

Deep neural architectures for structured output problems

Soufiane Belharbi

soufiane.belharbi@litislab.fr

Clément Chatelain

clement.chatelain@insa-rouen.fr



Joint work with: J.Lerouge, R.Herault, S.Adam, R.Modzelewski, F.Jardin, B.Labbe

May 19, 2015

- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)
- 3 Application of IODA to medical image labeling
- 4 Application of IODA to Facial Landmark Detection
- 5 Conclusion
- 6 Future Work on IODA

Traditional Machine Learning Problems

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

- Inputs $\mathcal{X} \in \mathbb{R}^d$: any type of input
- Outputs $\mathcal{Y} \in \mathbb{R}$ for the task: classification, regression, ...

Machine Learning for Structured Output Problems

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

- Inputs $\mathcal{X} \in \mathbb{R}^d$: any type of input
- Outputs $\mathcal{Y} \in \mathbb{R}^{d'}$, $d' > 1$ a structured object (dependencies)

See C. Lampert slides [3].

Traditional Machine Learning Problems

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

- Inputs $\mathcal{X} \in \mathbb{R}^d$: any type of input
- Outputs $\mathcal{Y} \in \mathbb{R}$ for the task: classification, regression, ...

Machine Learning for Structured Output Problems

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

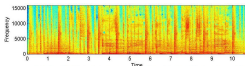
- Inputs $\mathcal{X} \in \mathbb{R}^d$: any type of input
- Outputs $\mathcal{Y} \in \mathbb{R}^{d'}$, $d' > 1$ a structured object (dependencies)

See C. Lampert slides [3].

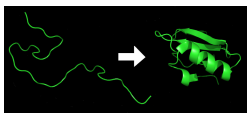
Data = *representation* (values) + *structure* (dependencies)

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi
 auctor lorem non justo. Nam lacus libero, pretium at, laboris vitae, ultricies et,
 tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna,
 vitae ornare odio metus a mi. Morbi ac ceri et nisl hendrerit mollis. Suspendisse
 ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et
 magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna.
 Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Text: part-of-speech
 tagging, translation



speech \leftrightarrow *text*



Protein folding



Image

Structured data

Approaches that Deal with Structured Output Data

- ▶ Kernel based methods: Kernel Density Estimation (KDE)
- ▶ Discriminative methods: Structure output SVM
- ▶ Graphical methods: HMM, CRF, MRF, ...

Drawbacks

- Perform one single data transformation
- Difficult to deal with *high dimensional* data

Ideal approach

- ▶ Structured output problems
- ▶ High dimension data
- ▶ Multiple data transformation (complex mapping functions)

Deep neural networks?

Approaches that Deal with Structured Output Data

- ▶ Kernel based methods: Kernel Density Estimation (KDE)
- ▶ Discriminative methods: Structure output SVM
- ▶ Graphical methods: HMM, CRF, MRF, ...

Drawbacks

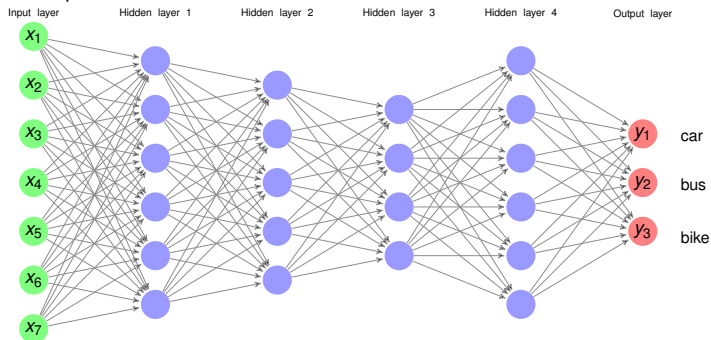
- Perform one single data transformation
- Difficult to deal with *high dimensional* data

Ideal approach

- ▶ Structured output problems
- ▶ High dimension data
- ▶ Multiple data transformation (complex mapping functions)

Deep neural networks?

