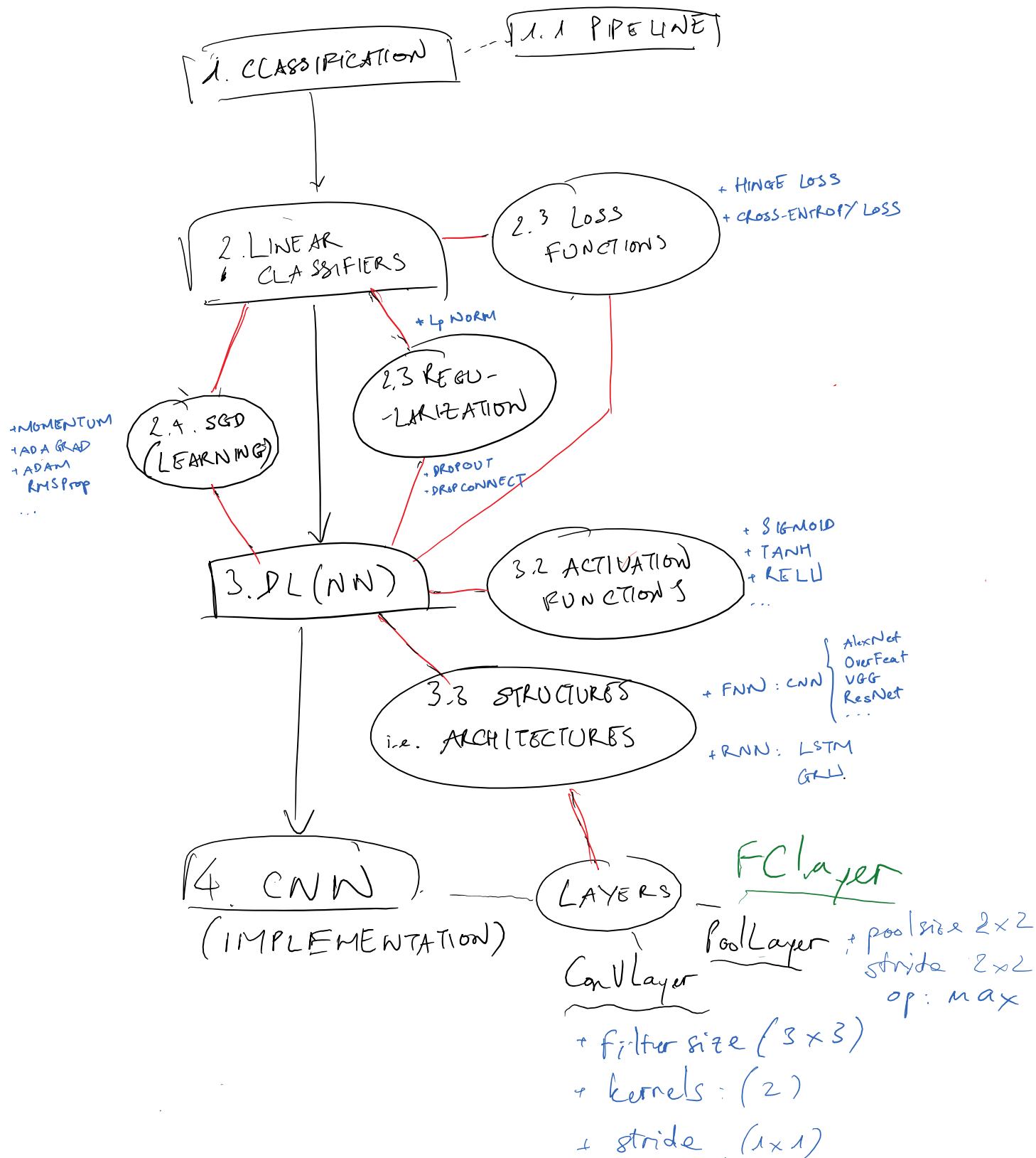


## OVERVIEW



+ zero padding (1)

