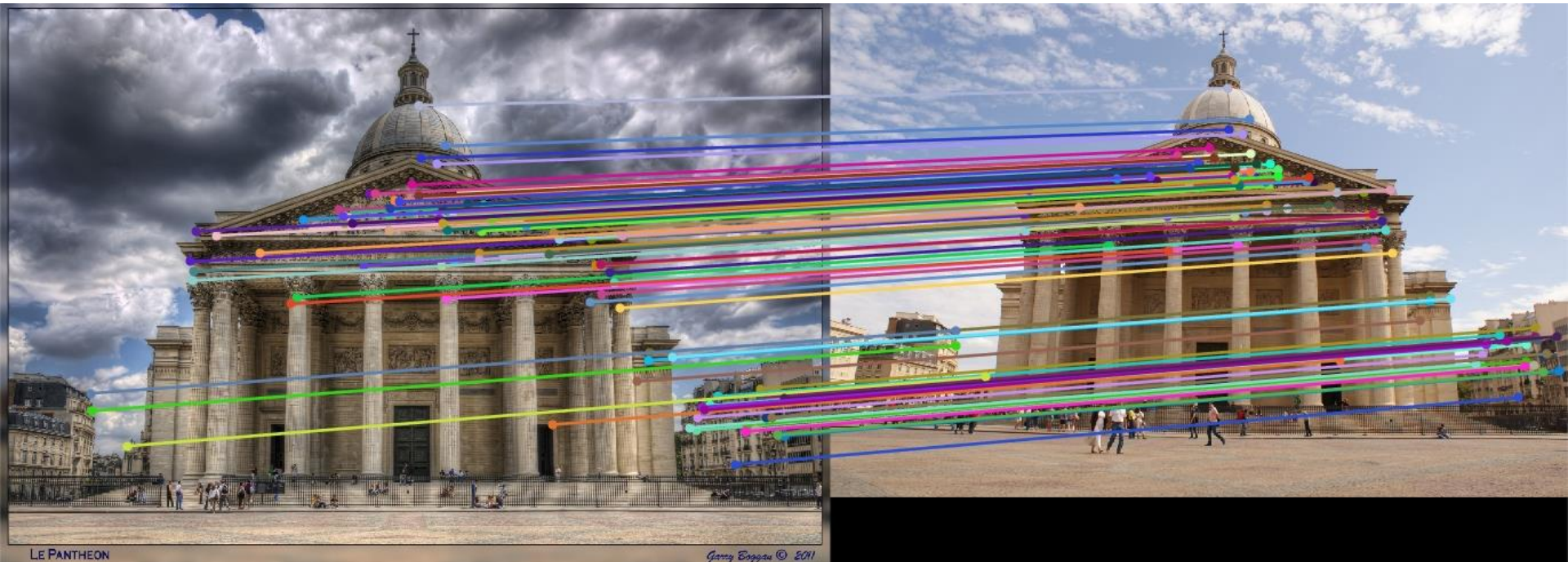# Project 2 highlights

# Project 2 highlights

- Jared Duncan
- Nathaniel Glaser
- Anh Thai
- Yuhui Zhao
- Jiaye Liu
- seunghwan lee
- yinglin li

# Deep Learning 1
# Neural Net Basics

Computer Vision

James Hays

# Outline

- Neural Networks
- *Convolutional* Neural Networks
- Variants
  - Detection
  - Segmentation
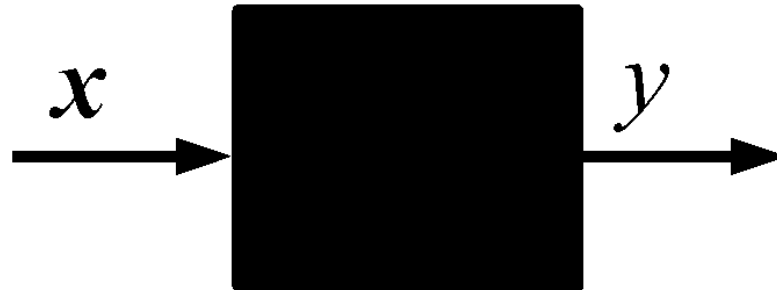  - Siamese Networks
- Visualization of Deep Networks

