

Deep Learning for Graphics

Neural Network Basics

Niloy Mitra UCL

lasonas Kokkinos UCL/Facebook

UCL

Paul Guerrero

Vladimir Kim Adobe Research Kostas Rematas U Washington

Tobias Ritschel

UCL



facebook Artificial Intelligence Research







Timetable

	Niloy	lasonas	Paul	Vova	Kostas	Tobias	Ersin
Introduction	Х	Х	Х			Х	
Theory	Х						
NN Basics		Х	Х			Х	
Supervised Applications							
Data						Х	
Unsupervised Applications			Х				
Beyond 2D	Х		Х	Х			Х
Outlook	Х	Х	Х	Х	Х	Х	Х



Introduction to Neural Networks



Goal: Learn a Parametric Function

$$f_{\theta}: \mathbb{X} \longrightarrow \mathbb{Y}$$

 θ : function parameters, $\ \mathbb{X}$: source domain $\ \mathbb{Y}$: target domain these are learned

Examples:

Image Classification:

Image Synthesis:

 $\begin{aligned} f_{\theta} : \mathbb{R}^{w \times h \times c} &\longrightarrow \{0, 1, \dots, k-1\} \\ w \times h \times c : \text{image dimensions} \quad k : \text{class count} \\ f_{\theta} : \mathbb{R}^{n} &\longrightarrow \mathbb{R}^{w \times h \times c} \\ n : \text{latent variable count} \quad w \times h \times c : \text{image dimensions} \end{aligned}$



Machine Learning 101: Linear Classifier



$$f_{\theta} : \mathbb{R}^{n} \longrightarrow \{0, 1\}$$
$$f_{\theta}(x) = \begin{cases} 1 & \text{if } wx + b \ge 0\\ 0 & \text{if } wx + b < 0 \end{cases}$$
$$\theta = \{w, b\}$$

Each data point has a class label:

 y^i

Nonlinear decision boundaries





Slide Credit: Marc'Aurelio Ranzato, Yann LeCun















'Neuron': Cascade of Linear and Nonlinear Function

basic building block







Activation functions

function



Image Credit: Olivier Grisel and Charles Ollion



Perceptrons (60's)





Multi-Layer Perceptrons (~1985)





EG Course Deep Learning for Graphics

Reminder: Non-linear decision boundaries





Nonlinear mapping $\mathbb{R}^2 o \mathbb{R}^2$

Evolution of isocontours as parameters change



$$y_1 = g(w_{1,1}x_1 + w_{1,2}x_2 + w_{1,3})$$

$$y_2 = g(w_{2,1}x_1 + w_{2,2}x_2 + w_{2,3})$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$
$$\mathbf{y} = g(\mathbf{W}\mathbf{x})$$

http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/



From non-separable to linearly separable



http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/



Linearizing a 2D classification task (4 hidden layers)



http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/

Linearization: may need higher dimensions





http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/

Linearization: may need higher dimensions



http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/

Linearization: may need higher dimensions



http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/



Hidden Layers: intuitively, what do they do?

Intuition: learn "dictionary" for objects

"Distributed representation":

represent (and classify) objects by mixing & mashing reusable parts

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



Deep Learning = Hierarchical Compositionality





Deep Learning = Hierarchical Compositionality





MLP Demo: playground.tensorflow.org





Training and Optimization



Neural Network Training: Old & New Tricks

Old:

Back-propagation algorithm

Stochastic Gradient Descent, Momentum, "weight decay"

```
New: (last 5-6 years)
```

Dropout

ReLUs

Batch Normalization

Residual Networks



Training Goal

Our network implements a parametric function:

$$f_{\theta}: \mathbb{X} \longrightarrow \mathbb{Y} \qquad \hat{y} = f(x; \theta)$$

During training, we search for parameters that minimize a loss: $\min_{\theta} L_f(\theta)$

Example: L2 regression loss given target (x^i, y^i) pairs :

$$L_{f}(\theta) = \sum_{i} \|f(x^{i};\theta) - y^{i}\|_{2}^{2}$$



Gradient Descent Minimization Method

Initialize:
$$\theta_0$$

Update: $\theta_{i+1} = \theta_i - \alpha \nabla f(\theta_i)$

We can always make it converge for a **convex** function





Multiple Local Minima, based on initialization



Empirically all are almost equally good

Central research topic: how can this happen?



