:

$$\bar{r} = \bar{g}_1^2 + \bar{g}_2^2 + \dots + \bar{g}_t^2$$

para's that have
changed a lot
already should
change less
in future



$$\bar{g}_1 = (g_{11}, g_{12}, \dots, g_{1,000,000})$$

$$\bar{g}_2 = (g_{21}, g_{22}, \dots,)$$

SGD: $\theta \leftarrow \theta - \varepsilon \bar{g}$ $\Delta \theta = -\varepsilon \bar{g}$

AdaGrad:

$$\Delta \theta = -\frac{1}{\delta + \sqrt{r}} \theta \varepsilon \bar{g}$$

τ
Smoothing constant

$$\bar{r} = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix}^T$$

1,000,000

$$\bar{g} = \begin{bmatrix} \cdot \\ \vdots \\ \cdot \\ \cdot \end{bmatrix}$$

RMS prop (Root mean square prop)

$$\bar{r} \leftarrow \rho \bar{r} + (-\rho) g^2$$

$0 < \rho < 1$

AdaGrad: $\bar{r} \leftarrow \bar{r} + g^2$

After t iteration:

$$\bar{r}_t = (1-\rho) (\underline{\rho^{t-1} g_1^2} + \rho^{t-2} g_2^2 + \dots + g_{t-1}^2)$$

Adam

Momentum

$$\bar{v} \leftarrow \alpha \bar{v} - \varepsilon \bar{g} \quad \Delta \bar{\theta} = \bar{v}$$

$$\bar{u} = -\bar{v}$$

$$\bar{u} \leftarrow \alpha \bar{u} + \varepsilon \bar{g} \quad \Delta \bar{\theta} = -\bar{u}$$

RMS prop

$$\bar{r} \leftarrow \rho \bar{r} + (1-\rho) \bar{g}^2$$

$$\Delta \bar{\theta} = -\varepsilon \bar{g} \frac{1}{\delta + \sqrt{\bar{r}}}$$

Adam (Adaptive Moments)

$$\bar{s} \leftarrow \rho_1 \bar{s} + (1-\rho_1) \bar{g}$$

$$\bar{r} \leftarrow \rho_2 \bar{r} + (1-\rho_2) \bar{g}^2$$

$$\bar{s} \leftarrow \frac{1}{1-\rho_1 t} \bar{s}$$

$$\bar{r} \leftarrow \frac{1}{1-\rho_2 t} \bar{r}$$

moment
correction

$$\Delta \theta = -\varepsilon \bar{s} \odot \frac{1}{\delta + \sqrt{\bar{r}}}$$

$$S \leftarrow p_i \bar{S} + ((-p_i) \bar{g})$$

$$\hat{S} = \frac{1}{1-p_i t} S$$

$$\frac{t}{1-p_i g_1}$$

$$2 \quad ((-p_i)(p_i g_1 + g_2))$$

$$3 \quad (1-p_i)(p_i^2 g_1 + p_i g_2 + g_3)$$

:

$$E[g_t] = \int p(t) \underline{g_t} dt = \frac{p_i^2}{z} g_1 + \frac{p_i}{z} g_2 + \frac{1}{z} g_3$$

