

Deep Learning 1

Neural Net Basics

Many slides by James Heys,
Marc'Aurelio Ranzato

1

Outline

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

2

1

Supervised Learning

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

x^i i-th input training example

y^i i-th target label

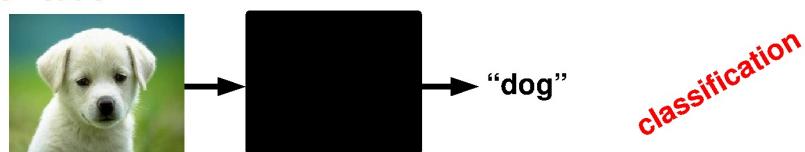
P number of training examples



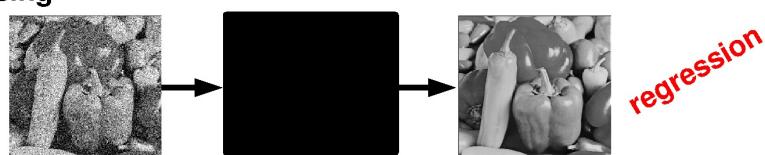
Goal: predict the target label of unseen inputs.

Supervised Learning: Examples

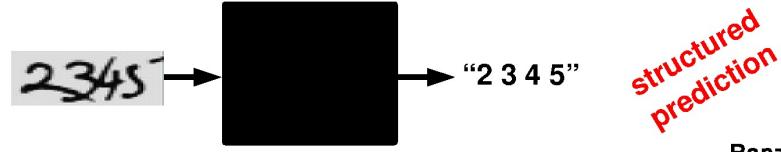
Classification

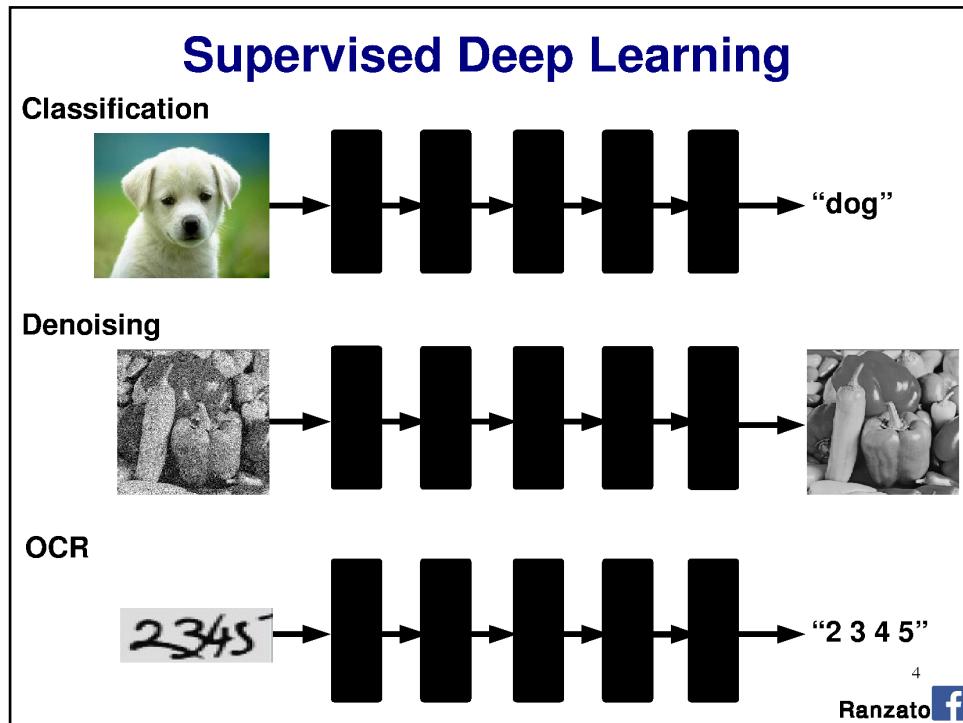


Denoising

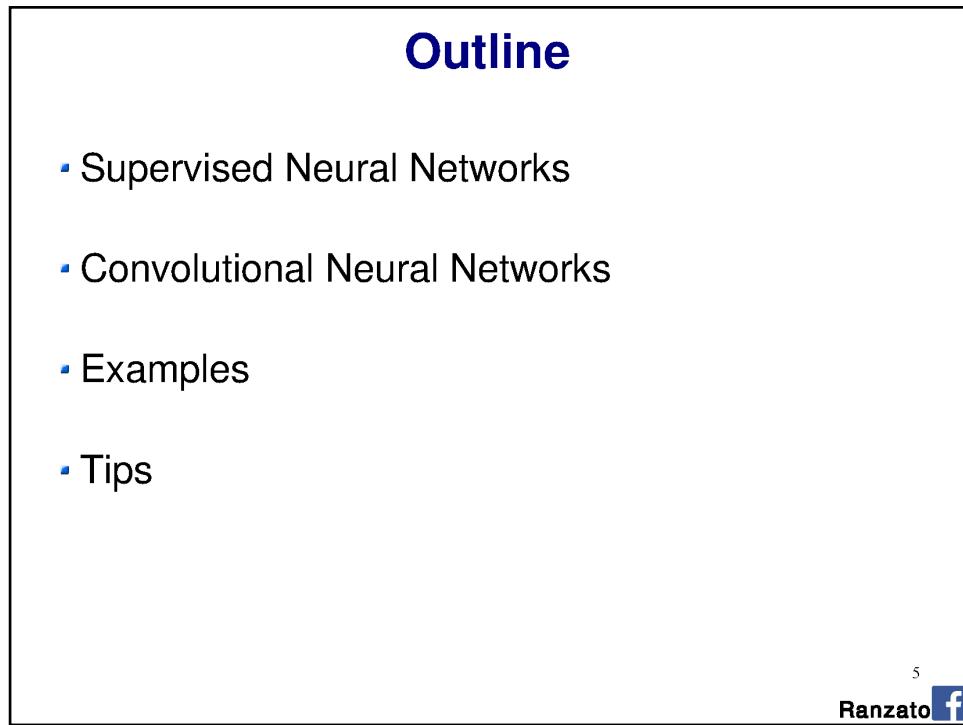


OCR





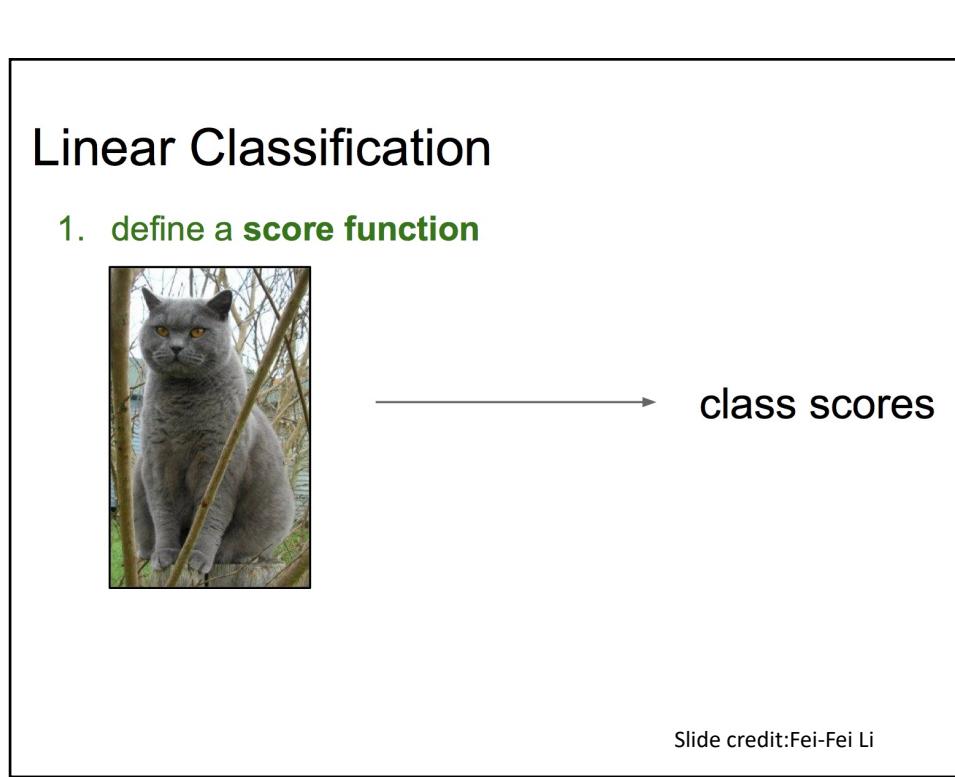
5



6



7



8

Linear Classification

1. define a **score function**

$f(x_i, W, b) = Wx_i + b$

data (image)

“weights”

“bias vector”

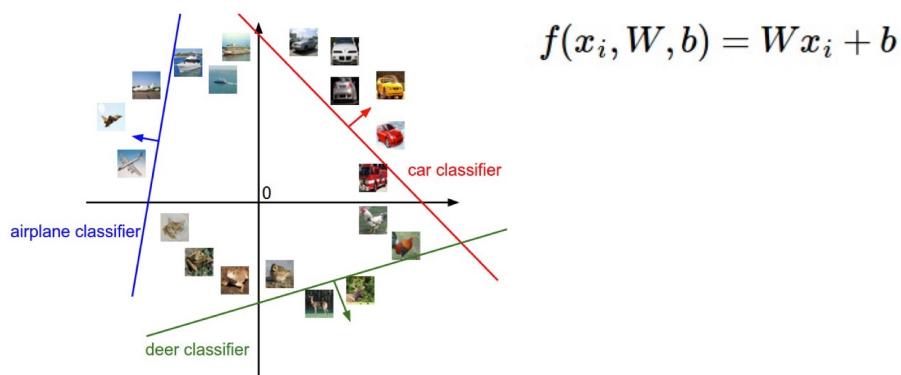
“parameters”

class scores

Slide credit:Fei-Fei Li

9

Interpreting a Linear Classifier



Slide credit:Fei-Fei Li

10

LOSS

2. Define a loss function (or cost function, or objective)

- scores, label \longrightarrow loss.
 $f(x_i, W)$ y_i L_i

Example:

$$f(x_i, W) = [13, -7, 11]$$

$$y_i = 0 \quad \text{_____}^{\uparrow}$$

Question: if you were to assign a single number to how “unhappy” you are with these scores, what would you do?

Slide credit:Fei-Fei Li

11

Multiclass SVM Loss

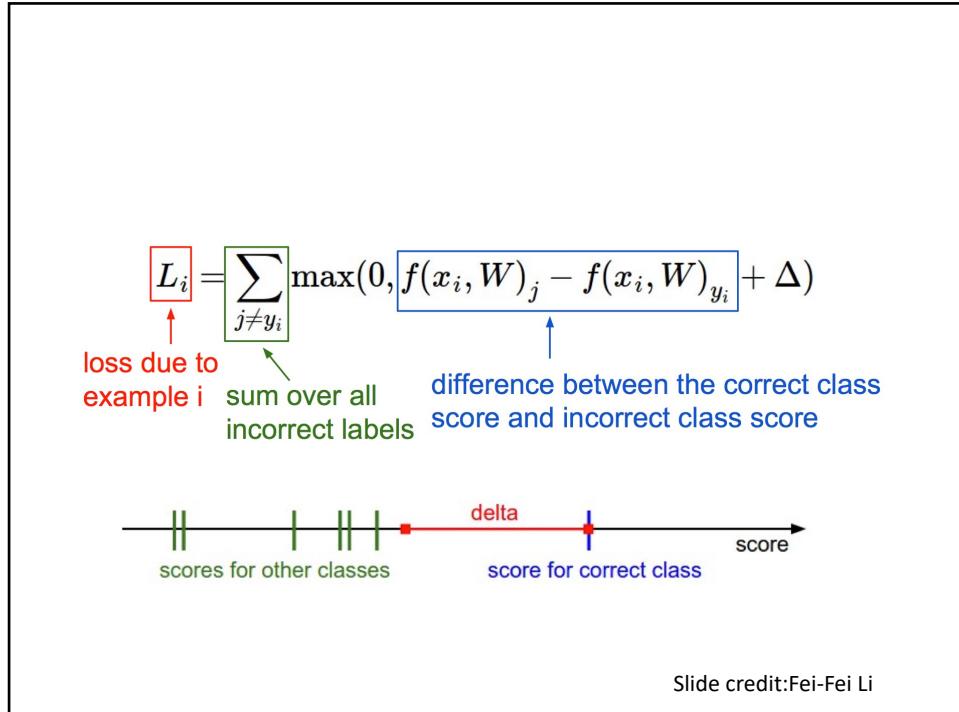
2. Define a loss function (or cost function, or objective)
One (of many ways) to do it: **Multiclass SVM Loss**

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

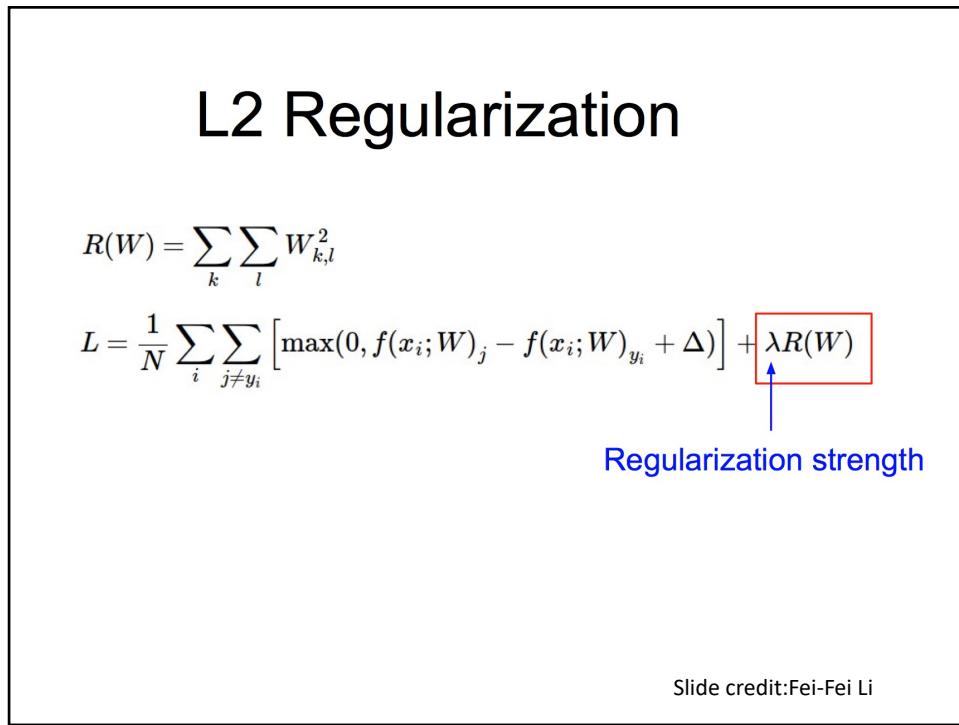
↑ loss due to example i
 ↑ sum over all incorrect labels
 ↑ difference between the correct class score and incorrect class score