Traditional Deep neural Network



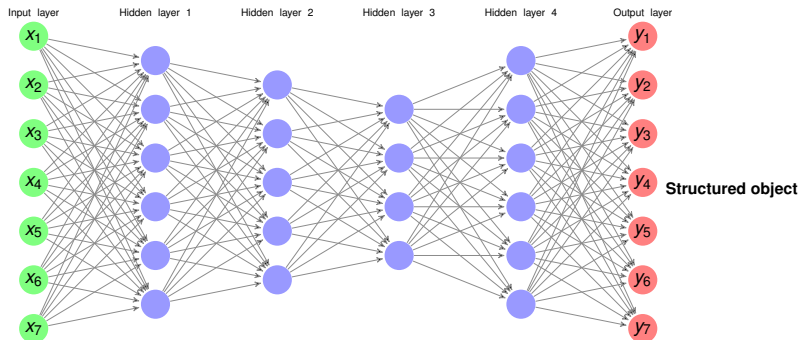
- ▶ High dimension data **OK**
- ▶ Multiple data transformation (complex mapping functions) **OK**
- ▶ Structured output problems **NO**

Plan

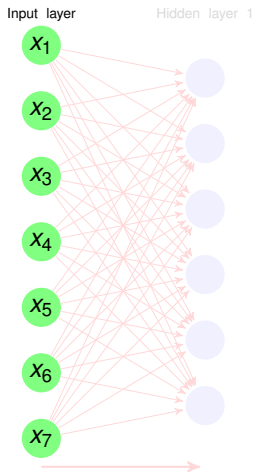
- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)**
- 3 Application of IODA to medical image labeling
- 4 Application of IODA to Facial Landmark Detection
- 5 Conclusion
- 6 Future Work on IODA

IODA:

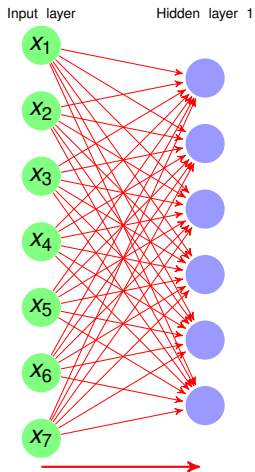
- ▶ Incorporate the output structure by learning
- ▶ Discover hidden dependencies in the outputs