1st order moment

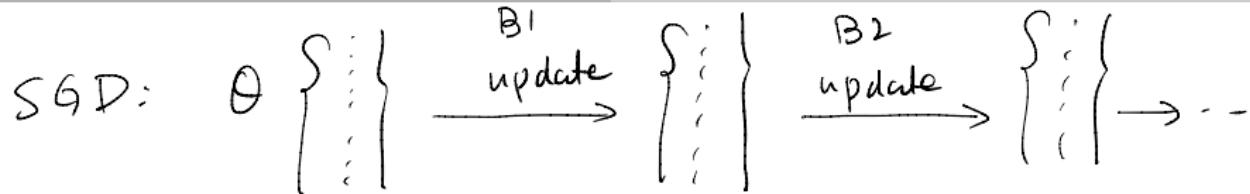
$$\oplus = 1$$

$$p(2)$$

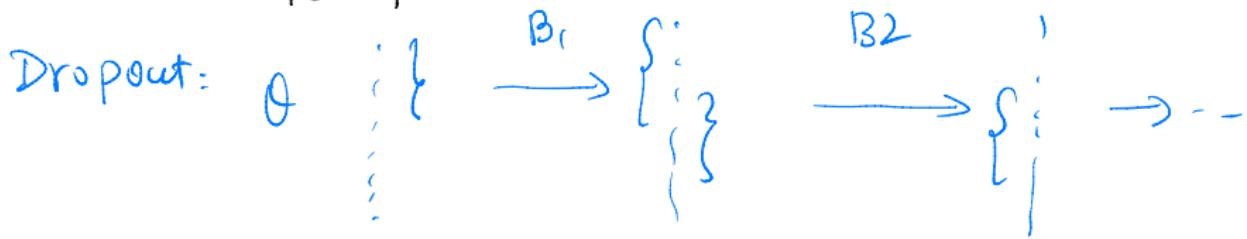
$$r: \int p(t) g_t^2 dt$$

2nd order moment

2022-03-11



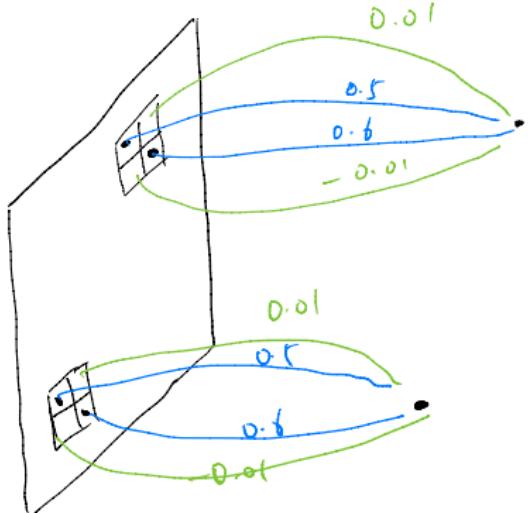
All parameters can co-adapt
to fit noise



Breaks the co-adaptation

2022-03-16

L07: CNN



- * One unit for detecting one feature

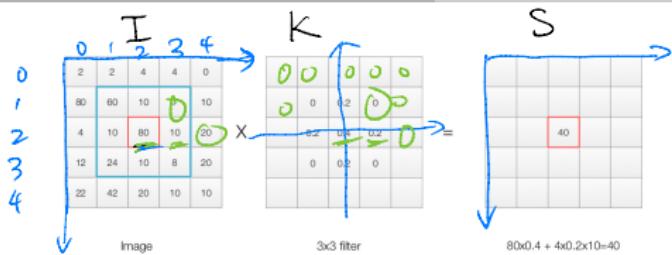
- * Detect the same feature elsewhere Share parameter

Another inductive bias:
shift invariant

Filter | Kernel:

$$\begin{matrix} 0.5 & 0.0 \\ 0.0 & 0.6 \end{matrix}$$

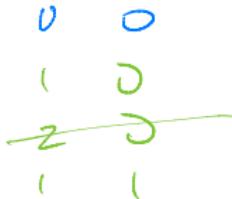
- * Detect multiple feature? Multiple filters