# Supervised Learning

$$\left\{(\boldsymbol{x}^i, y^i), i=1\dots P\right\}$$ training dataset

$\boldsymbol{x}^i$     i-th input training example

$y^i$     i-th target label

$P$     number of training examples



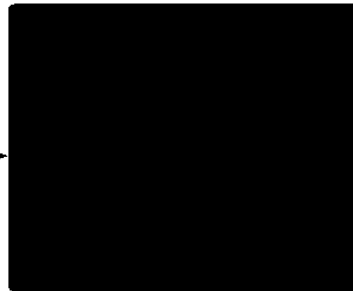Goal: predict the target label of unseen inputs.

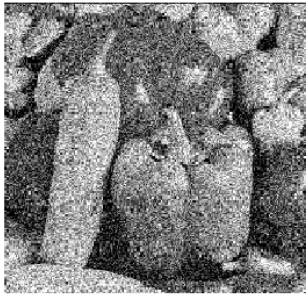**Ranzato**

# Supervised Learning: Examples

**Classification**



→ "dog"

*classification*

**Denoising**



*regression*

**OCR**



→ "2 3 4 5"

*structured prediction*

3

**Ranzato**

# Supervised Deep Learning

**Classification**



→ "dog"

**Denoising**



**OCR**



→ "2 3 4 5"

4

# Outline

- Supervised Neural Networks

- Convolutional Neural Networks

- Examples

- Tips

**Ranzato**

# Neural Networks

Assumptions (for the next few slides):
- The input image is vectorized (disregard the spatial layout of pixels)
- The target label is discrete (classification)

**Question:** what class of functions shall we consider to map the input into the output?

**Answer:** composition of simpler functions.

**Follow-up questions:** Why not a linear combination? What are the "simpler" functions? What is the interpretation?

**Answer:** later...

6

**Ranzato**

# Neural Networks: example

$$x \rightarrow \boxed{max(0, W^1 x)} \xrightarrow{h^1} \boxed{max(0, W^2 h^1)} \xrightarrow{h^2} \boxed{W^3 h^2} \xrightarrow{o}$$

$x$   input

$h^1$   1-st layer hidden units

$h^2$   2-nd layer hidden units

$o$   output

Example of a 2 hidden layer neural network (or 4 layer network, counting also input and output).

**Ranzato**

# Forward Propagation

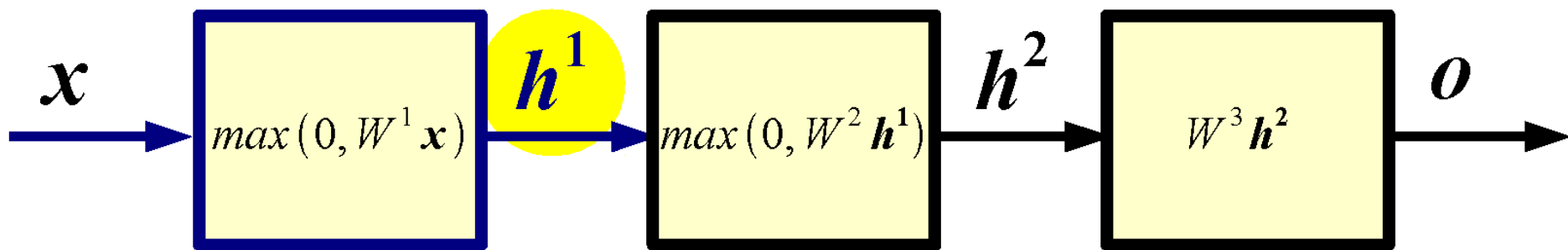**Def.:** Forward propagation is the process of computing the output of the network given its input.