On to the gradients!

All you need is gradients





Chain Rule



Given y(x) and dL/dy, What is dL/dx?

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun



Chain Rule



Given y(x) and dL/dy, What is dL/dx? $\frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dx}$

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun



'Another Brick in the Wall'



Given y(x) and dL/dy, What is dL/dx? $\frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dx}$

Slide Credit: Marc'Aurelio Ranzato, Yann LeCun





ep Learning for Graphics

Toy example: single sigmoidal unit

$$f(w,x) = \frac{1}{1 + \exp(-(w_0x_0 + w_1x_1 + w_2))}$$

Composition of differentiable blocks:

$$f(x) = \frac{1}{x} \quad \rightarrow \quad f'(x) = -\frac{1}{x^2}$$

$$f_c(x) = c + x \quad \rightarrow \quad f'(x) = 1$$

$$f(x) = e^x \quad \rightarrow \quad f'(x) = e^x$$

$$f_a(x) = ax \quad \rightarrow \quad f'(x) = a$$



Computation graph & automatic differentiation



Slide Credit: Justin Johnson


Multi-Layer Perceptrons



Slide Credit: G. Hinton



EG Course Deep Learning for Graphics

Multi-Layer Perceptrons





EG Course Deep Learning for Graphics

Back-propagation Algorithm





KeepCalmAndPosters.com

Training Goal

Our network implements a parametric function:

$$f_{\theta}: \mathbb{X} \longrightarrow \mathbb{Y} \qquad \hat{y} = f(x; \theta)$$

During training, we search for parameters that minimize a loss: $\min_{\theta} L_f(\theta)$

Example: L2 regression loss given target (x^i, y^i) pairs :

$$L_{f}(\theta) = \sum_{i} \|f(x^{i};\theta) - y^{i}\|_{2}^{2}$$



A Neural Network for Multi-way Classification





Hidden layer









EG Course Deep Learning for Graphics



Hidden layer



EG Course Deep Learning for Graphics



Hidden layer



Objective for linear regression



 $h(\mathbf{b}) = \mathbf{b}$



Objective for multi-class classification





Neural network in forward mode: recap

Network output:
$$\mathbf{\hat{y}} = f(\mathbf{x}; \mathbf{v}, \mathbf{w})$$

Loss (prediction error): $l(\mathbf{\hat{y}}, \mathbf{y})$

What we need to compute for gradient descent:

$$\frac{\partial l(\mathbf{\hat{y}}, \mathbf{y})}{\partial \mathbf{v}_i} \quad \frac{\partial l(\mathbf{\hat{y}}, \mathbf{y})}{\partial \mathbf{w}_j}$$







Hidden layer





Hidden layer





Linear Layer in Forward Mode: All For One

$$b_{m} = \sum_{h=1}^{H} z_{h} w_{h,m}$$

$$z_{1}$$

$$z_{h}$$

$$b_{m}$$

$$z_{H}$$



Linear Layer in Backward Mode: One From All

$$b_{m} = \sum_{h=1}^{H} z_{h} w_{h,m}$$

$$w_{h,m} \qquad \qquad \frac{\partial L}{\partial b_{m}}$$

$$\frac{\partial L}{\partial z_h} = \sum_{c=1}^C \frac{\partial L}{\partial b_c} \cdot \frac{\partial b_c}{\partial z_h} = \sum_{c=1}^C \frac{\partial L}{\partial b_c} w_{h,c}$$



Linear Layer Parameters in Backward: 1-to-1

$$\frac{\partial L}{\partial w_{h,m}} = \sum_{c=1}^{C} \frac{\partial L}{\partial b_c} \cdot \frac{\partial b_c}{\partial w_{h,m}} = \frac{\partial L}{\partial b_m} z_h$$



Hidden layer





Hidden layer







EUROGRAPHICS 2018





Neural Network Training: Old & New Tricks

Old:

Back-propagation algorithm

Stochastic Gradient Descent, Momentum, "weight decay"

```
New: (last 5-6 years)
```