# (Single-label) Classification, ANN

Sunday, September 25, 2016 10:50 PM

$x \Rightarrow \Phi$ : need feature extraction  $f(\cdot)$  and (linear) classifier  $c$

Linear classifier  $c$ : assume feature representation  $z$  given

	2-class	K-class
Example of $c$	2-way softmax classifier	K-way softmax classifier
Loss function $J(\cdot)$	ideal: $\sum_n \mathbb{I}(y^{(n)} \neq \hat{y}^{(n)})$ approximation: binary cross-entropy	ideal: $\sum_n \mathbb{I}(y^{(n)} \neq \hat{y}^{(n)})$ approximation: multinomial CE
score function $s(\cdot)$	sigmoid	softmax
Decision $d(\cdot)$	$\hat{y}^{(n)} = \mathbb{I}(s(x^{(n)}) > 0.5)$	$\hat{y}^{(n)} = \arg \max_k s_k(x^{(n)})$

Learning: minimize  $J(\cdot)$  with respect to each  $(W, b)$  by (stochastic) gradient descent

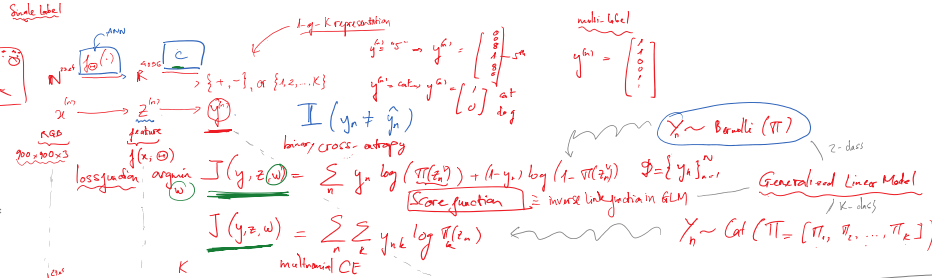
- Feature extractor  $f(\cdot)$ :  $z$  is actually not given
    - $z = f(x) \leftarrow f(\cdot) \equiv \text{NN}$
    - $z = f^{(2)}(f^{(1)}(x)) \leftarrow f^{(2)} \circ f^{(1)} \equiv \text{DNN}$
    - activation function  $\sigma^{(l)}(\cdot)$
- $\Rightarrow$  Learning: minimize  $J(\cdot)$  w.r.t. each  $(W, b)$  by (stochastic) gradient descent

$\sigma^{(2)} = \text{sigmoid}$   
 $\hat{y} = \text{ReLU}$   
 $\hat{y} = K \cdot \text{ReLU}$

$z_1 = \sigma(W_{11}^{(1)}x + b_1^{(1)})$   
 $z_2 = \sigma(W_{21}^{(1)}x + b_2^{(1)})$   
 $\vdots$   
 $z_{d_{\text{max}}} = \sigma(W_{d_{\text{max}}1}^{(1)}x + b_{d_{\text{max}}}^{(1)})$

$\hat{y} = \{W_{11}^{(1)}, b_{11}^{(1)}\}$   
 $z = f(x) = \sigma^{(2)}(W^{(1)}x + b^{(1)})$

$S = c(\hat{z}) = c(f^{(2)}(f^{(1)}(x)))$   
 $S = \sigma^{(2)}(W^{(2)}f^{(1)}(x) + b^{(2)})$



Deep Neural Network

2-class:  $S(x) = \frac{e^{w_1^T x + b_1}}{1 + e^{w_1^T x + b_1}} \text{ (sigmoid)}$   
 K-class:  $S_k(x) = \frac{e^{w_k^T x + b_k}}{\sum_k e^{w_k^T x + b_k}} \text{ (softmax)}$