Slide credit:Fei-Fei Li

12



13



14

Putting it all together:

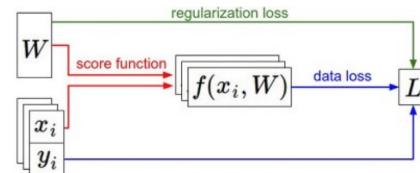
Linear Classification

SVM:

$$L = \frac{1}{N} \sum_i \sum_{j \neq y_i} [\max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + \Delta)] + \lambda \sum_k \sum_l W_{k,l}^2$$

Softmax:

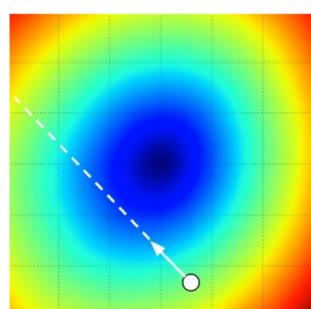
$$L = \frac{1}{N} \sum_i -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) + \lambda \sum_k \sum_l W_{k,l}^2$$



Slide credit:Fei-Fei Li

15

Optimization Landscape



Slide credit:Fei-Fei Li

16

Gradient Descent

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

Numerical gradient: slow :, approximate :, easy to write :)

Analytic gradient: fast :, exact :, error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient

Slide credit:Fei-Fei Li

17

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?

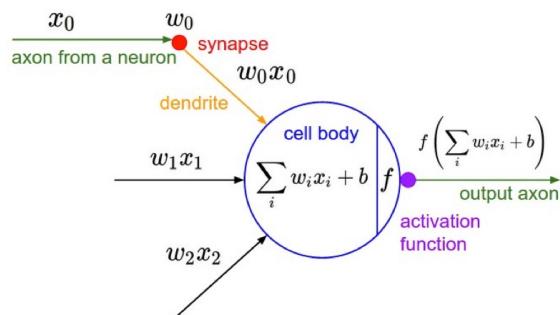
Answer: later...

6

Ranzato 

18

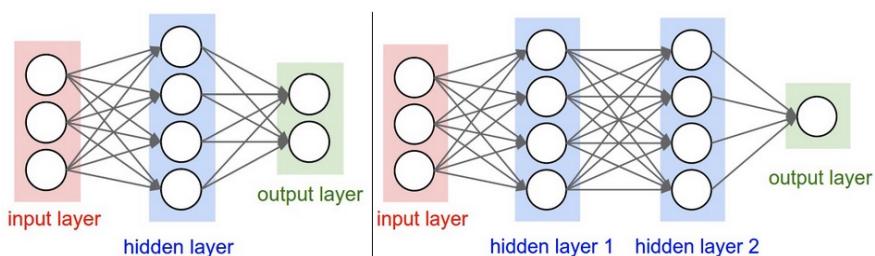
A Single Neuron Activation Functions



Slide credit: Fei-Fei Li

19

Neural Network Structure



Left: A 2-layer Neural Network (one hidden layer of 4 neurons (or units) and one output layer with 2 neurons), and three inputs.
 Right: A 3-layer neural network with three inputs, two hidden layers of 4 neurons each and one output layer. Notice that in both cases there are connections (synapses) between neurons across layers, but not within a layer.

20

Activation functions

Step

$$g(a) = \begin{cases} 0 & a < 0 \\ 1 & a \geq 0 \end{cases}$$

Sigmoidal

$$g(a) = \frac{1}{1 + \exp(-a)}$$

Hyperbolic tangent

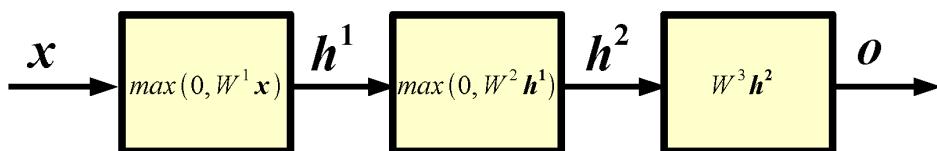
$$g(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$$

Rectified Linear Unit (ReLU)

$$g(a) = \max(0, a)$$

21

Neural Networks: example



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).

7

22

Forward Propagation

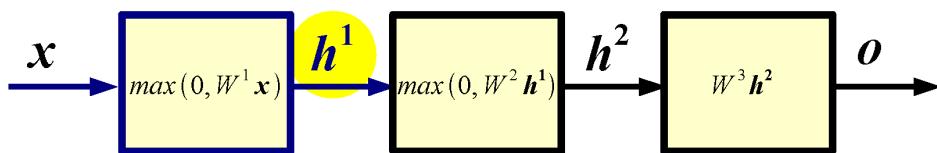
Def.: Forward propagation is the process of computing the output of the network given its input.

8

Ranzato 

23

Forward Propagation



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

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

W^1 1-st layer weight matrix or weights

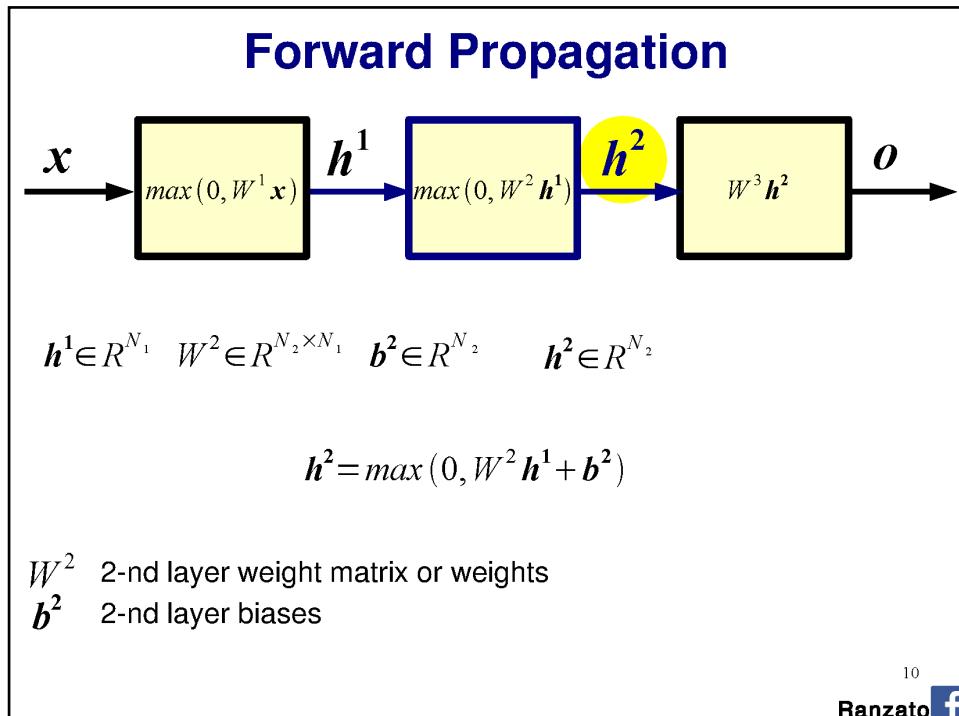
\mathbf{b}^1 1-st layer biases

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**”.

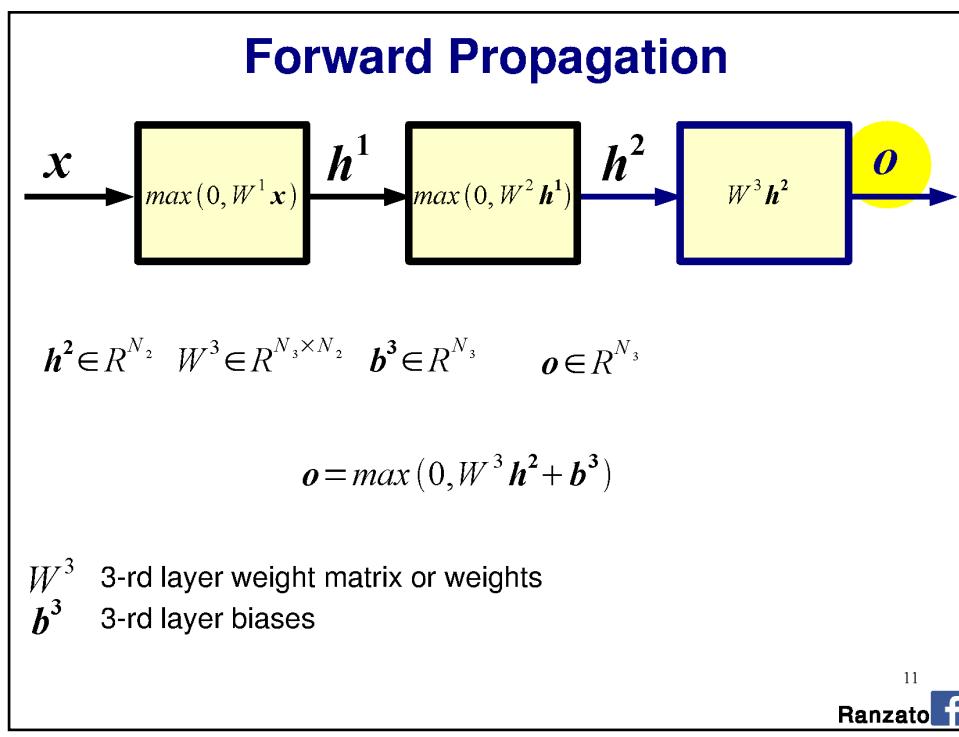
9

Ranzato 

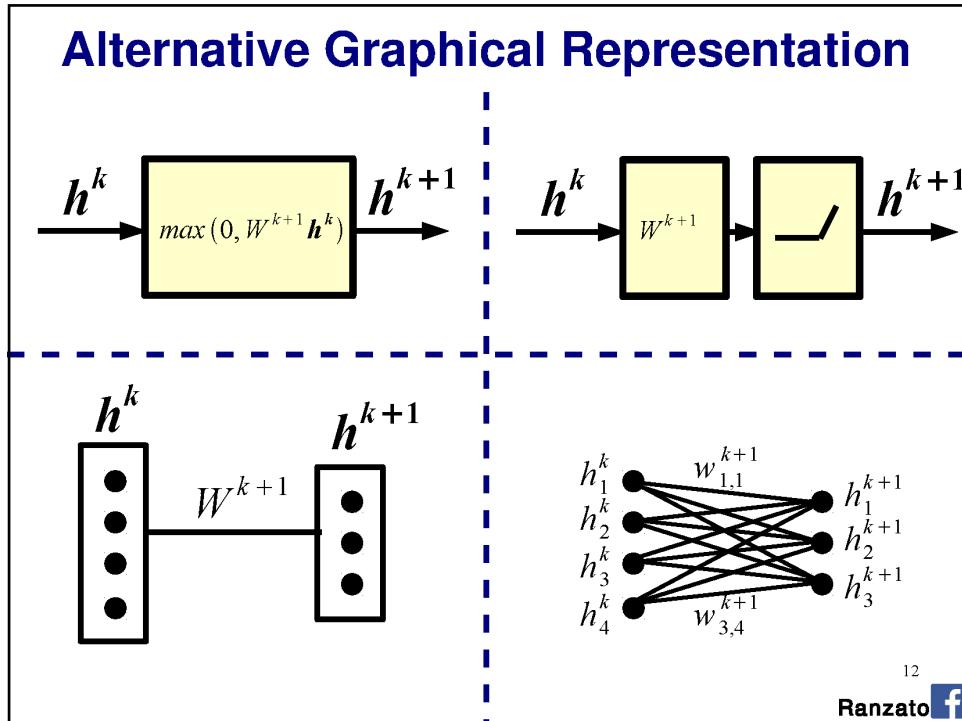
24



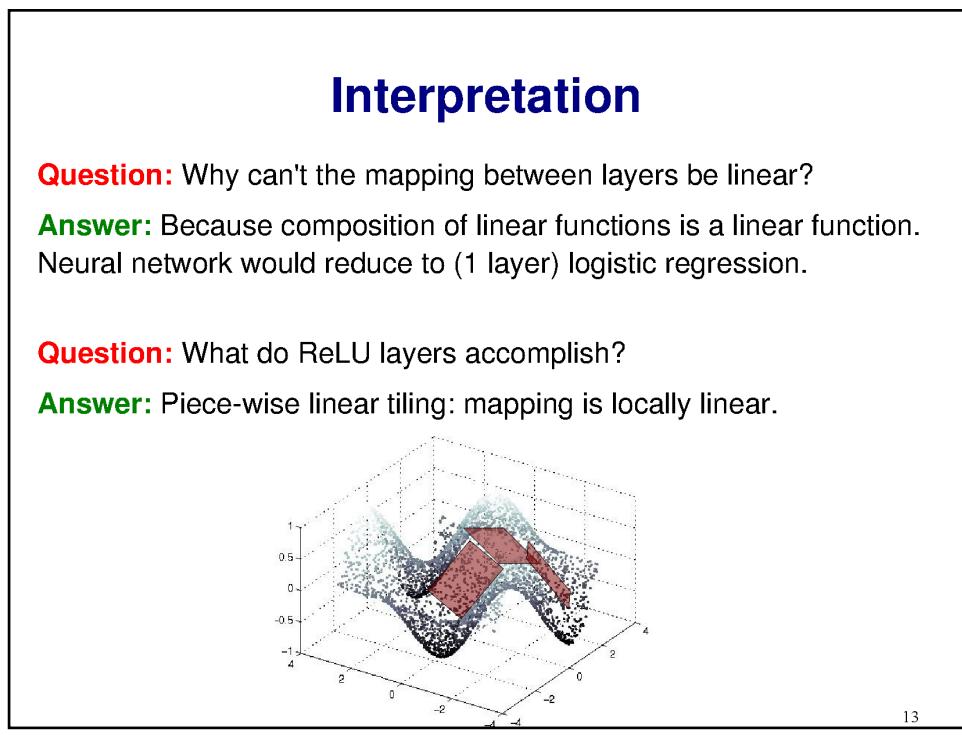
25



26



27



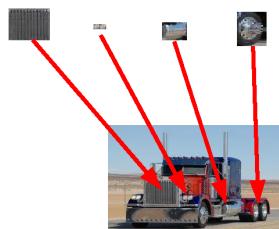
28

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)