Dropout

ReLUs

Batch Normalization



Training Objective for N training samples

$$L(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^{N} l(\mathbf{y}^{i}, \hat{\mathbf{y}}^{i}) + \sum_{l} \lambda_{l} \sum_{k,m} (\mathbf{W}_{k,m}^{l})^{2}$$
Per-sample loss
Per-layer regularization

Gradient descent: $\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon \nabla_{\mathbf{W}} L(\mathbf{W}_t)$

(I,k,m) element of gradient vector:

ROGRAPHICS

$$\frac{\partial L}{\partial \mathbf{W}_{k,m}^{l}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial l(\mathbf{y}^{i}, \mathbf{\hat{y}}^{i})}{\partial \mathbf{W}_{k,m}^{l}} + 2\lambda_{l} \mathbf{W}_{k,m}^{l}$$
Back-prop for i-th example

If N=10⁶, we will need to run back-prop 10⁶ times to update **W** once!

Stochastic Gradient Descent (SGD)

Gradient: **Batch:** [1..N]

$$\frac{\partial L}{\partial \mathbf{W}_{k,m}^{l}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial l(\mathbf{y}^{i}, \mathbf{\hat{y}}^{i})}{\partial \mathbf{W}_{k,m}^{l}} + 2\lambda_{l} \mathbf{W}_{k,m}^{l}$$

Noisy ('Stochastic') Gradient: Minibatch: B elements b(1), b(2),..., b(B): sampled from [1,N]

$$\frac{\partial L}{\partial \mathbf{W}_{k,m}^{l}} \simeq \frac{1}{B} \sum_{i=1}^{B} \frac{\partial l(\mathbf{y}^{b(i)}, \mathbf{\hat{y}}^{b(i)})}{\partial \mathbf{W}_{k,m}^{l}} + 2\lambda_{l} \mathbf{W}_{k,m}^{l}$$

Epoch: N samples, N/B batches



Code example

Gradient Descent vs Stochastic Gradient Descent

Regularization in SGD: Weight Decay

Gradient: Batch: [1..N]

$$\frac{\partial L}{\partial \mathbf{W}_{k,m}^{l}} = \frac{1}{N} \sum_{i=1}^{N} \frac{\partial l(\mathbf{y}^{i}, \mathbf{\hat{y}}^{i})}{\partial \mathbf{W}_{k,m}^{l}} + 2\lambda_{l} \mathbf{W}_{k,m}^{l}$$

Noisy ('Stochastic') Gradient: Minibatch: B elements b(1), b(2),..., b(B): sampled from [1,N]

$$\frac{\partial L}{\partial \mathbf{W}_{k,m}^{l}} \simeq \frac{1}{B} \sum_{i=1}^{B} \frac{\partial l(\mathbf{y}^{b(i)}, \hat{\mathbf{y}}^{b(i)})}{\partial \mathbf{W}_{k,m}^{l}} + 2\lambda_{l} \mathbf{W}_{k,m}^{l}$$

Back-prop on minibatch "Weight decay"

Epoch: N samples, N/B batches



Learning rate





Gradient Descent





(S)GD with adaptable stepsize





(S)GD with momentum



Main idea: retain long-term trend of updates, drop oscillations

(S)GD
$$\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \nabla_{\mathbf{W}} L(\mathbf{W})$$

(S)GD + momentum

$$\mathbf{V}_{t+1} = \mu \mathbf{V}_t + (1 - \mu) \nabla_{\mathbf{W}} L(\mathbf{W}_t)$$
$$\mathbf{W}_{t+1} = \mathbf{W}_t - \epsilon_t \mathbf{V}_{t+1}$$
EG



Code example

Multi-layer perceptron classification

Step-size Selection & Optimizers: research problem

- Nesterov's Accelerated Gradient (NAG)
- R-prop
- AdaGrad
- RMSProp
- AdaDelta
- Adam
- ...