**Ranzato**

# Forward Propagation

$$x \quad \boxed{max\,(0, W^1 x)} \quad h^1 \quad \boxed{max\,(0, W^2 h^1)} \quad h^2 \quad \boxed{W^3 h^2} \quad o$$

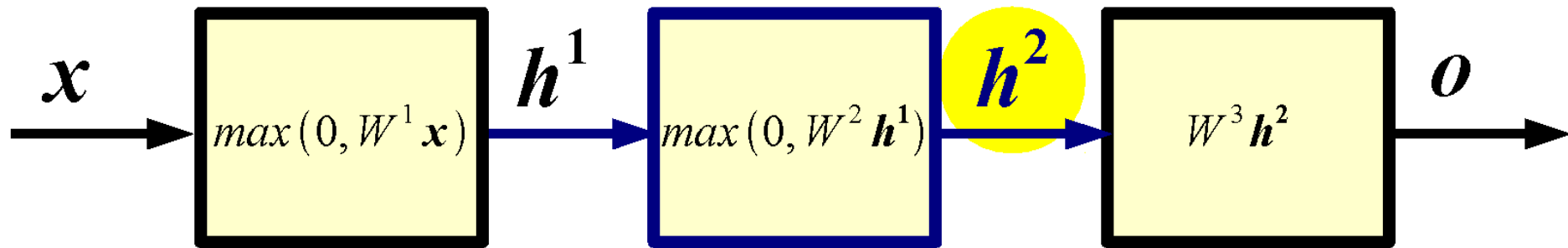$$x \in R^D \qquad W^1 \in R^{N_1 \times D} \qquad b^1 \in R^{N_1} \qquad h^1 \in R^{N_1}$$

$$h^1 = max\,(0, W^1 x + b^1)$$

$W^1$   1-st layer weight matrix or weights

$b^1$   1-st layer biases

The non-linearity $u = max\,(0, v)$ is called **ReLU** in the DL literature. Each output hidden unit takes as input all the units at the previous layer: each such layer is called "**fully connected**".

**Ranzato**

# Forward Propagation

$$x \quad \boxed{max(0, W^1 x)} \quad h^1 \quad \boxed{max(0, W^2 h^1)} \quad h^2 \quad \boxed{W^3 h^2} \quad o$$
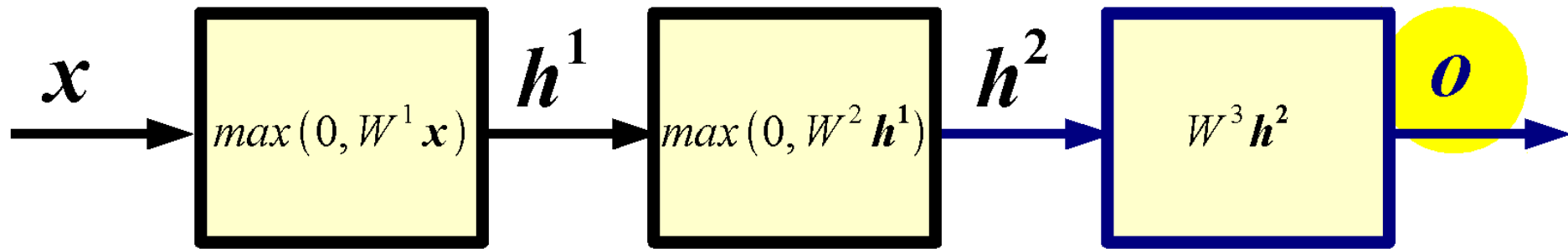
$$h^1 \in R^{N_1} \quad W^2 \in R^{N_2 \times N_1} \quad b^2 \in R^{N_2} \qquad h^2 \in R^{N_2}$$

$$h^2 = max(0, W^2 h^1 + b^2)$$

$W^2$  2-nd layer weight matrix or weights

$b^2$  2-nd layer biases

# Forward Propagation



$$h^2 \in R^{N_2} \quad W^3 \in R^{N_3 \times N_2} \quad b^3 \in R^{N_3} \qquad o \in R^{N_3}$$