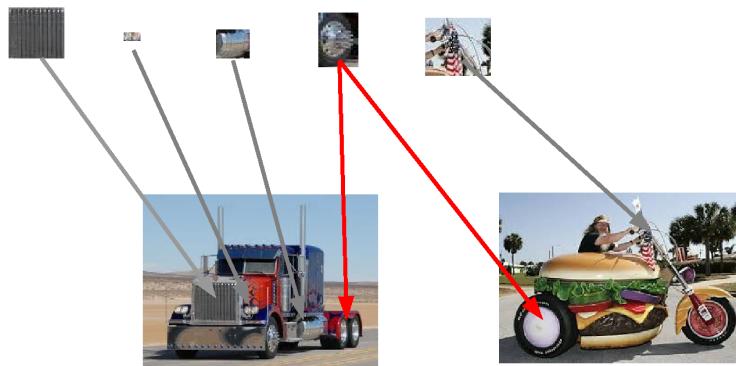
14
Ranzato

29

Interpretation

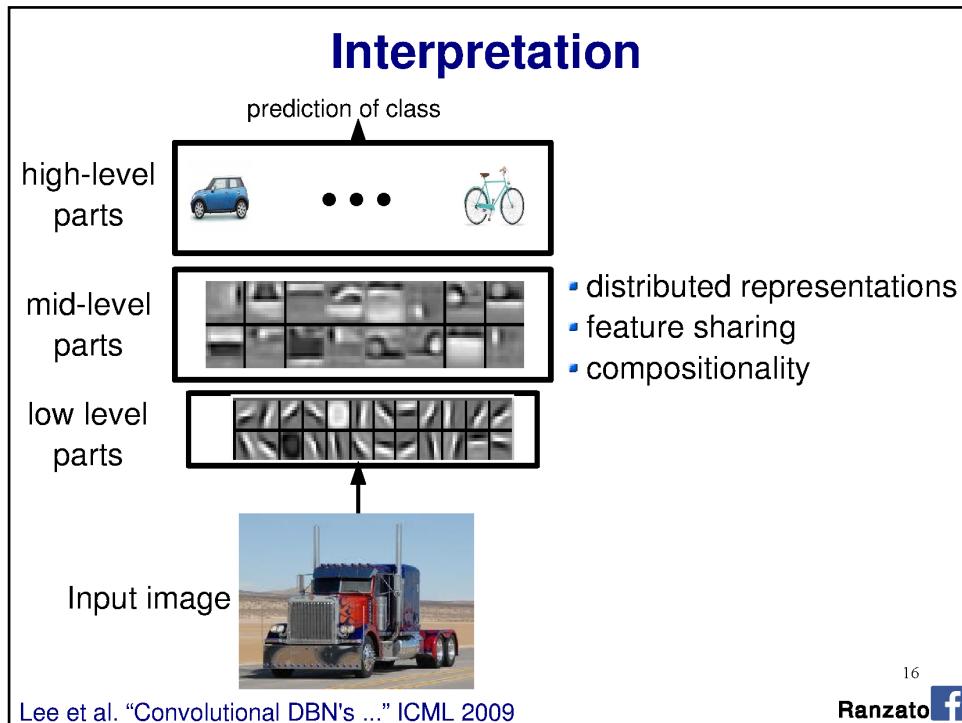
[1 1 0 0 0 1 0 1 0 0 0 0 1 1 0 1 ...] motorbike

[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ...] truck

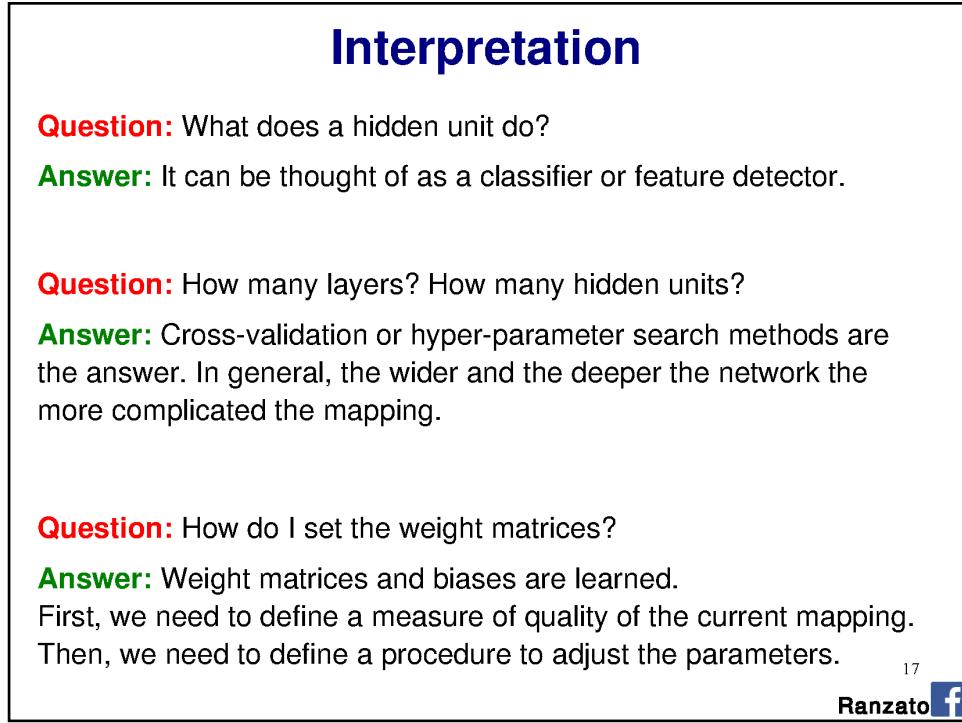


15
Ranzato

30



31



32

How Good is a Network?

$y = [0^1 0..0 1^k 0..0^c]$

Probability of class k given input (softmax):

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

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

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

18
Ranzato

33

Neural network properties

- **Theorem (Universal function approximators):** A two-layer network with a sufficient number of neurons can approximate any continuous function to any desired accuracy
- **Practical considerations:**
 - Can be seen as learning the features
 - Large number of neurons
 - Danger for overfitting
 - Hill-climbing procedure can get stuck in bad local optima

[Approximation by Superpositions of Sigmoidal Function](#), 1989

Slide credit: Pieter Abeel and Dan Klein

34

Training

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

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

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

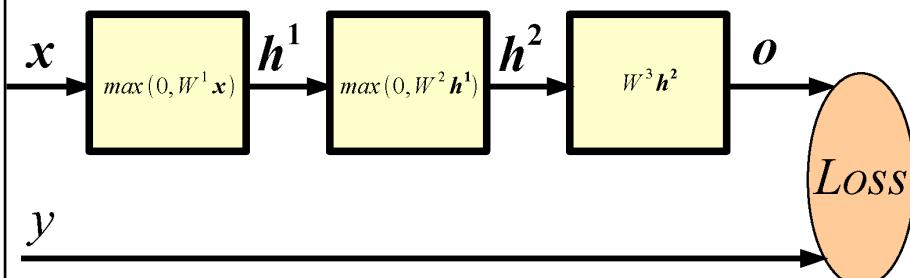
Answer: Chain rule, a.k.a. **Backpropagation**! That is the procedure to compute gradients of the loss w.r.t. parameters in a multi-layer neural network.

19

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

35

Key Idea: Wiggle To Decrease Loss



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

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

$$L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon)$$

Then, update:

$$W_{i,j}^1 \leftarrow W_{i,j}^1 + \epsilon \operatorname{sgn}(L(\mathbf{x}, y; \boldsymbol{\theta}) - L(\mathbf{x}, y; \boldsymbol{\theta} \setminus W_{i,j}^1, W_{i,j}^1 + \epsilon))$$

Ranzato f

20

36

Derivative w.r.t. Input of Softmax

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

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

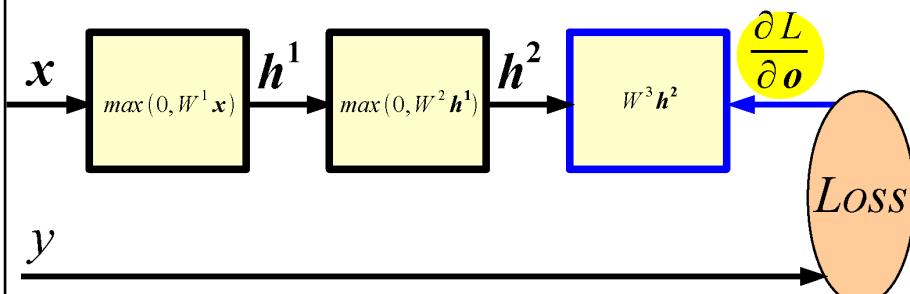
By substituting the first formula in the second, and taking the derivative w.r.t. θ we get:

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

21

37

Backward Propagation

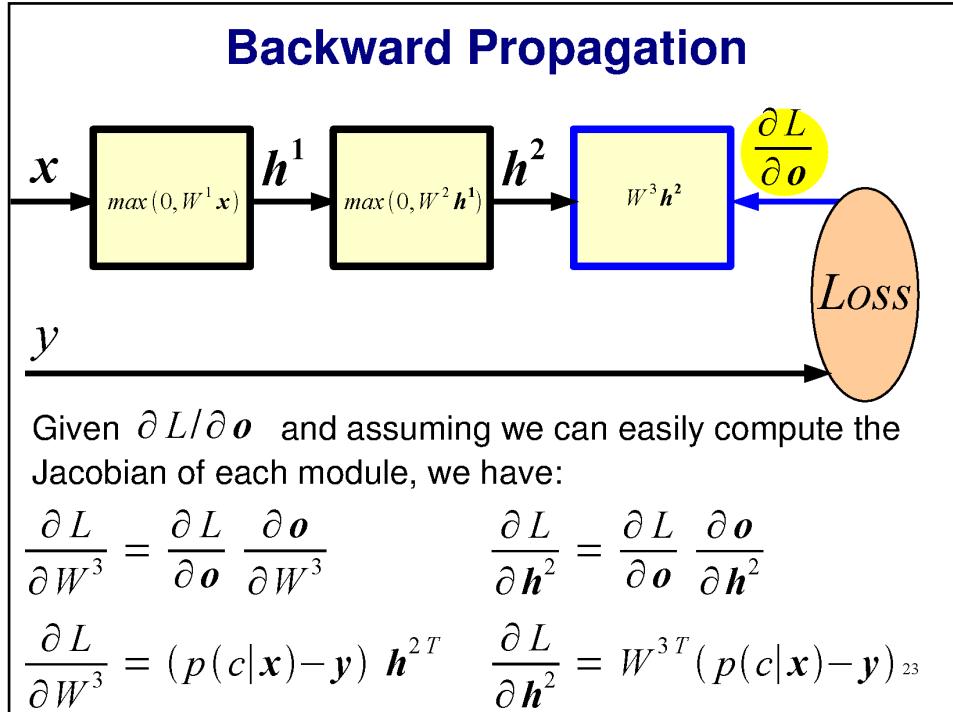


Given $\partial L / \partial \theta$ and assuming we can easily compute the Jacobian of each module, we have:

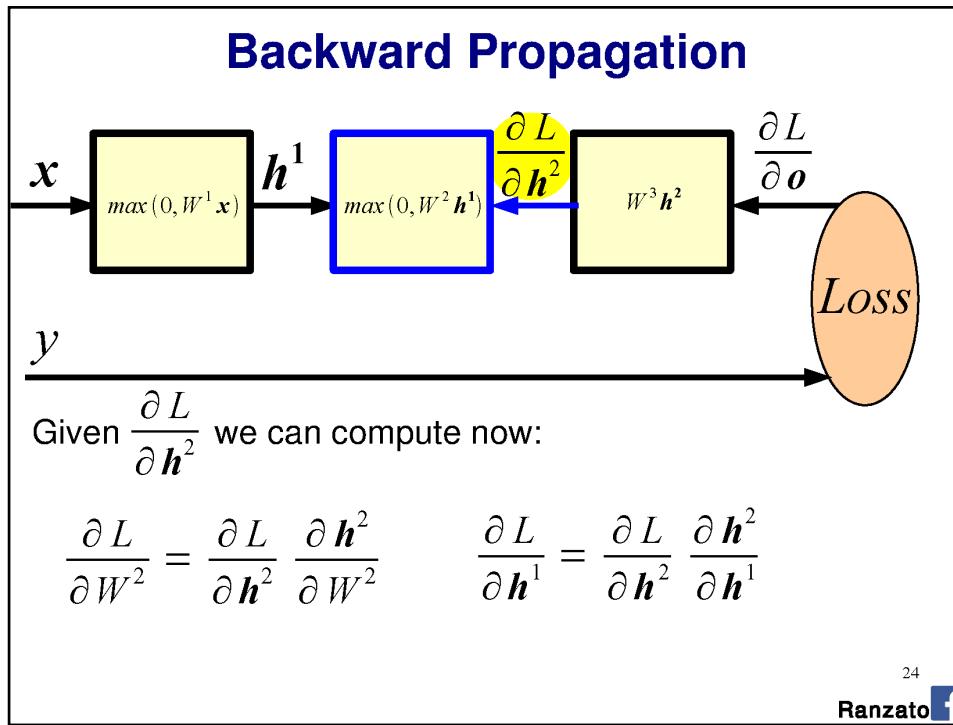
$$\frac{\partial L}{\partial W^3} = \frac{\partial L}{\partial \theta} \frac{\partial \theta}{\partial W^3} \quad \frac{\partial L}{\partial h^2} = \frac{\partial L}{\partial \theta} \frac{\partial \theta}{\partial h^2}$$