Neural Network Training: Old & New Tricks

Old: (80's)

Stochastic Gradient Descent, Momentum, "weight decay"

New: (last 5-6 years)

Dropout

ReLUs

Batch Normalization



Linearization: may need higher dimensions





http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/



Reminder: Overfitting, in images

Classification



Regression




Previously: 12 Regularization

$$\begin{split} L(\mathbf{W}) &= \frac{1}{N} \sum_{i=1}^{N} l(\mathbf{y}^{i}, \hat{\mathbf{y}}^{i}) + \sum_{l} \lambda_{l} \sum_{k,m} (\mathbf{W}_{k,m}^{l})^{2} \\ & \text{Per-sample loss} \quad \text{Per-layer regularization} \end{split}$$



Dropout



Each sample is processed by a 'decimated' neural net

Decimated nets: distinct classifiers But: they should all do the same job



Dropout block



$$\begin{array}{lll} z_i^{(l+1)} & = & \mathbf{w}_i^{(l+1)} \mathbf{y}^l + b_i^{(l+1)}, \\ y_i^{(l+1)} & = & f(z_i^{(l+1)}), \end{array}$$

$$\begin{array}{lll} r_{j}^{(l)} & \sim & \mathrm{Bernoulli}(p), \\ \widetilde{\mathbf{y}}^{(l)} & = & \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\ z_{i}^{(l+1)} & = & \mathbf{w}_{i}^{(l+1)} \widetilde{\mathbf{y}}^{l} + b_{i}^{(l+1)}, \\ y_{i}^{(l+1)} & = & f(z_{i}^{(l+1)}). \end{array}$$



'Feature noising'

EG Course Deep Learning for Graphics

Test time: Deterministic Approximation



EG Course Deep Learning for Graphics

Dropout Performance



Figure 4: Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.



Neural Network Training: Old & New Tricks

Old: (80's)

Stochastic Gradient Descent, Momentum, "weight decay"

New: (last 5-6 years)

Dropout

ReLUs

Batch Normalization



'Neuron': Cascade of Linear and Nonlinear Function





Reminder: a network in backward mode





EG Course Deep Learning for Graphics

Vanishing Gradients Problem

Gradient signal
from above

$$\frac{\partial l}{\partial a_k} = \sum_m \frac{\partial l}{\partial z_m} \frac{\partial z_m}{\partial a_k} = \begin{bmatrix} \partial l \\ \partial z_k \end{bmatrix} g'(a_k) = \frac{\partial l}{\partial z_k} g(a_k)(1 - g(a_k))$$

Do this 10 times: updates in the first layers get minimal Top layer knows what to do, lower layers "don't get it" Sigmoidal Unit: Signal is not getting through!





Vanishing Gradients Problem: ReLU Solves It



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

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



Neural Network Training: Old & New Tricks

Old: (80's)

Stochastic Gradient Descent, Momentum, "weight decay"

New: (last 5-6 years)

Dropout

ReLUs

Batch Normalization



External Covariate Shift: your input changes







"Whitening": Set Mean = 0, Variance = 1

Photometric transformation: $I \rightarrow a I + b$





Original Patch and Intensity Values



- Then make it have unit variance:

Brightness Decreased





```
\sigma^{2} = \frac{1}{N} \sum_{x,y} Z(x,y)^{2}ZN(x,y) = \frac{Z(x,y)}{\sigma}
```



EG Course **Deep Learning for Graphics**

• Make each patch have zero mean:

$$\mu = \frac{1}{N} \sum_{x,y} I(x,y)$$
$$Z(x,y) = I(x,y) - \mu$$

Internal Covariate Shift

Neural network activations during training: moving target





EG Course Deep Learning for Graphics

Batch Normalization

Whiten-as-you-go:

 Normalize the activations in each layer within a minibatch.

Learn the mean and variance
$$(\gamma, \beta)$$
 of each laver as parameters





Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift S loffe and C Szegedy (2015)



Batch Normalization: used in all current systems

- Multi-layer CNN's train faster with fewer data samples (15x).
- Employ faster learning rates and less network regularizations.
- Achieves state of the art results on ImageNet.



number of mini-batches



Convolutional Neural Networks



Fully-connected Layer



- we have not enough training samples anyway..