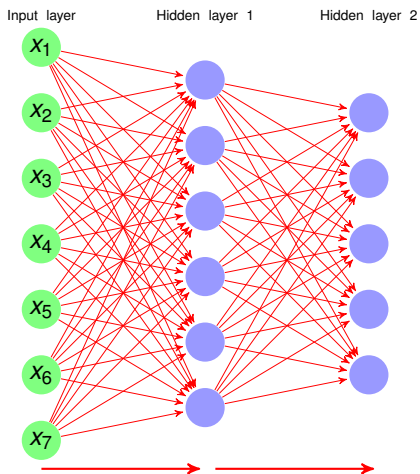
Training IODA



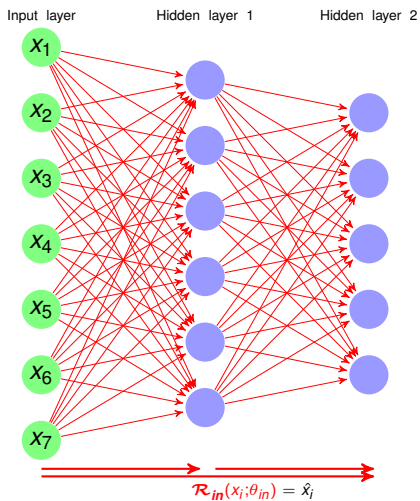
Training IODA



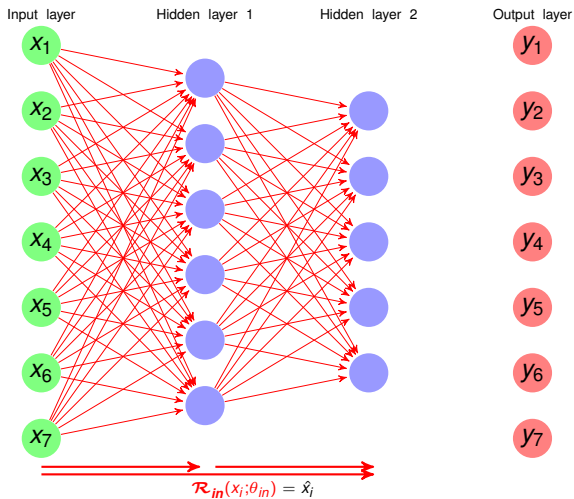
Training IODA



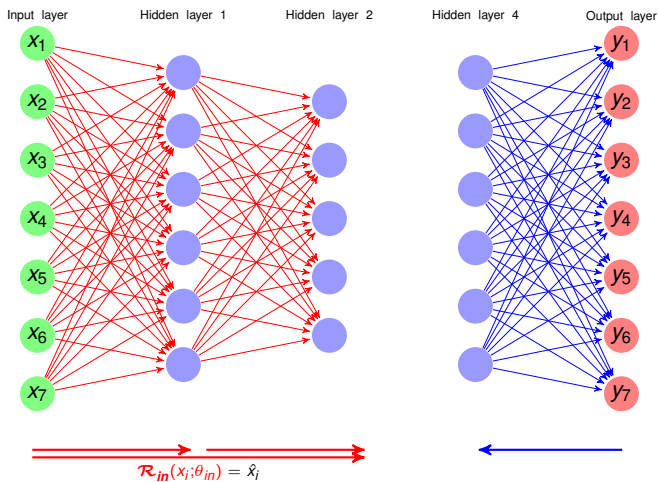
Training IODA



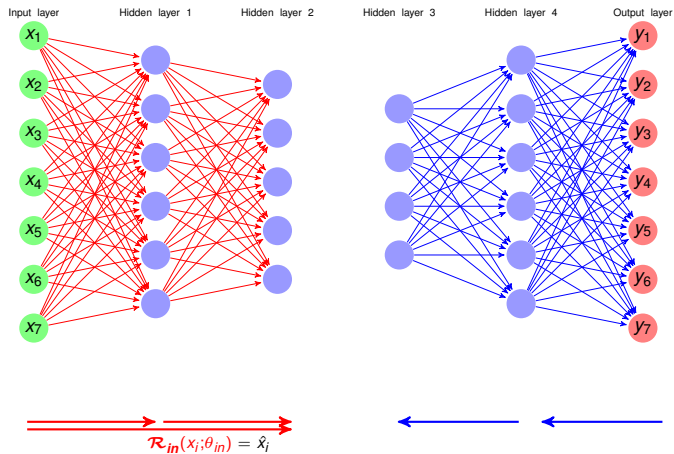
Training IODA



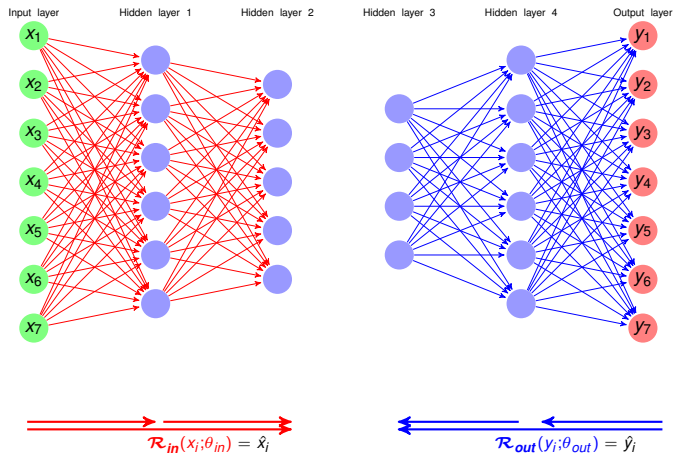
Training IODA



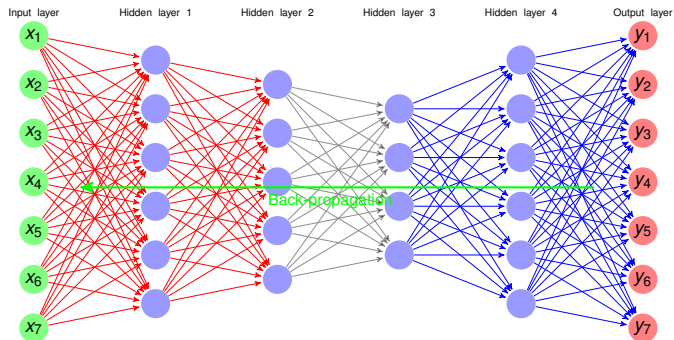
Training IODA



Training IODA



Training IODA



$$\mathcal{R}_{in}(x_j; \theta_{in}) = \hat{x}_j$$

$$\mathcal{R}_{out}(y_j; \theta_{out}) = \hat{y}_j$$

$$\mathcal{M}(x_j; \theta_{in}, \theta_{out}) = \hat{y}_j$$

IODA framework: $\min_{\theta} \mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$

$$\mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y})) = \frac{1}{n} \sum_{i=1}^n \left[\underbrace{\mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i)}_{\text{Learn (input} \rightarrow \text{output) dependencies}} \right. \\ \left. + \underbrace{\ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i)}_{\text{Learn input dependencies}} \right. \\ \left. + \underbrace{\ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i)}_{\text{Learn output dependencies}} \right]$$

$\mathcal{C}(\cdot), \ell_{in}(\cdot), \ell_{out}(\cdot)$: defined costs.

$\min_{\theta} \mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$ is hard to solve \Rightarrow split $\mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$

IODA framework: $\min_{\theta} \mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$