$$o = max\left(0, W^3 h^2 + b^3\right)$$

$W^3$  3-rd layer weight matrix or weights

$b^3$  3-rd layer biases

**Ranzato**

# Alternative Graphical Representation

$h^k$    $max(0, W^{k+1} h^k)$    $h^{k+1}$

$h^k$    $W^{k+1}$    $h^{k+1}$

$h^k$    $W^{k+1}$    $h^{k+1}$

$h^k_1$   $w^{k+1}_{1,1}$   $h^{k+1}_1$

$h^k_2$    $h^{k+1}_2$

$h^k_3$    $h^{k+1}_3$

$h^k_4$   $w^{k+1}_{3,4}$

# Interpretation

**Question:** Why can't the mapping between layers be linear?

**Answer:** Because composition of linear functions is a linear function. Neural network would reduce to (1 layer) logistic regression.

**Question:** What do ReLU layers accomplish?

**Answer:** Piece-wise linear tiling: mapping is locally linear.

**Ranzato**

# Interpretation

**Question:** Why do we need many layers?

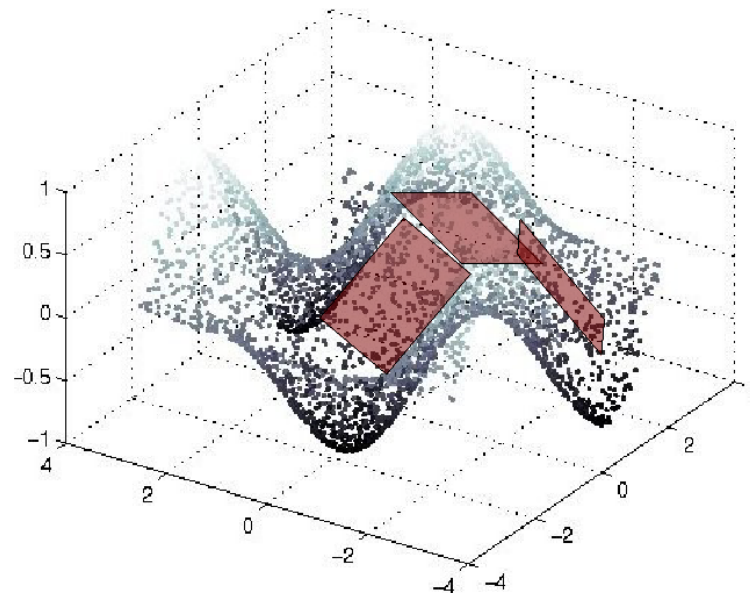**Answer:** When input has hierarchical structure, the use of a hierarchical architecture is potentially more efficient because intermediate computations can be re-used. DL architectures are efficient also because they use **distributed representations** which are shared across classes.

[0  0  **1**  0  0  0  0  **1**  0  0  **1**  **1**  0  0  **1**  0 ... ]  truck feature



Exponentially more efficient than a 1-of-N representation (a la k-means)

14

**Ranzato**

# Interpretation

[ 1 1 0 0 0 1 0 **1** 0 0 0 0 1 1 0 1… ] motorbike

[ 0 0 1 0 0 0 0 **1** 0 0 1 1 0 0 1 0 … ] truck

**Ranzato**

# Interpretation



prediction of class

high-level parts

mid-level parts

low level parts

Input image

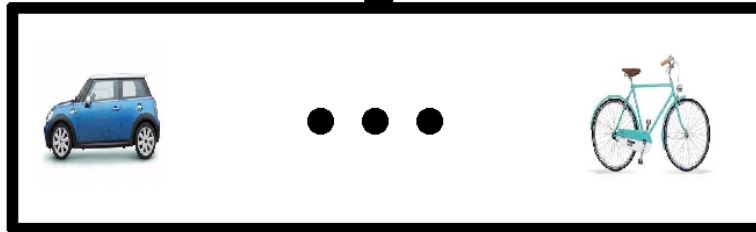- distributed representations
- feature sharing
- compositionality

Lee et al. "Convolutional DBN's ..." ICML 2009

Ranzato

# Interpretation

**Question:** What does a hidden unit do?

**Answer:** It can be thought of as a classifier or feature detector.

**Question:** How many layers? How many hidden units?

**Answer:** Cross-validation or hyper-parameter search methods are the answer. In general, the wider and the deeper the network the more complicated the mapping.

