# Course Notes

# CreativeAI: Deep Learning for Graphics

## SIGGRAPH Asia 2018
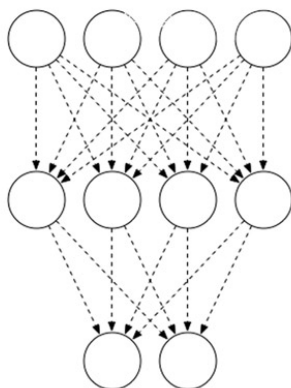
**Niloy J. Mitra**
UCL

**Iasonas Kokkinos**
UCL / Facebook

**Paul Guerrero**
UCL

**Nils Thuerey**
TU Munich

**Tobias Ritschel**
UCL

**Abstract**

In computer graphics, many traditional problems are now better handled by deep-learning based data-driven methods. In applications that operate on regular 2D domains, like image processing and computational photography, deep networks are state-of-the-art, beating dedicated hand-crafted methods by significant margins. More recently, other domains such as geometry processing, animation, video processing, and physical simulations have benefited from deep learning methods as well. The massive volume of research that has emerged in just a few years is often difficult to grasp for researchers new to this area. This tutorial gives an organized overview of core theory, practice, and graphics-related applications of deep learning.

# 1    Course Content and Syllabus

# 2   About the Lecturers

**Niloy J. Mitra**   leads the Smart Geometry Processing group in the Department of Computer Science at University College London. He received his PhD degree from Stanford University. His research interests include shape analysis, geometry processing, and computational design and fabrication. Niloy received the ACM Siggraph Significant New Researcher Award in 2013 and the BCS Roger Needham award in 2015. His work has twice been selected and featured as research highlights in the Communication of ACM, received best paper award at ACM Symposium on Geometry Processing 2014, best software SGP 2017, and Honourable Mention at Eurographics 2014.

**Iasonas Kokkinos**   obtained the Diploma of Engineering in 2001 and the Ph.D. Degree in 2006 from the School of Electrical and Computer Engineering of the National Technical University of Athens in Greece, and the Habilitation Degree in 2013 from Universit Paris-Est. He is currently a faculty at the University College London and Facebook AI Research (FAIR). His research activity is currently focused on deep learning for computer vision, focusing in particular on structured prediction for deep learning and multi-task learning architectures. He has been awarded a young researcher grant by the French National Research Agency, has served as associate editor for the Image and Vision Computing journal and the Computer Vision and Image Understanding journal, and serves regularly as a reviewer and area chair for all major computer vision conferences and journals.

**Paul Guerrero**   is a Post-Doc at University College London, working on shape analysis and image editing, combining methods from machine learning, optimization, and computational geometry. He received his PhD in computer science from Vienna University of Technology. Paul has published several research papers in high-quality journals, is a regular reviewer fo conferences and journals, and a conference IPC member.

**Nils Thuerey**   is an Associate Professor at the Technical University of Munich (TUM). He works in the field of computer graphics, with a particular emphasis on physics simulations and deep learning algorithms. After studying computer science, Nils Thuerey acquired a PhD on liquid simulations in 2006 (both at the University of Erlangen-Nuremberg). Until 2010 he held a position as a post-doctoral researcher at ETH Zurich. He received a tech-Oscar from the AMPAS in 2013 for his research on controllable smoke effects. Subsequently, he worked for three years as R&D lead at ScanlineVFX, before he started at TUM in October 2013.

**Tobias Ritschel**   is a Senior Lecturer at University College London. Previously he was a junior research group leader at the Max Planck Center for Visual Computing and Communication at Max Planck Institut Informatik. His interests include interactive and non-photorealistic rendering, human perception, and data-driven graphics. Ritschel received a PhD in computer graphics from Max Planck Institut Informatik. In 2011, he received the Eurographics PhD dissertation award and the Eurographics Young Researcher Award in 2014.

# CreativeAI:
# Deep Learning for Graphics

| Niloy Mitra | Iasonas Kokkinos | Paul Guerrero | Nils Thuerey | Tobias Ritschel |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

---

# People



| Niloy Mitra | Iasonas Kokkinos | Paul Guerrero | Nils Thuerey | Tobias Ritschel |
|---|---|---|---|---|

# Timetable

|  |  | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|:---:|:---:|:---:|:---:|:---:|
| **Theory and Basics** | Introduction | X | X | X | X | X |
|  | Theory | X |  |  | X |  |
|  | NN Basics | X | X |  |  |  |
|  | Alternatives to Direct Supervision |  |  | X |  |  |
|  | *15 min. break* |  |  |  |  |  |
| **State of the Art** | Feature Visualization |  |  |  |  | X |
|  | Image Domains |  | X |  |  | X |
|  | 3D Domains |  |  | X |  | X |
|  | Motion and Physics | X |  |  | X |  |

# Code Examples

**PCA/SVD basis**
**Linear Regression**
**Polynomial Regression**
**Stochastic Gradient Descent vs. Gradient Descent**
**Multi-layer Perceptron**
**Edge Filter 'Network'**
**Convolutional Network**
**Filter Visualization**
**Weight Initialization Strategies**
**Colorization Network**
**Autoencoder**
**Variational Autoencoder**
**Generative Adversarial Network**

http://geometry.cs.ucl.ac.uk/dl4g/

# Two-way Communication

- *This tutorial is given for the first time!*

- Our aim is to convey what we found to be relevant so far.

- You are invited/encouraged to give feedback
  - On-line form
  - Speakup. Please send us your criticism/comments/suggestions
  - Ask questions, please!

- **Thanks to many people who helped so far with slides/comments.**

# Course Overview

- **Part I: Introduction and ML Basics**

- Part II: Supervised Neural Networks: Theory and Applications

- Part III: Unsupervised Neural Networks: Theory and Applications

- Part IV: Beyond Image Data

# Representations in CG

• Images (e.g., pixel grid)

• Volume (e.g., voxel grid)

• Meshes (e.g., vertices/edges/faces)

• Animation (e.g., skeletal positions over time; cloth dynamics over time)

• Pointclouds (e.g., point arrays)

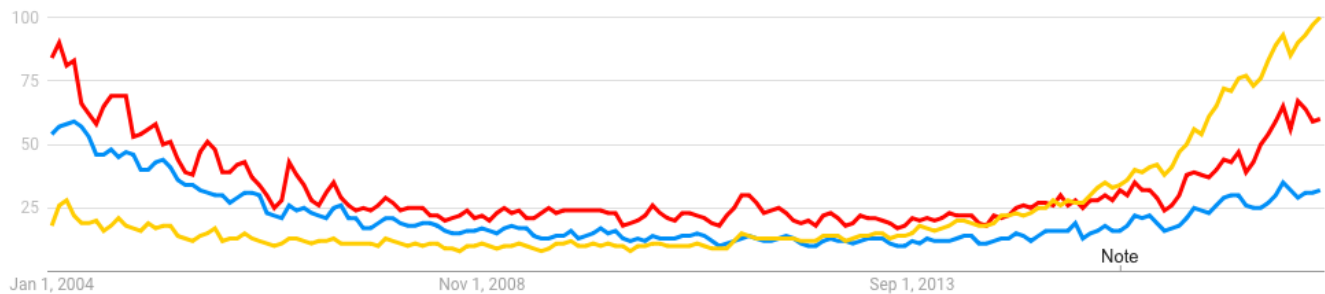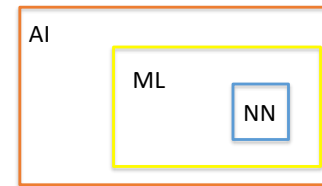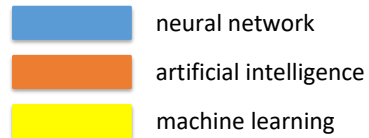• Physics simulations (e.g., fluid flow over space/time)

SIGGRAPH
ASIA 2018
T O K Y O                                                                                                7

# Problems in Computer Graphics

• Feature detection (image features, point features)    $\mathbb{R}^{m \times m} \to \mathbb{Z}$

• Denoising, Smoothing, etc.    $\mathbb{R}^{m \times m} \to \mathbb{R}^{m \times m}$

• Embedding, Distance computation    $\mathbb{R}^{m \times m, m \times m} \to \mathbb{R}^{d}$

• Rendering    $\mathbb{R}^{m \times m} \to \mathbb{R}^{m \times m}$

• Animation    $\mathbb{R}^{3m \times t} \to \mathbb{R}^{3m}$

• Physical simulation    $\mathbb{R}^{3m \times t} \to \mathbb{R}^{3m}$

• Generative models    $\mathbb{R}^{d} \to \mathbb{R}^{m \times m}$

SIGGRAPH
ASIA 2018
T O K Y O                                                                                                8

# Rise of Machine Learning



neural network

artificial intelligence

machine learning

AI

ML

NN

100

75

50

25

Jan 1, 2004          Nov 1, 2008          Sep 1, 2013

Note

9

# Data-driven Algorithms (Supervised)

**Labelled data
(supervision data)** → ML algorithm

↓

**Test data
(run-time data)** → Trained model → Prediction
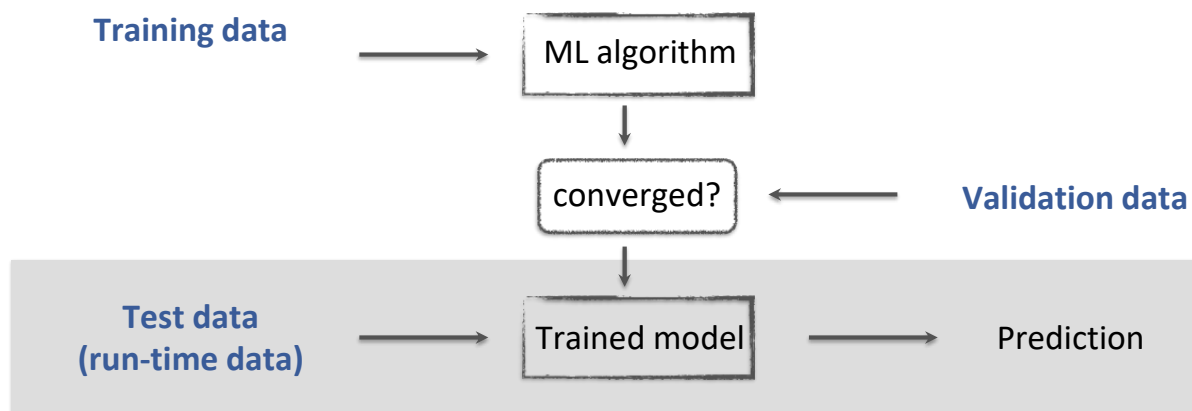
1
0

# Data-driven Algorithms (Supervised)

Labelled data
(supervision data) → ML algorithm

**Validation data
(supervision data)**

converged? ←

Test data
(run-time data) → Trained model → Prediction

**Implementation Practice: Training: 70%; Validation: 15%; Test 15%**

SIGGRAPH ASIA 2018 TOKYO

1
1

# Data-driven Algorithms (Unsupervised)

**Training data** → ML algorithm

converged? ← **Validation data**

**Test data
(run-time data)** → Trained model → Prediction

**Implementation Practice: Training: 70%; Validation: 15%; Test 15%**

SIGGRAPH ASIA 2018 TOKYO

1
2

# Various ML Approaches (Supervised approaches)

Input data

13

# Rise of Learning

- 1958:     Perceptron
- 1974:     Backpropagation
- 1981:     Hubel & Wiesel wins Nobel prize for 'visual system'
- 1990s:   SVM era
- 1998:     CNN used for handwriting analysis
- **2012:     AlexNet wins ImageNet**

14

# Rise of Machine Learning (in Graphics)

█████ machine learning

█████ neural network

**SIG+SA+EG+SGP+EGSR**

14%
12%
10%
8%
6%
4%
2%
0%

2013            2017

**Eurographics**

14%
12%
10%
8%
6%
4%
2%
0%

2013            2017

15

---

# What is Special about Graphics?

- **Image Processing** (image translation tasks)

- Many sources of input data — **model building**
  (e.g., images, scanners, motion capture)

- Many sources of **synthetic data** — can serve as supervision data
  (e.g., rendering, animation)

- Many problems in **generative models**

16

# End-to-end: **Features**

- *Old days*
  - First some handy features were extracted, e.g. edges or corners (hand-crafted)
  - Second, some AI was ran on that features (optimized)
- *Now*
  - End-to-end
  - Move away from hand-crafted representations

input image     edge image     $2^1/_2$-D sketch     3-D model



1
7

# End-to-end: **Loss**

- *Old days*
  - Evaluation came after
  - It was a bit optional:
    - You might still have a good algorithm without a good way of quantifying it
    - Evaluation helped publishing
- *Now*
  - It is essential and build-in
  - If the loss is not good, the result is not good
  - Evaluation happens automatically

- While still much is left to do, this makes graphics much more reproducable

1
8

# End-to-end: **Data**

- *Old days*
  - Test with some toy examples
  - Deploy on real stuff
  - Maybe collect some data later

- *Now*
  - Test and deploy need to be as identical as you can
  - Need to collect data first
  - No two steps



1
9

# **Examples in Graphics**

Geometry

Image
manipulation

Animation

Rendering

2
0

# Examples in Graphics

Geometry

Procedural
modelling

Mesh segmentation

Learning
deformations

Colorization

Sketch
simplification

## Image
## manipulation

Animation

BRDF estimation

Fluid

Boxification

Real-time rendering

## Animation

## Rendering

Denoising

Facial animation

PCD processing

2
1

# Examples in Graphics



Procedural
modelling

Mesh segmentation

Learning
deformations

Colorization

Sketch
simplification

BRDF estimation

Fluid

Animation

Boxification

Real-time rendering

Denoising

Facial animation

PCD processing

2
2

# Course Information (slides/code/comments)



http://geometry.cs.ucl.ac.uk/creativeai/

SIGGRAPH Asia Course **CreativeAI: Deep Learning for Graphics**

# CreativeAI: Deep Learning for Graphics

# Theory

| Niloy Mitra | Iasonas Kokkinos | Paul Guerrero | Nils Thuerey | Tobias Ritschel |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

# Timetable

| | | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|---|---|---|---|---|
| **Theory and Basics** | Introduction | X | X | X | X | X |
| | Theory | X | | | X | |
| | NN Basics | X | X | | | |
| | Alternatives to Direct Supervision | | | X | | |
| | *15 min. break* | | | | | |
| **State of the Art** | Feature Visualization | | | | | X |
| | Image Domains | | X | | | X |
| | 3D Domains | | | X | | X |
| | Motion and Physics | X | | | X | |

# Machine Learning

**Machine learning** is a field of **computer science** that uses statistical techniques to give computer systems the ability to *learn* (i.e., progressively improve performance on a specific task) with **data**, without being explicitly programmed.

*'ML' coined by Arthur Samuel, 1959.*

data  $\longrightarrow$  model building  $\longrightarrow$  prediction

3

# Machine Learning Variants

- **Supervised**
  - Classification
  - Regression
  - Data consolidation
- **Unsupervised**
  - Clustering
  - Dimensionality Reduction
- **Weakly supervised/semi-supervised**
  Some data supervised, some unsupervised
- **Reinforcement learning**
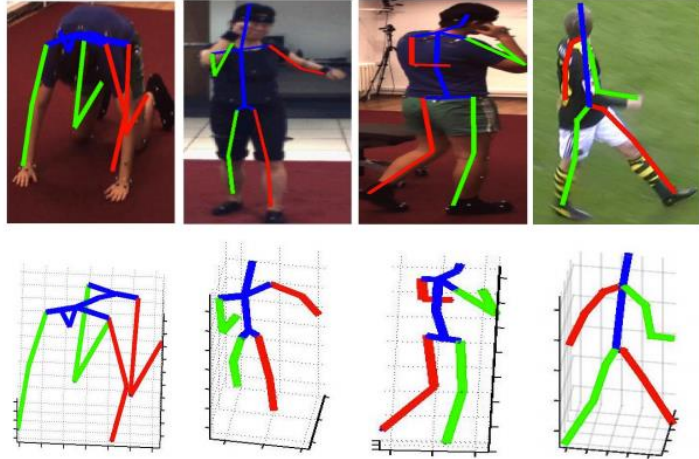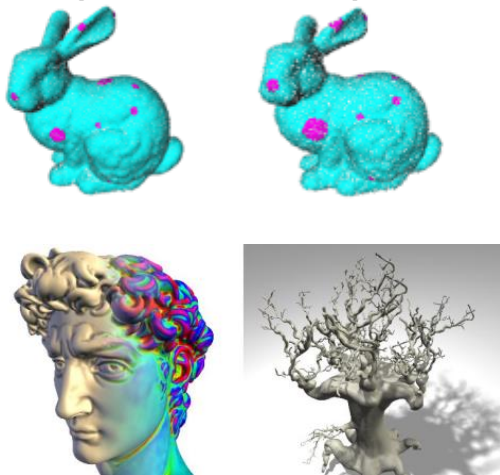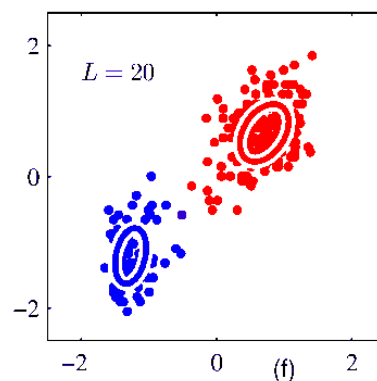  Supervision: sparse reward for a sequence of decisions

4

# Machine Learning Variants

- **Supervised**
  - **Classification**
  - Regression
  - Data consolidation
- **Unsupervised**
  - Clustering
  - Dimensionality Reduction
- **Weakly supervised/semi-supervised**
  Some data supervised, some unsupervised
- **Reinforcement learning**
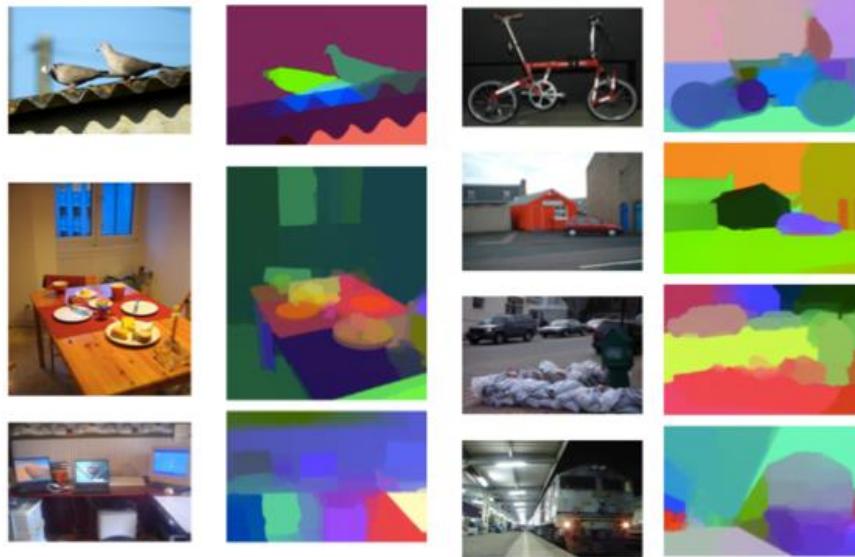  Supervision: sparse reward for a sequence of decisions