$$\mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y})) = \frac{1}{n} \sum_{i=1}^n \left[\underbrace{\mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i)}_{\text{Learn (input} \rightarrow \text{output) dependencies}} \right. \\ \left. + \underbrace{\ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i)}_{\text{Learn input dependencies}} \right. \\ \left. + \underbrace{\ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i)}_{\text{Learn output dependencies}} \right]$$

$\mathcal{C}(\cdot), \ell_{in}(\cdot), \ell_{out}(\cdot)$: defined costs.

$\min_{\theta} \mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$ is hard to solve \Rightarrow split $\mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y}))$

Relaxed-simplified version of IODA

1 Unsupervised training:

→ **Input** dependencies : $\min_{\theta_{in}} \frac{1}{n} \sum_{i=1}^n \ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i)$

→ **output** dependencies: $\min_{\theta_{out}} \frac{1}{n} \sum_{i=1}^n \ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i)$

2 Standard supervised learning:

$$\min_{\theta, \theta_{in}, \theta_{out}} \frac{1}{n} \sum_{i=1}^n \mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i)$$

Open source implementation

*Implemented using our library: **Crino** [1] [Python-Theano based].*

Relaxed-simplified version of IODA

1 Unsupervised training:

→ *Input* dependencies : $\min_{\theta_{in}} \frac{1}{n} \sum_{i=1}^n \ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i)$

→ *output* dependencies: $\min_{\theta_{out}} \frac{1}{n} \sum_{i=1}^n \ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i)$

2 Standard supervised learning:

$$\min_{\theta, \theta_{in}, \theta_{out}} \frac{1}{n} \sum_{i=1}^n \mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i)$$

Open source implementation

Implemented using our library: **Crino** [1] [Python-Theano based].