Score function:  $S_k(x) = \frac{e^{w_k^T x + b_k}}{\sum_k e^{w_k^T x + b_k}}$

Generalized Linear Model (GLM):  $\hat{y} = w_0 + w_1 z_1 + w_2 \log(z_2) + w_3$

Learning

optimal algo: SGD

SGD: 1. Init  $\{W_{dk}^{(0)}\}, \{b_k^{(0)}\}$   
 2. Loop until convergence (e.g.  $|J^{(i+1)} - J^{(i)}| < \epsilon$ )

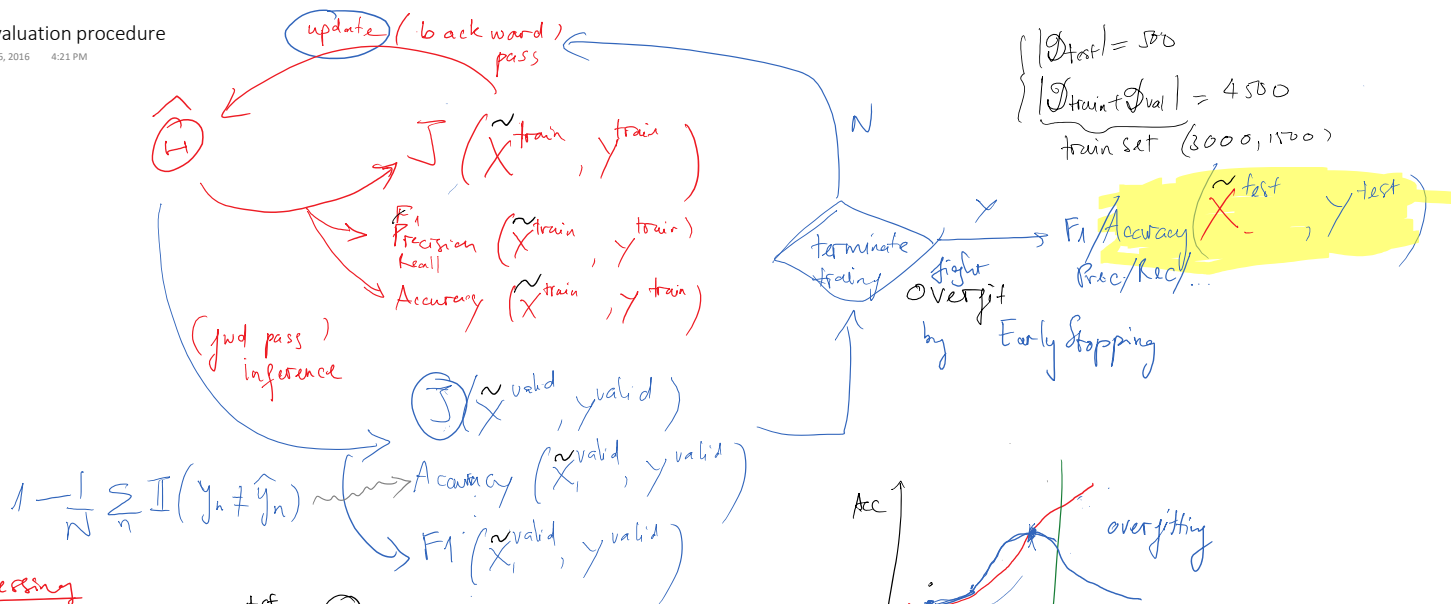
update expressions:

- Adam:  $+ \alpha, \beta$
- Ada grad:  $+ \alpha, \beta$
- RMS prop:  $+ \alpha, \beta$

weight decay

backprop

exploring / vanishing gradient



### Data preprocessing

Normalization  $X_n^{test} \rightarrow \tilde{X}_n^{test} = X_n^{test} - \bar{X}$