SIGGRAPH
ASIA 2018
TOKYO

5

# Classification Examples

- Digit Recognition



- Spam Detection



- Face detection



SIGGRAPH
ASIA 2018
TOKYO

6

# Segmentation + Classification in Real Images



**Evaluation measures:** Confusion matrix, ROC curve, precision, recall, etc.

7

# `Faceness' Function: Classifier



**background**

**decision boundary**

*face*

8

# Face Detection



CMU Science Lab

9

# Machine Learning Variants

- **Supervised**
  - Classification
  - **Regression**
  - Data consolidation
- **Unsupervised**
  - Clustering
  - Dimensionality Reduction
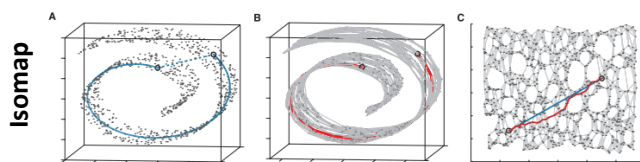- **Weakly supervised/semi-supervised**
  Some data supervised, some unsupervised
- **Reinforcement learning**
  Supervision: sparse reward for a sequence of decisions

10

# Human Face/Pose Estimation



[Blanz and Vetter, Siggraph, 1999]

11

# Regression: Model Estimation

[Mitra et al. SoCG, 2003]



[Guennebaud et al., Siggraph, 2007]                    [Zwicker et al., EGSR, 2005]

12

# Machine Learning Variants

- **Supervised**
  - Classification
  - Regression
  - Data consolidation
- **Unsupervised**
  - **Clustering**
  - Dimensionality Reduction
- **Weakly supervised/semi-supervised**
  Some data supervised, some unsupervised
- **Reinforcement learning**
  Supervision: sparse reward for a sequence of decisions

13

# Clustering: Color Points According to X



14

# Clustering Examples: Image Segmentation using NCuts

# Clustering Examples



[Chu et al., TVCG, 2009]

[Zheng et al., Eurographics, 2014]

# Machine Learning Variants

- **Supervised**
  - Classification
  - Regression
  - Data consolidation
- **Unsupervised**
  - Clustering
  - **Dimensionality Reduction**
- **Weakly supervised/semi-supervised**
  Some data supervised, some unsupervised
- **Reinforcement learning**
  Supervision: sparse reward for a sequence of decisions

17

# Dimensionality Reduction (Manifold Learning)

**Isomap**



[Tenenbaum et al., Science, 2000]

[Yang et al., TOG, 2011]

**Face Manifold**



parameterized embedding

[Averkiou et al., Eurographics, 2014]

18

# Example of Nonlinear Manifold: Faces



$$\mathbf{x}_1 \qquad \frac{1}{2}(\mathbf{x}_1 + \mathbf{x}_2) \qquad \mathbf{x}_2$$

19

# Morphing (Interpolation in Shape Space)



Input poses

[Kilian et al., Siggraph, 2007]      20

# Moving Along Learned Face Manifold



**Trajectory along the "male" dimension**

**Trajectory along the "young" dimension**

[Lample et. al. Fader Networks, NIPS 2017]

21

# Notations: Vectors and Matrices

- linear **independence**; **rank** of a matrix
- **span** of a matrix

| | |
|---|---|
| vector | $\mathbf{x}$ |
| matrix | $\mathbf{A}_{m \times n} = [\mathbf{a}_1 \ldots \mathbf{a}_n]$ |
| linear equation | $\mathbf{A}\mathbf{x} = \mathbf{b}$ |
| inner prod. | $< \mathbf{x}, \mathbf{y} >= \mathbf{x}^T \mathbf{y}$ $\qquad \|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}}$ |
| | $\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\|\|\mathbf{y}\| \cos(\theta)$ |

22

# Notations: Vectors and Matrices (cont.)

$$\|\mathbf{x}\|_p = (|x_1|^p + |x_2|^p + \dots)^{1/p}$$

$$L_1, L_2, L_p, L_\infty$$

$$\|\mathbf{x}\|_p = \max\{|x_1|, |x_2|, \dots\} \qquad p = \infty$$

range $\qquad \mathcal{R}(\mathbf{A}) = \{\mathbf{A}\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\}$

null space $\quad \mathcal{N}(\mathbf{A}) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A}\mathbf{x} = 0\}$

23

# Eigenvectors and Eigenvalues

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

$$\mathbf{T} = [\mathbf{v}_1 \ \mathbf{v}_2 \dots]$$

$$\mathbf{A}\mathbf{e}_i = \lambda_i \mathbf{e}_i$$

$$\mathbf{T}^{-1} A \mathbf{T} = \mathrm{diag}(\lambda_1, \lambda_2, \dots)$$

- All eigenvalues of symmetric matrices are real.
- Any real symmetric nxn matrix has a set of **n** mutually orthogonal eigenvectors.

24

# Code Example



```python
rng = np.random.RandomState(10)
X = np.dot(rng.rand(2, 2), rng.randn(2, 500)).T

mean_vec = np.mean(X, axis=0)
cov_mat = (X - mean_vec).T.dot((X - mean_vec)) / (X.shape[0]-1)
eig_vals, eig_vecs = np.linalg.eig(cov_mat)
```

25

# Morphable Faces



26

# Singular Value Decomposition (SVD)

- Very useful for matrix manipulation.
- Used for robust numerical computation.

$$A = U\Sigma V^T$$

scaling            rotation

$$A = A^T = U\Sigma U^T$$

# Code Example



```
mean_vec = np.mean(X, axis=0)
cov_mat = (X - mean_vec).T.dot((X - mean_vec)) / (X.shape[0]-1)
matU, sigma, matV = np.linalg.svd(cov_mat)
```

# Differentiation (chain rule recap)

$z = f \circ g(x) = f(g(x))$

$z = f(y)$

$y = g(x)$

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} = f'(y)g'(x) = f'(g(x))g'(x)$$

$z = \sin(5x)$

$\quad = \dfrac{d\sin(5x)}{d(5x)}\dfrac{d(5x)}{dx}$

$\quad = 5\cos(5x)$

$f(x_0)$

slope: $f'(x_0)$

$x_0$

SIGGRAPH
ASIA 2018
T O K Y O

29

# Derivative Matrix

$$f : \mathbb{R}^n \to \mathbb{R}^m$$

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \qquad \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_j} = \begin{bmatrix} \frac{\partial \mathbf{f}_1(\mathbf{x})}{\partial x_j} \\ \vdots \\ \frac{\partial \mathbf{f}_m(\mathbf{x})}{\partial x_j} \end{bmatrix}_{m \times 1}$$

*Jacobian matrix*

$$\mathbf{L} = D\mathbf{f} = \begin{bmatrix} \dfrac{\partial f(\mathbf{x})}{\partial x_1} & \cdots & \dfrac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix}_{m \times n}$$

SIGGRAPH
ASIA 2018
T O K Y O

30

# Regression: Continuous Output

31

# Learning a Function

$$y = f_w(x)$$

32

# Learning a Function

**prediction**     **method**

$$y = f_w(x)$$

**parameters**     **input**

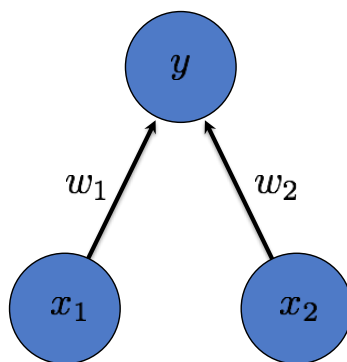| | | | |
|---|---|---|---|
| Calculus | $x \in \mathbb{R}$ | **Classification:** | $y \in \{0, 1\}$ |
| Vector calculus | $\mathbf{x} \in \mathbb{R}^d$ | **Regression:** | $y \in \mathbb{R}$ |

*Machine learning: can work also for discrete inputs, strings, images, meshes, animations, …*

SIGGRAPH
ASIA 2018
TOKYO

33

# Learning a Simple Separator/Classifier



separating hyperplane

$$y = f(w_1 x_1 + w_2 x_2)$$

SIGGRAPH
ASIA 2018
TOKYO

34

# Combining Simple Functions/Classifiers

2 layers of
trainable
weights

convex region

# Combining Simple Functions/Classifiers

3 layers of
trainable
weights

complex polygons

# Learning a Function: Modeling

method

prediction

$$y = f_w(x) = f(x; w)$$

parameters       input

$$w \in \mathbb{R}$$
$$\mathbf{w} \in \mathbb{R}^K$$

37

# Regression

1. Least Squares fitting

2. Nonlinear error function and gradient descent

3. Perceptron training

38

# Regression

**1. Least Squares fitting**

2. Nonlinear error function and gradient descent

3. Perceptron training

---

# Assumption: Linear Function

$$y = f_{\mathbf{w}}(\mathbf{x}) = f(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w}^T \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{d=1}^{D} \mathbf{w}_d \mathbf{x}_d$$

$$\mathbf{x} \in \mathbb{R}^D, \mathbf{w} \in \mathbb{R}^D$$

# Reminder: Linear Classifier

$\mathbf{x}_i$ positive: $\mathbf{x}_i \cdot \mathbf{w} \geq 0$

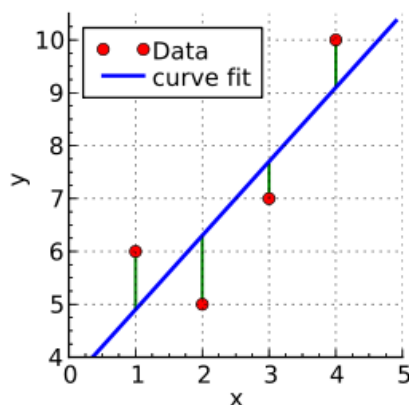$\mathbf{x}_i$ negative: $\mathbf{x}_i \cdot \mathbf{w} < 0$

*supervised setting*

labelled input

$$y_t = \begin{cases} +1 & \bullet \\ -1 & \bullet \end{cases}$$

feature coordinate

feature coordinate

41

# Which Line to Pick?

$\mathbf{x}_i$ positive: $\mathbf{x}_i \cdot \mathbf{w} \geq 0$

$\mathbf{x}_i$ negative: $\mathbf{x}_i \cdot \mathbf{w} < 0$

*supervised setting*

labelled input

$$y_t = \begin{cases} +1 & \bullet \\ -1 & \bullet \end{cases}$$

feature coordinate

feature coordinate

42

# Linear Regression in 1D



$$\mathcal{S} = \{(x^i, y^i)\}, \quad i = 1 \ldots, N$$

Training set: input–output pairs

$$x^i \in \mathbb{R}, \quad y^i \in \mathbb{R}$$
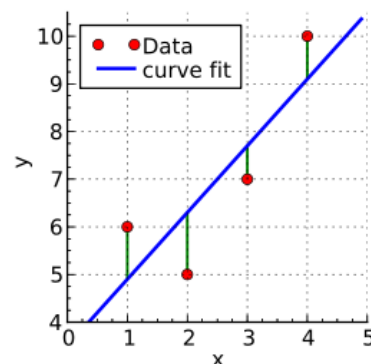
43

# Linear regression in 1D

$$y^i = w_0 + w_1 x_1^i + \epsilon^i \qquad w_0 \text{ bias}$$

$$= w_0 x_0^i + w_1 x_1^i + \epsilon^i, \quad x_0^i = 1, \quad \forall i$$

$$= \mathbf{w}^T \mathbf{x}^i + \boxed{\epsilon^i}$$
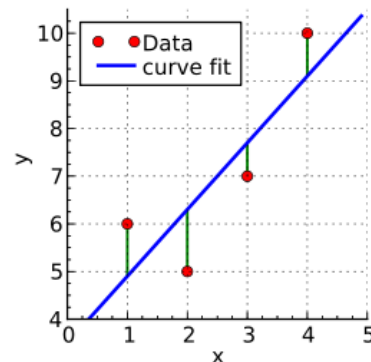$$\text{noise}$$



44

## Sum of Square Errors *(MSE without the mean)*

$$y^i = \mathbf{w}^T \mathbf{x}^i + \epsilon^i$$

Loss function: sum of squared errors

$$L(\mathbf{w}) = \sum_{i=1}^{N} (\epsilon^i)^2$$

In two variables:

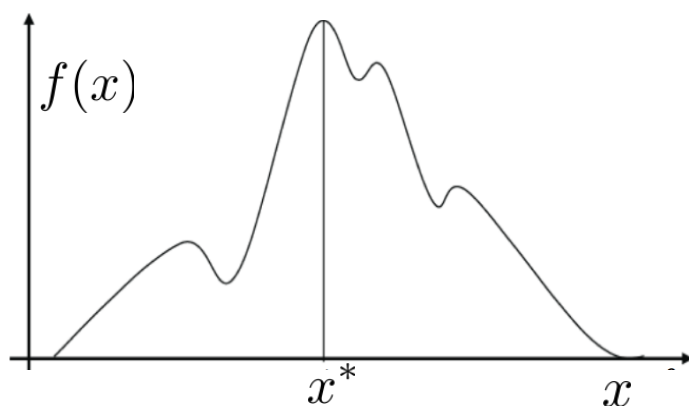$$L(w_0, w_1) = \sum_{i=1}^{N} \left[ y^i - \left( w_0 x_0^i + w_1 x_1^i \right) \right]^2$$

*Question: what is the best (or least bad) value of **w**?*

SIGGRAPH
ASIA 2018
TOKYO

45

---

## Calculus 101

$f(x)$

$x^*$      $x$

$$x^* = \operatorname{argmax}_x f(x)$$

SIGGRAPH
ASIA 2018
TOKYO

46

# Local Extrema Condition



$$x^* = \operatorname{argmax}_x f(x) \quad \rightarrow \quad f'(x^*) = 0$$

47

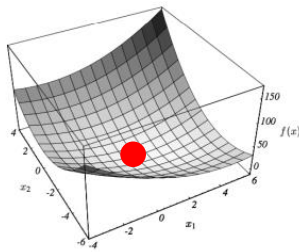# Local Extrema Condition



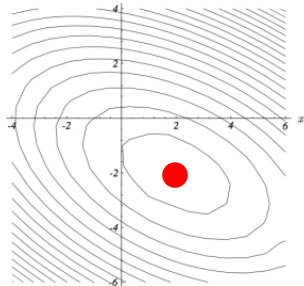$$x^* = \operatorname{argmax}_x f(x) \quad \rightarrow \quad f'(x^*) = 0$$
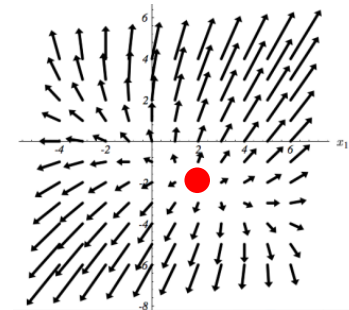
48

# Vector Calculus 101



$$f(\mathbf{x})$$

**2D function graph**

$$f(\mathbf{x}) = c$$

**isocontours**

$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

**gradient field**

at minimum of function: $\quad \nabla f(\mathbf{x}) = \mathbf{0}$

49

# Line Fitting

$$L(w_0, w_1) = \sum_{i=1}^{N} \left[ y^i - \left( w_0 x_0^i + w_1 x_1^i \right) \right]^2$$

$$\frac{\partial L(w_0, w_1)}{\partial w_0} = \sum_{i=1}^{N} \frac{\partial \left[ y^i - \left( w_0 x_0^i + w_1 x_1^i \right) \right]^2}{\partial w_0}$$

$$\frac{\partial L(w_0, w_1)}{\partial w_0} = \sum_{i=1}^{N} \frac{\partial [z^i]^2}{\partial z^i} \frac{\partial z^i}{\partial w_0} = \sum_{i=1}^{N} (2z^i)(-x_0^i)$$

$$= -2 \sum_{i=1}^{N} (y^i - (w_0 x_0^i + w_1 x_1^i)) x_0^i$$

$$\boxed{z^i = y^i - (w_0 x_0^i + w_1 x_1^i)}$$

50

# Line Fitting (continued)

$$\frac{\partial L(w_0, w_1)}{\partial w_0} = \sum_{i=1}^{N} \frac{\partial \left[ y^i - \left( w_0 x_0^i + w_1 x_1^i \right) \right]^2}{\partial w_0}$$

$$= -2 \sum_{i=1}^{N} \left( y^i x_0^i - w_0 x_0^i x_0^i - w_1 x_1^i x_0^i \right)$$

$$\frac{\partial L(w_0, w_1)}{\partial w_0} = 0 \qquad \sum_{i=1}^{N} y^i x_0^i = w_0 \sum_{i=1}^{N} x_0^i x_0^i + w_1 \sum_{i=1}^{N} x_1^i x_0^i$$

# Line Fitting (continued)

$$\sum_{i=1}^{N} y^i x_0^i = w_0 \sum_{i=1}^{N} x_0^i x_0^i + w_1 \sum_{i=1}^{N} x_1^i x_0^i$$

$$\sum_{i=1}^{N} y^i x_1^i = w_0 \sum_{i=1}^{N} x_0^i x_1^i + w_1 \sum_{i=1}^{N} x_1^i x_1^i$$

2x2 system
of equations

$$\begin{bmatrix} \sum_{i=1}^{N} y^i x_0^i \\ \sum_{i=1}^{N} y^i x_1^i \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{N} x_0^i x_0^i & \sum_{i=1}^{N} x_0^i x_1^i \\ \sum_{i=1}^{N} x_0^i x_1^i & \sum_{i=1}^{N} x_1^i x_1^i \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

# Line Fitting (continued)

$$\begin{bmatrix} \sum_{i=1}^{N} y^i x_0^i \\ \sum_{i=1}^{N} y^i x_1^i \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{N} x_0^i x_0^i & \sum_{i=1}^{N} x_0^i x_1^i \\ \sum_{i=1}^{N} x_0^i x_1^i & \sum_{i=1}^{N} x_1^i x_1^i \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \mathbf{w}$$

$$\mathbf{y} = \begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix} \qquad \mathbf{X} = \begin{bmatrix} x_0^1 & x_1^1 \\ \vdots & \vdots \\ x_0^N & x_2^N \end{bmatrix}$$

*Normal Equation*

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

SIGGRAPH
ASIA 2018
TOKYO

53

# Code Example

```python
import numpy as np
from numpy import array
from numpy import matmul
from numpy.linalg import inv
from numpy.random import rand
from matplotlib import pyplot

# generate data on a line perturbed with some noise
noise_margin= 2
w = rand(2,1) # w[0] is random constant term (offset from origin), w[1] is random linear term (slope)
x  = np.linspace(-5,5,20)
y = w[0] + w[1]*x + noise_margin*rand(len(x))

# create the design matrix: the x data, and add a column of ones for the constant term
X = np.column_stack( [np.ones([len(x), 1]), x.reshape(-1, 1)] )

# These are the normal equations in matrix form: w = (X' X)^-1 X' y
w_est = matmul(inv(matmul(X.transpose(),X)),X.transpose()).dot(y)

# For ridge regression, use regularizer
#weight = 0.01
#w_est = matmul(inv(matmul(X.transpose(),X) + weight*np.identity(2)),X.transpose()).dot(y)

# evaluate the x values in the fitted model to get estimated y values
y_est = w_est[0] + w_est[1]*x

# visualize the fitted model
pyplot.scatter(x, y, color='red')
pyplot.plot(x, y_est, color='blue')
pyplot.show()
```
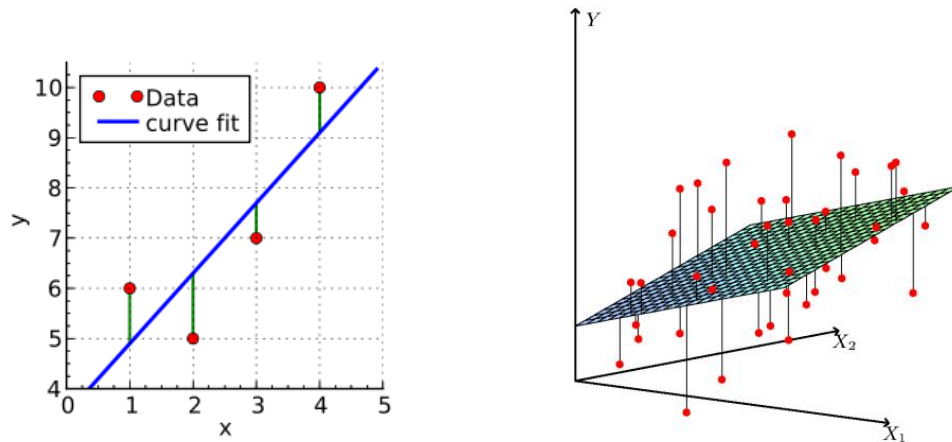
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