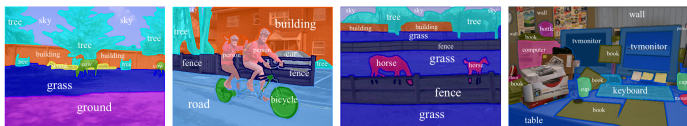
Plan

- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)
- 3 Application of IODA to medical image labeling**
- 4 Application of IODA to Facial Landmark Detection
- 5 Conclusion
- 6 Future Work on IODA

Image labeling problems

Definition

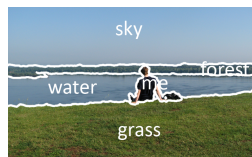
Assigning a label to each pixel of an image
(AKA "semantic segmentation")



Various applications in:

- Document image analysis (text, image, tables, etc.)
- Computer vision (road safety, natural scene understanding)
- Medical imaging (organ, tumour segmentation)

Image labeling problems



Output dependencies

- Local dependencies (neighbouring labels are correlated)
- Structural dependencies (sky is generally above grass)

→ Image labeling can be considered as a structured output problems

Application of IODA on a medical Image labeling problem

Collaboration with the Henri Becquerel Center (Quantif team)

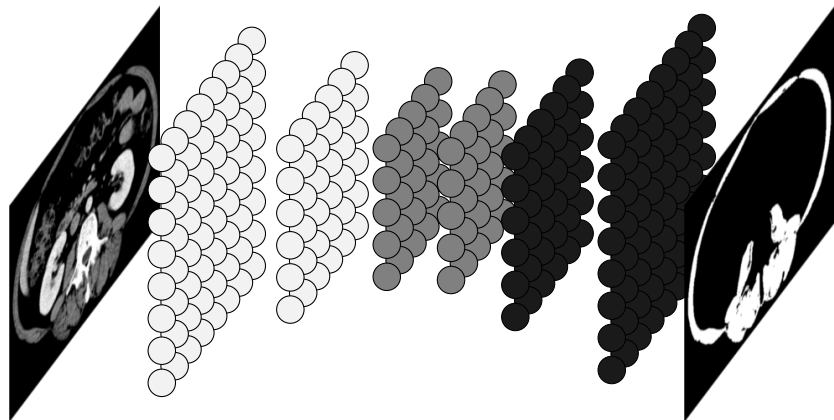
- Sarcopenia is a critical indication for lymphoma treatment
- Can be measured on scanner images by labeling skeletal muscle at L3 (third vertebra)
- 4 min/patient for a senior radiologist



Dataset

- 128 labeled L3 scanner images 512*512 pix
- Reference method from Chung (based on registration)

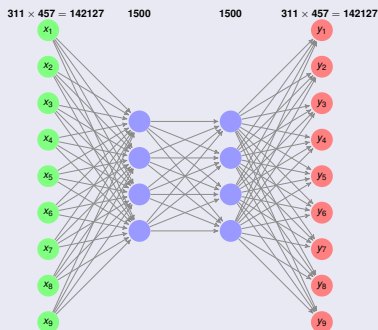
Input/Output Deep Architecture (IODA) for Image Labeling



IODA architecture for skeletal muscle segmentation

Implementation

Architecture (optimized on validation set)



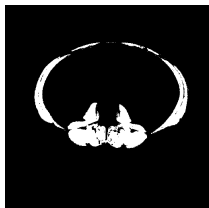
A few figures:

- 428 M parameters (!!)
- Less than an hour for training (GPU, 4Go)
- 201.2 ms for decision

