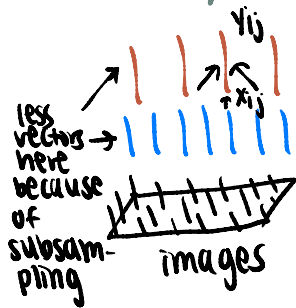# LECTURE 14
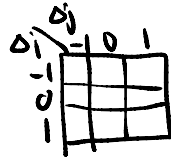
– CNN/conv. Net

*each layer corresponds to a box*
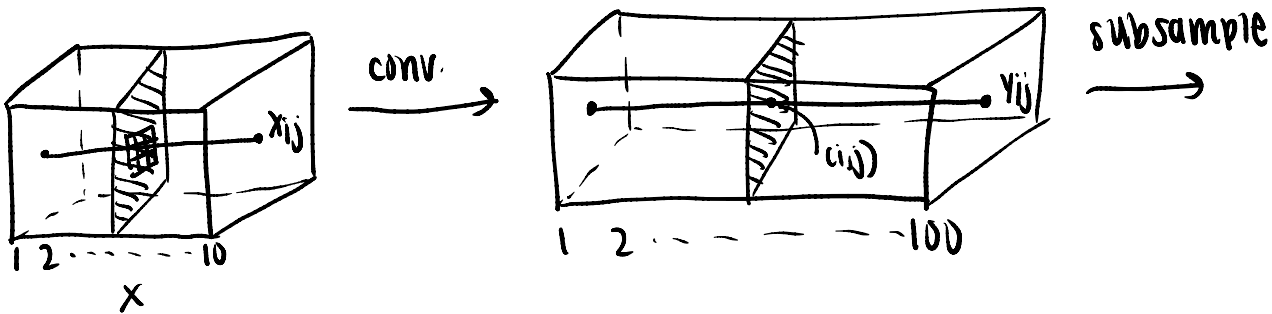
thought vectors compose
vectors at lower layer,
cover bigger range (more complex)



less vectors here because of subsampling    images

$$y_{ij} = \sum w_{\Delta i, \Delta j} x_{i+\Delta i, j+\Delta j} + b + \text{rectify}$$



composition



conv.    subsample

$x_{ij}$

1 2 ... 10    X

$y_{ij}$    (i,j)

1 2 ... 100
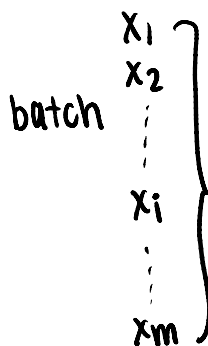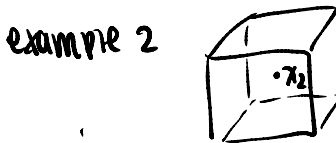
Le Net (80s) MNIST
Dan Net (2010) GPU, bigger model
✓ Alex Net (2012) Image Net (1000 categories, 1 mil. examples)
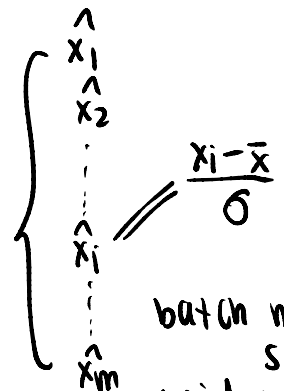        also a bigger version exists
            ↳ Amazon Turks

Google Le Net    1×1, 3×3, 5×5 filters
VGG            3×3, 16-layer, 19-layer, commonly used as <u>pre-trained</u> features
            (use learned weights to get features)
            – perceptual
            – 64 million parameters

## Batch Normalization (BN)

example 1        $x_1$

batch    $\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_i \\ \vdots \\ x_m \end{matrix}$

example 2        $x_2$

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

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

insert BN after each
weight layer

$\begin{matrix} \hat{x}_1 \\ \hat{x}_2 \\ \vdots \\ \hat{x}_i \\ \vdots \\ \hat{x}_m \end{matrix}$    $= \frac{x_i - \bar{x}}{\sigma}$
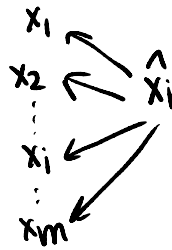
batch mean = 0
    s.d. = 1
avoid covariate
shift, stabilize
learning