SIGGRAPH
ASIA 2018
TOKYO

54

# Linear Regression (Line/Plane Fitting)

# LS Solution for Regression

$$L(\mathbf{w}) = \sum_{i=1}^{N}(y^i - \mathbf{w}^T\mathbf{x}^i)^2 = \sum_{i=1}^{N}(\epsilon^i)^2$$

$$L(\mathbf{w}) = \begin{bmatrix} \epsilon^1 & \epsilon^2 & \dots & \epsilon^N \end{bmatrix} \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \vdots \\ \epsilon^N \end{bmatrix}$$
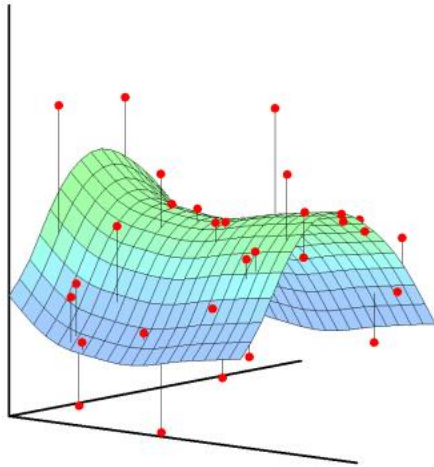
$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T\boldsymbol{\epsilon} \qquad \mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

# Generalized Linear Regression



$$\mathbf{x} \rightarrow \boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \vdots \\ \phi_M(\mathbf{x}) \end{bmatrix}$$

known nonlinearity

57

# 1D Example: k-th Degree Polynomial Fitting

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} 1 \\ x \\ \vdots \\ (x)^K \end{bmatrix}$$



$$\langle \mathbf{w}, \boldsymbol{\phi}(x) \rangle = w_0 + w_1 x + \ldots + w_k (x)^K$$

58

# Generalized Linear Regression

$$L(\mathbf{w}) = \sum_{i=1}^{N}(y^i - \mathbf{w}^T\boldsymbol{\phi}(\mathbf{x}^i))^T = \sum_{i=1}^{N}(\epsilon^i)^2$$

$$\begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^N \end{bmatrix} = \begin{bmatrix} \boldsymbol{\phi}(\mathbf{x}^1)^T \\ \boldsymbol{\phi}(\mathbf{x}^2)^T \\ \hline \vdots \\ \hline \boldsymbol{\phi}(\mathbf{x}^N)^T \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix} + \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \vdots \\ \epsilon^N \end{bmatrix}$$

**Nx1**      **NxM**      **Mx1**      **Nx1**

$$\boldsymbol{\phi}(\mathbf{x}) : \mathbb{R}^D \to \mathbb{R}^M$$

SIGGRAPH
ASIA 2018
T O K Y O

59

# LS Solution for Linear Regression

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T\boldsymbol{\epsilon}$$

$$\mathbf{X} = \begin{bmatrix} (\mathbf{x}^1)^T \\ (\mathbf{x}^2)^T \\ \hline \vdots \\ \hline (\mathbf{x}^N)^T \end{bmatrix}$$

$$\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

SIGGRAPH
ASIA 2018
T O K Y O

60

# LS Solution for Generalized Linear Regression

$$\mathbf{y} = \mathbf{\Phi}\mathbf{w} + \boldsymbol{\epsilon}$$

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon}$$

$$\mathbf{\Phi} = \begin{bmatrix} \phi(\mathbf{x}^1)^T \\ \phi(\mathbf{x}^2)^T \\ \vdots \\ \phi(\mathbf{x}^N)^T \end{bmatrix}$$

$$\mathbf{w}^* = (\mathbf{\Phi}^T\mathbf{\Phi})^{-1}\mathbf{\Phi}\mathbf{y}$$

SIGGRAPH
ASIA 2018
TOKYO

61

# Code Example



```python
import numpy as np
from numpy import array
from numpy import matmul
from numpy.linalg import inv
from numpy.random import rand
from matplotlib import pyplot

# generate data on a line perturbed with some noise
noise_margin= 3
w = 2*rand(3,1) # w[0] is random constant term (offset from origin), w[1] is random linear term, w[2] is random quadratic term
x = np.linspace(-5,5,20)
y = w[0] + w[1]*x + w[2]*x**2 + noise_margin*rand(len(x))

# create the design matrix: the x data, and add a column of ones for the constant term
X = np.column_stack( [np.ones([len(x), 1]), x.reshape(-1, 1), (x**2).reshape(-1, 1)] )

# These are the normal equations in matrix form: w = (X' X)^-1 X' y
w_est = matmul(inv(matmul(X.transpose(),X)),X.transpose()).dot(y)

# evaluate the x values in the fitted model to get estimated y values
y_est = w_est[0] + w_est[1]*x + w_est[2]*x**2

# visualize the fitted model
pyplot.scatter(x, y)
pyplot.plot(x, y_est, color='red')
pyplot.show()
```

$$\mathbf{w}^* = (\mathbf{\Phi}^T\mathbf{\Phi})^{-1}\mathbf{\Phi}\mathbf{y}$$

SIGGRAPH
ASIA 2018
TOKYO

62

# Underfitting vs. Overfitting



---

# Tuning Model's Complexity

A *flexible model* approximates the target function well in the training set

but can "**overtrain**" and have poor performance on the test set ("variance").

A *rigid model*'s performance is more predictable in the test set

but the model may not be good even on the training set ("bias").

# Regularized Linear Regression

$$\boldsymbol{\epsilon} = \mathbf{y} - \boldsymbol{\Phi}\mathbf{w}$$    residual vector

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T\boldsymbol{\epsilon}$$    linear regression: minimize model error

Complexity term:    $$R(\mathbf{w}) \doteq \|\mathbf{w}\|_2^2 = \mathbf{w}^T\mathbf{w}$$
(regularizer)

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T\boldsymbol{\epsilon} + \lambda\mathbf{w}^T\mathbf{w}$$

⬆ "data fidelity"    ⬆ complexity

minimum remains to be determined    scalar, remains to be determined

SIGGRAPH
ASIA 2018
T O K Y O

66

# Least Squares Solution

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T\boldsymbol{\epsilon}$$
$$= (\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$= \mathbf{y}^T\mathbf{y} - 2\mathbf{y}^T\mathbf{X}\mathbf{w} + \mathbf{w}^T\mathbf{X}^T\mathbf{X}\mathbf{w}$$

**Condition for minimum:**

$$\nabla L(\mathbf{w}^*) = \mathbf{0}$$

$$-2\mathbf{X}^T\mathbf{y} + 2\mathbf{X}^T\mathbf{X}\mathbf{w}^* = \mathbf{0}$$

$$\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

SIGGRAPH
ASIA 2018
T O K Y O

67

# Ridge regression: L2-regularized Linear Regression

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} + \lambda \mathbf{w}^T \mathbf{w}$$

$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w}^T \mathbf{I} \mathbf{w}$$

as before, for linear regression                                   identity matrix

$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right) \mathbf{w}$$

**Condition for minimum:**

$$\nabla L(\mathbf{w}^*) = \mathbf{0}$$

$$-2\mathbf{X}^T \mathbf{y} + 2(\mathbf{X}^T \mathbf{X} + \lambda I) \mathbf{w}^* = \mathbf{0}$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$\lambda$: "hyperparameter"

68

# Bias-Variance Tradeoff (function of λ)



69

# Selecting λ with Cross-validation

- Cross validation technique
  - Exclude part of the training data from parameter estimation
  - Use them only to predict the test error

- K-fold cross validation:
  - K splits, average K errors

- Use cross-validation for different values of λ
  - pick value that minimizes cross-validation error

  **Least glorious, most effective of all methods**



70

---

## Form of posterior distribution

Bernoulli-type conditional distribution

$$P(Y = 1|X = \mathbf{x}; \mathbf{w}) = \quad f(\mathbf{x}, \mathbf{w})$$
$$P(Y = 0|X = \mathbf{x}; \mathbf{w}) = \quad 1 - f(\mathbf{x}, \mathbf{w}) \quad \Big\} \rightarrow$$
$$P(Y = y|X = \mathbf{x}; \mathbf{w}) = f(\mathbf{x}, \mathbf{w})^y (1 - f(\mathbf{x}, \mathbf{w}))^{1-y}$$

Particular choice of form of f:

$$P(Y = 1|X = \mathbf{x}; \mathbf{w}) = g(\mathbf{w}^T \mathbf{x})$$

Sigmoidal:     $g(\alpha) = \dfrac{1}{1 + \exp(-a)}$

"squashing function":
$$-\infty \rightarrow 0$$
$$+\infty \rightarrow 1$$

71

# Logistic vs Linear Regression



Logistic Regression          Linear Regression

# From Two to Many

• How about multi-class classification?

# Multiple Classes & Linear Regression

C classes: one-of-c coding (or one-hot encoding)

4 classes, i-th sample is in 3$^{rd}$ class:
$$\mathbf{y}^i = (0, 0, 1, 0)$$

Matrix notation:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}^1 \\ \vdots \\ \mathbf{y}^N \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 \mid \ldots \mid \mathbf{y}_C \end{bmatrix} \qquad \text{where} \quad \mathbf{y}_c = \begin{bmatrix} y_c^1 \\ \vdots \\ y_c^N \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1 \mid \ldots \mid \mathbf{w}_C \end{bmatrix}$$

Loss function:

$$L(\mathbf{W}) = \sum_{c=1}^{C} (\mathbf{y}_c - \mathbf{X}\mathbf{w}_c)^T (\mathbf{y}_c - \mathbf{X}\mathbf{w}_c)$$

Least squares fit (decouples per class):

$$\mathbf{w}_c^* = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_c$$

SIGGRAPH ASIA 2018 TOKYO

74

---

# Linear Regression Masking Problem

**Class 1**          **Class 2**          **Class 3**

$$y^1 = \mathbf{x}\mathbf{w}^1$$

$$y^2 = \mathbf{x}\mathbf{w}^2$$

$$y^3 = \mathbf{x}\mathbf{w}^3$$

$x$

One linear discriminant per class:  $s_c(\mathbf{x}) = \mathbf{w}_c^T \mathbf{x}$

Nothing ever gets assigned to class 2!

2D version:

SIGGRAPH ASIA 2018 TOKYO

75

# Multiple classes & Logistic regression

Soft maximum (softmax) of competing classes:

$$P(y = c | \mathbf{x}; \mathbf{W}) = \frac{\exp(\mathbf{w}_c^T \mathbf{x})}{\sum_{c'=1}^{C} \exp(\mathbf{w}_{c'}^T \mathbf{x})} \doteq g_c(\mathbf{x}, \mathbf{W})$$



Discriminants (inputs)    Softmax (outputs)

76

# Logistic vs Linear Regression, n>2 classes



Linear regression    Logistic regression

Logistic regression does not exhibit the masking problem

77

# LS Solution (in vector form)

$$L(\mathbf{w}) = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon}$$
$$= (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$
$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w}$$

**Condition for minimum:**

$$\nabla L(\mathbf{w}^*) = \mathbf{0}$$
$$-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\mathbf{w}^* = \mathbf{0}$$
$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

SIGGRAPH
ASIA 2018
T O K Y O

78

# Gradient of Cross-entropy Loss

$$L(\mathbf{w}) = -\sum_{i=1}^{N} y^i \log g(\mathbf{w}^T \mathbf{x}^i) + (1 - y^i) \log(1 - g(\mathbf{w}^T \mathbf{x}^i))$$

$$\frac{\partial L(\mathbf{w})}{\partial w_k} = -\sum_{i=1}^{N} \left[ y^i \frac{1}{g(\mathbf{w}^T \mathbf{x}^i)} \frac{\partial g(\mathbf{w}^T \mathbf{x}^i)}{\partial w_k} + (1 - y^i) \frac{1}{1 - g(\mathbf{w}^T \mathbf{x}^i)} \left( -\frac{\partial g(\mathbf{w}^T \mathbf{x}^i)}{\partial w_k} \right) \right]$$

using
$$g(x) = \frac{1}{1 + \exp(-x)} \rightarrow \frac{dg}{dx} = g(x)(1 - g(x))$$

$$= -\sum_{i=N}^{N} \left[ y^i \frac{1}{g(\mathbf{w}^T \mathbf{x}^i)} - (1 - y^i) \frac{1}{1 - g(\mathbf{w}^T \mathbf{x}^i)} \right] g(\mathbf{w}^T \mathbf{x}^i)(1 - g(\mathbf{w}^T \mathbf{x}^i)) \frac{\partial \mathbf{w}^T \mathbf{x}^i}{\partial w_k}$$

$$= -\sum_{i \neq 1} \left[ y^i (1 - g(\mathbf{w}^T \mathbf{x}^i)) - (1 - y^i) g(\mathbf{w}^T \mathbf{x}^i) \right] x_k^i$$

$$= -\sum_{i=1} \left[ y^i - g(\mathbf{w}^T \mathbf{x}^i) \right] \mathbf{x}_k^i \qquad\qquad \nabla L(\mathbf{w}^*) = \mathbf{0}$$

**nonlinear system of equations!!**

SIGGRAPH
ASIA 2018
T O K Y O

79

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Fact: gradient at any point gives direction of fastest increase

SIGGRAPH
ASIA 2018
T O K Y O

80

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Fact: gradient at any point gives direction of fastest increase

Idea: start at a point and move in the direction opposite to the gradient

SIGGRAPH
ASIA 2018
T O K Y O

81

# Gradient Descent Minimization

$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Fact: gradient at any point gives direction of fastest increase

Idea: start at a point and move in the direction opposite to the gradient

SIGGRAPH
ASIA 2018
TOKYO

82

# Gradient Descent Minimization

$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Fact: gradient at any point gives direction of fastest increase

Idea: start at a point and move in the direction opposite to the gradient

SIGGRAPH
ASIA 2018
TOKYO

83

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Fact: gradient at any point gives direction of fastest increase

Idea: start at a point and move in the direction opposite to the gradient

Initialize:      $\mathbf{x}_0$

Update:      $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$      **i=0**

84

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Update:      $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$      **i=1**

85

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Update:     $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$     **i=1**

SIGGRAPH
ASIA 2018
T O K Y O

86

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Update:     $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$     **i=2**

SIGGRAPH
ASIA 2018
T O K Y O

87

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Update:           $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$          **i=2**

SIGGRAPH
ASIA 2018
TOKYO

88

# Gradient Descent Minimization



$$\nabla f(\mathbf{x}) = \left[ \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right]$$

Update:           $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$          **i=3**

SIGGRAPH
ASIA 2018
TOKYO

89

# Gradient Descent Minimization

Initialize:     $\mathbf{x}_0$

Update:         $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha \nabla f(\mathbf{x}_i)$

We can always make it converge for a convex function.

**convex**

$f(x)$

$g(x)$

$f(x_2)$

$tf(x_1) + (1-t)f(x_2)$

$f(tx_1 + (1-t)x_2)$

$f(x_1)$

$x_1$         $tx_1 + (1-t)x_2$    $x_2$

**non-convex**

120
100
80
60
40
20
0
-20
-10        -5         0         5         10

90

---

# XOR Problem                   $y = f(x_1, x_2)$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0     | 0     | 0   |
| 0     | 1     | 1   |
| 1     | 0     | 1   |
| 1     | 1     | 0   |

$x_2$

**X**

$x_1$

$$y = f(w_0, w_1, w_2) = \mathcal{H}(w_0 + w_1 x_1 + w_2 x_2)$$

91

## XOR Problem $y = f(z_1, z_2) = f(g_1(x_1, x_2), g_2(x_1, x_2))$

| $x_1$ | $x_2$ | $z_1$ |
|-------|-------|-------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

| $x_1$ | $x_2$ | $z_2$ |
|-------|-------|-------|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

| $z_1$ | $z_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



92

## XOR Problem

$$y = f(z_1, z_2) = f(g_1(x_1, x_2), g_2(x_1, x_2))$$
$$= f\left(\mathcal{H}(\mathbf{w}^1, x_1, x_2), \mathcal{H}(\mathbf{w}^2, x_1, x_2)\right)$$
$$= \mathcal{H}(\mathbf{w}^3, \mathcal{H}(g_1(\mathbf{w}^1, x_1, x_2)), \mathcal{H}(g_2(\mathbf{w}^2, x_1, x_2))$$



93

# XOR Problem

$$y = f(z_1, z_2) = f(g_1(x_1, x_2), g_2(x_1, x_2))$$
$$= f\left(\mathcal{H}(\mathbf{w}^1, x_1, x_2), \mathcal{H}(\mathbf{w}^2, x_1, x_2)\right)$$
$$= \mathcal{H}(\mathbf{w}^3, \mathcal{H}(g_1(\mathbf{w}^1, x_1, x_2)), \mathcal{H}(g_2(\mathbf{w}^2, x_1, x_2)))$$



**94**

# Course Information (slides/code/comments)



**http://geometry.cs.ucl.ac.uk/creativeai/**

SIGGRAPH Asia Course CreativeAI: Deep Learning for Graphics

# SIGGRAPH ASIA 2018 TOKYO

## CreativeAI: Deep Learning for Graphics

# Neural Network Basics

| **Niloy Mitra** | **Iasonas Kokkinos** | **Paul Guerrero** | **Nils Thuerey** | **Tobias Ritschel** |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

**UCL**          **facebook** Artificial Intelligence Research          **TUM** Technische Universität München

---

## Timetable

| | | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|---|---|---|---|---|
| **Theory and Basics** | Introduction | X | X | X | X | X |
| | Theory | X | | | X | |
| | NN Basics | X | X | | | |
| | Alternatives to Direct Supervision | | | X | | |
| | 15 min. break | | | | | |
| **State of the Art** | Feature Visualization | | | | | X |
| | Image Domains | | X | | | X |
| | 3D Domains | | | X | | X |
| | Motion and Physics | X | | | X | |

# Introduction to Neural Networks

SIGGRAPH
ASIA 2018
T O K Y O

---

# Goal: Learn a Parametric Function

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

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

Examples:

   Image Classification:    $f_\theta : \mathbb{R}^{w \times h \times c} \longrightarrow \{0, 1, \ldots, k-1\}$

                 $w \times h \times c$ : image dimensions     $k$ : class count

   Image Synthesis:     $f_\theta : \mathbb{R}^n \longrightarrow \mathbb{R}^{w \times h \times c}$

            $n$ : latent variable count     $w \times h \times c$ : image dimensions

SIGGRAPH
ASIA 2018
T O K Y O

# Machine Learning 101: Linear Classifier



$$f_\theta : \mathbb{R}^n \longrightarrow \{0, 1\}$$

$$f_\theta(x) = \begin{cases} 1 & \text{if } wx + b \geq 0 \\ 0 & \text{if } wx + b < 0 \end{cases}$$

$$\theta = \{w, b\}$$

Each data point has a class label:

$$y^i = \begin{cases} 1 & (\textcolor{red}{\bullet}) \\ 0 & (\textcolor{blue}{\bullet}) \end{cases}$$