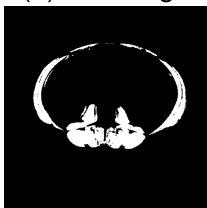
Qualitative results 1/2



(a) CT image



(b) Ground truth



(c) Chung



(d) IODA

Non-sarcopenic patient

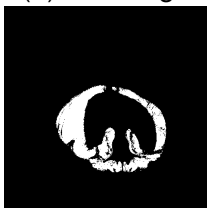
Qualitative results 2/2



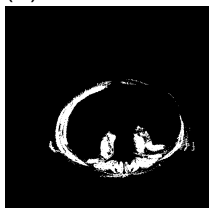
(a) CT image



(b) Ground truth



(c) Chung



(d) IODA

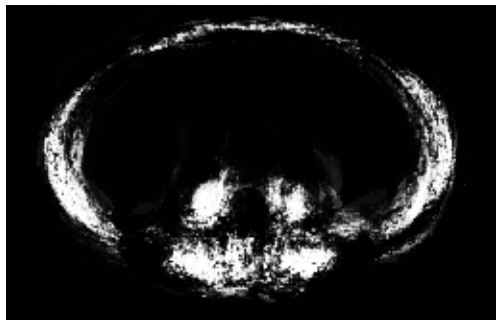
Sarcopenic patient

Quantitative results

Method	Diff. (%)	Jaccard (%)
Chung (reference method)	-10.6	60.3
No pre-train DA	0.12	85.88
Input pre-train DA	0.15	85.91
Input/Output pre-train DA (IODA)	3.37	88.47

The "blank test image"

Feed the network with a blank image



Published in *pattern recognition* [4]

Plan

- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)
- 3 Application of IODA to medical image labeling
- 4 Application of IODA to Facial Landmark Detection**
- 5 Conclusion
- 6 Future Work on IODA

► **Facial landmarks:**

set of **facial key points** with **coordinates** (x,y)

Task → predict the **shape**(set of points) given a facial image



⇒ Geometric dependencies ⇒ structured output problem

⇒ Apply IODA (regression task)

► **Facial landmarks:**

set of **facial key points** with **coordinates** (x,y)

Task → predict the **shape**(set of points) given a facial image



⇒ **Geometric dependencies** ⇒ **structured output problem**

⇒ Apply **IODA** (regression task)

► **Facial landmarks:**

set of **facial key points** with **coordinates** (x,y)

Task → predict the **shape** (set of points) given a facial image



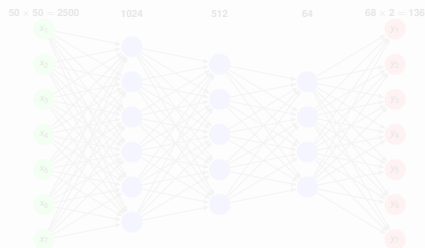
⇒ **Geometric dependencies** ⇒ **structured output problem**

⇒ **Apply IODA (regression task)**

Datasets & Performance Measures

- ▶ Datasets: LFPW(~1000 samples), HELEN(~2300 samples)
- ▶ Performance Measure:
 - ▶ Normalized Root Mean Square Error (**NRMSE**)
 - ▶ Cumulative Distribution Function: **CDF_{NRMSE}**
 - ▶ Area Under the CDF Curve (**AUC**) ****new****

Architecture (optimized on validation set)

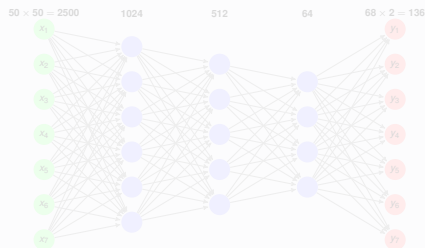


⇒ Total training on GPU takes less than 30mins.

Datasets & Performance Measures

- ▶ Datasets: LFPW(~1000 samples), HELEN(~2300 samples)
- ▶ Performance Measure:
 - ▶ Normalized Root Mean Square Error (**NRMSE**)
 - ▶ Cumulative Distribution Function: **CDF**_{NRMSE}
 - ▶ Area Under the CDF Curve (**AUC**) ****new****

Architecture (optimized on validation set)

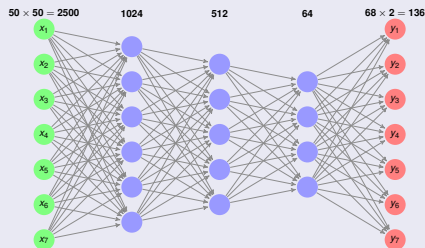


⇒ Total training on GPU takes less than 30mins.