$c \times w \times h \rightarrow c \times \tilde{w} \times \tilde{h}$

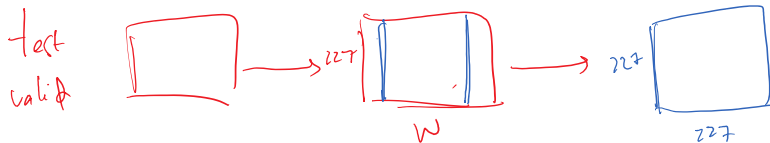
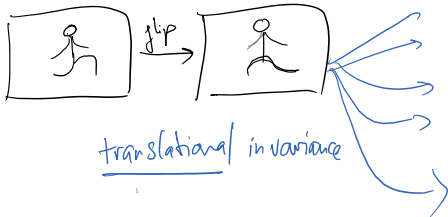
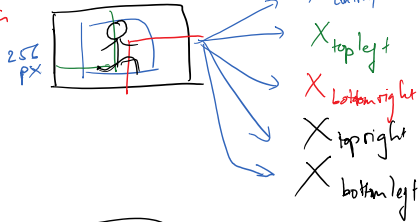
$\bar{X} \leftarrow \{X^{train}\}$

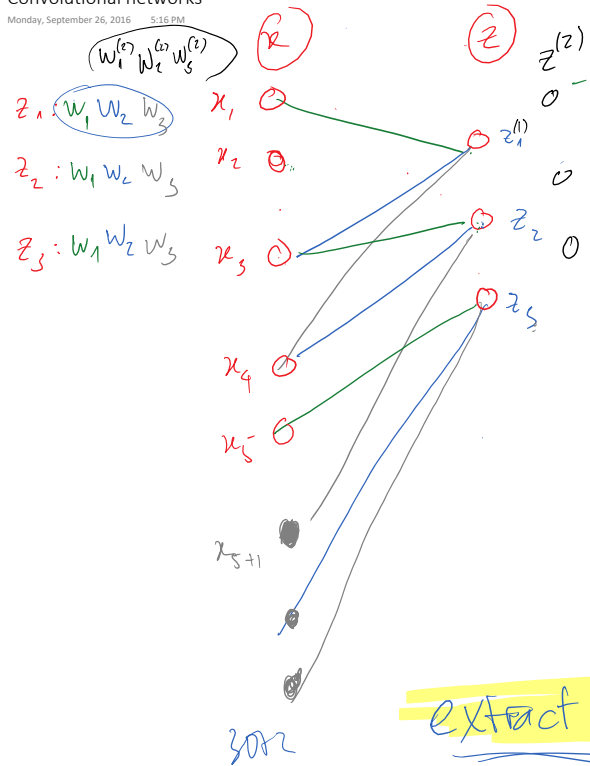
$w, h \in [0, 255]$   
etc

$mean(\tilde{w}) \approx 0$   
 $mean(\tilde{h}) \approx 0$

crop size:  $227 \times 227$   
 $X^{train}$

### Augmentation





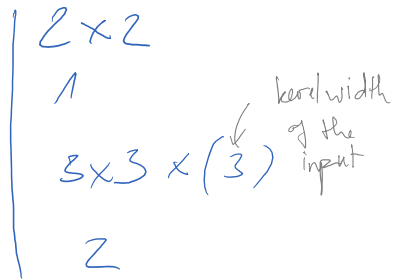
## Convolutional layer

- + locally connected  $\gg$  fully-connected (FC)
- + weight sharing

+ hyperparam of the structure