Feature coordinate $x_2$

Feature coordinate $x_1$

# Nonlinear decision boundaries

$\mathbb{X}$

$\mathbb{X}'$

$$g : \mathbb{X} \longrightarrow \mathbb{X}'$$



$$f_\theta(x) = \begin{cases} 1 & \text{if } w \, g(x) + b \geq 0 \\ 0 & \text{if } w \, g(x) + b < 0 \end{cases}$$

# Building A Complicated Function

Given a library of simple functions

$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a

complicated function

SIGGRAPH
ASIA 2018
T O K Y O

# Building A Complicated Function

Given a library of simple functions

$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a

complicated function

Idea 1: Linear Combinations

• Boosting

• Kernels

• …

$$f(x) = \sum_i \alpha_i g_i(x)$$

SIGGRAPH
ASIA 2018
T O K Y O

# Building A Complicated Function

Given a library of simple functions

$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a

complicated function

Idea 2: Compositions

• Decision Trees

• Deep Learning

$$f(x) = g_1(g_2(\ldots(g_n(x)\ldots)))$$

# Building A Complicated Function

Given a library of simple functions

$\sin(x)$

$\log(x)$

$\cos(x)$

$x^3$

$\exp(x)$

Compose into a

complicated function

Idea 2: Compositions

• Decision Trees

• Grammar models

• Deep Learning

$$f(x) = \log(\cos(\exp(\sin^3(x))))$$

$\sin(x)$ → $x^3$ → $\exp(x)$ → $\cos(x)$ → $\log(x)$

# 'Neuron': Cascade of Linear and Nonlinear Function

basic building block



$x_0$      $w_0$

axon from a neuron    synapse

$w_0 x_0$

dendrite

cell body    $f\left(\sum_i w_i x_i + b\right)$

$w_1 x_1$    $\sum_i w_i x_i + b \;\; f$    output axon

activation function

$w_2 x_2$

**Sigmoidal activation**

TOKYO

---

# Activation functions

—— function
—— derivative



| Step ("perceptron") | Sigmoidal ("logistic") | Hyperbolic tangent | Rectified Linear Unit (RELU) |
|---|---|---|---|
| $g(a) = \begin{cases} 0 & a < 0 \\ 1 & a \geq 0 \end{cases}$ | $g(a) = \dfrac{1}{1 + \exp(-a)}$ | $g(a) = \dfrac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$ | $g(a) = \max(0, a)$ |

Image Credit: Olivier Grisel and Charles Ollion

SIGGRAPH
ASIA 2018
TOKYO

# Perceptrons (60's)

**Apple**  **Orange**

**output units  e.g.
class labels**

**non-adaptive
hand-coded
features**

**Fixed
mapping**

**input units
e.g. pixels**

## XOR: perceptron killer

$x_2$

$x_1$

# Multi-Layer Perceptrons (~1985)

$$u_i = g\left(\sum_{k \in \mathcal{N}(i)} w_{k,i} g\left(\sum_{m \in \mathcal{N}(k)} w_{m,k} u_m + b_k\right) + b_i\right)$$

**outputs**

**hidden
layers**

**input vector**

# Reminder: Non-linear decision boundaries

**This is what the hidden layers should be doing!**

$\mathbb{X}$

$g : \mathbb{X} \longrightarrow \mathbb{X}'$

$\mathbb{X}'$

$$f_\theta(x) = \begin{cases} 1 & \text{if } w \ g(x) + b \geq 0 \\ 0 & \text{if } w \ g(x) + b < 0 \end{cases}$$

# Nonlinear mapping $\mathbb{R}^2 \to \mathbb{R}^2$

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

**Non-linearly separable data**

**Data mapped to learned space**

**Decision function**



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



**Points in 1D,
Decision in 2D**

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



Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

72

# Deep Learning = Hierarchical Compositionality

"car"

# Deep Learning = Hierarchical Compositionality

Low-Level Feature → Mid-Level Feature → High-Level Feature → Trainable Classifier → "car"

# MLP Demo: [playground.tensorflow.org](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:

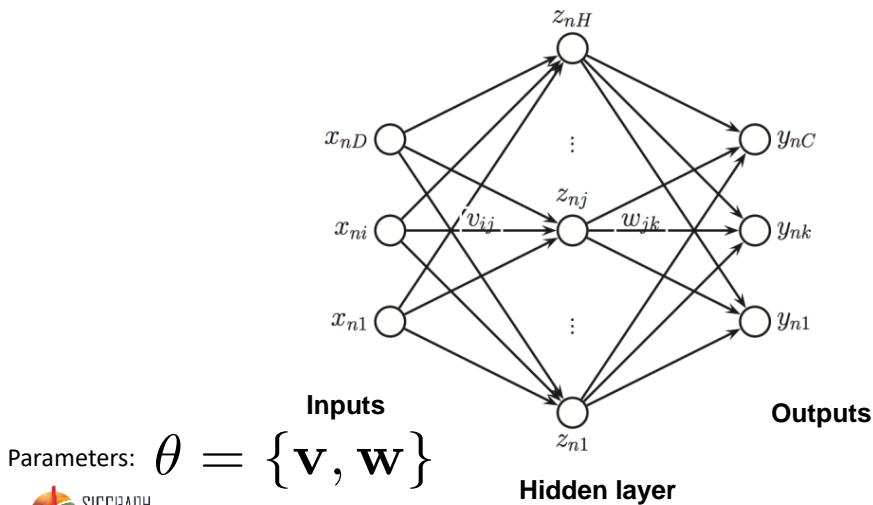$$\min_\theta L_f(\theta)$$

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

**Forward**

$X$          $Z$

$\theta$

**Backward**

$\dfrac{\partial L}{\partial X}$    $\left\{ \dfrac{\partial Z}{\partial X}, \dfrac{\partial Z}{\partial \theta} \right\}$    $\dfrac{\partial L}{\partial Z}$

$\dfrac{\partial L}{\partial \theta}$

# Chain Rule

$x$          $y$

$\dfrac{dL}{dx}$          $\dfrac{dL}{dy}$

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

## Chain Rule

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

$$\frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dx}$$

SIGGRAPH ASIA 2018 TOKYO

## 'Another Brick in the Wall'

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

$$\frac{dL}{dx} = \frac{dL}{dy} \cdot \frac{dy}{dx}$$

SIGGRAPH ASIA 2018 TOKYO

# Toy example: single sigmoidal unit

$$f(w, x) = \frac{1}{1 + \exp(-(w_0 x_0 + w_1 x_1 + w_2))}$$

**Composition of differentiable blocks:**

$$f(x) = \frac{1}{x} \quad \rightarrow \quad f'(x) = -\frac{1}{x^2}$$
$$f_c(x) = c + x \quad \rightarrow \quad f'(x) = 1$$
$$f(x) = e^x \quad \rightarrow \quad f'(x) = e^x$$
$$f_a(x) = ax \quad \rightarrow \quad f'(x) = a$$

SIGGRAPH
ASIA 2018
T O K Y O

# Computation graph & automatic differentiation



Slide Credit: Justin Johnson

SIGGRAPH
ASIA 2018
T O K Y O

# Multi-Layer Perceptrons

$$u_i = g\left(\sum_{k \in \mathcal{N}(i)} w_{k,i}\, g\left(\sum_{m \in \mathcal{N}(k)} w_{m,k} u_m + b_k\right) + b_i\right)$$

outputs

hidden layers

input vector

Slide Credit: G. Hinton

SIGGRAPH ASIA 2018 TOKYO

# Multi-Layer Perceptrons

Compare outputs with **correct answer** to get error signal

**Back-propagate error signal to get derivatives for learning**

outputs

hidden layers

input vector

Slide Credit: G. Hinton

SIGGRAPH ASIA 2018 TOKYO

## Back-propagation Algorithm



## 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

$$\mathbf{x}_n \xrightarrow{\mathbf{V}} \mathbf{a}_n \xrightarrow{g} \mathbf{z}_n \xrightarrow{\mathbf{W}} \mathbf{b}_n \xrightarrow{h} \hat{\mathbf{y}}_n$$



Parameters: $\theta = \{\mathbf{v}, \mathbf{w}\}$

# A Neural Network in Forward Mode ▶▶

$$\mathbf{a} = \mathbf{V}\mathbf{x}$$

# A Neural Network in Forward Mode ⏩

$$\mathbf{z} = g(\mathbf{a})$$

$$z_k = \frac{1}{1 + \exp(-a_k)}$$



**Inputs**

**Outputs**

**Hidden layer**

# A Neural Network in Forward Mode ⏩

$$\mathbf{b} = \mathbf{W}\mathbf{z}$$



**Inputs**

**Outputs**

**Hidden layer**

# A Neural Network in Forward Mode ▶▶



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

# Objective for linear regression

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



$$\hat{\mathbf{y}} = h(\mathbf{b}) \quad \mathbf{y}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Ground truth**

$$l(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{c=1}^{C} (\mathbf{y}_c - \hat{\mathbf{y}}_c)^2$$

# Objective for multi-class classification

Softmax unit    $\hat{y}_k = \dfrac{\exp(b_k)}{\sum_{c=1}^{C} \exp(b_c)}$



$$\hat{\mathbf{y}} = h(\mathbf{b}) \quad \mathbf{y}$$

$z_{nH}$

$x_{nD}$    $z_{nj}$    $y_{nC} \Longleftrightarrow$

$x_{ni}$   $v_{ij}$   $w_{jk}$   $y_{nk} \Longleftrightarrow$

$x_{n1}$   $y_{n1} \Longleftrightarrow$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

**Ground truth**

$z_{n1}$

$$l(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{c=1}^{C} \mathbf{y}_c \log(\hat{\mathbf{y}}_c) \quad \text{`Cross-entropy' loss}$$

---

# Neural network in forward mode: recap

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

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

What we need to compute for gradient descent:    $\dfrac{\partial l(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{v}_i} \quad \dfrac{\partial l(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{w}_j}$

# A Neural Network in Backward Mode ⏪



# A Neural Network in Backward Mode ⏪

**Hidden layer**

$$\mathbf{b} = \mathbf{Wz}$$



**This we want**

$$\frac{\partial l}{\partial w_{jk}} = \quad ?$$

# A Neural Network in Backward Mode ⏪

**Hidden layer**

$z_{nH}$

$$\mathbf{b} = \mathbf{W}\mathbf{z}$$

$\mathbf{y}$

$x_{nD}$    $y_{nC}$

$z_{nj}$

$x_{ni}$    $v_{ij}$    $w_{jk}$    $y_{nk}$

$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

$x_{n1}$    $y_{n1}$

**Outputs**

**This we want**    $z_{n1}$

$$\boxed{\frac{\partial l}{\partial z_j}} = \quad ?$$

---

# 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

$$b_m = \sum_{h=1}^{H} z_h w_{h,m}$$



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

## Linear Layer Parameters in Backward: 1-to-1

$$b_m = \sum_{h=1}^{H} z_h w_{h,m}$$



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

# A Neural Network in Backward Mode ⏪

**Hidden layer**

$z_{nH}$

$$\mathbf{b} = \mathbf{W}\mathbf{z}$$

$\mathbf{y}$

$x_{nD}$

$y_{nC}$

$\vdots$

$z_{nj}$    $\boxed{w_{jk}}$

$x_{ni}$    $v_{ij}$

$y_{nk}$

$x_{n1}$

$y_{n1}$

$\vdots$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

**This we want**   $z_{n1}$  **This we have**  **This we computed**

$$\boxed{\frac{\partial l}{\partial w_{jk}}} = \sum_m \boxed{\frac{\partial l}{\partial b_m}}\boxed{\frac{\partial b_m}{\partial w_{jk}}} = \frac{\partial l}{\partial b_m} z_j$$

---

# A Neural Network in Backward Mode ⏪

**Hidden layer**

$z_{nH}$

$$\mathbf{b} = \mathbf{W}\mathbf{z}$$

$\mathbf{y}$

$x_{nD}$

$y_{nC}$

$\vdots$

$\boxed{z_{nj}}$

$x_{ni}$    $v_{ij}$    $w_{jk}$

$y_{nk}$

$x_{n1}$

$y_{n1}$

$\vdots$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

**This we want**   $z_{n1}$  **This we have**  **This we computed**

$$\boxed{\frac{\partial l}{\partial z_j}} = \sum_m \boxed{\frac{\partial l}{\partial b_m}}\boxed{\frac{\partial b_m}{\partial z_j}} = \sum_m \frac{\partial l}{\partial b_m} w_{j,m}$$

# A Neural Network in Backward Mode ⏪

**Hidden layer**

$$z_k = \frac{1}{1 + \exp(-a_k)}$$

$$\mathbf{y}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

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

---

# A Neural Network in Backward Mode ⏪

**Hidden layer**

$$\mathbf{a} = \mathbf{V}\mathbf{x}$$

$$z_k = \frac{1}{1 + \exp(-a_k)}$$

$$\mathbf{y}$$

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

$$\frac{\partial l}{\partial v_{ij}} = \sum_k \frac{\partial l}{\partial a_k} \frac{\partial a_k}{\partial v_{ij}} = \frac{\partial l}{\partial a_j} x_i$$

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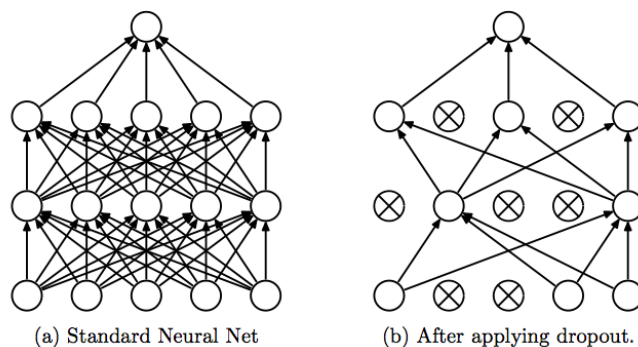
$$L(\mathbf{W}) = \frac{1}{N}\sum_{i=1}^{N} l(\mathbf{y}^i, \hat{\mathbf{y}}^i) + \sum_{l}\lambda_l \sum_{k,m}(\mathbf{W}_{k,m}^l)^2$$

Per-sample loss          Per-layer regularization

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

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

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

Back-prop for
i-th example

If N=$10^6$ , we will need to run back-prop $10^6$ times to update **W** once!

# Stochastic Gradient Descent (SGD)

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

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

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

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

**Epoch:** N samples, N/B batches

SIGGRAPH
ASIA 2018
T O K Y O

# Regularization in SGD: Weight Decay

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

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

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

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

Back-prop on minibatch      ''Weight decay''

**Epoch:** N samples, N/B batches

SIGGRAPH
ASIA 2018
T O K Y O

# Learning rate



loss

very high learning rate

low learning rate

high learning rate

good learning rate

epoch

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

# Gradient Descent



1    2    3    4

# (S)GD with adaptable stepsize



Too small: converge very slowly

Too big: overshoot and even diverge

Reduce size over time

**e.g.** $\quad \epsilon_t = \dfrac{c}{t}$

# (S)GD with momentum



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

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

(S)GD + momentum

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

## Step-size Selection & Optimizers: research problem

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



SIGGRAPH
ASIA 2018
TOKYO

---

# Code example

## Multi-layer perceptron classification

68

# 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



$$\mathbf{a} = \mathbf{V}\mathbf{x} \qquad \mathbf{b} = \mathbf{W}\mathbf{z}$$

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

# Reminder: Overfitting, in images

**Classification**

| Underfitting | just right $\longleftrightarrow$ | Overfitting |

**Regression**

# Previously: l2 Regularization

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

Per-sample loss                              Per-layer regularization

# Dropout



(a) Standard Neural Net          (b) After applying dropout.

**Each sample is processed by a 'decimated' neural net**

**Decimated nets: distinct classifiers**

**But: they should all do the same job**

# Dropout block



(a) Standard network          (b) Dropout network

Figure 3: Comparison of the basic operations of a standard and dropout network.

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

$$
\begin{aligned}
r_j^{(l)} &\sim \text{Bernoulli}(p), \\
\widetilde{\mathbf{y}}^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\
z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)}\widetilde{\mathbf{y}}^l + b_i^{(l+1)}, \\
y_i^{(l+1)} &= f(z_i^{(l+1)}).
\end{aligned}
$$

**'Feature noising'**

# Test time: Deterministic Approximation



(a) At training time                    (b) At test time

# 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



**Sigmoidal ("logistic")**

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

**Rectified Linear Unit (RELU)**

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

# Reminder: a network in backward mode

$$z_k = \frac{1}{1 + \exp(-a_k)}$$

$z_{nH}$

$x_{nD}$

$z_{nj}$

$v_{ij}$    $w_{jk}$

$x_{ni}$

$x_{n1}$

$y_{nC}$

$y_{nk}$

$y_{n1}$

$$\mathbf{y} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

**Outputs**

$z_{n1}$

**Gradient signal from above**

**scaling: <1  (actually <0.25)**

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

---

# Vanishing Gradients Problem

**Gradient signal from above**

**scaling: <1  (actually <0.25)**

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

**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!**

Sigmoid s(z)
Derivative s'(z)

Weighted sum

## Vanishing Gradients Problem: ReLU Solves It

**Gradient signal from above**     **Scaling: {0,1}**

$$\frac{\partial l}{\partial a_k} = \sum_m \frac{\partial l}{\partial z_m} \frac{\partial z_m}{\partial a_k} = \frac{\partial l}{\partial z_k} g'(a_k)$$

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

$$g'(a) = \left\{ \begin{array}{ll} 1 & a > 0 \\ 0 & a < 0 \end{array} \right.$$

# 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

10 am                          2pm                          7pm



---

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

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



Original Patch and Intensity Values

Brightness Decreased

Contrast increased,

- Make each patch have zero mean:

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

$$Z(x, y) = I(x, y) - \mu$$

- Then make it have unit variance:

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

$$ZN(x, y) = \frac{Z(x, y)}{\sigma}$$

# Internal Covariate Shift

**Neural network activations during training: moving target**

---

# Batch Normalization

**Whiten-as-you-go:**

- Normalize the activations in each layer within a mini-batch.