22

38



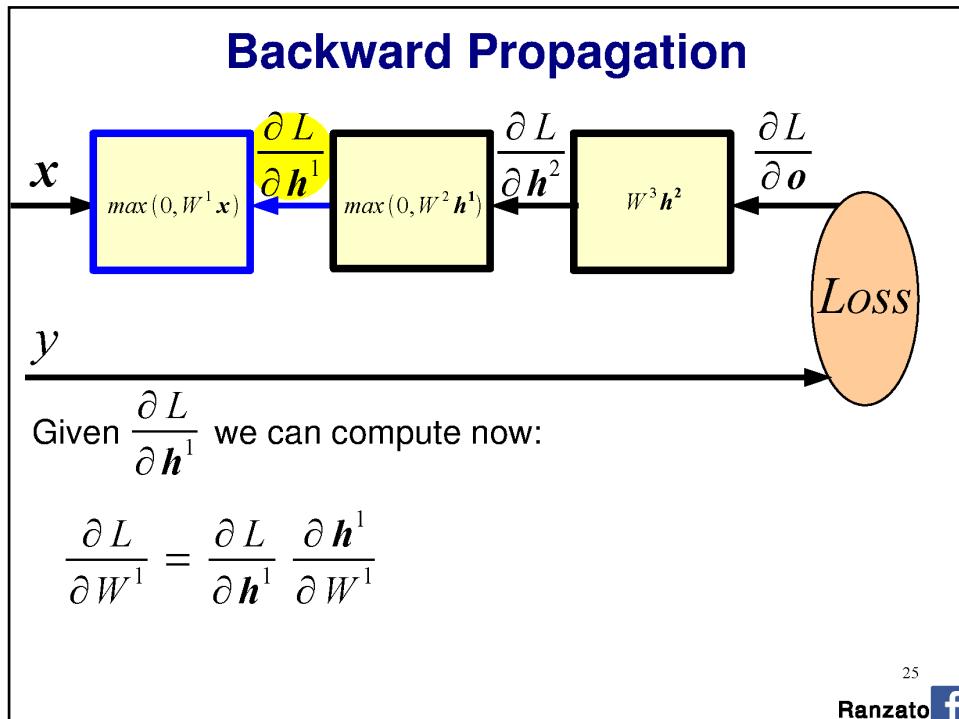
39



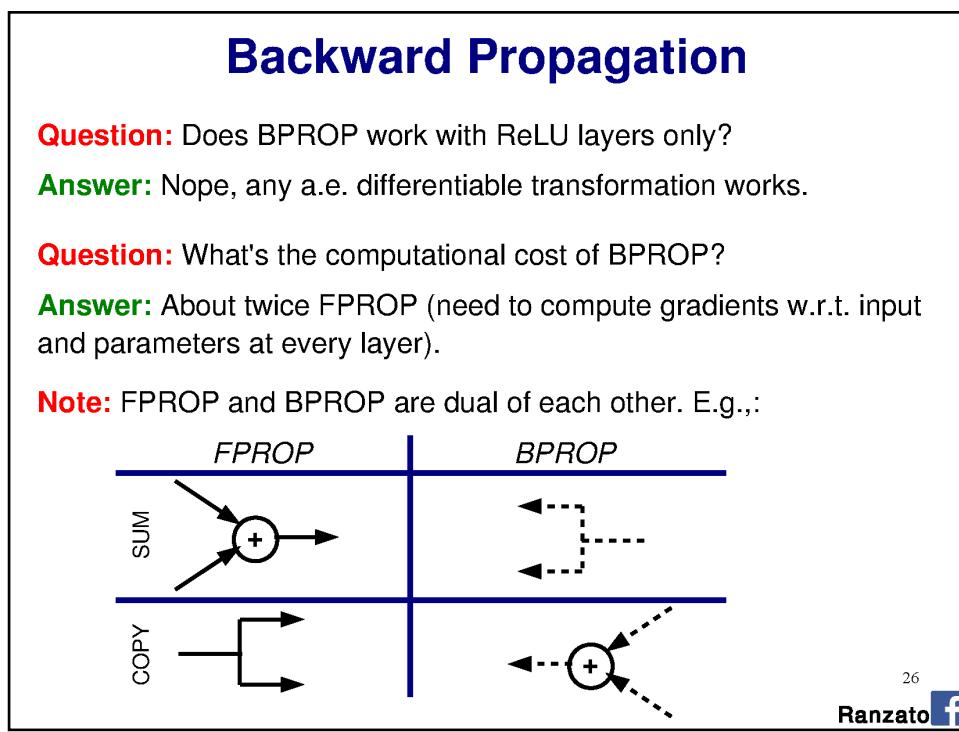
40

24

Ranzato



41



42

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:

$$\begin{aligned}\boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} - \eta \Delta \\ \Delta &\leftarrow 0.9 \Delta + \frac{\partial L}{\partial \boldsymbol{\theta}}\end{aligned}$$

Note: there are many other variants...

27



43

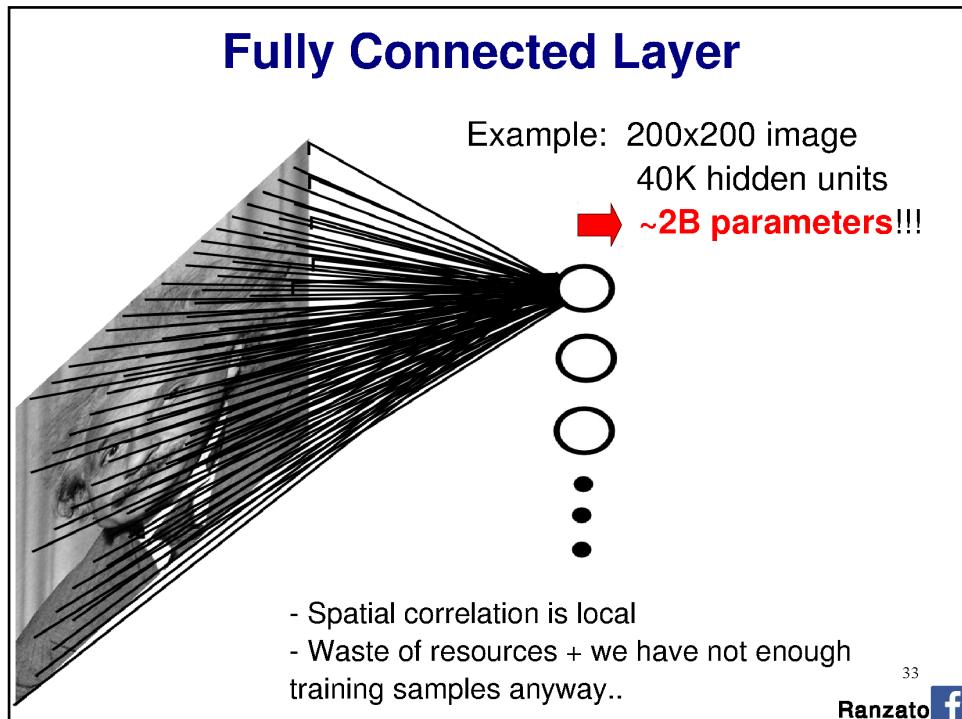
Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