EUROGRAPHICS

EG Course Deep Learning for Graphics



EUROGRAPHICS 2018

EG Course Deep Learning for Graphics

Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels


































































Fully-connected layer



#of parameters: K²

y_1		$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$	•••	$w_{1,K}$ -] [x_1
y_2		$w_{2,1}$	$w_{2,2}$	$w_{2,3}$	$w_{2,4}$	• • •	$w_{2,K}$		x_2
y_3		$w_{3,1}$	$w_{3,2}$	$w_{3,3}$	$w_{3,4}$	• • •	$w_{3,K}$		x_3
y_4	=	$w_{4,1}$	$w_{4,2}$	$w_{4,3}$	$w_{4,4}$	• • •	$w_{4,K}$		x_4
•				• •					• •
y_K		$w_{K,1}$	$w_{K,2}$	$w_{K,3}$	$w_{K,4}$	• • •	$w_{K,K}$		x_K





#of parameters: size of window

$\begin{bmatrix} y_1 \end{bmatrix}$		w_0	w_1	w_2	0	• • •	0	x_1
y_2		0	w_0	w_1	w_2	• • •	0	x_2
y_3		0	0	w_0	w_1	• • •	0	x_3
y_4	=	0	0	0	w_0	• • •	0	x_4
				• •				•
$\begin{bmatrix} y_K \end{bmatrix}$		0	0	0	0	• • •	w_0	x_K







Code example

Learning an edge filter



EG Course Deep Learning for Graphics

















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?



Pooling layer

By "pooling" (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.



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:





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)x(P+K-1)





Receptive field

-0-

K

A DESIGNATION OF THE OWNER

- 0-

K

A DESIGNATION OF TAXABLE PARTY.

k

10 C

H







Modern Architectures



CNNs, late 1980's: LeNet



Gradient-based learning applied to document recognition. Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. 1998

https://www.youtube.com/watch?v=FwFduRA_L6Q





What happened in between?



deep learning = neural networks (+ big data + GPUs) + a few more recent tricks!



EG Course Deep Learning for Graphics

Code example

Convolutional Network and Filter Visualizations

Parameter Initialization

 All zero initialization: All parameters get the same gradient and same updates

w = 0

• Random initialization: Is sometimes used in practice, but variance of output depends on number of inputs, which may cause instability early on

$$w \sim \mathcal{N}(\mu = 0, \sigma = 0.01)$$

• Kaiming initialization: divide standard deviation by number of inputs

 $w \sim \mathcal{N}(\mu = 0, \sigma = rac{0.01}{n})$ n : number of inputs (fan-in)



Code examples

Parameter initialization

CNNs, 2012



AlexNet

Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton: ImageNet classification with deep convolutional neural networks. Commun. ACM 60(6): 84-90 (2017)



CNNs, 2014: VGG



Karen Simonyan, Andrew Zisserman (=Visual Geometry Group) Very Deep Convolutional Networks for Large-Scale Image Recognition, arxiv, 2014.



ning for Graphics

CNNs, 2014: GoogLeNet



Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich Going Deeper with Convolutions, CVPR 2015



CNNs, 2015: ResNet



ResNet

Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, Deep Residual Learning for Image Recognition CVPR 2016



The Deeper, the Better

- Deeper networks can cover more complex problems
 - Increasingly large receptive field size & rich patterns





Going Deeper

- From 2 to 10: 2010-2012
 - ReLUs

...

- Dropout

Revolution of Depth 28.2 25.8 152 layers 16.4 11.7 22 layers 19 layers 7.3 **\ 6.7** 3.57 8 layers 8 layers shallow ILSVRC'15 ILSVRC'14 ILSVRC'14 ILSVRC'13 ILSVRC'12 ILSVRC'11 ILSVRC'10 ResNet GoogleNet VGG AlexNet



Going Deeper

- From 10 to 20: 2015
 - Batch Normalization





Going Deeper

- From 20 to 100/1000
 - Residual networks




Plain network: deeper is not necessarily better

- Plain nets: stacking 3x3 conv layers
- 56-layer net has higher training error and test error than 20-layer net





Residual Network

- Naïve solution
 - If extra layers are an identity mapping, then training errors can not increase





EG Course Deep Learning for Graphics

Residual Modelling: Basic Idea in Image Processing

• Goal: estimate update between an original image and a changed image



Residual Network

- Plain block
 - Difficult to make identity mapping because of multiple non-linear layers