**Question:** How do I set the weight matrices?

**Answer:** Weight matrices and biases are learned.
First, we need to define a measure of quality of the current mapping. Then, we need to define a procedure to adjust the parameters.

17

**Ranzato**

# How Good is a Network?



Probability of class k given input (softmax):

$$p(c_k = 1 | \boldsymbol{x}) = \frac{e^{o_k}}{\sum_{j=1}^{C} e^{o_j}}$$

(Per-sample) **Loss**; e.g., negative log-likelihood (good for classification of small number of classes):

$$L(\boldsymbol{x}, y; \boldsymbol{\theta}) = -\sum_j y_j \log p(c_j | \boldsymbol{x})$$

**Ranzato**

# Training

**Learning** consists of minimizing the loss (plus some regularization term) w.r.t. parameters over the whole training set.

$$\boldsymbol{\theta}^* = arg\ min_{\boldsymbol{\theta}} \sum_{n=1}^{P} L(\boldsymbol{x}^n, y^n; \boldsymbol{\theta})$$

**Question:** How to minimize a complicated function of the parameters?

**Answer:** Chain rule, a.k.a. **Backpropagation**! That is the procedure to compute gradients of the loss w.r.t. parameters in a multi-layer neural network.

Rumelhart et al. "Learning internal representations by back-propagating.." Nature 1986

# Key Idea: Wiggle To Decrease Loss



$$x \rightarrow \boxed{max(0, W^1 x)} \xrightarrow{h^1} \boxed{max(0, W^2 h^1)} \xrightarrow{h^2} \boxed{W^3 h^2} \xrightarrow{o} Loss$$

$$y \longrightarrow Loss$$

Let's say we want to decrease the loss by adjusting $W^1_{i,j}$.
We could consider a very small $\epsilon = 1e\text{-}6$ and compute:

$$L(x, y; \boldsymbol{\theta})$$

$$L(x, y; \boldsymbol{\theta} \backslash W^1_{i,j}, W^1_{i,j} + \epsilon)$$

Then, update:

$$W^1_{i,j} \leftarrow W^1_{i,j} + \epsilon\, sgn(L(x, y; \boldsymbol{\theta}) - L(x, y; \boldsymbol{\theta} \backslash W^1_{i,j}, W^1_{i,j} + \epsilon))$$

**Ranzato**

# Derivative w.r.t. Input of Softmax

$$p(c_k=1|\boldsymbol{x})=\frac{e^{o_k}}{\sum_j e^{o_j}}$$

$$L(\boldsymbol{x},y;\boldsymbol{\theta})=-\sum_j y_j \log p(c_j|\boldsymbol{x}) \qquad \boldsymbol{y}=[\overset{1}{0}\,0\,..\,0\,\overset{k}{1}\,0\,..\,\overset{c}{0}]$$

By substituting the fist formula in the second, and taking the derivative w.r.t. $\boldsymbol{o}$ we get:

$$\frac{\partial L}{\partial o}=p(c|\boldsymbol{x})-\boldsymbol{y}$$

**Ranzato**

# Backward Propagation



Given $\partial L / \partial \boldsymbol{o}$ and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \boldsymbol{o}} \frac{\partial \boldsymbol{o}}{\partial W^3} \qquad \frac{\partial L}{\partial \boldsymbol{h}^2} = \frac{\partial L}{\partial \boldsymbol{o}} \frac{\partial \boldsymbol{o}}{\partial \boldsymbol{h}^2}$$