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



$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

(b) Without BN          (c) With BN

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift
S Ioffe 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

**Example:  200x200 image**
**40K hidden units**
**~2B parameters!!!**

- **Spatial correlation is local**
- **Waste of resources**
- **we have not enough training samples anyway..**

# Locally-connected Layer

**Example: 200x200 image**
**40K hidden units**
**Filter size: 10x10**
**4M parameters**

**Note:** This parameterization is good when input image is registered (e.g., face recognition).

# Locally-connected Layer



**Example: 200x200 image**
**40K hidden units**
**Filter size: 10x10**
**4M parameters**

**Note:** **This parameterization is good when input image is registered (e.g., face recognition).**

# Convolutional Layer



**Share the same parameters across different locations (assuming input is stationary):**
**Convolutions with learned kernels**

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Convolutional Layer



# Convolutional Layer

# Fully-connected layer

**#of parameters: K²**

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ \vdots \\ y_K \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & \dots & w_{1,K} \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & \dots & w_{2,K} \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & \dots & w_{3,K} \\ w_{4,1} & w_{4,2} & w_{4,3} & w_{4,4} & \dots & w_{4,K} \\ & & \vdots & & & \\ w_{K,1} & w_{K,2} & w_{K,3} & w_{K,4} & \dots & w_{K,K} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \\ x_K \end{bmatrix}$$

# Convolutional layer

**#of parameters: size of window**

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

# Convolutional layer



$$* \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} =$$

# Code example

Learning an edge filter

112

# Convolutional layer

**Learn multiple filters.**

**E.g.: 200x200 image**
**100 Filters**
**Filter size: 10x10**
**10K parameters**

# Convolutional layer

$$h_i^n = \max\left\{ 0, \sum_{j=1}^{\#\text{input channels}} h_j^{n-1} * w_{ij}^n \right\}$$

**output feature map**

**input feature map**

**kernel**

$h_1^{n-1}$

$h_2^{n-1}$

$h_3^{n-1}$

**Conv. layer**

$h_1^n$

$h_2^n$

# Convolutional layer

$$h_i^n = \max \left\{ 0, \sum_{j=1}^{\#\text{input channels}} h_j^{n-1} * w_{ij}^n \right\}$$

**output feature map**

**input feature map**

**kernel**

$h_1^{n-1}$

$h_2^{n-1}$

$h_3^{n-1}$

$h_1^n$

$h_2^n$

# Convolutional layer

$$h_i^n = \max \left\{ 0, \sum_{j=1}^{\#\text{input channels}} h_j^{n-1} * w_{ij}^n \right\}$$

**output feature map**

**input feature map**

**kernel**

$h_1^{n-1}$

$h_2^{n-1}$

$h_3^{n-1}$

$h_1^n$

$h_2^n$

# 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



$h^{n-1}$        $h^n$        $h^{n+1}$

Conv. layer    Pool. layer

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)

# Pooling layer: receptive field size



$h^{n-1}$        $h^n$        $h^{n+1}$

Conv. layer    Pool. layer

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**



**Receptive field: layer 1**

Receptive field: layer 2



Receptive field: layer 3

**Receptive field: layer 4**



**Receptive field: layer 5**

Receptive field: layer 6



Receptive field: layer 7

# Receptive field: layer 8

# Modern Architectures

# CNNs, late 1980's: LeNet

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



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

---

# What happened in between?



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

# 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.

# 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
  - …

# 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

# Residual Modelling: Basic Idea in Image Processing

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

**Preserving base information**



**Some Network**

**residual**

**can treat perturbation**

# Residual Network

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



$x$

any two stacked layers

weight layer

relu

weight layer

relu

$H(x)$

# 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

# 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

# Fully-convolutional Neural Networks



**FCNN**

# Fully-convolutional Neural Networks



**FCNN**

# Fully-convolutional Neural Networks

FCNN

Fast     (shared convolutions)
Simple (dense)



# Fully Convolutional Neural Networks in Practice

32-fold decimation
224x224 to 7x7

FCNN

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

# Atrous convolution

**downsample x 2**      **convolve**      **'implant' in image coordinates**



**filter 'atrous'**

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
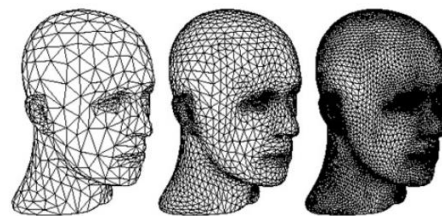
# Atrous convolution = Dilated Convolution



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 $3 \times 3$. (b) $F_2$ is produced from $F_1$ by a 2-dilated convolution; each element in $F_2$ has a receptive field of $7 \times 7$. (c) $F_3$ is produced from $F_2$ by a 4-dilated convolution; each element in $F_3$ has a receptive field of $15 \times 15$. The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.
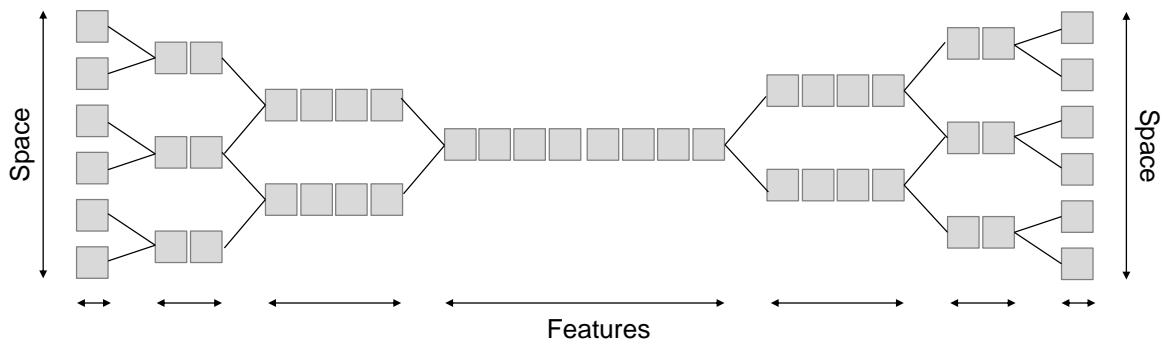
**F. Yu, V. Koltun, Multi-Scale Context Aggregation by Dilated Convolutions, ICLR 2016**

# Graphics: Multiresolution
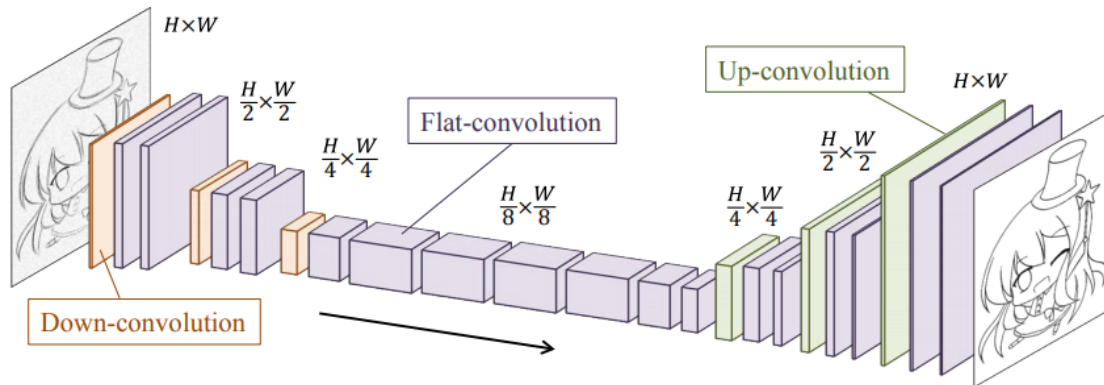
# Encoder-decoder



# Interpretation

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

# Encoder-decoder



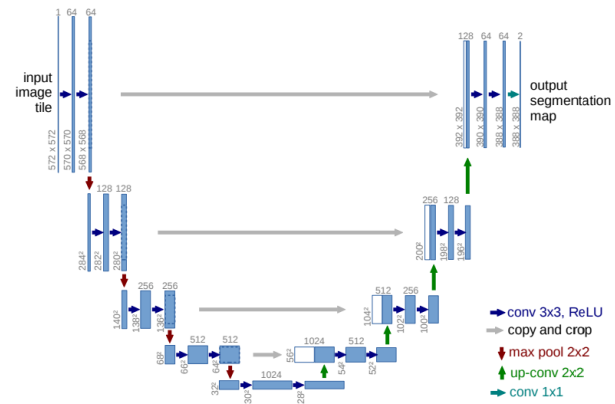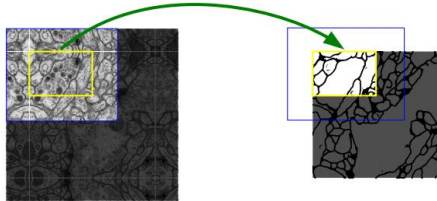Learning to simplify. Simo-Serra et al. 2016

# Up-sampling

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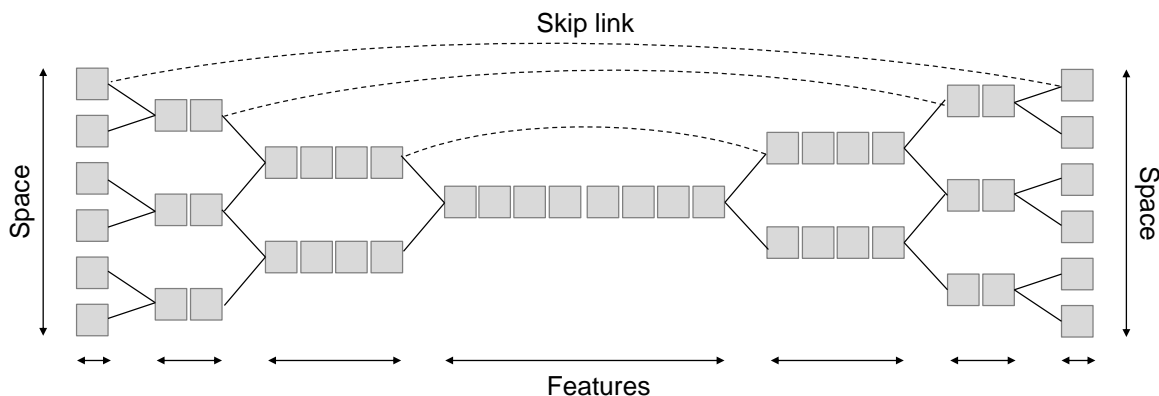  - Padding (insert zeros)
  - Transpose convolutions

# Encoder-decoder + Skip connections

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



conv 3x3, ReLU
copy and crop
max pool 2x2
up-conv 2x2
conv 1x1

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

# Encoder-decoder with skip connections
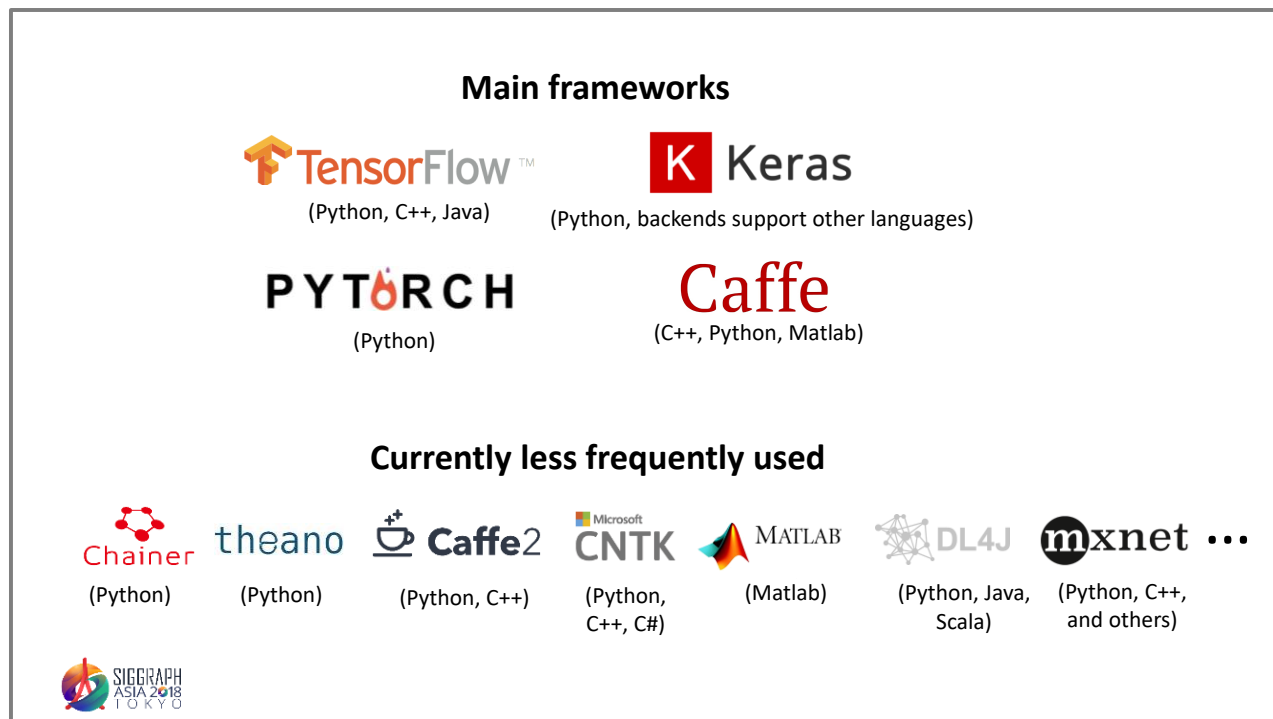


Skip link

Space

Space

Features

# 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
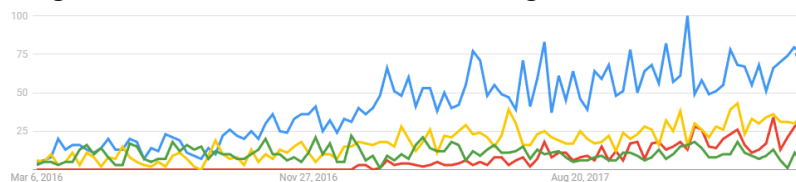  - Residual coding
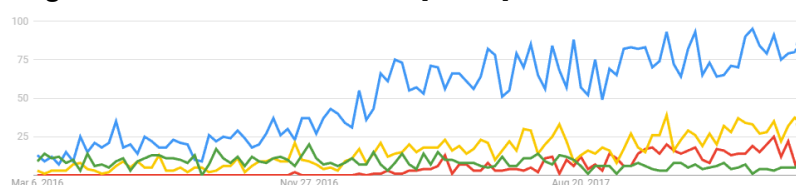  - Pyramidal schemes (Laplacian pyramid, etc.)

# Deep Learning Frameworks

**Main frameworks**

TensorFlow™
(Python, C++, Java)

K Keras
(Python, backends support other languages)

PYTORCH
(Python)

Caffe
(C++, Python, Matlab)

**Currently less frequently used**

Chainer
(Python)

theano
(Python)

Caffe2
(Python, C++)

Microsoft CNTK
(Python, C++, C#)

MATLAB
(Matlab)

DL4J
(Python, Java, Scala)

mxnet
(Python, C++, and others)

· · ·



# Popularity

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

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

— TensorFlow™
K Keras
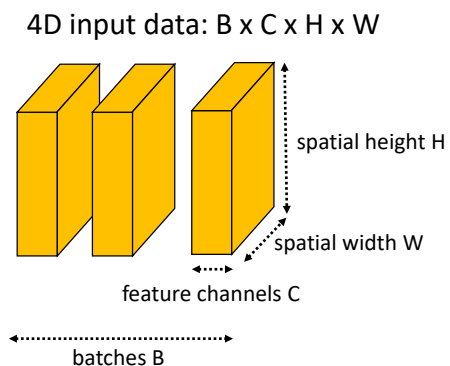PYTORCH
Caffe

# 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)
```
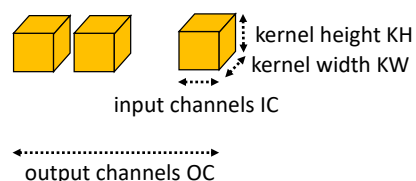
# Tensors

- Frameworks typically represent data as tensors
- Examples:

4D input data: B x C x H x W

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

spatial height H

spatial width W

feature channels C

batches B

kernel height KH
kernel width KW

input channels IC

output channels OC
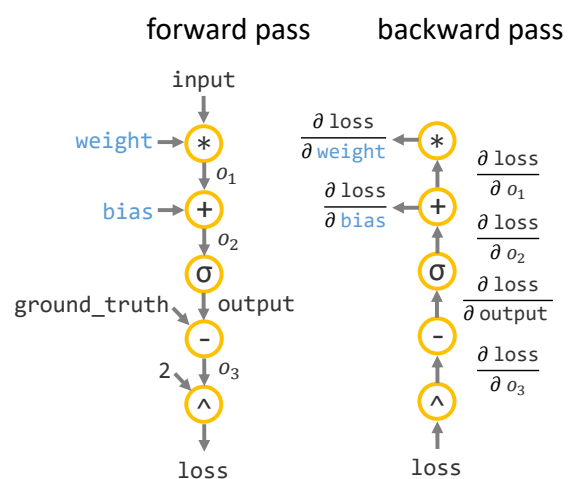
# What Does a Deep Learning Framework Do?

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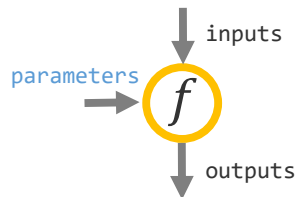- sometimes: distributed computing

# Automatic Differentiation & the Computation Graph

```
parameters = (weight, bias)

output = σ(weight * input + bias)

loss = (output - ground_truth)^2

# gradients of loss with respect to parameters
gradients = backpropagate(loss, parameters)
```
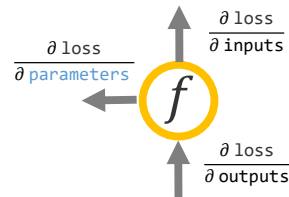
Since loss is a scalar, the gradients
are the same size as the parameters

forward pass

input

weight $\rightarrow$ *

$o_1$

bias $\rightarrow$ +

$o_2$

σ

ground_truth    output

−

$2$    $o_3$

^

loss

backward pass

$\frac{\partial\, loss}{\partial\, weight}$ $\leftarrow$ *    $\frac{\partial\, loss}{\partial\, o_1}$

$\frac{\partial\, loss}{\partial\, bias}$ +    $\frac{\partial\, loss}{\partial\, o_2}$

σ    $\frac{\partial\, loss}{\partial\, output}$

−    $\frac{\partial\, loss}{\partial\, o_3}$

^

loss

# Automatic Differentiation & the Computation Graph

inputs

parameters $f$

outputs

$$\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{inputs}}$$

$$\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{parameters}}$$ $f$

$$\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{outputs}}$$

outputs = forward(inputs, parameters)

$$\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{inputs}}\,,\quad\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{parameters}} = \mathrm{backward}(\frac{\partial\,\mathrm{loss}}{\partial\,\mathrm{outputs}})$$

SIGGRAPH
ASIA 2018
TOKYO

# 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)
```