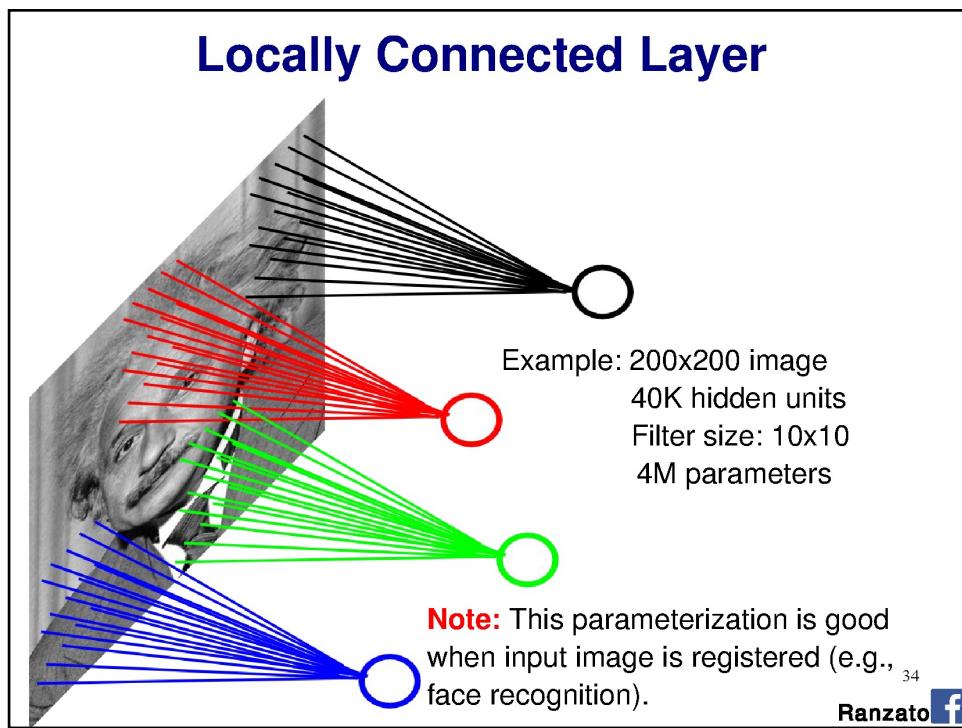
32



44

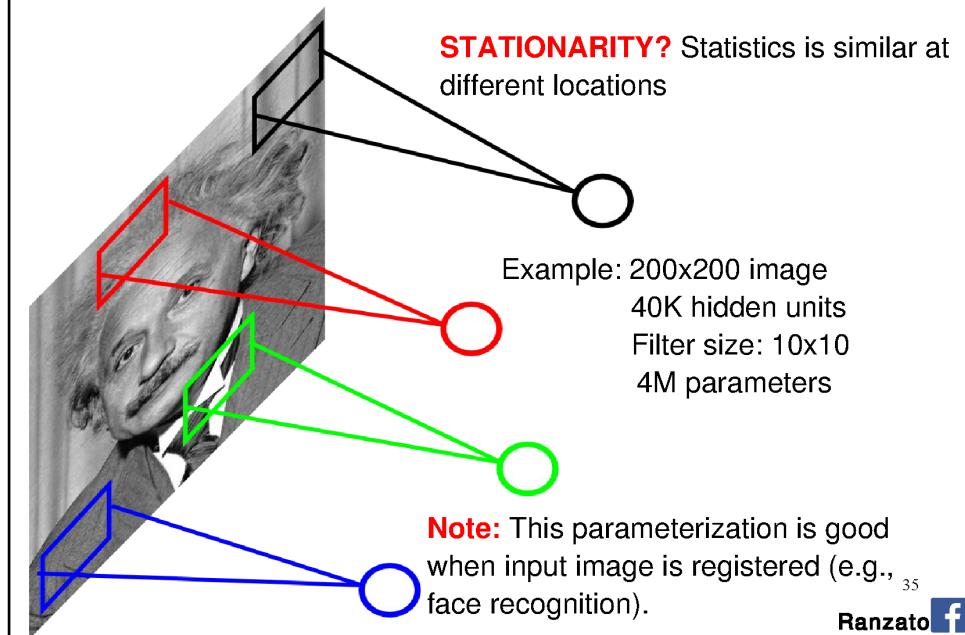


45



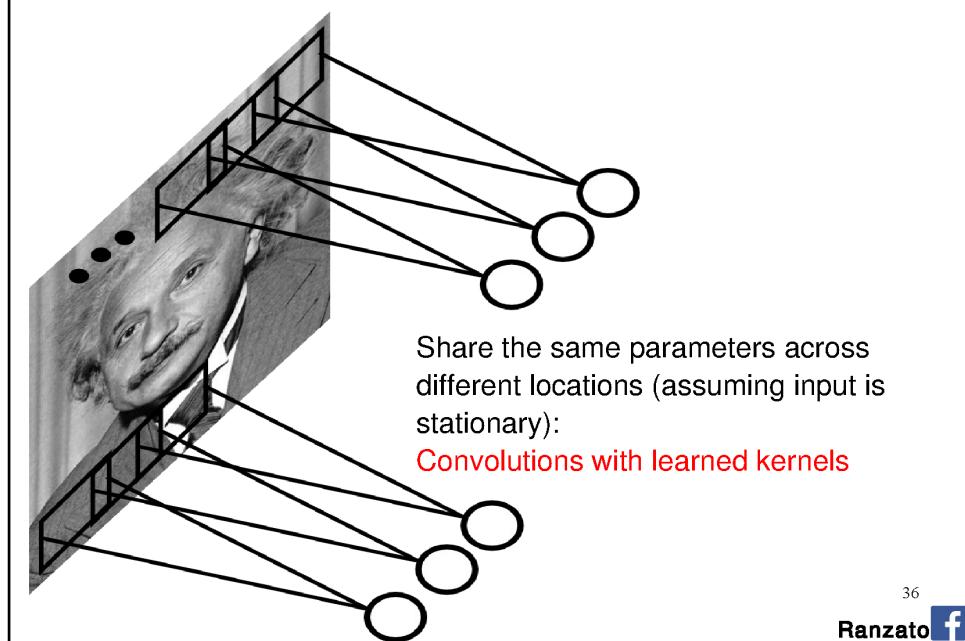
46

Locally Connected Layer



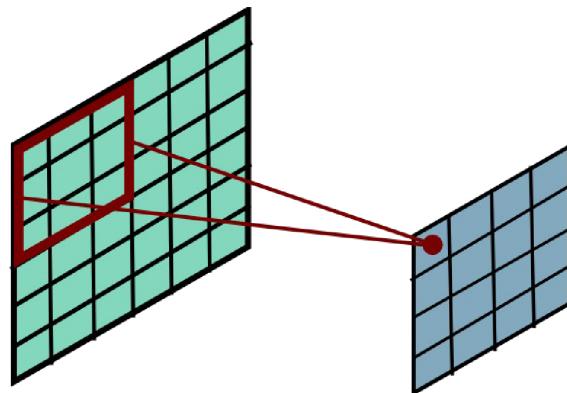
47

Convolutional Layer



48

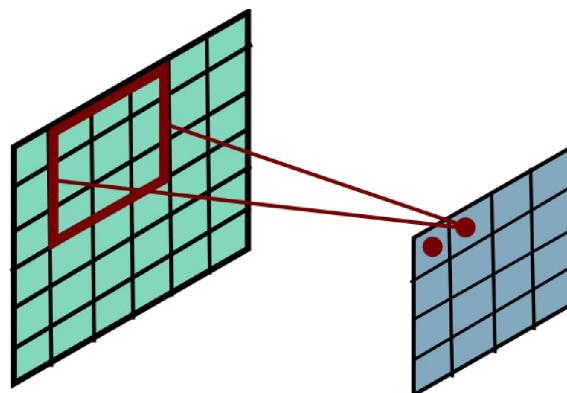
Convolutional Layer



Ranzato

49

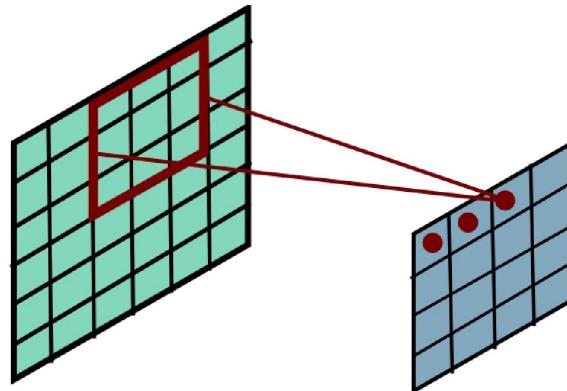
Convolutional Layer



Ranzato

50

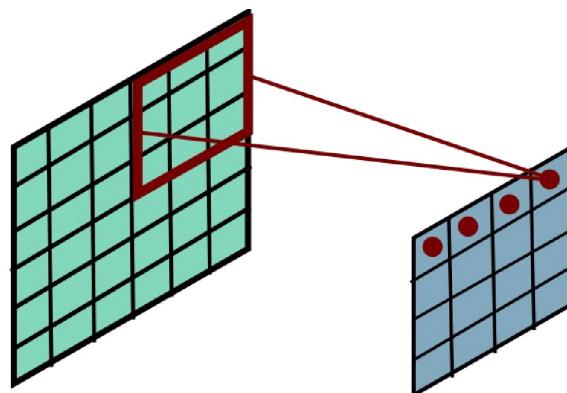
Convolutional Layer



Ranzato

51

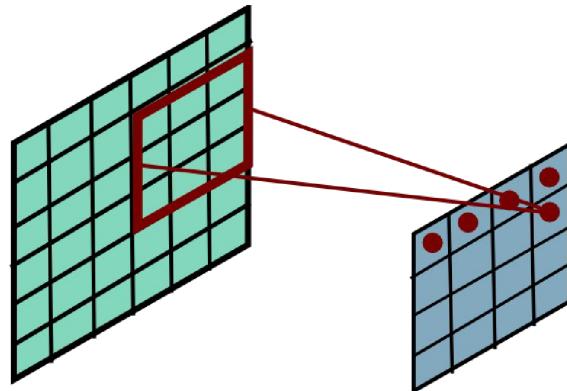
Convolutional Layer



Ranzato

52

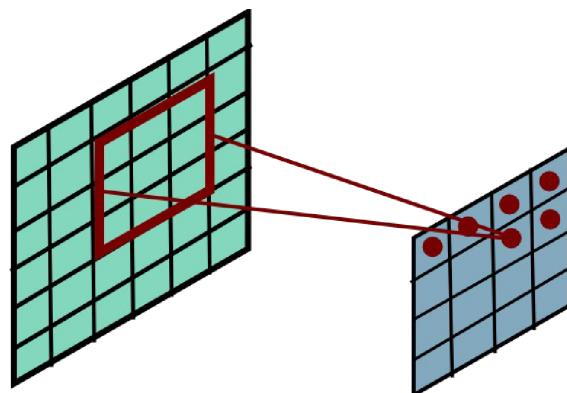
Convolutional Layer



Ranzato

53

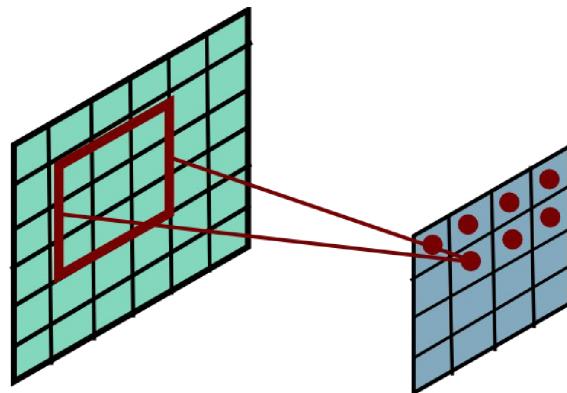
Convolutional Layer



Ranzato

54

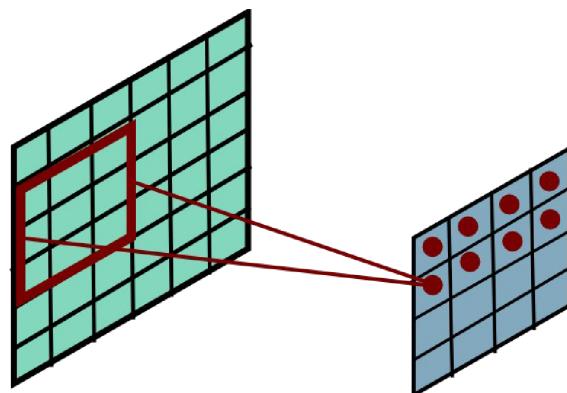
Convolutional Layer



Ranzato

55

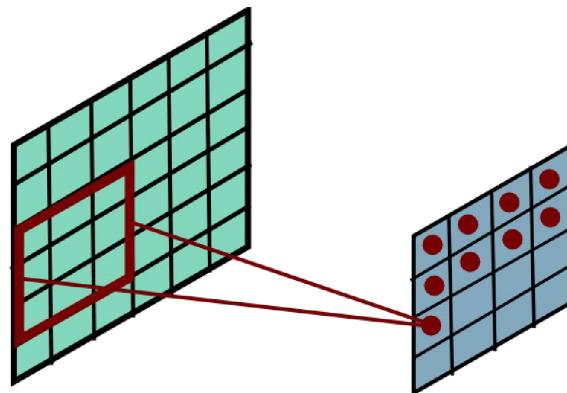
Convolutional Layer



Ranzato

56

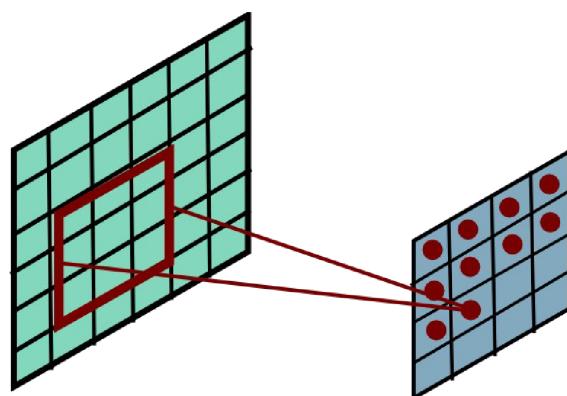
Convolutional Layer



Ranzato

57

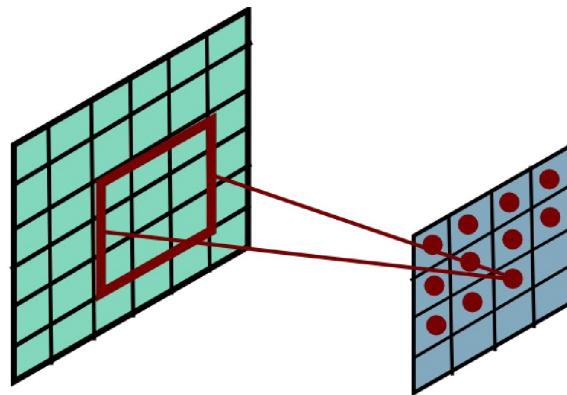
Convolutional Layer



Ranzato

58

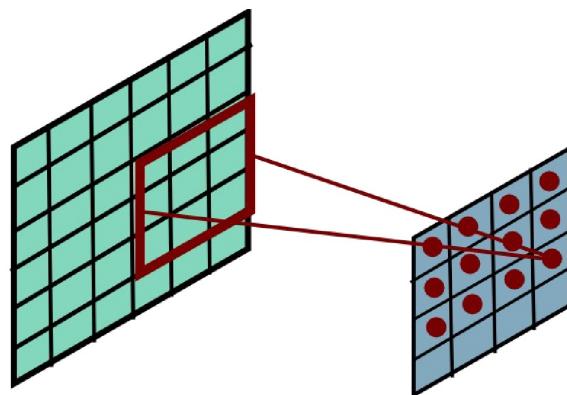
Convolutional Layer



Ranzato

59

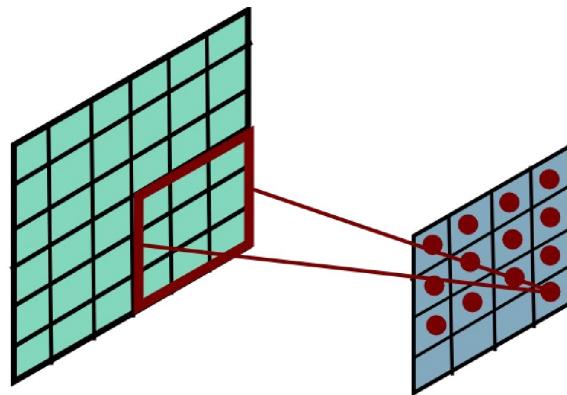
Convolutional Layer



Ranzato

60

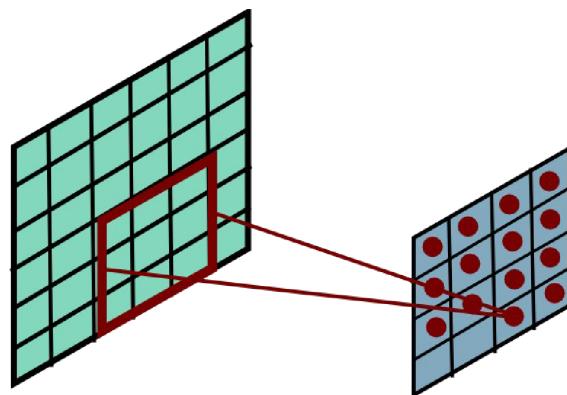
Convolutional Layer



Ranzato

61

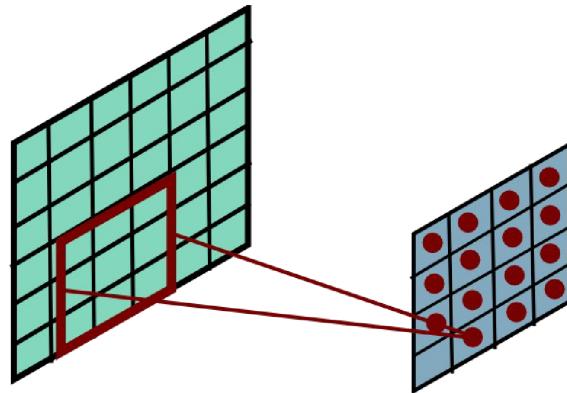
Convolutional Layer



Ranzato

62

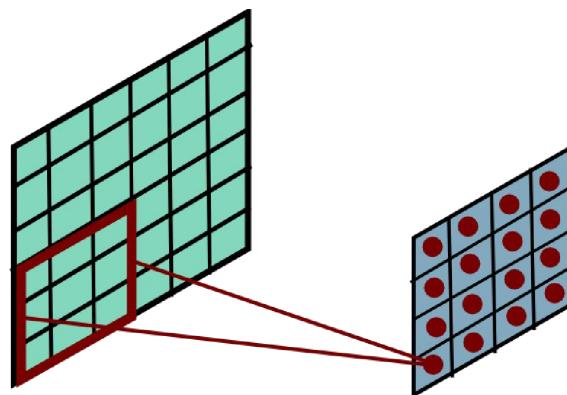
Convolutional Layer



Ranzato

63

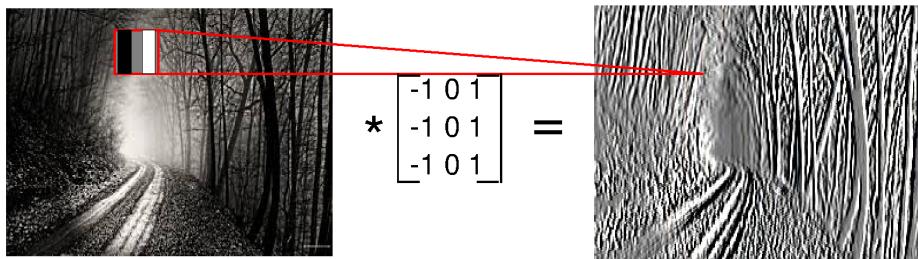
Convolutional Layer



Ranzato

64

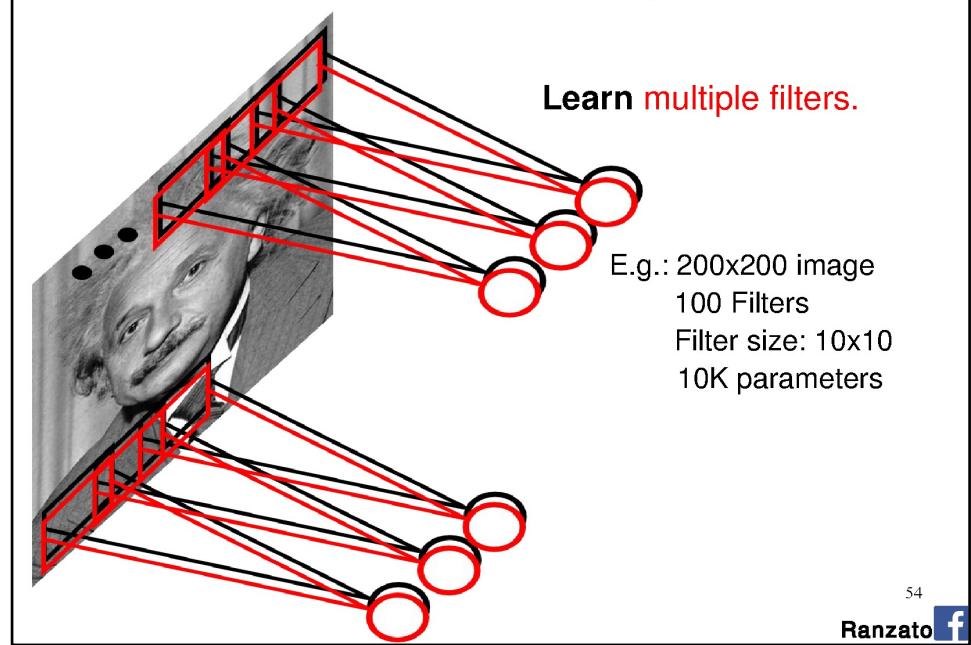
Convolutional Layer



Ranzato

65

Convolutional Layer

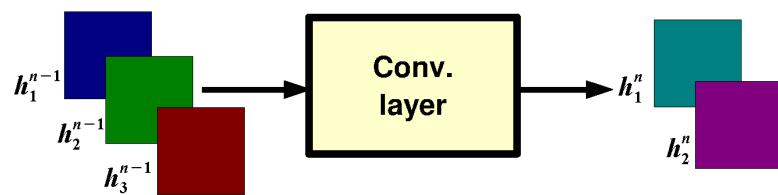


66

Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output feature map input feature map kernel



55

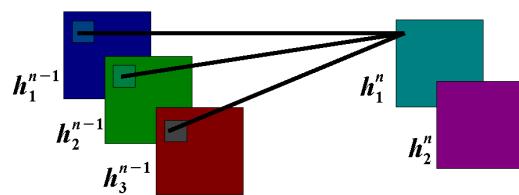
Ranzato

67

Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output feature map input feature map kernel



56

Ranzato

68

Convolutional Layer

$$h_j^n = \max(0, \sum_{k=1}^K h_k^{n-1} * w_{kj}^n)$$

output feature map input feature map kernel

The diagram illustrates the computation of output feature maps h_1^n and h_2^n from input feature maps h_1^{n-1} , h_2^{n-1} , and h_3^{n-1} . A kernel is applied to each input feature map to produce the output feature maps. The output feature map h_1^n is the result of applying the kernel to the top-left block of h_1^{n-1} and the top-right block of h_2^{n-1} . The output feature map h_2^n is the result of applying the kernel to the bottom-left block of h_2^{n-1} and the bottom-right block of h_3^{n-1} .

57
Ranzato

69

Convolutional Layer

Question: What is the size of the output? What's the computational cost?

Answer: It is proportional to the number of filters and depends on the stride. If kernels have size $K \times K$, input has size $D \times D$, stride is 1, and there are M input feature maps and N output feature maps then:

- the input has size $M @ D \times D$
- the output has size $N @ (D-K+1) \times (D-K+1)$
- the kernels have $M \times N \times K \times K$ coefficients (which have to be learned)
- cost: $M \times K \times K \times N \times (D-K+1) \times (D-K+1)$

Question: How many feature maps? What's the size of the filters?

Answer: Usually, there are more output feature maps than input feature maps. Convolutional layers can increase the number of hidden units by big factors (and are expensive to compute). The size of the filters has to match the size/scale of the patterns we want to detect (task dependent).

58
Ranzato

70

Key Ideas

A standard neural net applied to images:

- scales quadratically with the size of the input
- does not leverage stationarity

Solution:

- connect each hidden unit to a small patch of the input
- share the weight across space

This is called: **convolutional layer**.

A network with convolutional layers is called **convolutional network**.

59

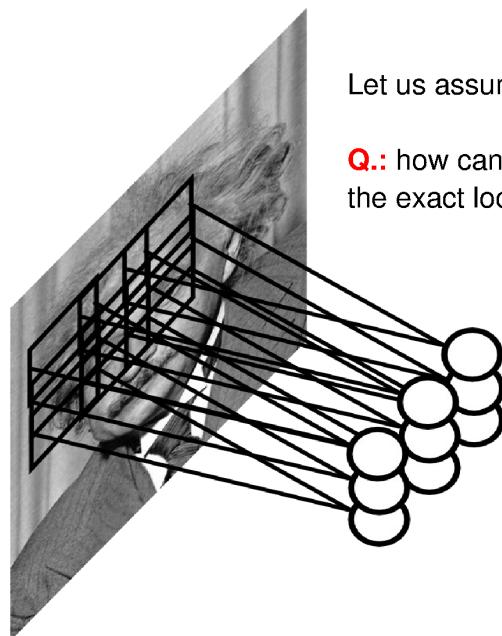
LeCun et al. "Gradient-based learning applied to document recognition" IEEE 1998

71

Pooling Layer

Let us assume filter is an “eye” detector.

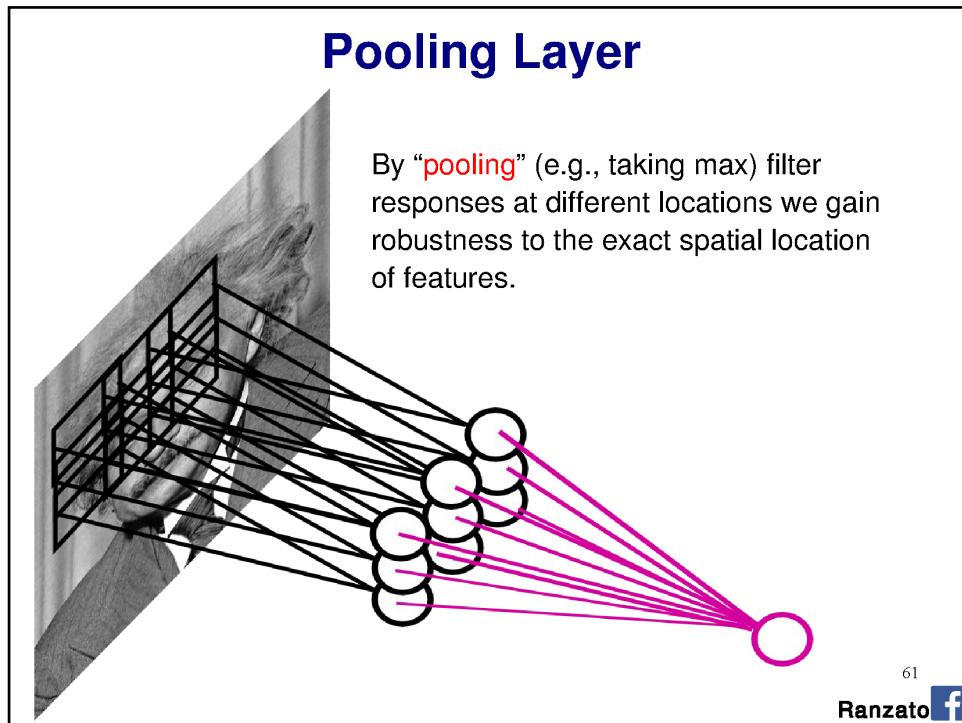
Q.: how can we make the detection robust to the exact location of the eye?



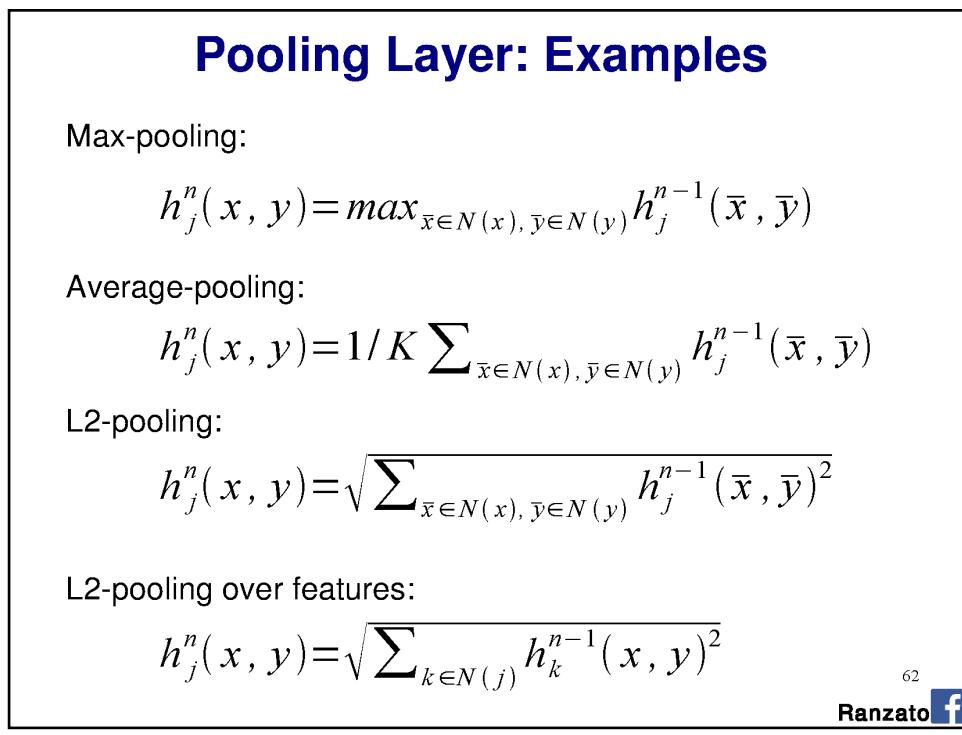
60

Ranzato

72



73



74

Pooling Layer

Question: What is the size of the output? What's the computational cost?

Answer: The size of the output depends on the stride between the pools. For instance, if pools do not overlap and have size KxK, and the input has size DxD with M input feature maps, then:

- output is $M @ (D/K) \times (D/K)$
- the computational cost is proportional to the size of the input (negligible compared to a convolutional layer)

Question: How should I set the size of the pools?

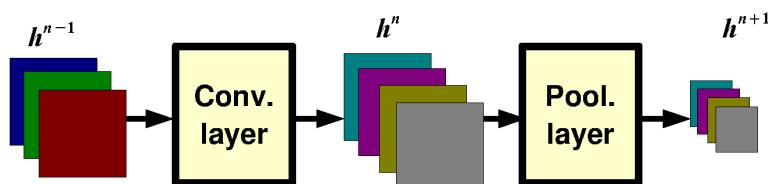
Answer: It depends on how much “invariant” or robust to distortions we want the representation to be. It is best to pool slowly (via a few stacks of conv-pooling layers).

63

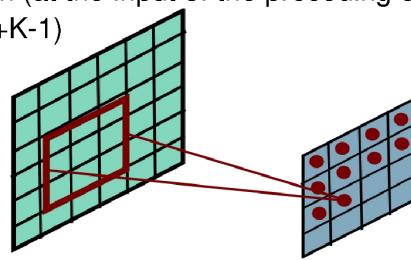
Ranzato 

75

Pooling Layer: Receptive Field Size



If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: $(P+K-1) \times (P+K-1)$

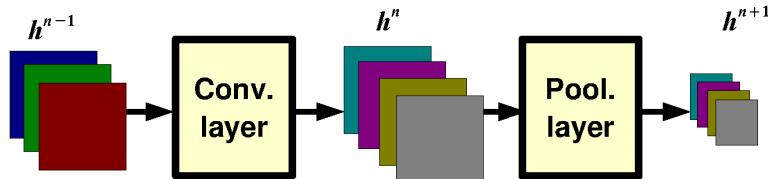


66

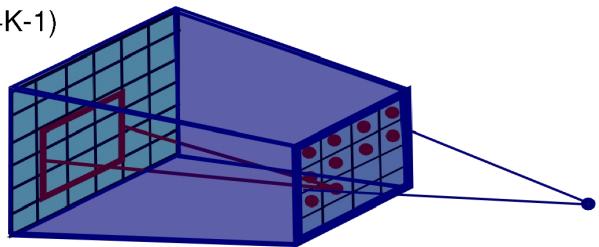
Ranzato 

76

Pooling Layer: Receptive Field Size



If convolutional filters have size $K \times K$ and stride 1, and pooling layer has pools of size $P \times P$, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:
 $(P+K-1) \times (P+K-1)$

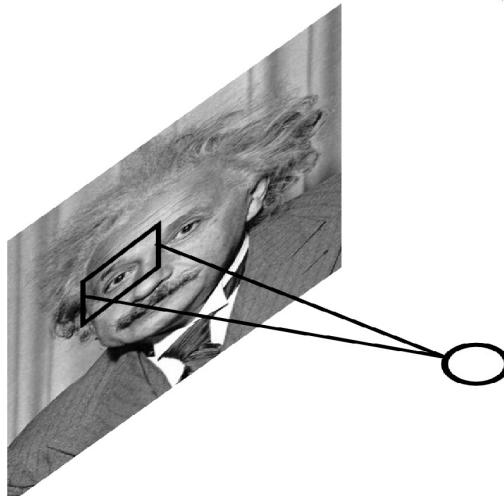


67
Ranzato

77

Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$

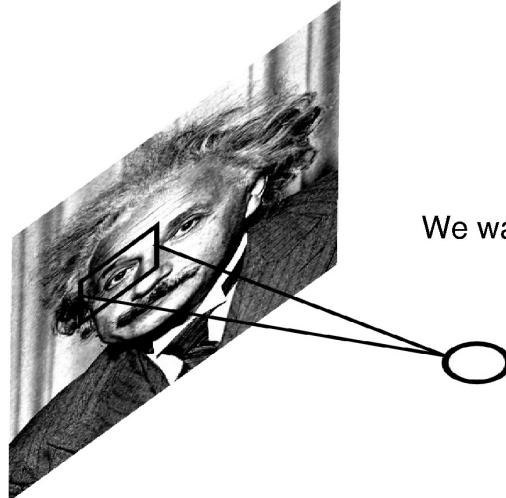


68
Ranzato

78

Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



We want the same response.