$$\Rightarrow \quad y_i = \beta \hat{x}_i + r$$

$$y_1$$
$$\vdots$$
$$y_m$$

back-prop of BN

$$x_1 \quad \hat{x}_i$$
$$x_2 \leftarrow$$
$$\vdots$$
$$x_i$$
$$\vdots$$
$$x_m$$

## layer normalization



normalize within the
same layer of the
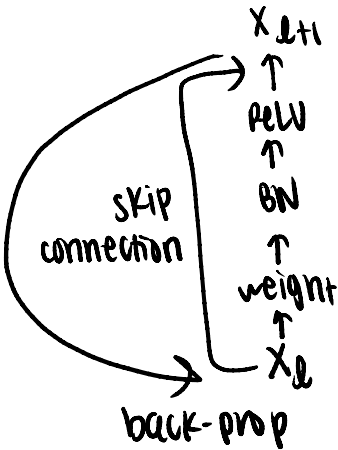same example

## residual network (Res Net)

conventional CNN:
$$x_{\ell+1} = f_\ell(x_\ell)$$
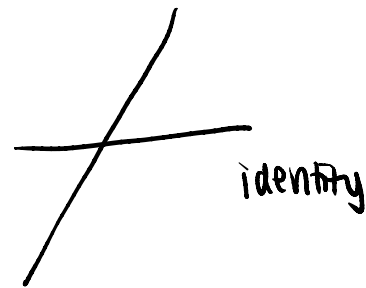↓
weight + BN + ReLU

Res Net:
$$x_{\ell+1} = x_\ell + \boxed{f_\ell(x_\ell)}$$

— residual:
correction, refinement,
update of previous layer

$$x_{\ell+1}$$
↑
ReLU
↑
BN
↑
weight
↑
$$x_\ell$$

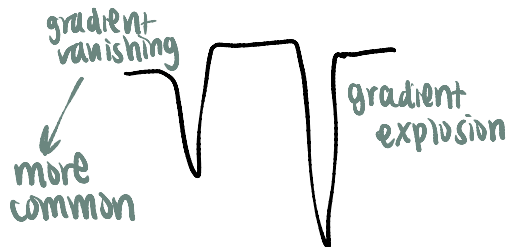SKIP
connection

back-prop

why does this perform better than conventional?

① ReLU has difficulty to learn identity mapping



ReLU          identity
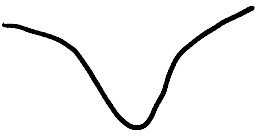
② loss function
conventional CNN

gradient
vanishing

more
common

gradient
explosion

- even w/o gradient vanishing,
have complex loss function / landscape
(sub-optimal local modes)

In ResNet

- closer to convex
- loss function better
  behaved
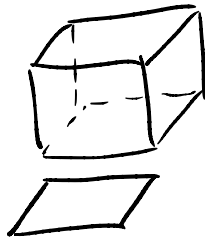  ↳ less multi-modal

③ Iterative Algorithm
   recurrent network

$$X_{t+1} = X_t + f_t(X_t)$$
   (e.g. gradient descent)

9 res. blocks

↑ subsample

3 res. blocks

Res 50/200

gradient

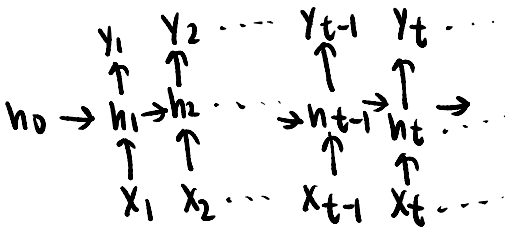$$\frac{\partial X_{\ell+1}}{\partial X_\ell} = \mathbb{1} + \frac{\partial f_\ell}{\partial X_\ell}$$

avoid gradient vanishing

## Recurrent Neural Network (RNN)

output y

hidden h

input x

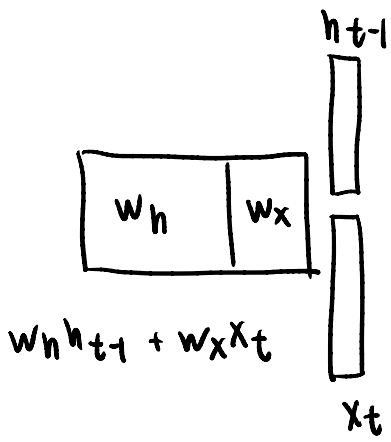$y_1 \quad y_2 \quad \cdots \quad y_{t-1} \quad y_t \cdots$