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)
```

SIGGRAPH
ASIA 2018
TOKYO

# 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

- 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

# Keras

K Keras

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

# Caffe

Caffe

- 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 | - | pytorch-tf/ MMdnn | model-converters/ nn_toolsconvert-to-tensorflow/MMdnn | MMdnn/ nn_tools | None | crosstalk/MMdnn | None | MMdnn |
| pytorch | pytorch2keras (over Keras) | - | Pytorch2keras/ nn-transfer | Pytorch2caffe/pytorch-caffe-darknet-convert | onnx-caffe2 | ONNX | None | None |
| keras | nn_tools /convert-to-tensorflow/keras_to_tensorflow/keras_to_tensorflow/MMdnn | MMdnn/ nn-transfer | - | MMdnnnn_tools | None | MMdnn | None | MMdnn |
| caffe | MMdnn/nn_tools/caffe-tensorflow | MMdnn/ pytorch-caffe-darknet-convert/ pytorch-resnet | caffe_weight_converter/ caffe2keras/nn_tools/ kerascaffe2keras/ Deep_Learning_Model_Converter/MMdnn | - | CaffeToCaffe2 | crosstalkcaffe/CaffeConverterMMdnn | None | mxnet/tools/caffe_converter/ResNet_caffe2mxnet/ MMdnn |
| 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 | MMdnn | MMdnn | MMdnn | MMdnn/MXNet2Caffe/ Mxnet2Caffe | None | MMdnn | None | - |

---

# ONNX

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

# MMdnn



- Converters available for several frameworks
- Common intermediate representation, but no clear standard

# Course Information (slides/code/comments)



**http://geometry.cs.ucl.ac.uk/creativeai/**

CreativeAI: Deep Learning for Graphics

# Alternatives to Direct Supervision

**Niloy Mitra**    **Iasonas Kokkinos**    **Paul Guerrero**    **Nils Thuerey**    **Tobias Ritschel**

UCL       UCL/Facebook       UCL       TU Munich       UCL

---

# Timetable

|  |  | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|:---:|:---:|:---:|:---:|:---:|
| Theory and Basics | Introduction | X | X | X | X | X |
|  | Theory | X |  |  | X |  |
|  | NN Basics | X | X |  |  |  |
|  | Alternatives to Direct Supervision |  |  | X |  |  |
|  | — 15 min. break — |  |  |  |  |  |
| State of the Art | Feature Visualization |  |  |  |  | X |
|  | Image Domains |  | X |  |  | X |
|  | 3D Domains |  |  | X |  | X |
|  | Motion and Physics | X |  |  | X |  |

# Unsupervised Learning

- There is no direct ground truth for the quantity of interest

- Autoencoders
- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)

---

# Autoencoders

Goal: Meaningful features that capture the main factors of variation in the dataset
- These are good for classification, clustering, exploration, generation, …
- We have no ground truth for them

Features $z$

Encoder

Input data $x$

# Autoencoders

Goal: Meaningful features that capture the main factors of variation

Features that can be used to reconstruct the image

$\hat{x}$

Decoder

Features
(Latent variables)  $z$

L2 Loss function:
$\|x - \hat{x}\|^2$

Encoder

Input data  $x$

# Autoencoders

Linear Transformation for Encoder and Decoder give result close to PCA

Deeper networks give better reconstructions, since basis can be non-linear

Original

Autoencoder

PCA

Decoder

$W_1^T$
2000
$W_2^T$
1000
$W_3^T$
500
$W_4^T$
30   Code layer
$W_4$
500
$W_3$
1000
$W_2$
2000
$W_1$

Encoder

# Example: Document Word Prob. → 2D Code

LSA (based on PCA)             Autoencoder



Image Credit: Reducing the Dimensionality of Data with
Neural Networks, Hinton and Salakhutdinov

# Example: Semi-Supervised Classification

- Many images, but few ground truth labels

start unsupervised
train autoencoder on many images

supervised fine-tuning
train classification network on labeled images



Loss function
(Softmax, etc)

Predicted Label   $\hat{y}$

GT Label   $y$

Decoder

L2 Loss function:
$\|x - \hat{x}\|^2$

Classifier

Features
(Latent Variables)   $z$

Encoder

Features   $z$

Encoder

Input data   $x$

$x$

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Code example

### Autoencoder
### (autoencoder.ipynb)

9

# Generative Models

- Assumption: the dataset are samples from an unknown distribution $p_{\text{data}}(x)$
- Goal: create a new sample from $p_{\text{data}}(x)$ that is not in the dataset

Dataset

?

Generated

Image credit: *Progressive Growing of GANs for Improved Quality, Stability, and Variation,* Karras et al.

# Generative Models

- Assumption: the dataset are samples from an unknown distribution $p_{\text{data}}(x)$
- Goal: create a new sample from $p_{\text{data}}(x)$ that is not in the dataset



Dataset                                    Generated

Image credit: *Progressive Growing of GANs for Improved Quality, Stability, and Variation,* Karras et al.

# Generative Models



$$p_\theta(x) \approx p_{\text{data}}(x)$$

Generator with parameters $\theta$

$$p(z)$$

known and
easy to sample from

# Generative Models

$$p_\theta(x) \approx p_{\mathrm{data}}(x)$$

Generator with parameters $\theta$

$$p(z)$$

known and
easy to sample from

**How to measure similarity of $p_\theta(x)$ and $p_{\mathrm{data}}(x)$?**

1) Likelihood of data in $p_\theta(x)$

**Variational Autoencoders (VAEs)**

2) Adversarial game:
*Discriminator* distinguishes $p_\theta(x)$ and $p_{\mathrm{data}}(x)$  **vs**  *Generator* makes it hard to distinguish

**Generative Adversarial Networks (GANs)**

---

# Autoencoders as Generative Models?

$\hat{x}$

Decoder = Generator?

$z$

random

Feature space / latent space

- A trained decoder transforms some features $z$ to approximate samples from $p_{\mathrm{data}}(x)$
- What happens if we pick a random $z$?
- We do not know the distribution $p(z)$ of features that decode to likely samples

Image Credit: *Reducing the Dimensionality of Data with Neural Networks*, Hinton and Salakhutdinov

# Variational Autoencoders (VAEs)

- Pick a parametric distribution $p(z)$ for features
- The generator maps $p(z)$ to an image distribution $p_\theta(x)$ (where $\theta$ are parameters)

$$p_\theta(x) = \int p_\theta(x|z) \; p(z) \; dz$$

- Train the generator to maximize the likelihood of the data in $p_\theta(x)$:

$$\max_\theta \sum_{x \in \text{data}} \log p_\theta(x)$$

$p_\theta(x|z)$

Generator with parameters $\theta$

$z$

sample

$p(z)$

# Outputting a Distribution

Normal distribution
$$p_\theta(x|z) = N(x; \mu(z), \Sigma(z))$$

$\mu$ $\Sigma$

Generator with parameters $\theta$

$z$

sample

$p(z)$

Bernoulli distribution
$$p_\theta(x|z) = Bern(x; r(z))$$

$r$

Generator with parameters $\theta$

$z$

sample

$p(z)$

# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

$$\theta^* = \arg\max_\theta \sum_{x \in \text{data}} \log \int p_\theta(x|z) \; p(z) \; dz$$

$$\theta^* \approx \arg\max_\theta \mathbb{E}_{x_i \sim p_{\text{data}}(x)} \mathbb{E}_{z \sim p(z)} \log p_\theta(x_i|z)$$

- SGD approximates the expected values over $(z, x_i)$ samples
- In each training iteration, sample $z$ from $p(z)$ …
- … and $x_i$ randomly from the dataset, and maximize:

$$\max_\theta \log p_\theta(x_i|z)$$

---

# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

Loss function:
$-\log p_\theta(x_i|z)$

$p_\theta(x|z)$    $x_i$

Generator with parameters $\theta$    Random from dataset

$z$

sample

$p(z)$

- In each training iteration, sample $z$ from $p(z)$ …
- … and $x_i$ randomly from the dataset
- SGD approximates the expected values over $(z, x_i)$ samples

# Variational Autoencoders (VAEs): Naïve Sampling (Monte-Carlo)

Loss function:
$$-\log p_\theta(x_i|z)$$

$p_\theta(x|z)$      $x_i$

Generator with parameters $\theta$     Random from dataset

$z$

sample

$p(z)$

$z$ with non-zero loss gradient for $x_i$

- In each training iteration, sample $z$ from $p(z)$ ...
- ... and $x_i$ randomly from the dataset
- SGD approximates the expected values over $(z, x_i)$ samples
- Few $(z, x_i)$ pairs have non-zero gradients



# Variational Autoencoders (VAEs): The Encoder

Loss function:
$$-\log p_\theta(x_i|z)$$

$$p_\theta(x) = \int p_\theta(x|z) \; p(z) \; dz$$

$p_\theta(x|z)$

Generator with parameters $\theta$

$z$

sample

$q_\phi(z|x_i)$

Encoder with parameters $\phi$

$x_i$

- During training, another network can guess a good $z$ for a given $x_i$
- $q_\phi(z|x_i)$ should be much smaller than $p(z)$
- This also gives us the data point $x_i$

# Variational Autoencoders (VAEs): The Encoder

Loss function:
$$-\log p_\theta(x_i|z) + KL(\ q_\phi(z|x_i)\ \|\ p(z)\ )$$

$p_\theta(x|z)$

Generator with parameters $\theta$

sample

$z$

$q_\phi(z|x_i)$

Encoder with parameters $\phi$

$x_i$

- Can we still easily sample a new $z$?
- Need to make sure $q_\phi(z|x_i)$ approximates $p(z)$
- Regularize with KL-divergence
- Negative loss can be shown to be a lower bound for the likelihood, and equivalent if
$$q_\phi(z|x) = p_\theta(z|x)$$

SIGGRAPH
ASIA 2018
TOKYO

# Reparameterization Trick

$p_\theta(x|z)$

Generator with parameters $\theta$

sample

$z$

Backprop?

$q_\phi(z|x_i)$

Encoder with parameters $\phi$

$x_i$

Example when $q_\phi(z|x_i) = N(z; \mu(x_i), \sigma(x_i))$:

$z = \sigma + \mu \cdot \epsilon$ , where $\epsilon \sim N(0,1)$

$$\frac{\partial z}{\partial \phi} = \frac{\partial \mu}{\partial \phi} + \frac{\partial \sigma}{\partial \phi} \cdot \epsilon$$

$z$

Backprop

sample

$\epsilon$         $\mu$         $\sigma$

$N(0,1)$

Encoder with parameters $\phi$

Does not depend on parameters $\phi$

$x_i$

SIGGRAPH
ASIA 2018
TOKYO

# Generating Data

MNIST

Frey Faces

sample

$\hat{x}$

$p_\theta(x|z)$

Generator with parameters $\theta$

$z$

sample

$p(z)$

Image Credit: *Auto-Encoding Variational Bayes*, Kingma and Welling

# Demos

**VAE on MNIST**
http://dpkingma.com/sgvb_mnist_demo/demo.html

**VAE on Faces**
http://vdumoulin.github.io/morphing_faces/online_demo.html

24

# Code example

**Variational Autoencoder
(variational_autoencoder.ipynb)**

25

# Generative Adversarial Networks



**Player 1: generator**
Scores if discriminator
can't distinguish output
from real image

**Player 2: discriminator** → real/fake
Scores if it can distinguish
between real and fake

from dataset

# Naïve Sampling Revisited

Loss function:
$$-\log p_\theta(x_i|z)$$

$p_\theta(x|z)$          $x_i$

Generator with         Random from dataset
parameters $\theta$

$z$

sample

$p(z)$

$z$ with non-zero
loss gradient for $x_i$

- Few $(z, x_i)$ pairs have non-zero gradients
- This is a problem of the maximum likelihood
- Use a different loss: Train a discriminator network to measure similarity $p_\theta(x) \approx p_{\text{data}}(x)$

SIGGRAPH
ASIA 2018
TOKYO

# Why Adversarial?

$\approx p_{\text{data}}(\hat{x})$

$D_\psi$: discriminator
with parameters $\psi$

$\hat{x}$

$G_\theta$: generator
with parameters $\theta$

$z$

sample

$p(z)$

- If discriminator approximates $p_{\text{data}}(x)$:
- $x^*$ at maximum of $p_{\text{data}}(x)$ has lowest loss
- Optimal $p_\theta(x)$ has single mode at $x^*$, small variance

$$D_\psi \approx p_{\text{data}}(\hat{x})$$

$x^*$    $p_{\text{data}}(x)$         $p_\theta(x)$

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

SIGGRAPH
ASIA 2018
TOKYO

# Why Adversarial?

$$\approx \frac{p_{\text{data}}(\hat{x})}{p_{\text{data}}(\hat{x}) + p_\theta(\hat{x})}$$

$D_\psi$: discriminator
with parameters $\psi$

$\hat{x}$

$G_\theta$ : generator
with parameters $\theta$

$z$

sample

$p(z)$

• For GANs, the discriminator instead approximates:

$$\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)}$$

→ depends on the generator

$D_\psi \approx p_{\text{data}}(\hat{x})$    $D_\psi \approx \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)}$

$x^*$    $p_{\text{data}}(x)$    $p_\theta(x)$    $p_\theta(x)$

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

---

# Why Adversarial?

$p_{\text{data}}(x)$

$p_\theta(x)$

$p_\theta(x)$

$p_\theta(x)$

VAEs:
Maximize likelihood of
**data samples** in $p_\theta(x)$

GANs:
Adversarial game

Maximize likelihood of
**generator samples** in
approximate $p_{\text{data}}(x)$

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

# Why Adversarial?

$$\approx KL(p_{\text{data}} \parallel p_\theta) \quad \approx JS(p_{\text{data}} \parallel p_\theta) \quad \approx KL(p_\theta \parallel p_{\text{data}})$$



$p_{\text{data}}(x)$     $p_\theta(x)$     $p_\theta(x)$     $p_\theta(x)$

| $p_{\text{data}}(x)$ | VAEs: Maximize likelihood of **data samples** in $p_\theta(x)$ | GANs: Adversarial game | Maximize likelihood of **generator samples** in approximate $p_{\text{data}}(x)$ |

Image Credit: *How (not) to Train your Generative Model: Scheduled Sampling, Likelihood, Adversary?*, Ferenc Huszár

# GAN Objective



probability that $\hat{x}$ is not fake

$D_\psi(\hat{x})$

$D_\psi$ :discriminator

$\hat{x}$

$G_\theta$ :generator

$z$

sample

$p(z)$

fake/real classification loss (BCE):
$$L(\theta, \psi) = -0.5 \, \mathbb{E}_{x \sim p_{\text{data}}} \, \log D_\psi(x)$$
$$- 0.5 \, \mathbb{E}_{x \sim p_\theta} \, \log(1 - D_\psi(x))$$

Discriminator objective:
$$\min_\psi L(\theta, \psi)$$

Generator objective:
$$\max_\theta L(\theta, \psi)$$

# Non-saturating Heuristic

$$L(\theta, \psi) = -0.5 \, \mathbb{E}_{x \sim p_{\text{data}}} \, \log D_\psi(x)$$
$$- 0.5 \, \mathbb{E}_{x \sim p_\theta} \, \log(1 - D_\psi(x))$$

Generator loss is negative binary cross-entropy:

$$L_G(\theta, \psi) = 0.5 \, \mathbb{E}_{x \sim p_\theta} \, \log(1 - D_\psi(x)) \quad \text{poor convergence}$$



Image Credit: NIPS 2016 Tutorial: Generative Adversarial Networks, Ian Goodfellow

---

# Non-saturating Heuristic

Generator loss is negative binary cross-entropy:

$$L_G(\theta, \psi) = 0.5 \, \mathbb{E}_{x \sim p_\theta} \, \log(1 - D_\psi(x)) \quad \text{poor convergence}$$

Flip target class instead of flipping the sign for generator loss:

$$L_G(\theta, \psi) = -0.5 \, \mathbb{E}_{x \sim p_\theta} \, \log D_\psi(x) \quad \text{good convergence – like BCE}$$



Image Credit: NIPS 2016 Tutorial: Generative Adversarial Networks, Ian Goodfellow

# GAN Training

Generator training

Loss:
$$L_G(\theta, \psi) = -\log D_\psi(\hat{x})$$

$D_\psi(\hat{x})$

$\uparrow$ $D_\psi$ :discriminator



$\hat{x}$

$\uparrow$ $G_\theta$ :generator



$z$

Sample $\curvearrowleft$

$p(z)$

Discriminator training

Loss:
$$L_D(\theta, \psi) = -0.5\log(1 - D_\psi(\hat{x})) - 0.5\log D_\psi(x_i)$$

$D_\psi(\hat{x})$ $D_\psi(x_i)$

$\uparrow$ $D_\psi$ :discriminator $\uparrow$

 

$\hat{x}$ $x_i$

from dataset

Interleave in each training step

# DCGAN

- First paper to successfully use CNNs with GANs
- Due to using novel components (at that time) like batch norm., ReLUs, etc.



man
with glasses  $-$  man
without glasses  $+$  woman
without glasses  $=$  woman with glasses

Image Credit: *Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks*, Radford et al.

# InfoGAN

varying $c_i$



$c_1$

$c_2$

$c_3$

$D_\psi(\hat{x})$

$D_\psi$ :discriminator

$\hat{x}$

$G_\theta$ :generator

maximize mutual information

$I(c; G_\theta(z, c))$

sample

$z$  $c$

$p(z)$

Image Credit: *InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets*, Chen et al.

# Code example

**Generative Adversarial Network (gan.ipynb)**

38

# Conditional GANs (CGANs)

- ≈ learn a mapping between images from example pairs