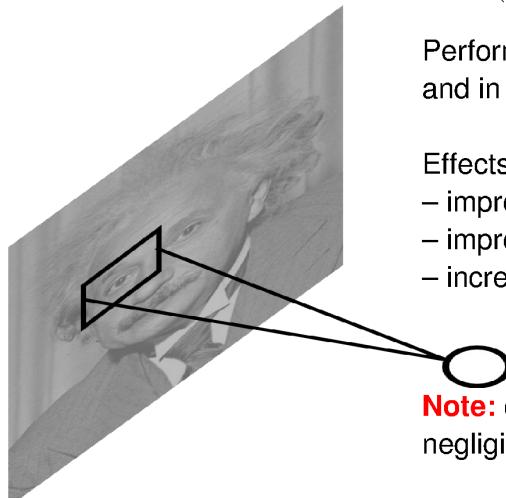
69

Ranzato

79

Local Contrast Normalization

$$h^{i+1}(x, y) = \frac{h^i(x, y) - m^i(N(x, y))}{\sigma^i(N(x, y))}$$



Performed also across features
and in the higher layers..

Effects:

- improves invariance
- improves optimization
- increases sparsity

Note: computational cost is
negligible w.r.t. conv. layer.

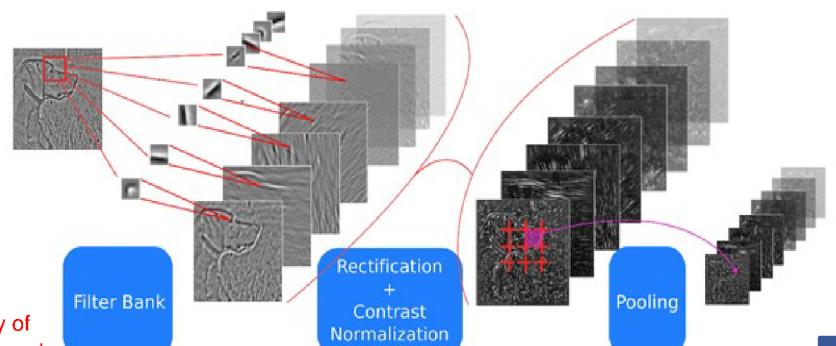
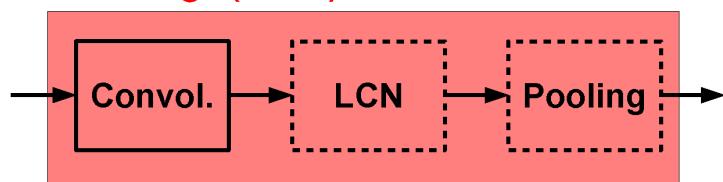
70

Ranzato

80

ConvNets: Typical Stage

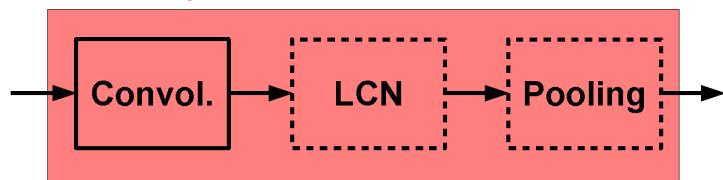
One stage (zoom)



81

ConvNets: Typical Stage

One stage (zoom)

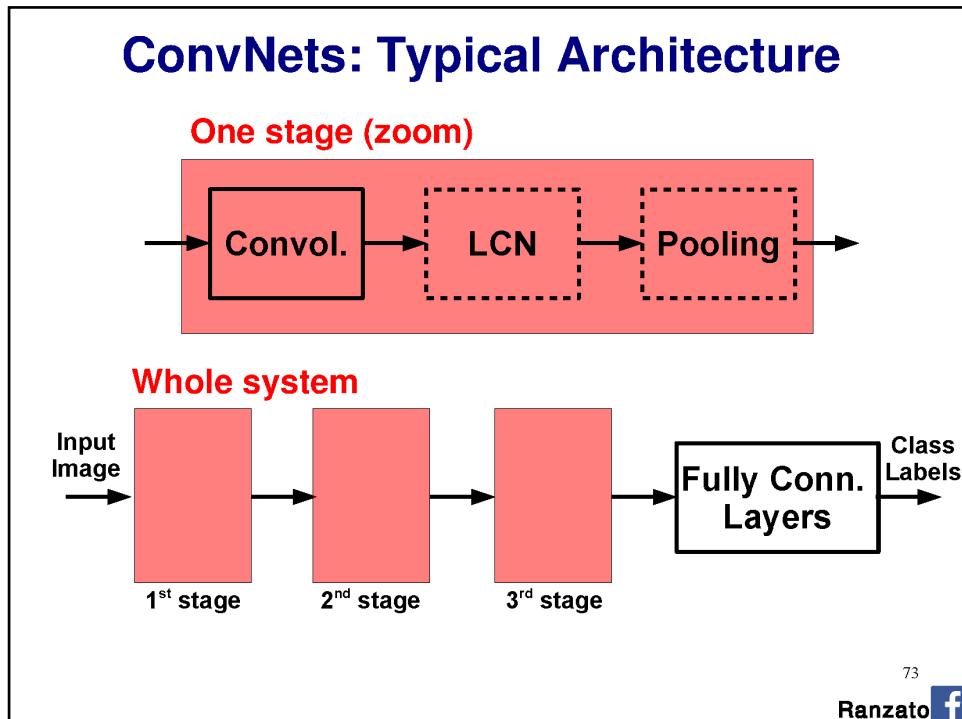


Conceptually similar to: SIFT, HoG, etc.

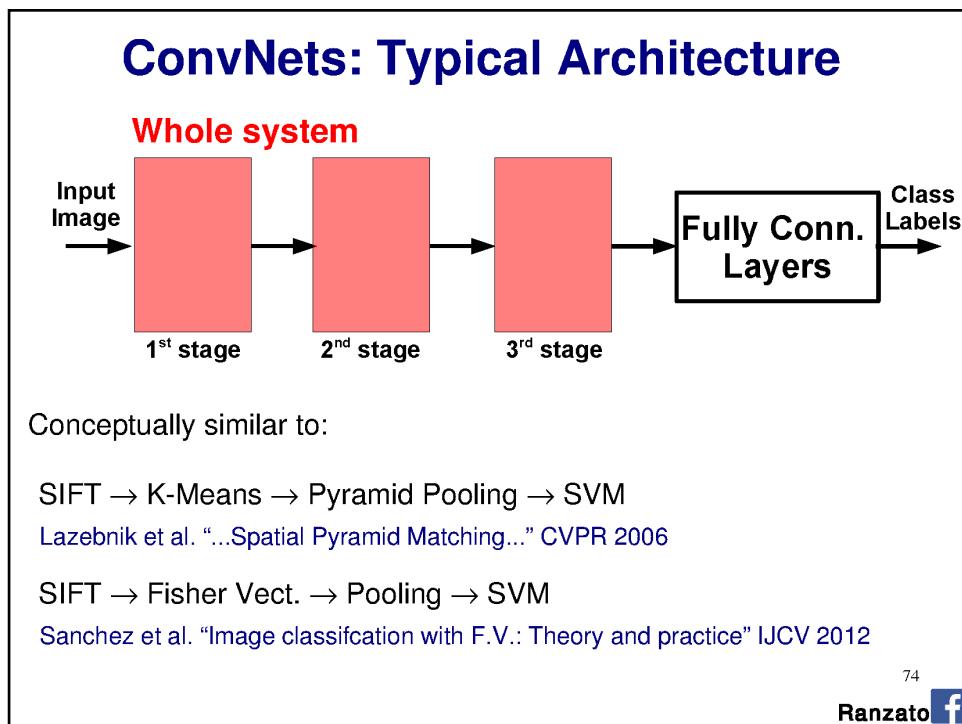
72

Ranzato

82



83



84

ConvNets: Training

All layers are differentiable (a.e.).
We can use standard back-propagation.

Algorithm:

- Given a small mini-batch**
- F-PROP
 - B-PROP
 - PARAMETER UPDATE

75



85

Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples
- Tips

81



86

CONV NETS: EXAMPLES

- OCR / House number & Traffic sign classification



Ciresan et al. "MCDNN for image classification" CVPR 2012

Wan et al. "Regularization of neural networks using dropconnect" ICML 2013

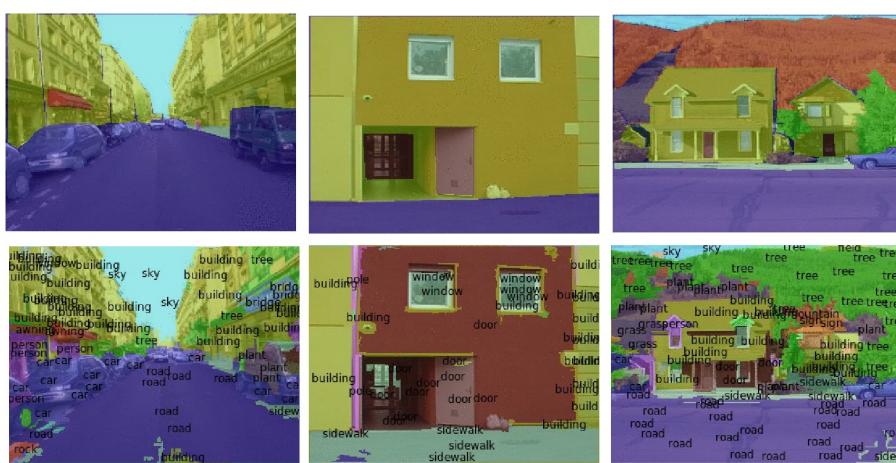
82

Jaderberg et al. "Synthetic data and ANN for natural scene text recognition" arXiv 2014

87

CONV NETS: EXAMPLES

- Scene Parsing



Farabet et al. "Learning hierarchical features for scene labeling" PAMI 2013

85

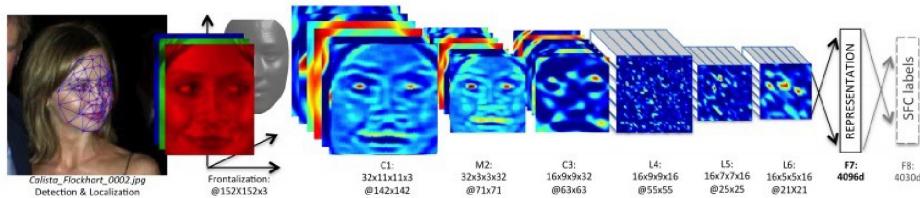
Pinheiro et al. "Recurrent CNN for scene parsing" arxiv 2013