Residual Network

• Residual block

- If identity were optimal, easy to set weights as 0
- If optimal mapping is closer to identity, easier to find small fluctuations

Appropriate for treating perturbation as keeping a base information





Residual Network: deeper is better

• Deeper ResNets have lower training error





EG Course Deep Learning for Graphics

Residual Network: deeper is better





CNNs, 2017: DenseNet

Densely Connected Convolutional Networks, CVPR 2017 Gao Huang, Zhuang Liu, Laurens van der Maaten, Kilian Q. Weinberger



Recently proposed, better performance/parameter ratio





Image-to-Image



Image-to-image

- So far we mapped an image image to a number or label
- In graphics, output often is "richer":
 - An image
 - A volume
 - A 3D mesh
 - ...
- Architectures
 - Encoder-Decoder
 - Skip connections







EG Course Deep Learning for Graphics





EG Course Deep Learning for Graphics









EG Course Deep Learning for Graphics



Fast (shared convolutions) Simple (dense)



EG Course Deep Learning for Graphics

Fully Convolutional Neural Networks in Practice



32-fold decimation 224x224 to 7x7



Fast (shared convolutions) Simple (dense) Low resolution



Receptive field arithmetic



https://medium.com/mlreview/a-guide-to-receptive-field-arithmetic-for-convolutional-neural-networks-e0f514068807



EG Course Deep Learning for Graphics

Atrous convolution



S. Mallat, An introduction to wavelets, 1989

DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, Alan L. Yuille

EUROGRAPHICS

EG Course Deep Learning for Graphics

Atrous convolution = Dilated Convolution



Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a) F_1 is produced from F_0 by a 1-dilated convolution; each element in F_1 has a receptive field of 3×3 . (b) F_2 is produced from F_1 by a 2-dilated convolution; each element in F_2 has a receptive field of 7×7 . (c) F_3 is produced from F_2 by a 4-dilated convolution; each element in F_3 has a receptive field of 15×15 . The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.

F. Yu, V. Koltun, Multi-Scale Context Aggregation by Dilated Convolutions, ICLR 2016



Graphics: Multiresolution









Encoder-decoder





Interpretation

- Turns image into vector
- This vector is a very compact and abstract "code"
- Turns code back into image



Encoder-decoder



Learning to simplify. Simo-Serra et al. 2016



Code example

Colorization Network

Up-sampling

- We saw
 - ... how to keep resolution
 - … how to reduce it with pooling
- But how to increase it again?
- Options
 - Interpolation
 - Padding (insert zeros)
 - Transpose convolutions



Encoder-decoder + Skip connections

- 1st: Reduce resolutions as before
- 2nd: Increase resolution
- Transposed convolutions





U-Net: Convolutional Networks for Biomedical Image Segmentatio. Ronneberger et al. 2015



Encoder-decoder with skip connections





Interpretation

- Turns image into vector
- Turns vector back into image
- At every step of increasing the resolution, check back with the input to preserve details
- Familiar trick to graphics people
 - (Haar) wavelet
 - Residual coding
 - Pyramidal schemes (Laplacian pyramid, etc.)



Deep Learning Frameworks



Main frameworks



(Python, C++, Java)



(Python, backends support other languages)



(Python)



Currently less frequently used



Popularity

Google Trends for search terms: "[name] github"



Google Trends for search terms: "[name] tutorial"





PYTORCH

Caffe

Typical Training Steps

```
for i = 1 .. max_iterations
```

```
input, ground_truth = load_minibatch(data, i)
```

```
output = network_evaluate(input, parameters)
```

```
loss = compute_loss(output, ground_truth)
```

```
# gradients of loss with respect to parameters
gradients = network_backpropagate(loss, parameters)
```

```
parameters = optimizer_step(parameters, gradients)
```



Tensors

- Frameworks typically represent data as tensors
- Examples:



4D convolution kernel: OC x IC x KH x KW



output channels OC



What Does a Deep Learning Framework Do?