# Backward Propagation

$$x \rightarrow \boxed{max(0, W^1 \boldsymbol{x})} \xrightarrow{\boldsymbol{h}^1} \boxed{max(0, W^2 \boldsymbol{h}^1)} \xrightarrow{\boldsymbol{h}^2} \boxed{W^3 \boldsymbol{h}^2}$$

$$\frac{\partial L}{\partial \boldsymbol{o}}$$
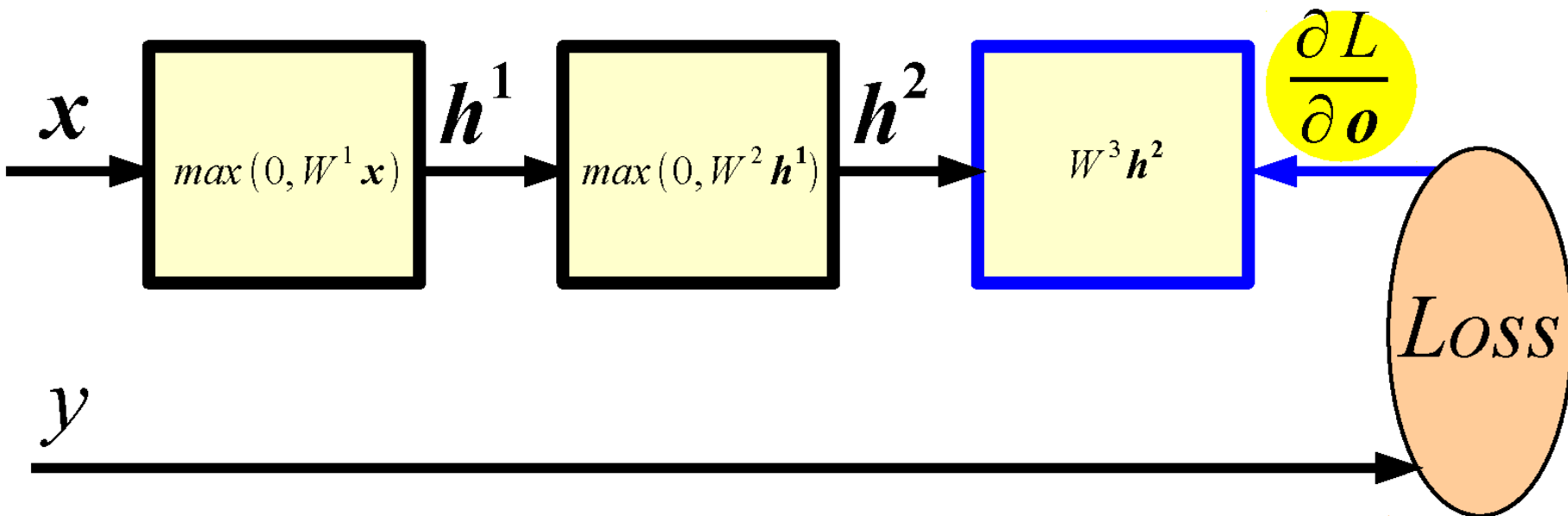
Loss

$y$

Given $\partial L / \partial \boldsymbol{o}$ and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \boldsymbol{o}} \frac{\partial \boldsymbol{o}}{\partial W^3} \qquad \frac{\partial L}{\partial \boldsymbol{h}^2} = \frac{\partial L}{\partial \boldsymbol{o}} \frac{\partial \boldsymbol{o}}{\partial \boldsymbol{h}^2}$$

$$\frac{\partial L}{\partial W^3} = (p(c|\boldsymbol{x}) - \boldsymbol{y}) \, \boldsymbol{h}^{2T} \qquad \frac{\partial L}{\partial \boldsymbol{h}^2} = W^{3T} (p(c|\boldsymbol{x}) - \boldsymbol{y})$$

23

# Backward Propagation



Given $\dfrac{\partial L}{\partial \boldsymbol{h}^2}$ we can compute now:

$$\frac{\partial L}{\partial W^2} = \frac{\partial L}{\partial \boldsymbol{h}^2}\, \frac{\partial \boldsymbol{h}^2}{\partial W^2} \qquad\qquad \frac{\partial L}{\partial \boldsymbol{h}^1} = \frac{\partial L}{\partial \boldsymbol{h}^2}\, \frac{\partial \boldsymbol{h}^2}{\partial \boldsymbol{h}^1}$$