Ranzato

88

CONV NETS: EXAMPLES

- Face Verification & Identification



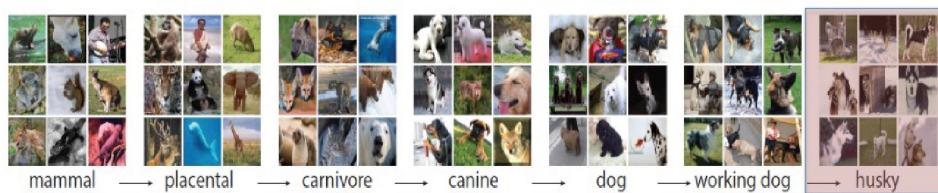
92

Taigman et al. "DeepFace..." CVPR 2014

Ranzato

89

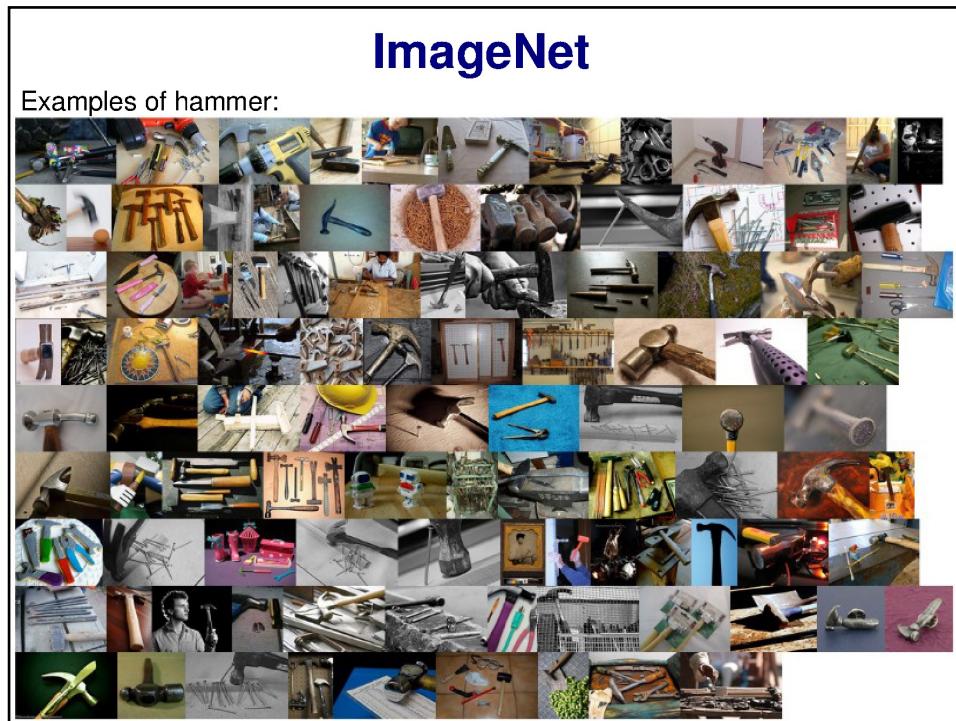
Dataset: ImageNet 2012



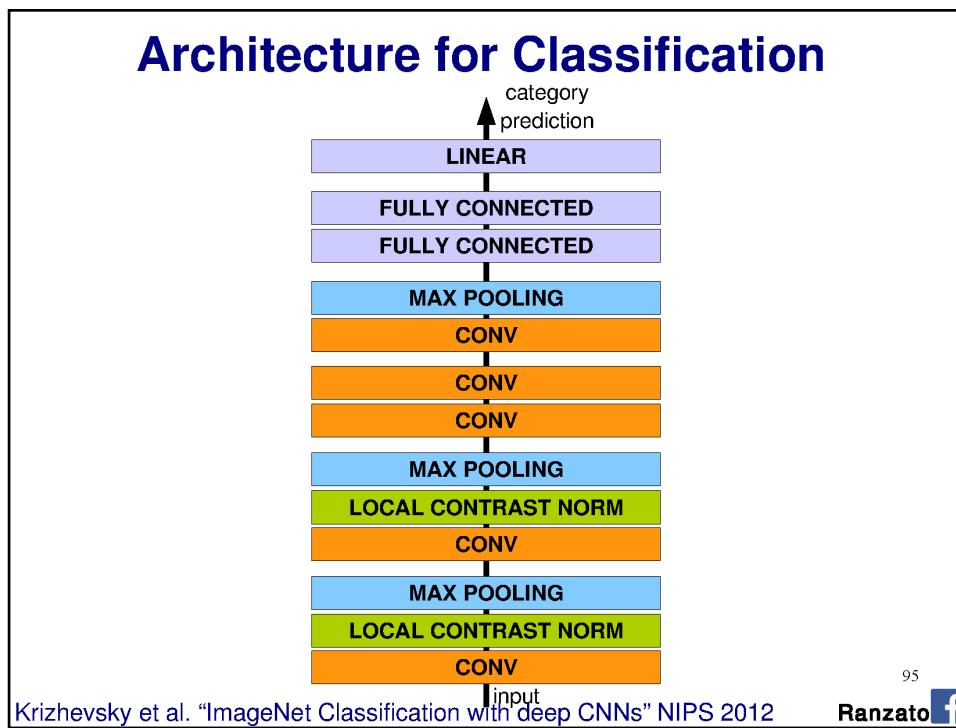
- S: (a) **Eskimo dog**: husky (breed of heavy-coated Arctic sled dog)
 - direct hypernym / inherited hypernym / user term
 - S: (a) **working dog** (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
 - S: (a) **domestic dog, Canis familiaris** (a member of the genus *Canis* (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) "the dog barked all night"
 - S: (a) **canine, canid** (any of various fissiped mammals with nonretractile claws and typically long muzzles)
 - S: (a) **carnivore** (a terrestrial or aquatic flesh-eating mammal) "terrestrial carnivores have four or five clawed digits on each limb"
 - S: (a) **placental, placental mammal, eutherian, eutherian mammal** (mammals having a placenta; all mammals except monotremes and marsupials)
 - S: (a) **mammal, mammalian** (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - S: (a) **vertebrate, craniate** (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
 - S: (a) **chordate** (any animal of the phylum Chordata having a notochord or spinal column)
 - S: (a) **animal, animate being, beast, brute, creature, fauna** (a living organism characterized by voluntary movement)
 - S: (a) **organism, being** (a living thing that has (or can develop) the ability to act or function independently)
 - S: (a) **living thing, animate thing** (a living (or once living) entity)
 - S: (a) **whole, unit** (an assemblage of parts that is regarded as a single entity) "how big is that part compared to the whole?"; "the team is a unit"
 - S: (a) **object, physical object** (a tangible and visible entity, an entity that can cast a shadow) "it was full of rackets, balls and other objects"
 - S: (a) **physical entity** (an entity that has physical existence)
 - S: (a) **entity** (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

Deng et al. "Imagenet: a large scale hierarchical image database" CVPR 2009

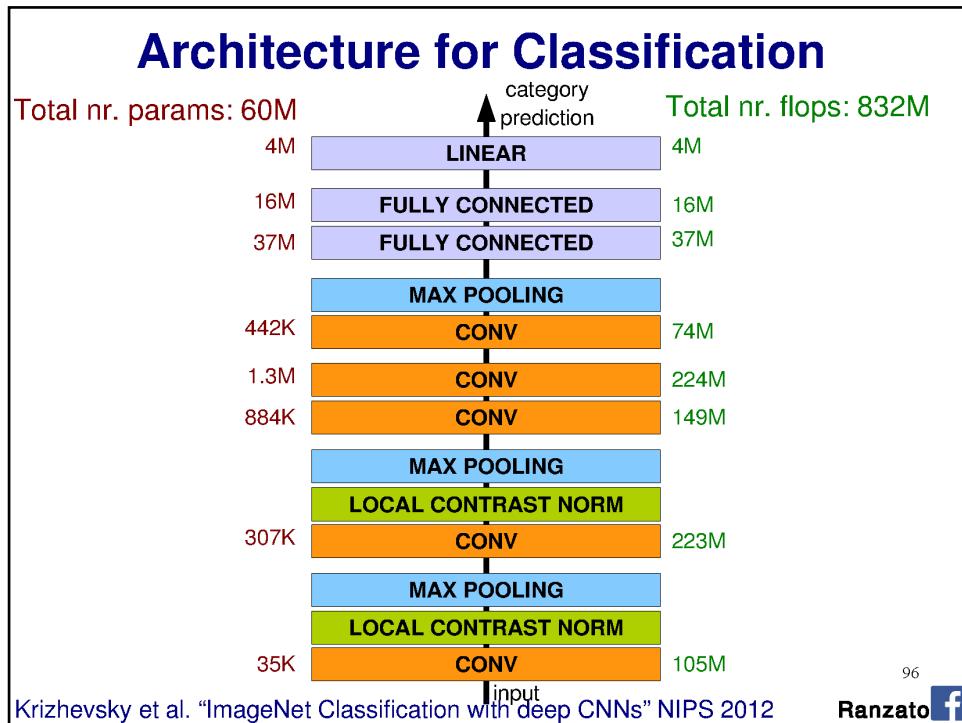
90



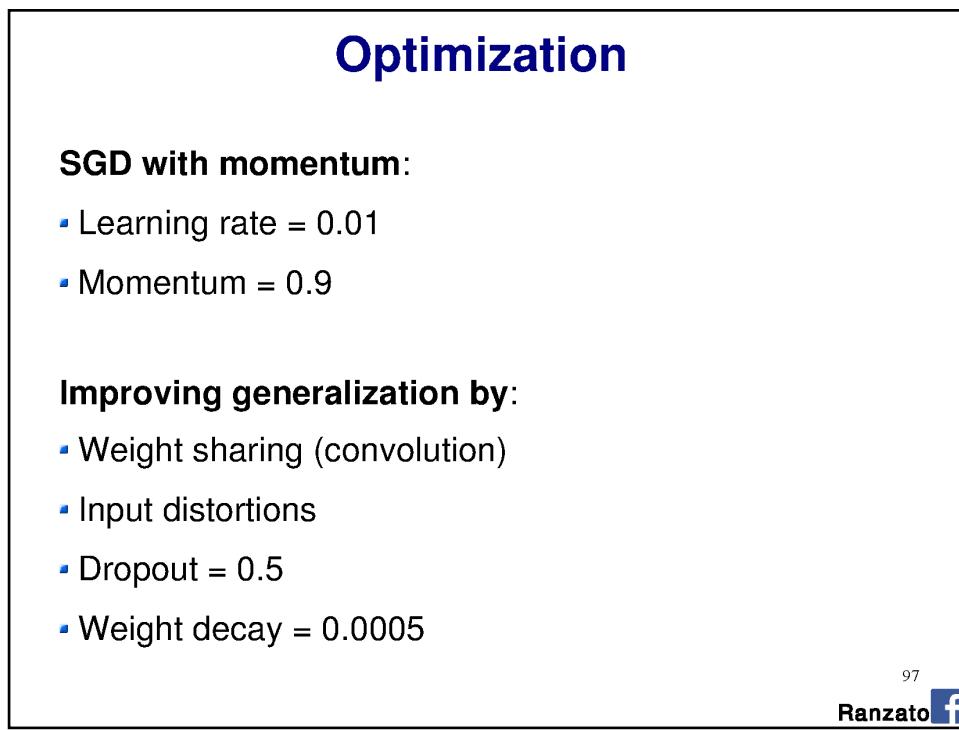
91



92

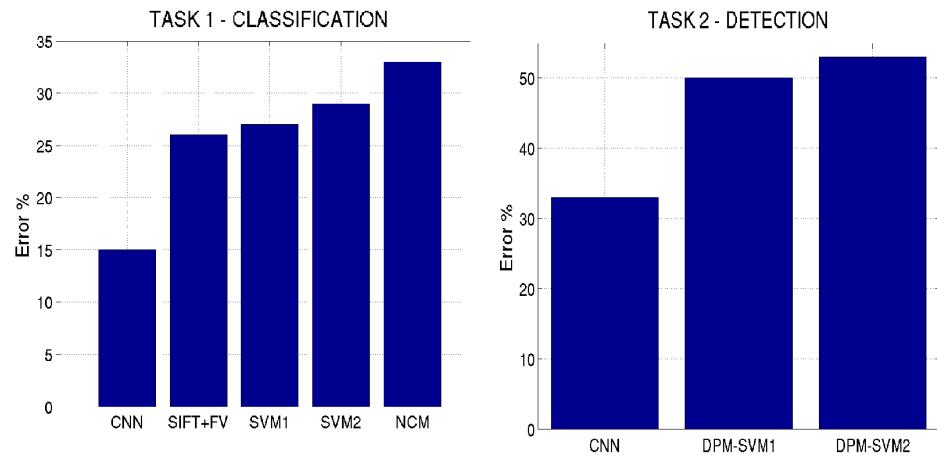


93



94

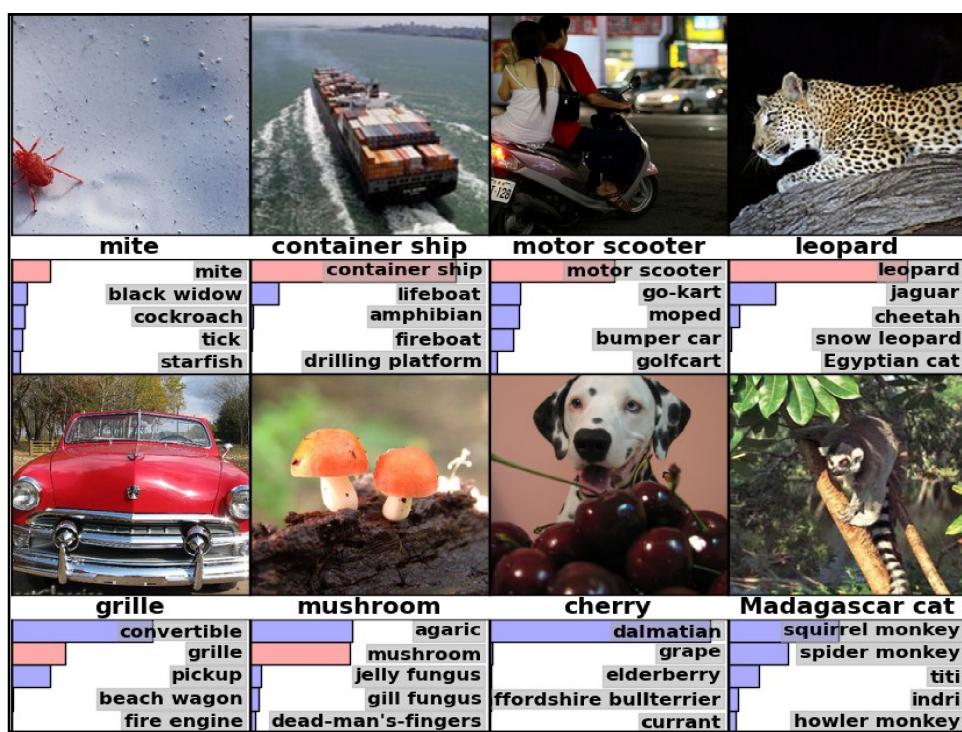
Results: ILSVRC 2012



Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

Ranzato

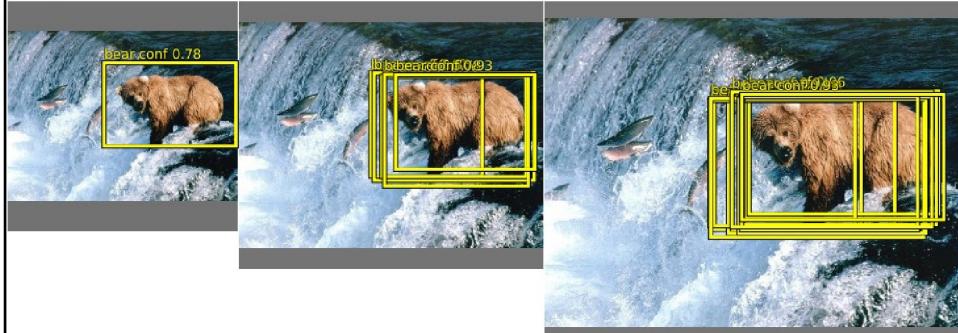
98



96

CONV NETS: EXAMPLES

- Object detection



Sermanet et al. "OverFeat: Integrated recognition, localization, ..." arxiv 2013

Girshick et al. "Rich feature hierarchies for accurate object detection..." arxiv 2013 91

Szegedy et al. "DNN for object detection" NIPS 2013

Ranzato