- Approximate sampling from a conditional distribution $p_{\text{data}}(x \mid c)$

$$c \qquad x \qquad\qquad c \qquad x \qquad\qquad c \qquad x$$



Image Credit: *Image-to-Image Translation with Conditional Adversarial Nets*, Isola et al.

# Conditional GANs

**Generator training**

Loss:
$$L_G(\theta, \psi) = -\log D_\psi(\hat{x}, c)$$

$D_\psi(\hat{x})$

$D_\psi$ :discriminator

$\hat{x}$

$G_\theta$ :generator

$z$

$c$

sample $p(z)$

**Discriminator training**

Loss:
$$L_D(\theta, \psi) = -0.5 \log(1 - D_\psi(\hat{x}, c)) - 0.5 \log D_\psi(x_i, c)$$

$D_\psi(\hat{x})$      $D_\psi(x_i)$

$D_\psi$ :discrim.

$\hat{x}$      $x_i$

from dataset

$c$      $c$

Image Credit: *Image-to-Image Translation with Conditional Adversarial Nets*, Isola et al.

# Conditional GANs: Low Variation per Condition

**Generator training**

Loss:
$$L_G(\theta, \psi) = -\log D_\psi(\hat{x}, c)$$

$D_\psi(\hat{x})$

$D_\psi$ :discriminator

$\hat{x}$

$G_\theta$ :generator

$z$ is often omitted
in favor of dropout
in the generator

$c$

**Discriminator training**

Loss:
$$L_D(\theta, \psi) = -0.5\log(1 - D_\psi(\hat{x}, c)) - 0.5\log D_\psi(x_i, c)$$

$D_\psi(\hat{x})$      $D_\psi(x_i)$

$D_\psi$ :discrim.

$\hat{x}$      $x_i$

from dataset

$c$      $c$

Image Credit: *Image-to-Image Translation with Conditional Adversarial Nets*, Isola et al.

SIGGRAPH
ASIA 2018
T O K Y O

---

# Demos

**CGAN**
https://affinelayer.com/pixsrv/index.html

42

# CycleGANs

- Less supervision than CGANs: mapping between unpaired datasets
- Two GANs + cycle consistency



Image Credit: *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks*, Zhu et al.

# CycleGAN: Two GANs …

- Not conditional, so this alone does not constrain generator input and output to match



Image Credit: *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks*, Zhu et al.

# CycleGAN: … and Cycle Consistency



L1 Loss function:
$$\|y - \tilde{y}\|_1$$

$\tilde{y}$

$F_{\theta_2}$ :generator2

$\hat{x}$

$G_{\theta_1}$ :generator1

$y$

L1 Loss function:
$$\|x - \tilde{x}\|_1$$

$\tilde{x}$

$G_{\theta_1}$ :generator1

$\hat{y}$

$F_{\theta_2}$ :generator2

$x$

Image Credit: *Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks*, Zhu et al.

# Unstable Training

GAN training can be unstable

Three current research problems (may be related):

• Reaching a Nash equilibrium (the gradient for both $L_G$ and $L_D$ is 0)

• $p_\theta$ and $p_{\text{data}}$ initially don't overlap

• Mode Collapse

# GAN Training

- Vector-valued loss: $\mathbf{L}(\theta, \psi) = \begin{pmatrix} L_G(\theta, \psi) \\ L_D(\theta, \psi) \end{pmatrix}$

- In each iteration, gradient descent approximately follows this vector over the parameter space $(\theta, \psi)$:

$$\mathbf{V}(\theta, \psi) = \begin{pmatrix} \frac{\partial}{\partial \theta} L_G(\theta, \psi) \\ \frac{\partial}{\partial \psi} L_D(\theta, \psi) \end{pmatrix}$$

$\psi$

$\theta$

$\longrightarrow \ \frac{\partial}{\partial \theta} L_G(\theta, \psi)$

$\longrightarrow \ \frac{\partial}{\partial \psi} L_D(\theta, \psi)$

$\longrightarrow \ \mathbf{V}(\theta, \psi)$

SIGGRAPH
ASIA 2018
T O K Y O

# Reaching Nash Equilibrium

Gradient field example

$\mathbf{V}(\theta, \psi)$ Example

Nash equilib.

$\psi$

$\theta$

SIGGRAPH
ASIA 2018
T O K Y O

Image Credit: *GANs are Broken in More than One Way: The Numerics of GANs*, Ferenc Huszár

# Reaching Nash Equilibrium

Solution attempt: relaxation with term: $-\nabla L = \frac{\partial}{\partial \theta} \left\| \mathbf{V}(\theta, \psi) \right\|_2^2$

Non-conservative field $v$     Conservative field $-\nabla L$     Combined field $v - 0.6 \nabla L$

saddle points

equilibria

no relaxation has cycles     full relaxation introduces bad Nash equilibria     mixture works sometimes

Image Credit: *GANs are Broken in More than One Way: The Numerics of GANs*, Ferenc Huszár

# Generator and Data Distribution Don't Overlap

Standard

$p_\theta(x)$
$p_{\text{data}}(x)$
$D_\psi(x)$

Instance noise: adding noise to generated and real images

$p_\theta(x) * N(0, \sigma)$
$p_{\text{data}}(x) * N(0, \sigma)$
$D_\psi(x)$

Wasserstein GANs: EMD as distance between $p_\theta$ and $p_{\text{data}}$

$p_\theta(x)$
$p_{\text{data}}(x)$
$D_\psi(x)$

Roth et al. suggest an analytic convolution with a gaussian:
*Stabilizing Training of Generative Adversarial Networks through Regularization, Roth et al. 2017*

Image Credit: *Amortised MAP Inference for Image Super-resolution*, Sønderby et al.

# Mode Collapse

Optimal $D_\psi(x)$:   $\dfrac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_\theta(x)}$

$p_\theta$ only covers one or a few modes of $p_{\text{data}}$



$p_{\text{data}}$            $p_\theta$ after n training steps



Image Credit: *Wasserstein GAN*, Arjovsky et al.
*Unrolled Generative Adversarial Networks*, Metz et al.

# Mode Collapse

Solution attempts:

- Minibatch comparisons:  Discriminator can compare instances in a minibatch (*Improved Techniques for Training GANs*, Salimans et al.)
- Unrolled GANs: Take k steps with the discriminator in each iteration, and backpropagate through all of them to update the generator



$p_{\text{data}}$    Standard GAN         $p_\theta$ after n training steps

Unrolled GAN with k=5

$p_\theta$ after n training steps

Image Credit: *Wasserstein GAN*, Arjovsky et al.
*Unrolled Generative Adversarial Networks*, Metz et al.

# Progressive GANs

- Resolution is increased progressively during training
- Also other tricks like using minibatch statistics and normalizing feature vectors



Image Credit: *Progressive Growing of GANs for Improved Quality, Stability, and Variation*, Karras et al.

53

# Disentanglement

$z$

$z_a$ $z_b$ $\cdots$

Entangled: different properties may be mixed up over all dimensions

Disentangled: different properties are in different dimensions

specified property: number

other properties



specified property: character

other properties



Image Credit: *Disentangling factors of variation in deep representations using adversarial training*, Mathieu et al.

# Summary

- Autoencoders
  - Can infer useful latent representation for a dataset
  - Bad generators
- VAEs
  - Can infer a useful latent representation for a dataset
  - Better generators due to latent space regularization
  - Lower quality reconstructions and generated samples (usually blurry)
- GANs
  - Can not find a latent representation for a given sample (no encoder)
  - Usually better generators than VAEs
  - Currently unstable training (active research)

---

# Course Information (slides/code/comments)



http://geometry.cs.ucl.ac.uk/creativeai/

# CreativeAI: Deep Learning for Graphics

# Feature Visualization

| Niloy Mitra | Iasonas Kokkinos | Paul Guerrero | Nils Thuerey | Tobias Ritschel |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

facebook
Artificial Intelligence Research

---

# Timetable

| | | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|---|---|---|---|---|
| Theory and Basics | Introduction | X | X | X | X | X |
| | Theory | X | | | X | |
| | NN Basics | X | X | | | |
| | Alternatives to Direct Supervision | | | X | | |
| | *15 min. break* | | | | | |
| State of the Art | **Feature Visualization** | | | | | **X** |
| | Image Domains | | X | | | X |
| | 3D Domains | | | X | | X |
| | Motion and Physics | X | | | X | |

# What to Visualize

- Features (activations)
- Weights (filter kernels in a CNN)
- Inputs that maximally activate some class probabilities or features
- Inputs that maximize the error (adversarial examples)

# Feature Samples

- In good training, features are usually sparse
- Can find "dead" features that never activate



Images from: http://cs231n.github.io/understanding-cnn/

# Feature Distribution using t-SNE

- Low-dimensional embedding of the features for visualization



t-SNE embedding of image features
in a CNN layer

before training                                    after training
t-SNE embedding of MNIST (images of digits) features in a CNN layer, colored by class

Images from: https://cs.stanford.edu/people/karpathy/cnnembed/ and
Rauber et al. *Visualizing the Hidden Activity of Artificial Neural Networks*. TVCG 2017

# Weights

- Useful for CNN kernels, not useful for fully connected layers
- Kernels are typically smooth and diverse after a successful training



first layer filters of AlexNet

Images from: http://cs231n.github.io/understanding-cnn/

# Inputs that Maximize Feature Response

Local maxima of the response for class:

| Indian Cobra | Pelican | Ground Beetle |



Images from: Yosinski et al. *Understanding Neural Networks Through Deep Visualization*. ICML 2015

# Inputs that Maximize the Error

$$\max_{\delta \in \Delta} \mathcal{L}(x + \delta, y; \theta) \qquad \Delta = \{\delta \in \mathbb{R}^d \mid \|\delta\|_p \leq \varepsilon\}$$



$+.007 \times$

$x$ — "Panda" 55.7% conf.

$\delta$

$x + \delta$ — "Gibbon" 99.3% conf.

Images from: Goodfellow et al. *Explaining and Harnessing Adversarial Examples*. ICLR 2015

# Course Information (slides/code/comments)

http://geometry.cs.ucl.ac.uk/creativeai/

SIGGRAPH ASIA 2018 TOKYO

SIGGRAPH Asia Course **CreativeAI: Deep Learning for Graphics**

# CreativeAI: Deep Learning for Graphics

# Image Domains

| Niloy Mitra | Iasonas Kokkinos | Paul Guerrero | Nils Thuerey | Tobias Ritschel |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

---

# Timetable

| | | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|---|---|---|---|---|
| Theory and Basics | Introduction | X | X | X | X | X |
| | Theory | X | | | X | |
| | NN Basics | X | X | | | |
| | Alternatives to Direct Supervision | | | X | | |
| | 15 min. break | | | | | |
| State of the Art | Feature Visualization | | | | | X |
| | Image Domains | | X | | | X |
| | 3D Domains | | | X | | X |
| | Motion and Physics | X | | | X | |

# Sketch Simplification

- *Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, Simon-Serra et al., 2016

- *Deep Extraction of Manga Structural Lines*, Li et al., 2017



# Sketch Simplification: *Learning to Simplify*



*Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, Simo-Serra et al.

# Sketch Simplification: *Learning to Simplify*

- Loss for thin edges saturates easily
- Authors take extra steps to align input and ground truth edges



Pencil: input
Red: ground truth

*Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, **Simo-Serra et al.**

5

# Image Decomposition

- A selection of methods:
- *Direct Instrinsics*, Narihira et al., 2015
- *Learning Data-driven Reflectance Priors for Intrinsic Image Decomposition,* Zhou et al., 2015
- *Decomposing Single Images for Layered Photo Retouching*, Innamorati et al. 2017



6

# Image Decomposition: Decomposing
# *Single Images for Layered Photo Retouching*



→ Stride-2 Convolution    ⇒ Stride-1 Convolution    → Resize-convolution

7

# Colorization

- Concurrent methods:
    - *Let there be Color!*, Iizuka et al., 2016
    - *Colorful Image Colorization*, Zhang et al. 2016
    - *Learning Representations for Automatic Colorization,* Larsson et al., 2016
    - *Real-Time User-Guided Image Colorization with Learned Deep Priors*, Zhang et al. 2017



8

# Colorization: *Let There Be Color!*



*Let there be Color!:* Iizuka et al.

# Colorization: *Colorful Image Colorization*



input output
direct regression  probability distr.

Image Credit: *Colorful Image Colorization*, Zhang et al.

# Sketch Simplification

- *Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, Simon-Serra et al., 2016
- *Deep Extraction of Manga Structural Lines*, Li et al., 2017



11

# Sketch Simplification: *Learning to Simplify*



*Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, Simo-Serra et al.

12

# Sketch Simplification: *Learning to Simplify*

- Loss for thin edges saturates easily
- Authors take extra steps to align input and ground truth edges



Pencil: input
Red: ground truth

*Learning to Simplify: Fully Convolutional Networks for Rough Sketch Cleanup*, **Simo-Serra et al.**

13

# Image Decomposition

- A selection of methods:
- *Direct Instrinsics*, Narihira et al., 2015
- *Learning Data-driven Reflectance Priors for Intrinsic Image Decomposition,* Zhou et al., 2015
- *Decomposing Single Images for Layered Photo Retouching*, Innamorati et al. 2017



14

# Image Decomposition: Decomposing
## *Single Images for Layered Photo Retouching*



→ Stride-2 Convolution    ⇒ Stride-1 Convolution    → Resize-convolution

15

# Colorization

- Concurrent methods:
    - *Let there be Color!*, Iizuka et al., 2016
    - *Colorful Image Colorization*, Zhang et al. 2016
    - *Learning Representations for Automatic Colorization,* Larsson et al., 2016
    - *Real-Time User-Guided Image Colorization with Learned Deep Priors*, Zhang et al. 2017



16

# Colorization: *Let There Be Color!*



*Let there be Color!:* Iizuka et al.

# Colorization: *Colorful Image Colorization*



input

output

direct regression    probability distr.

Image Credit: *Colorful Image Colorization*, Zhang et al.

# LDR to HDR Image Reconstruction:

- Concurrently:
- *Deep Reverse Tone Mapping*, Endo et al. 2017
- *HDR image reconstruction from a single exposure using deep CNNs*, Eilertsen et al. 2017

19

# Reflectance Maps

- Paint a sphere as if it is made of a material under a certain illumination of another object in a photo



*Deep Reflectance Maps*. Rematas et al. **CVPR 2015**

20

# DeLight

- Factor BRDF and (HDR) Illumination



*Reflectance and Natural Illumination from Single-Material Specular Objects Using Deep Learning.* Georgoulis et al. **PAMI 2017**

21

# Denoising Renderings

- Concurrent:

- *Kernel-Predicting Convolutional Networks for Denoising Monte Carlo Renderings,* Bako et al. 2017

- *Interactive Reconstruction of Monte Carlo Image Sequences using a Recurrent Denoising Autoencoder,* Chaitanya et al. 2017 (more on Autoencoders later)



TRAINING
*Kernel-Predicting Convolutional Networks for Denoising Monte Carlo Renderings*, **Bako et al.**

22

# Denoising Renderings:



*Kernel-Predicting Convolutional Networks for Denoising Monte Carlo Renderings*, Bako et al. **SIGGRAPH 2017**

23

# 3D Pose Estimation: *VNECT*



*VNect: Real-time 3D Human Pose Estimation with a Single RGB Camera*, Mehta et al., **SIGGRAPH 2017**

# Object Detection: Fast(er)-RCNN

- Fast/Faster R-CNN
  - ✓Good speed
  - ✓Good accuracy
  - ✓Intuitive
  - ✓Easy to use



Ross Girshick. "Fast R-CNN". ICCV 2015.
Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

# Mask R-CNN

- Mask R-CNN = **Faster R-CNN** with **FCN** on RoIs



FCN on RoI

# Mask R-CNN results on COCO



# Mask R-CNN for Human Keypoint Detection

- 1 keypoint = 1-hot "mask"
- Human pose = 17 masks
- Softmax over spatial locations
  - e.g. $56^2$-way softmax on 56x56

Mask R-CNN frame-by-frame



Mask R-CNN frame-by-frame

# UberNet: a "universal" network for all tasks



https://github.com/jkokkin/UberNet

I. Kokkinos, UberNet: Training a Universal CNN for *Low- Mid- and High-Level* Vision, CVPR 2017

# What is the ultimate vision task?

"Inverse graphics": understand how an image was generated from a scene

If we focus on a single object category: surface-based models