$$S(2,2) = I(2,2) K(0,0) + I(3,2) K(1,0) + \dots$$

$$S(i,j) = \sum_{m,n} I(i+m, j+n) K(m,n)$$

Cross-correlation



$$k'(m, n) = k(-m, -n)$$

$$\Rightarrow k(m, n) = k'(-m, -n)$$

$$s(i, j) = \sum_{m, n} I(i+m, j+n) k(m, n)$$

$$= \sum_{m, n} I(i+m, j+n) k'(-m, -n)$$

$$a = -m$$

$$b = -n$$

$$= \sum_{a, b} I(i-a, j-b) k'(a, b)$$

Convolution

Input: $W_1 \times H_1 \times D_1$

K $F \times F$ filters

stride: S

padding: P

Output: $W_2 \times H_2 \times D_2$

$$D_2 = K$$

$$W_2 = \frac{1}{S}(W_1 + 2P - F) + 1$$

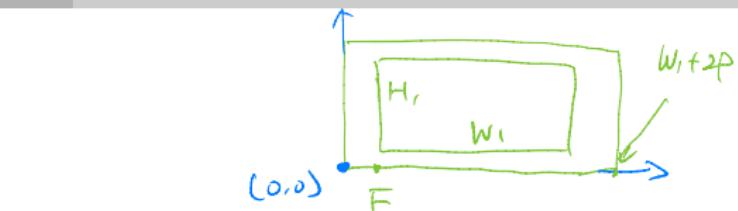
$$H_2 = \frac{1}{S}(H_1 + 2P - F) + 1$$

of pairs:

$$(F^2 + 1)K$$

$$\# \text{ of FLOPS} = 2 \underbrace{F^2 D_1}_{\text{floating point operations}} \underbrace{W_2 H_2 D_2}_{\text{operations}} = \frac{W_1 + 2P - \frac{F}{2} - \frac{F}{2}}{S} + 1$$

Floating point
operations



on X-axis:

$$(0,0)$$

$$\frac{F}{2}$$

center

1st possible location for filter

$$\frac{F}{2}$$

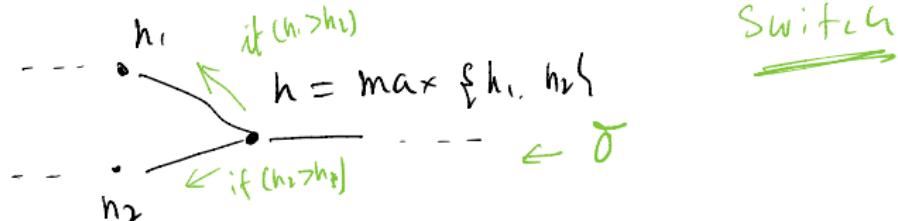
last possible location - - -

$$W_1 + 2P - \frac{F}{2}$$

of possible location - - -

$$\frac{W_1 + 2P - \frac{F}{2} - \frac{F}{2}}{S} + 1$$

Max Pooling . Back prop



$$\frac{\partial h}{\partial h_1} = \begin{cases} 1 & h_1 > h_2 \\ 0 & h_1 < h_2 \end{cases}$$

$$\frac{\partial h}{\partial h_2} = \begin{cases} 1 & h_2 > h_1 \\ 0 & h_2 < h_1 \end{cases}$$

Data Normalization

$$x_1 \in [0, 1], \quad x_2 \in [0, 1,000]$$

$$\hat{y} = w_0 + w_1 x_1 + \underline{w_2 x_2}$$

$$\Rightarrow J(\bar{w}) = E[L(y, \hat{y})] + \lambda (w_1^2 + w_2^2)$$

Regularization affects which parameter more:

$$w_1^2 \downarrow \qquad \qquad w_2^2 \downarrow$$

Impact on \hat{y} less more

Impact on L : less more

Ridge regression : * Reduce $w_1^2 + w_2^2$
* Minimize L

What to do: Reduce w_1^2 more \Rightarrow bias

$$x_1 \in [0, 1], \quad x_2 \in [0, 1,000]$$

$$\hat{y} = w_0 + w_1 x_1 + \underline{w_2 x_2}$$

$$J(\bar{w}) = E[L(y, \hat{y})] + \lambda (w_1^2 + w_2^2)$$

$$\frac{\partial J}{\partial w_1} = E \left[\frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_1} \right] + 2\lambda w_1$$