- Tensor math
- Common network operations/layers
- Gradients of common operations
- Backpropagation
- Optimizers
- GPU implementations of the above
- usually: data loading, network parameter saving/loading
- sometimes: distributed computing



Automatic Differentiation & the Computation Graph



are the same size as the parameters



loss

loss

Automatic Differentiation & the Computation Graph




Static vs Dynamic Computation Graphs

- Static analysis allows optimizations and distributing workload
- Dynamic graphs make data-driven control flow easier
- In static graphs, the graph is usually defined in a separate 'language'
- Static graphs have less support for debugging

```
define once,
evaluate during training
Static
```

```
x = Variable()
loss = if_node(x < parameter[0],
    x + parameter[0],
    x - parameter[1])
for i = 1 .. max_iterations
    x = data()
    run(loss)
    backpropagate(loss, parameters)</pre>
```

define implicitly by running operations, a new graph is created in each evaluation **Dynamic**

```
for i = 1 .. max_iterations
    x = data()
    if x < parameter[0]
        loss = x + parameter[0]
    else
        loss = x - parameter[1]
    backpropagate(loss, parameters)</pre>
```



Tensorflow



- Currently the largest community
- Static graphs (dynamic graphs are in development: Eager Execution)
- Good support for deployment
- Good support for distributed computing
- Typically slower than the other three main frameworks on a single GPU



PyTorch

PYTÖRCH

- Fast growing community
- Dynamic graphs
- Distributed computing is in development (some support is already available)
- Intuitive code, easy to debug and good for experimenting with less traditional architectures due to dynamic graphs
- Very Fast







- A high-level interface for various backends (Tensorflow, CNTK, Theano)
- Intuitive high-level code
- Focus on optimizing time from idea to code
- Static graphs







- Created earlier than Tensorflow, PyTorch or Keras
- Less flexible and less general than the other three frameworks
- Static graphs
- Legacy to be replaced by Caffe2: focus is on performance and deployment
 - Facebook's platform for Detectron (Mask-RCNN, DensePose, ...)



Converting Between Frameworks

- Example: develop in one framework, deploy in another
- Currently: a large range of converters, but no clear standard
- Standardized model formats are in development

from https://github.com/ysh329/deep-learning-model-convertor

convertor	tensorflow	pytorch	keras	caffe	caffe2	CNTK	chainer	mxnet
tensorflow	-	<u>pytorch-tf</u> / <u>MMdnn</u>	model-converters/ nn_toolsconvert-to- tensorflow/MMdnn	MMdnn/ nn_tools	None	<u>crosstalk/MMdnn</u>	None	<u>MMdnn</u>
pytorch	<u>pytorch2keras</u> (over Keras)	-	Pytorch2keras/ nn-transfer	Pytorch2caffe/pytorch- caffe-darknet-convert	onnx-caffe2	ONNX	None	None
keras	nn tools /convert-to- tensorflow/keras to ten sorflow/keras to tensorf low/MMdnn	<u>MMdnn</u> / <u>nn-transfer</u>	-	<u>MMdnnnn tools</u>	None	<u>MMdnn</u>	None	<u>MMdnn</u>
caffe	MMdnn/nn_tools/caffe- tensorflow	MMdnn/ pytorch- caffe-darknet- convert/ pytorch- resnet	caffe weight converter/ caffe2keras/nn tools/ kerascaffe2keras/ Deep Learning Model <u>Converter/MMdnn</u>	-	CaffeToCaffe2	<u>crosstalkcaffe/CaffeCon</u> verterMMdnn	None	<u>mxnet/tools/caffe_conve</u> <u>rter/ResNet_caffe2mxne</u> <u>t/MMdnn</u>
caffe2	None	ONNX	None	None	-	ONNX	None	None
CNTK	MMdnn	ONNX MMdnn	MMdnn	MMdnn	ONNX	-	None	MMdnn
chainer	None	chainer2pytorch	None	None	None	None	-	None
mxnet	<u>MMdnn</u>	<u>MMdnn</u>	MMdnn	MMdnn/MXNet2Caffe/ Mxnet2Caffe	None	<u>MMdnn</u>	None	-



ONNX

- Standard format for models
- Native support in development for Pytorch, Caffe2, Chainer, CNTK, and MxNet
- Converter in development for Tensorflow

 Converters available for several frameworks

MMdnn



 Common intermediate representation, but no clear standard



The end