Datasets & Performance Measures

- ▶ Datasets: LFPW(~1000 samples), HELEN(~2300 samples)
- ▶ Performance Measure:
 - ▶ Normalized Root Mean Square Error (**NRMSE**)
 - ▶ Cumulative Distribution Function: **CDF_{NRMSE}**
 - ▶ Area Under the CDF Curve (**AUC**) ****new****

Architecture (optimized on validation set)



⇒ Total **training** on **GPU** takes less than **30mins**.

No pre-train DA



Input pre-train DA

Input/Output pre-train
DA (IODA)

Visual results LFPW

No pre-train DA



Input pre-train DA

Input/Output pre-train
DA (IODA)

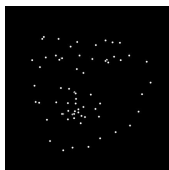
Visual results HELEN

	LFPW		HELEN	
	AUC	CDF _{0.1}	AUC	CDF _{0.1}
Mean shape	66.15%	18.30%	63.30%	16.97%
No pre-train DA 0-0-0	77.60%	50.89%	80.91%	69.69%
Input pre-train DA 1-0-0	79.25%	62.94%	82.13%	76.36%
2-0-0	79.10%	58.48%	82.39%	75.75%
3-0-0	79.51%	65.62%	82.25%	77.27%
Input/Output pre-train DA 1-0-1	80.66%	68.30%	83.95%	83.03%
1-1-1	81.50%	72.32%	83.51%	80.90%
1-0-2	81.00%	71.42%	83.91%	82.42%
1-1-2	81.06%	70.98%	83.81%	83.03%
1-0-3	81.91%	74.55%	83.72%	80.30%
2-0-1	81.32%	72.76%	83.61%	80.00%
2-1-1	81.47%	70.08%	84.11%	83.33%
2-0-2	81.35%	71.87%	83.88%	82.12%
3-0-1	81.62%	72.76%	83.38%	78.48%

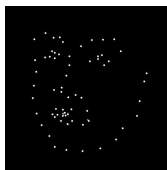
Performance of mean shape, NDA, IDA and IODA on LFPW and HELEN.

Blank Image Test

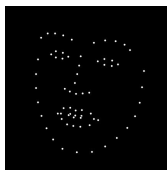
Feed a blank image to a trained network
⇒ what is the output?



No pre-train
DA 0-0-0



Input pre-train DA
3-0-0



Input/Output pre-train
DA 1-0-3

The outputs on LFPW

Paper submitted to ECML 2015 (arXiv [2]).

Plan

- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)
- 3 Application of IODA to medical image labeling
- 4 Application of IODA to Facial Landmark Detection
- 5 Conclusion**
- 6 Future Work on IODA

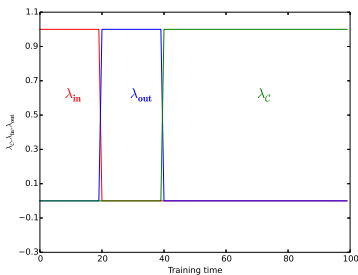
- ▶ Fully **neural** based approach
- ▶ Able to learn the **output dependencies** in **high dimension**
- ▶ Efficient on two real world problems

Plan

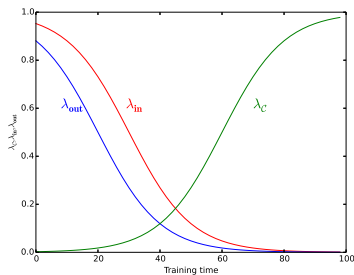
- 1 Structured Output Problems
- 2 Input/Output Deep Architecture (IODA)
- 3 Application of IODA to medical image labeling
- 4 Application of IODA to Facial Landmark Detection
- 5 Conclusion
- 6 Future Work on IODA**