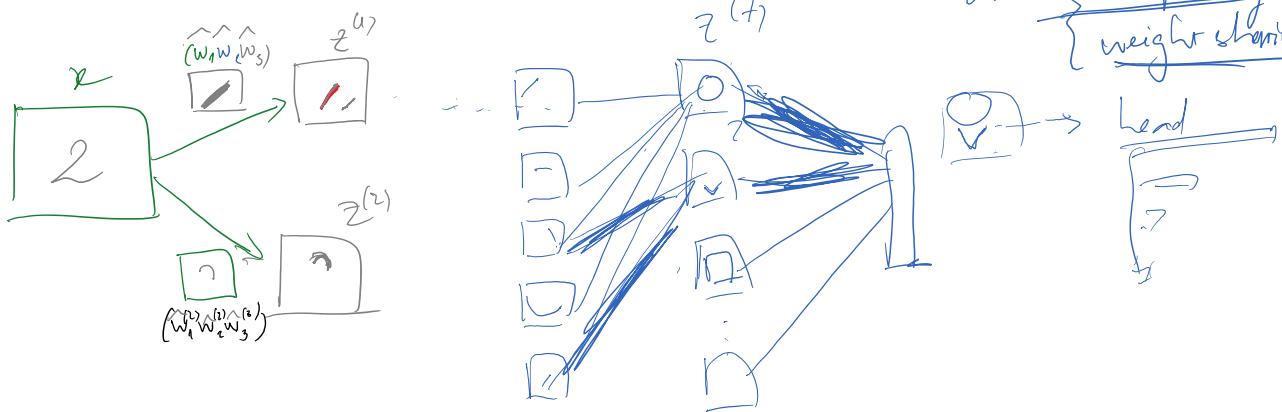
- stride: 2
- zero-padding: 3
- filter size: 3
- # filters (# kernels): 2

eg  $2 \times 2 \times 3 \times 3 \times 3$



extract invariant features w.r.t translation

due to maxpooling  
weight sharing

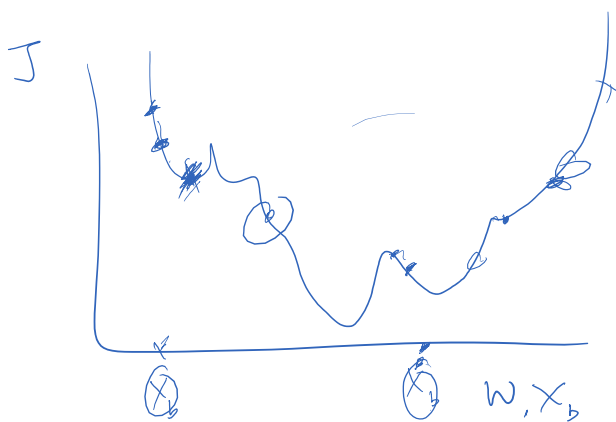


POOL layer

$z_{conv}$

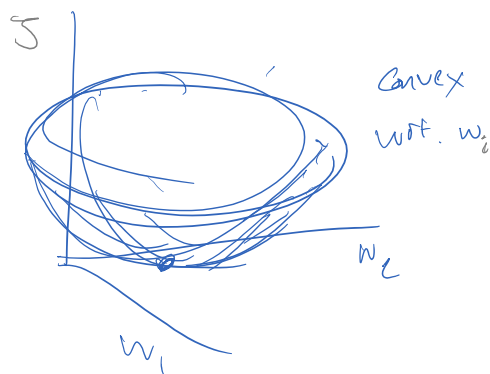
$z_{pool}$

SGD



Stochastic optian (w)

non-convex loss fn



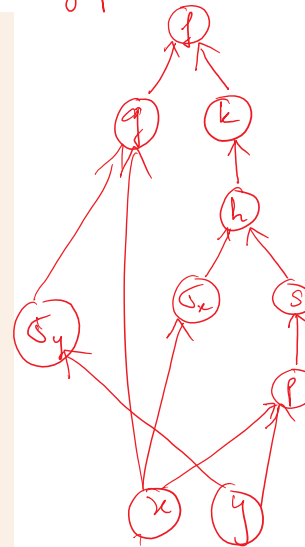
# backprop

Monday, May 2, 2016 5:55 PM

ALGO IS THE SAME FOR ALL  $\textcircled{1}$  STEP 1 have a variable dependency graph

ANALOGY TO FNN:

- $f(\cdot) \equiv$  loss function
- other variables  $\equiv W, b$  at each layer



eg: compute  $\left. \frac{\partial f}{\partial x} \right|_{x=x_0}$  and  $\left. \frac{\partial f}{\partial x} \right|_{y=y_0}$  when  $f(x,y) = \frac{x + \sigma(y)}{\sigma(x) + (x+y)^2}$

forward		backward	
		DONE	
1	$\sigma_y = \sigma(y) = \frac{1}{1+e^{-y}}$	8 (*)	$dy \leftarrow dy + d\sigma_y \cdot \frac{\partial \sigma_y}{\partial y} = dy + (1 - \sigma_y) \sigma_y$
2	$g(x, \sigma_y) = x + \sigma_y$	7 (*)	$dx \leftarrow dx + dg \cdot \frac{\partial g}{\partial x} = dx + dg$ $d\sigma_y = \frac{\partial f}{\partial \sigma_y} = dg \cdot \frac{\partial g}{\partial \sigma_y} = dg$
3 (*)	$\sigma_x = \sigma(x) = \frac{1}{1+e^{-x}}$	6 (*)	$dx \leftarrow dx + d\sigma_x \cdot \frac{\partial \sigma_x}{\partial x} = dx + (1 - \sigma_x) \sigma_x$
4 (*)	$p(x,y) = x + y$	5 (*)	$dx = \frac{\partial f}{\partial x} \leftarrow dp \cdot \frac{\partial p}{\partial x} = dp$ $dy = \frac{\partial f}{\partial y} \leftarrow dp \cdot \frac{\partial p}{\partial y} = dp$
5	$s(p) = p^2$	4	$dp = \frac{\partial f}{\partial p} = ds \cdot \frac{\partial s}{\partial p} = ds \cdot 2p$
6	$h(\sigma_x, s) = \sigma_x + s$	3	$d\sigma_x = \frac{\partial f}{\partial \sigma_x} = dh \cdot \frac{\partial h}{\partial \sigma_x} = dh$ $ds = \frac{\partial f}{\partial s} = dh \cdot \frac{\partial h}{\partial s} = dh$
7	$k(h) = \frac{1}{h}$	2	$dh = \frac{\partial f}{\partial h} = dk \cdot \frac{\partial k}{\partial h} = dk \cdot \left(-\frac{1}{h^2}\right)$
8	$f(g,k) = g \cdot k$	1	$dg = \frac{\partial f}{\partial g} = k$ $dk = \frac{\partial f}{\partial k} = g$
DONE			

STEP 2 Compute derivative  $\frac{\partial f}{\partial \eta}$  (denoted as  $dy$ )

w.r.t. direct child  $\theta_i$  of  $\eta$

$$\frac{\partial f}{\partial \eta} := \frac{\partial f}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \eta} \quad \text{i.e. } d\theta_i$$

STEP 3: if a parent  $\eta$  has  $\geq 2$  children  
say  $\theta_1, \theta_2, \dots$ , ADD UP

$$\frac{\partial f}{\partial \eta} += d\theta_1 \frac{\partial \theta_1}{\partial \eta}$$

$$\frac{\partial f}{\partial \eta} += d\theta_2 \frac{\partial \theta_2}{\partial \eta}$$

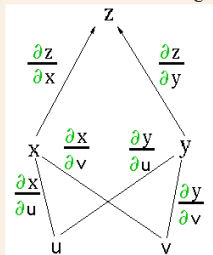
(multivariate calculus chain-rule: add up at fork (parent who has 2+ children))

Note

Backprop intuitions: passing the messages

<http://cs231n.github.io/optimization-2/>

multivariable chain-rule: gradients add up at forks (\*)



$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial u} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial u}$$

$$\frac{\partial z}{\partial v} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial v} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial v}$$

Multi-label learning strategy

Monday, September 26, 2016 6:47 PM