$h_0 \to h_1 \to h_2 \cdots \to h_{t-1} \to h_t \to \cdots$

$x_1 \quad x_2 \cdots \quad x_{t-1} \quad x_t \cdots$

→ time

$y_t$
↑
$h_t$   thought
vector
↑
$x_t$

thought process

## vanilla RNN

$$h_t = \text{Rectify}\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

weight concatenate
(generic notation)
differ in each occurrence

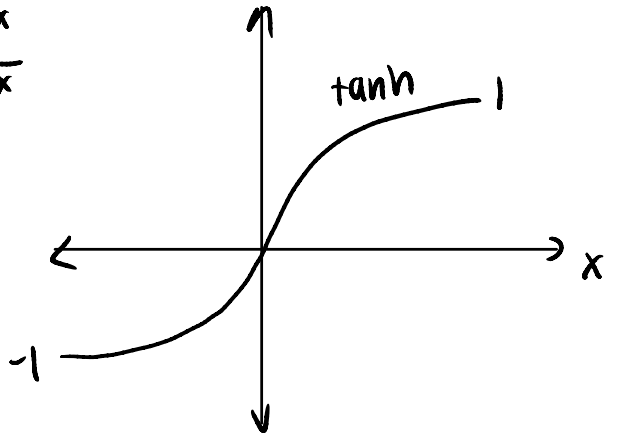$w_h$ | $w_x$

$h_{t-1}$

$x_t$

$w_h h_{t-1} + w_x x_t$

Rectify: RELU sigmoid

$$tanh = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$
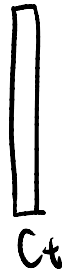
$y_t = $ Rectify $(w, h_t)$

softmax → another $w$

tanh → 1

x

-1
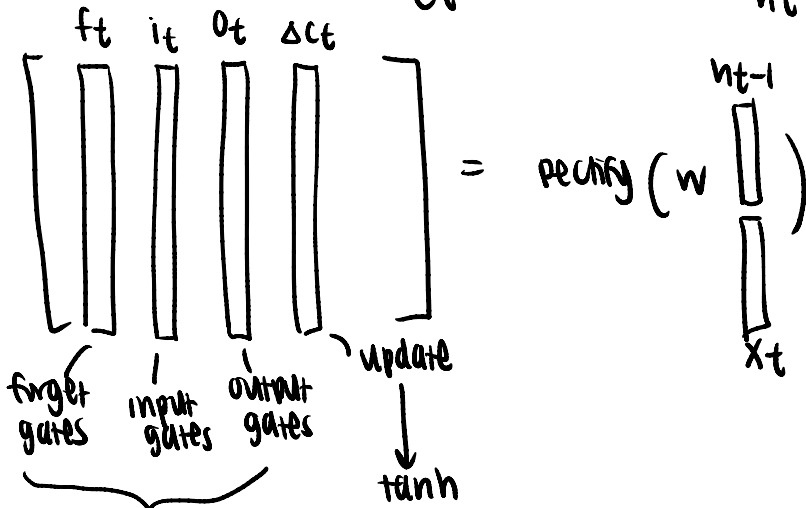
suffer gradient vanishing in back-prop through time
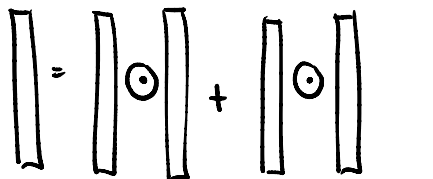solution: long short term memory (LSTM)

memory cells $C_t$

same dimension as $h_t$

$f_t$ $i_t$ $o_t$ $\Delta c_t$

$= $ Rectify $(w \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix})$

forget gates, input gates, output gates, update

IF-then o/i

tanh

sigmoid    tanh

element wise product

$C_t = f_t \odot C_{t-1} + i_t \odot \Delta C_t$

$[\ ] = [\ ] \odot [\ ] + [\ ] \odot [\ ]$

special case $f_t = 1$ $i_t = 1$

$C_t = C_{t-1} + \Delta C_t$   (RES Net)

$h_t = o_t \odot$ Rectify $(C_t)$

$[\ ] = [\ ] \odot [\ ]$

$y_t = $ Rectify $(w h_t)$