$$= E \left[\frac{\partial L}{\partial \hat{y}} \right] x_1 + 2\lambda w_1$$

Smaller

$$\frac{\partial J}{\partial w_2} = E \left[\frac{\partial L}{\partial \hat{y}} \right] x_2 + 2\lambda w_2$$

larger

Unnormalized Data

$$x_{11} \quad x_{12} \quad \dots \quad x_{1D}$$

$$x_{22} \quad x_{22}, \quad \dots \quad x_{21})$$

— — — —

$$x_{N1}, x_{N2}, \dots, x_{NQ}$$

$$\mu_1 \quad \mu_2 \quad \mu_D$$

$$\Gamma: \sigma_1 \quad \sigma_2 \quad \dots \quad \sigma_p$$

Normalized Data

$$x_{11}^{\gamma} \quad x_{12}^{\gamma} \quad - \quad x_{10}^{\gamma}$$

•  

$$\hat{x}_{N1} \quad \hat{x}_{N2} \quad \hat{x}_N^D$$

$$\mu = \overbrace{0 \quad 0}^0$$

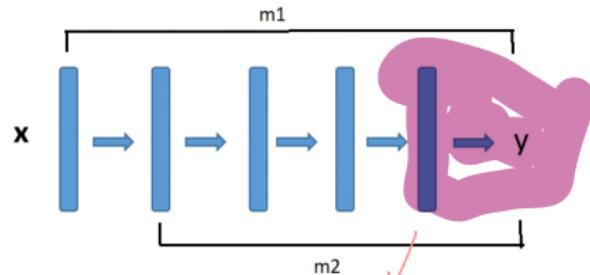
$$\sigma: \quad 1 \quad 1 \quad \underline{1}$$

$$\mu_j = \sum_{i=1}^N x_{ij} \quad \text{mean of c}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{ij} - \mu_j)^2$$

$$\hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sigma_i}$$

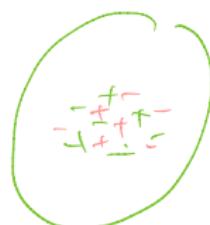
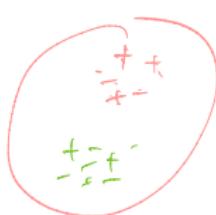
Batch Normalization



Without
Normalization

iteration 10

iteration 15



Covariate shift

Inception Module

Input: $H_i \times W_i \times D_i$ # of para's

(a) K $F \times F$ filters $(F^2 D_i + 1) K$

\Rightarrow Output 1

(b) D 1×1 filters

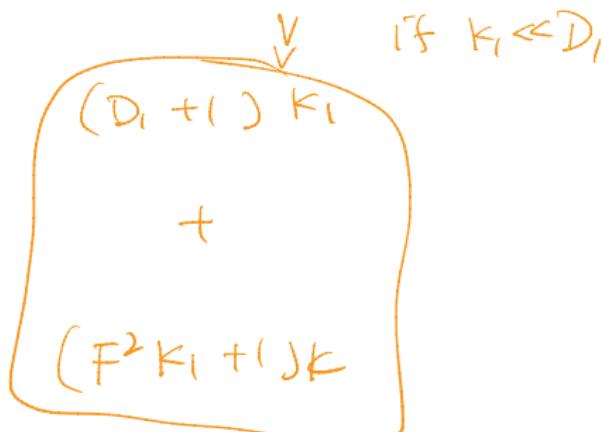
\Rightarrow intermediate output:

$H_i \times W_i \times K_1$

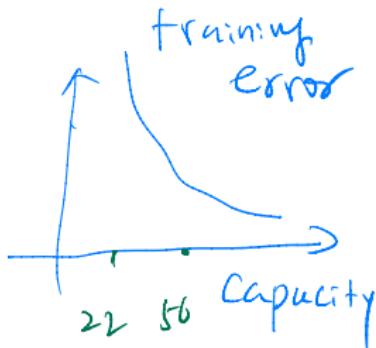
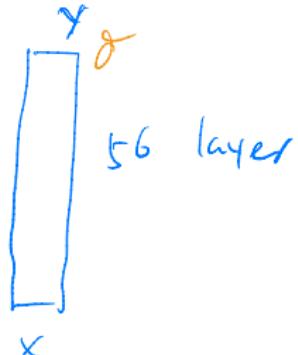
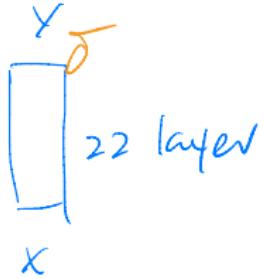
② K $F \times F$ filters

\Rightarrow Output 2:

Same shape as Output 1



Gradient



Training
error

Theory high

lower

layer layer

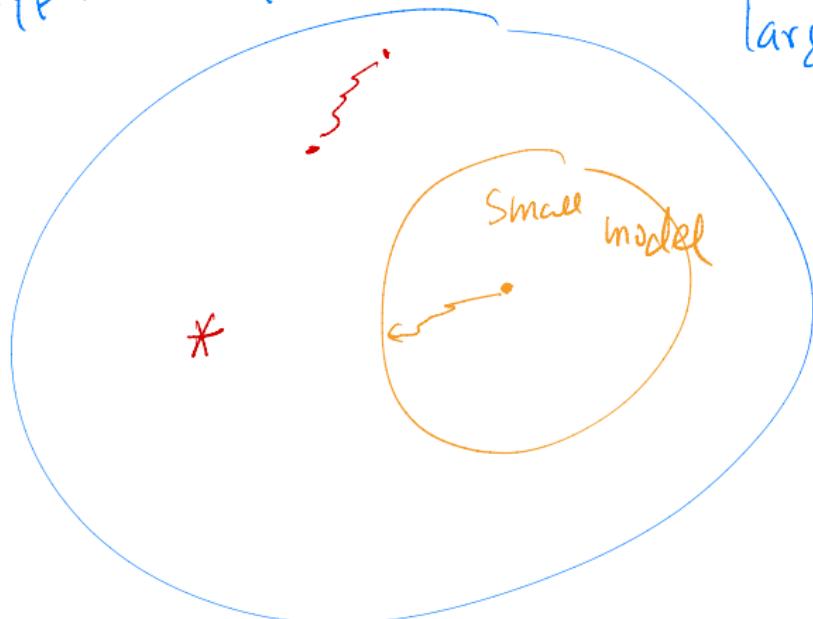
Practice lower

higher

model model

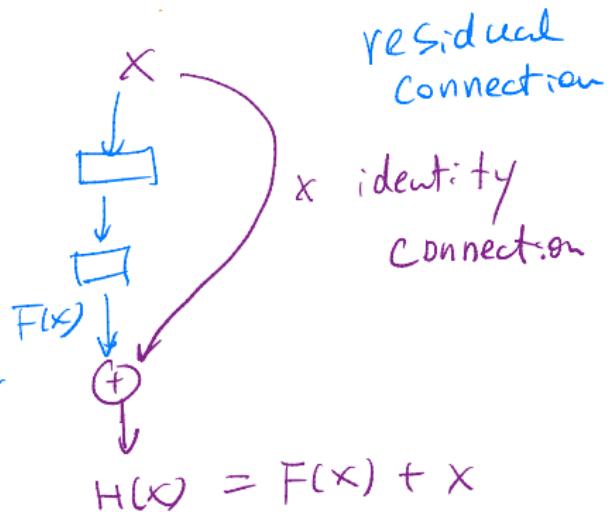
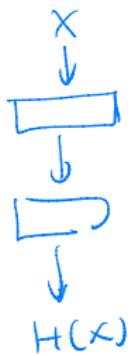
hypothesis space

larger model



Residual Module

plain net



$$F(x) = \frac{H(x) - x}{\text{residual}}$$