**UberNet:**
**Universal Network**



**DensePose:**
**Unified model**

# DenseReg: dense image-to-face regression



R. A. Guler, G. Trigeorgis, E. Antonakos, P. Snape, S. Zafeiriou, I. Kokkinos,

DenseReg: Fully Convolutional Dense Shape Regression In-the-Wild, CVPR 2017

# DensePose: dense image-to-body correspondence



DensePose-RCNN: ~25 FPS

R. A. Guler, N. Neverova, I. Kokkinos "DensePose: Dense Human Pose Estimation In The Wild", CVPR'18

# SFSNet: incorporating image formation in model



SfSNet: Learning Shape, Reflectance and Illuminance of Faces 'in the wild' Soumyadip Sengupta Angjoo Kanazawa Carlos D. Castillo David W. Jacobs, CVPR 2018

# Beyond single frames: end-to-end optical flow



FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks

Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, Thomas Brox

University of Freiburg, Germany

——— Supplementary Material ———

# End-to-end Structure From Motion



- DeMoN: Depth and Motion Network for Learning Monocular Stereo, B. Ummenhofer, et al, CVPR 2017
- Unsupervised learning of depth and ego-motion from video, T Zhou, M Brown, N Snavely, DG Lowe, CVPR 2017

# Monocular depth & normal estimation



Input Image          Depth          Normals          Labels

- D. Eigen and R. Fergus, Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-Scale Convolutional Architecture, ICCV 2015

---

# Course Information (slides/code/comments)



http://geometry.cs.ucl.ac.uk/creativeai/

# CreativeAI: Deep Learning for Graphics

# 3D Domains

| **Niloy Mitra** | **Iasonas Kokkinos** | **Paul Guerrero** | **Nils Thuerey** | **Tobias Ritschel** |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

# Timetable

|  |  | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|:---:|:---:|:---:|:---:|:---:|
| **Theory and Basics** | Introduction | X | X | X | X | X |
|  | Theory | X |  |  | X |  |
|  | NN Basics | X | X |  |  |  |
|  | Alternatives to Direct Supervision |  |  | X |  |  |
|  | — 15 min. break — |  |  |  |  |  |
| **State of the Art** | Feature Visualization |  |  |  |  | X |
|  | Image Domains |  | X |  |  | X |
|  | **3D Domains** |  |  | **X** |  | **X** |
|  | Motion and Physics | X |  |  | X |  |

# Motivating Applications

Deep neural network predicts the **next best part** to add and its **position** to enable non-expert users to create novel shapes.

[Sung et al. 2017]

3

# CrossLink: Linking Images and 3D Models

[Heuting et al. 2015]

4

# Motivating Applications

understanding 3D shapes can benefit image understanding



**Physically based Rendering**

[Zhang et al. 2017]

5

---

# Motivating Applications: Semantic Scene Understanding



(a) depth

(b) visible surface

floor
wall
bed
window
sofa
objects
furniture

(c) output

[Song et al. 2017]

6

# Motivating Applications: Semantic Scene Understanding

[Kelly et al. 2017]

# Representation for 3D

- Image-based

- Volumetric

- Point-based

- Surface-based

- Parametric

9

# Representation for 3D

- **Image-based**

- Volumetric

- Point-based

- Surface-based

- Parametric

1
0

# Representation for 3D: Multi-view CNN



**regular image analysis networks**

3D shape model
rendered with
different virtual cameras

2D rendered
images

our multi-view CNN architecture

output class
predictions

[Kalogerakis et al. 2015]

1
1

# Representation for 3D: Local Multi-view CNN



view based convolutional network

point descriptor

Segmentation
Correspondence
Feature matching
Predicting semantic functions

**localized renderings for point-wise features**        [Huang et al. 2018]

1
2

# 3D-R$^2$N$^2$ (3D Recurrent Reconstruction Neural Network)

- Multiple views are treated as image sequence
- An LSTM controls what part of the latent representation is updated by each view



Choy et al. *3d-r2n2: A unified approach for single and multi-view 3d object reconstruction*. ECCV 2016

---

# Representation for 3D

- **Image-based**
  - **PROS**: directly use image networks, good performance
  - **CONS**: rendering is slow and memory-heavy, not very geometric

- Volumetric
- Point-based
- Surface-based
- Parametric

# Representation for 3D

- Image-based

- **Volumetric**

- Point-based

- Surface-based

- Parametric

1
5

# Representation for 3D: Volumetric



4000

object label 10          1200

512 filters of
stride 1

160 filters of
stride 2

48 filters of
stride 2

3D voxel input

- Add one dimension to kernels and intermediate outputs:
batches x channels x w x h          batches x channels x d x
w x h

- Does not scale well to high resolutions

[Xiao et al. 2014]

16

214

# Representation for 3D: Volumetric Deformation



[Yumer et al. 2014]

# Efficient Volumetric Datastructures



[Wang et al. 2017]

# Efficient Volumetric Datastructures

Generator / Decoder

Encoder

Volumetric (Up-) Convolutions

Cropping

Wang et al. 2017

*only generate non-empty voxels*

3 Labels (free space / boundary / occupied space)

[Hane et al. 2018]

19

# Efficient Volumetric Datastructures

$64^3$          $128^3$          $256^3$

[Hane et al. 2018]

20

# Octree Generating Networks



Tatarchenko et al. *Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs.*
ICCV 2017

# Deep Marching Cubes

• Input Domain: images, volumetric grids, point clouds
• Output Domain: Meshes



Yiyi Liao et al. *Deep Marching Cubes. CVPR 2018*

# Learning to Complete 3D Scans    (slide credit: Matthias Niessner)



Input Partial Scan      Completed Scan      Predicted Semantics

[Dai et al. 2018]

# State-of-the-art 3D Reconstructions    (slide credit: Matthias Niessner)



TOG'17 [Dai et al.]: BundleFusion

# Problem: Incomplete Scan Geometry

# Problem: Incomplete Scan Geometry

# Learning from Synthetic Data

# Recall: Semantic Scene Understanding



(a) depth

(b) visible surface

floor
wall
bed
window
sofa
objects
furniture

(c) output

[Song et al. 2017]

# Learning to Complete 3D Scans  (slide credit: Matthias Niessner)



Scenes from SUNCG [Song et al. 17]

2
9

# Dependent Predictions: Autoregressive Neural Networks

- PixelCNN [van den Oord 2015, van den Oord 2016, Reed 2017]



- WaveNet [van den Oord 2016]



3
0

# Dependent Predictions: Autoregressive Neural Networks



1

1→2

1,2→3

1,2,3→4

1,2,3,4→5

1,2,3,4,5→6

1,2,3,4,5,6→7

1,2,3,4,5,6,7→8

[Dai et al. 2018]

31

# ScanComplete: Completing 3D Scans

Input          Completion          Ground Truth



[Dai et al. 2018]

32

# ScanComplete: Completing 3D Scans



[Dai et al. 2018]

# Geometry Abstraction / Simplification



*Learning Shape Abstractions by Assembling Volumetric Primitives*, Tulsiani et al. 2016

# Geometry Abstraction / Simplification:



*Learning Shape Abstractions by Assembling Volumetric Primitives*, Tulsiani et al. 2016

35

# SplatNet



Hang Su et al. *Splatnet: Sparse lattice networks for point cloud processing.* CVPR 2018

# Representation for 3D

- Image-based
- Volumetric
  - **PROS**: modify image networks
  - **CONS**: special layers for hierarchical datastructures, still too coarse
- Point-based
- Surface-based
- Parametric

# Representation for 3D

- Image-based

- Volumetric

- **Point-based**

- Surface-based

- Parametric

# Point Clouds

- Common representation
- Easy to obtain from meshes, depth scans, laser scans
- Difficulty: invariance to point order



point cloud $\mathbb{P}$

$$\begin{array}{cc} p_1 & p_3 \\ p_2 & p_1 \\ p_3 & p_2 \\ \vdots & \vdots \end{array} \quad \longrightarrow \quad F(\mathbb{P})$$

invariance to all $n!$ cases

# Point Interpretation

Samples from
a probability distribution

(Irregular) samples of
a continuous function

# PointNet



point cloud $\mathbb{P}$

order-independent

$$p_1 \quad p_3 - f(p_3)$$
$$p_2 \quad p_1 - f(p_1) \longrightarrow \sum_{p_j \in \mathbb{P}} g(f(p_j)) = F(\mathbb{P})$$
$$p_3 \quad p_2 - f(p_2)$$

symmetric

order-independent

feature vector for a **point**

feature vector for the **point cloud**

Qi et al. *Pointnet: Deep learning on point sets for 3d classification and segmentation*. CVPR 2017

---

# PointNet for Point Cloud Analysis



PointNet

mug?
table?
car?

Classification | Part Segmentation | Semantic Segmentation

# PointNet for Point Cloud Analysis: PointNet++



N points in (x,y)

N₁ points in (x,y,**f**)

N₂ points in (x,y,**f'**)

[Qi et al. 2018]

4
3

# PointNet for Local Point Cloud Analysis



[Guerrero et al. 2018]

4
4

# MCCNN



Hermosilla et al. *Monte Carlo Convolution for Learning on Non-Uniformly Sampled Point Clouds*. SIGGRAPH Asia 2018

# PointNet for Point Cloud Synthesis

generated output needs to be compare to some true shape



Input

Earth Mover Distance as loss function

$$d_{EMD}(S_1, S_2) = \min_{\phi:S_1 \to S_2} \sum_{x \in S_1} \|x - \phi(x)\|_2$$

[Su et al. 2017]

46

# Representation for 3D

- Image-based
- Volumetric
- Point-based
- **Surface-based**

Surface models used in engineering (i.e., CAD) and computer graphics (i.e., meshes)



**Image**      **Generated Volume**      **Generated Points**      **Generated Surface**

47

# AtlasNet for Surface Generation

*condition decoded points on 2D patches*



Latent shape representation

Sampled 2D point

MLP

Generated 3D point

[Groueix et al. 2018]

48

# AtlasNet for Surface Generation

*condition decoded points on 2D patches*

2D Image

3D Point Cloud

Latent representation can be inferred from images or point clouds

Sampled 2D point

MLP

Generated 3D point

[Groueix et al. 2018]

49



# AtlasNet for Surface Generation

2D Image

3D Point Cloud

Latent representation can be inferred from images or point clouds

Sampled 2D point

MLP

Quad Mesh is generated by mapping a regular grid in 2D domain to 3D points

[Groueix et al. 2018]

50

# AtlasNet for Surface Generation

BONUS: natural space to store
textures for CG

Latent representation can be
inferred from images or point clouds

2D Image

3D Point Cloud

Sampled
2D point

MLP

5
1

# Texture Transfer

input image

image-to-shape texture transfer

shape-to-shape texture transfer

edited original image

[Wang et al. 2016]

5
2

# Parameterization for Surface Analysis

map 3D surface to 2D domain



[Maron et al. 2017]

5
3

# Parameterization for Surface Analysis

map 3D surface to 2D domain



$\mathcal{T}$

[Maron et al. 2017]

5
4

# Parameterization for Surface Analysis

• Map 3D surface to 2D domain

• One such mapping: flat torus (seamless => translation-invariant)

• Many mappings exists: sample a few and average result

• Which functions to map?
XYZ, normals, curvature, …

[Maron et al. 2017]
5
5

# Parameterization for Surface Analysis



[Maron et al. 2017]
5
6

# Other Parameterizations

**Geometry Image**

**Metric Alignment**



[Sinha et al. 2017]

[Ezuz et al. 2017]

5
7

# Other Parameterizations

geodesic discs

parameterize in spectral domain



Spatial domain

**Spectral domain**

5
8

# Other Parameterizations



[Masci et al. 2015]

59

# Discrete Laplacian

(slide credit: Michael Bronstein)



**Undirected graph** $(\mathcal{V}, \mathcal{E})$

$(\Delta f)_i \approx \sum_{(i,j)\in\mathcal{E}} w_{ij}(f_i - f_j)$

**Triangular mesh** $(\mathcal{V}, \mathcal{E}, \mathcal{F})$

$(\Delta f)_i \approx \dfrac{1}{a_i} \sum_{(i,j)\in\mathcal{E}} \dfrac{\cot\alpha_{ij}+\cot\beta_{ij}}{2}(f_i - f_j)$

$a_i = $ local area element

$f = \hat{f}_1 \phi_1 + \hat{f}_2 \phi_2 + \hat{f}_3 \phi_3 + \ldots$

60

# Transferring Correspondence



Reference

Texture transferred from reference to query shapes

[Monti et al. 2016]

6
1

# Spectral Methods      (slide credit: Michael Bronstein)



GEOMETRIC DEEP LEARNING    ABOUT   WORKSHOPS   TUTORIALS   PAPERS & CODE   CONTACTS

## GEOMETRIC DEEP LEARNING

Geometric Deep Learning is one of the most emerging fields of the Machine Learning community. This website represents a collection of materials of this particular research area.

FIND OUT MORE

6
2

# SyncSpecCNN



before synchronization                    after synchronization

[Yi et al. 2017]

63

# 3D volumes form Xrays



*Single-Image Tomography: 3D Volumes from 2D Cranial X-Rays.* Henzler et al. **EG 2018**

64

# Representation for 3D

- Image-based

- Volumetric

- Point-based

- Surface-based

  **Parametric**

# Procedural Parameter Estimation



a) Sketch    b) Suggested snippets    c) Generated building    d) Rendering of city corner

*Interactive Sketching of Urban Procedural Models*, Nishida et al. 2016

# Procedural Parameter Estimation:
## *Interactive Sketching of Urban Procedural Models*



*Interactive Sketching of Urban Procedural Models*, **Nishida et al.**

67

# Course Information (slides/code/comments)



http://geometry.cs.ucl.ac.uk/creativeai/

SIGGRAPH Asia Course **CreativeAI: Deep Learning for Graphics**

# CreativeAI: Deep Learning for Graphics

# Motion and Physics

| **Niloy Mitra** | **Iasonas Kokkinos** | **Paul Guerrero** | **Nils Thuerey** | **Tobias Ritschel** |
|---|---|---|---|---|
| UCL | UCL/Facebook | UCL | TU Munich | UCL |

**facebook**
Artificial Intelligence Research

---

# Timetable

|  |  | Niloy | Iasonas | Paul | Nils | Tobias |
|---|---|---|---|---|---|---|
| Theory and Basics | Introduction | X | X | X | X | X |
|  | Theory | X |  |  | X |  |
|  | NN Basics | X | X |  |  |  |
|  | Alternatives to Direct Supervision |  |  | X |  |  |
|  | *15 min. break* |  |  |  |  |  |
| State of the Art | Feature Visualization |  |  |  |  | X |
|  | Image Domains |  | X |  |  | X |
|  | 3D Domains |  |  | X |  | X |
|  | **Motion and Physics** | **X** |  |  | **X** |  |

# Deep Learning for Fluids

Tompson et. al 2017



Long et. al 2017



Schenck et. al 2017

3

# High Resolution Simulation of Liquids

Latent-space encoding

Volumetric decoding

Temporal prediction

[Latent-space Physics: Towards Learning the Temporal Evolution of Fluid Flow, arXiv 2018]
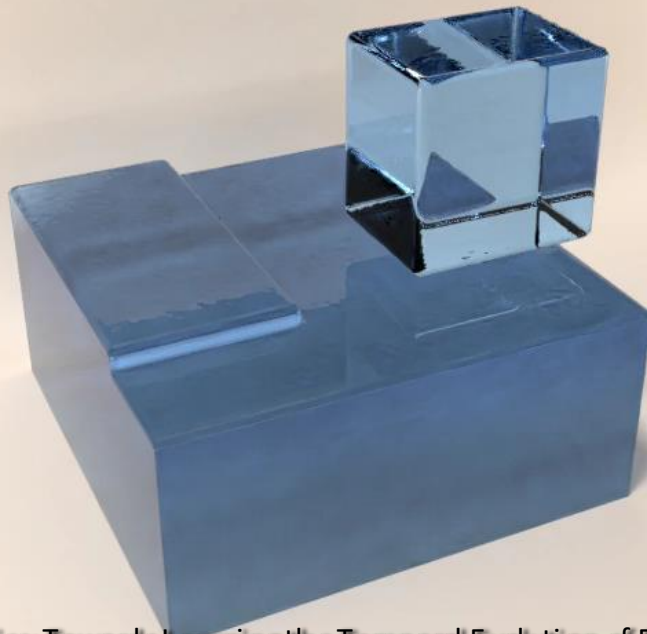
4

(slide credit: Nils Thuerey)

Examples from training data set, 64³
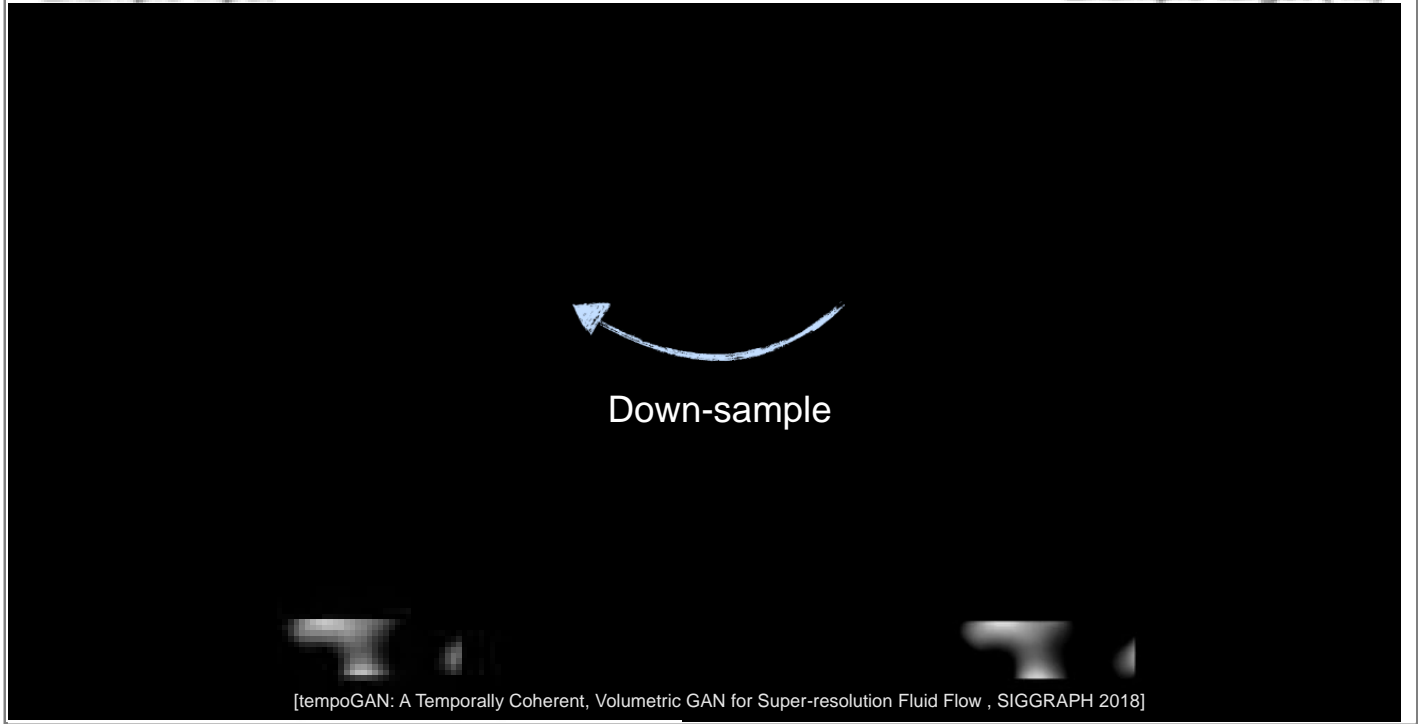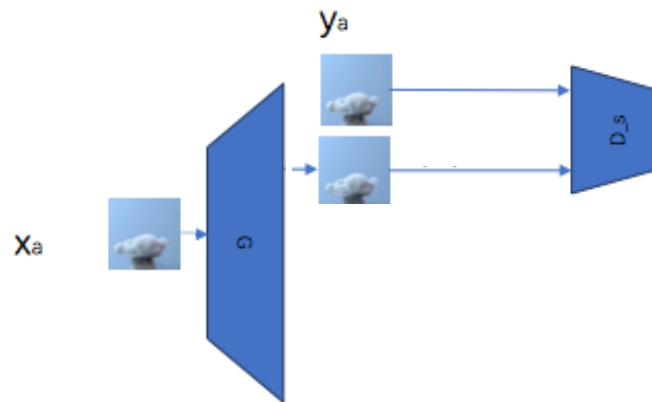
5



(slide credit: Nils Thuerey)

Further Examples, 128³

[Latent-space Physics: Towards Learning the Temporal Evolution of Fluid Flow, arXiv 2018]

Example input          (slide credit: Nils Thuerey)          Example target (4x)

Down-sample

[tempoGAN: A Temporally Coherent, Volumetric GAN for Super-resolution Fluid Flow , SIGGRAPH 2018]

# Architecture Overview          (slide credit: Nils Thuerey)



Xa          G

8

# Architecture Overview

**(slide credit: Nils Thuerey)**



# Architecture Overview

**(slide credit: Nils Thuerey)**
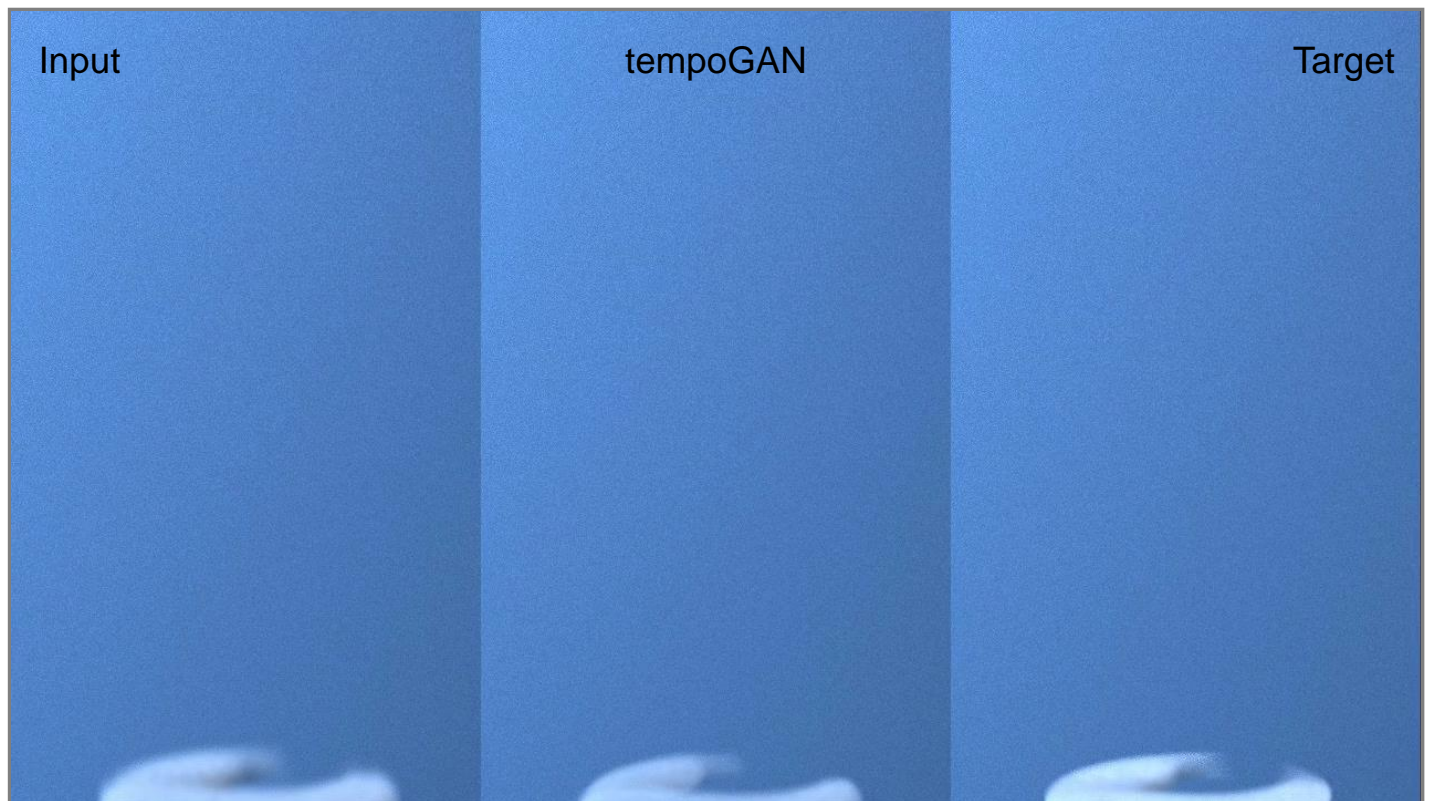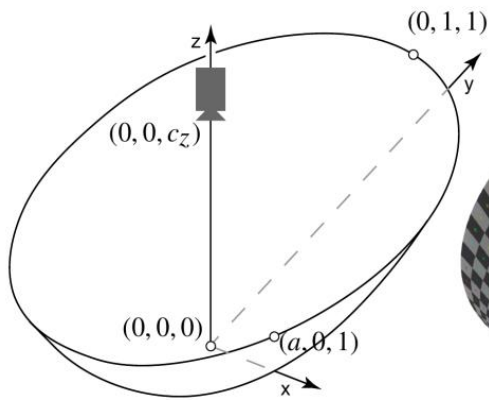
"Loss" for generator

# Architecture Overview

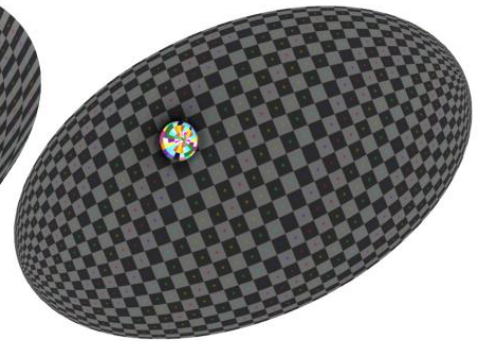**(slide credit: Nils Thuerey)**



Advection encoded in loss for G

11



Input | tempoGAN | Target

# Learning Rolling Motion



# Learning Rolling Motion

Nicholas Watters, Andrea Tacchetti, Theophane Weber, Razvan Pascanu, Peter Battaglia, Daniel Zoran (DeepMind): **Visual Interaction Networks**, NIPS 2017

Image resolution: 128x128

ours

# Course Information (slides/code/comments)



http://geometry.cs.ucl.ac.uk/creativeai/



SIGGRAPH Asia Course **CreativeAI: Deep Learning for Graphics**