**Ranzato** f

# Backward Propagation



Given $\dfrac{\partial L}{\partial \boldsymbol{h}^1}$ we can compute now:

$$\frac{\partial L}{\partial W^1} = \frac{\partial L}{\partial \boldsymbol{h}^1} \frac{\partial \boldsymbol{h}^1}{\partial W^1}$$

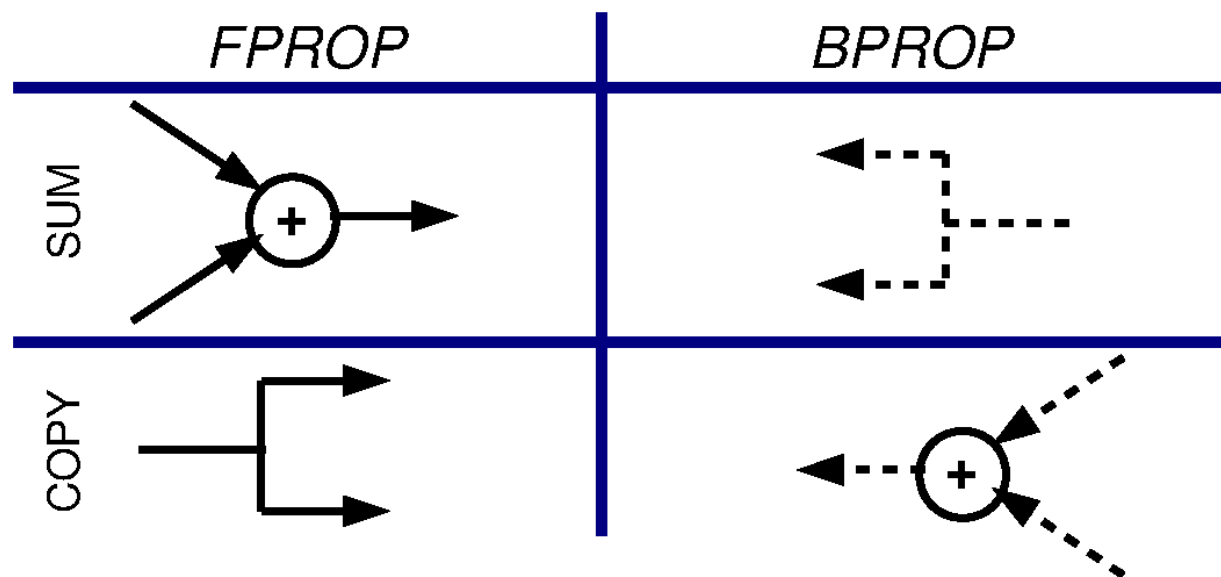**Ranzato** f

# Backward Propagation

**Question:** Does BPROP work with ReLU layers only?

**Answer:** Nope, any a.e. differentiable transformation works.

**Question:** What's the computational cost of BPROP?

**Answer:** About twice FPROP (need to compute gradients w.r.t. input and parameters at every layer).

**Note:** FPROP and BPROP are dual of each other. E.g.,:

| | *FPROP* | *BPROP* |
|---|---|---|
| SUM | | |
| COPY | | |



26

**Ranzato**

# Optimization

**Stochastic Gradient Descent** **(on mini-batches):**

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \frac{\partial L}{\partial \boldsymbol{\theta}} , \eta \in (0, 1)$$

**Stochastic Gradient Descent with Momentum:**

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \boldsymbol{\Delta}$$

$$\boldsymbol{\Delta} \leftarrow 0.9 \boldsymbol{\Delta} + \frac{\partial L}{\partial \boldsymbol{\theta}}$$

**Note: there are many other variants...**

27

**Ranzato**

# Outline

- Supervised Neural Networks

- **Convolutional Neural Networks**

- Examples

- Tips

**Ranzato**