1 Embedded Pre-training (draft on arXiv):

$$\mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y})) = \frac{1}{n} \sum_{i=1}^n \left[\lambda_c \mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i) \right. \\ \left. + \lambda_{in} \ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i) \right. \\ \left. + \lambda_{out} \ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i) \right]$$



Relaxed IODA



Embedded IODA

2 Use of **unlabeled** data:

$$\begin{aligned}
 \mathcal{L}(\theta, \mathcal{D}(\mathbf{x}, \mathbf{y})) &= \frac{1}{n} \sum_{i=1}^n \lambda_{\mathcal{C}} \mathcal{C}(\mathcal{M}(x_i; \theta, \theta_{in}, \theta_{out}), y_i) \\
 &\quad + \frac{1}{n + n_{in}} \sum_{i=1}^{n+n_{in}} \lambda_{in} \ell_{in}(\mathcal{R}_{in}(x_i; \theta_{in}), x_i) \\
 &\quad + \frac{1}{n + n_{out}} \sum_{i=1}^{n+n_{out}} \lambda_{out} \ell_{out}(\mathcal{R}_{out}(y_i; \theta_{out}), y_i)
 \end{aligned}$$

n_{in}, n_{out} potentially **huge unlabeled** input, output data.

3 Convolutional IODA:

Convolutional layers are efficient in feature extraction

⇒ Use convolutional layers instead of auto-encoders in the input-layers

- [1] Crino, a neural-network library based on Theano. <https://github.com/jlerouge/crino>, 2014.
- [2] S. Belharbi, C. Chatelain, R. Hérault, and S. Adam.
Input/Output Deep Architecture for Structured Output Problems.
ECML, 2015.
- [3] CH. Lampert.
Slides: Learning with Structured Inputs and Outputs, http://www.di.ens.fr/willow/events/cvml2010/materials/INRIA_summer_school_2010_Christoph.pdf, 2010.
- [4] J. Lerouge, R. Herault, C. Chatelain, F. Jardin, and R. Modzelewski.
loda: An input output deep architecture for image labeling.
Pattern Recognition, 48(9):2847–2858, 2015.

Thank you for your attention.

Sets	Train samples	Test samples
LFPW	811	224
HELEN	2000	330

Number of samples in datasets.

Normalized Root Mean Square Error (NRMSE)

$$NRMSE(s_p, s_g) = \frac{1}{n * D} \sum_{i=1}^n \|s_{pi} - s_{gi}\|_2,$$

s_p, s_g predicted, ground truth shape. D inter-ocular distance of s_g

Cumulative Distribution Function: CDF_{NRMSE}

$$CDF_x = \frac{CARD(NRMSE \leq x)}{N}$$

$CARD(.)$ cardinal of a set. N number of images.

e.g. $CDF_{0.1} = 0.4$ means that 40% of images have an NRMSE error less or equal than 0.1

Area Under the CDF Curve (AUC) ****new****: more numerical precision

- Plot a CDF_{NRMSE} curve by varying NRMSE in $[0, 0.5]$.
- Calculate the area under this curve.

Normalized Root Mean Square Error (NRMSE)

$$NRMSE(s_p, s_g) = \frac{1}{n * D} \sum_{i=1}^n \|s_{pi} - s_{gi}\|_2,$$

s_p, s_g predicted, ground truth shape. D inter-ocular distance of s_g

Cumulative Distribution Function: CDF_{NRMSE}

$$CDF_x = \frac{CARD(NRMSE \leq x)}{N}$$

$CARD(.)$ cardinal of a set. N number of images.

e.g. $CDF_{0.1} = 0.4$ means that 40% of images have an NRMSE error less or equal than 0.1

Area Under the CDF Curve (AUC) ****new****: more numerical precision

- Plot a CDF_{NRMSE} curve by varying NRMSE in $[0, 0.5]$.
- Calculate the area under this curve.

Normalized Root Mean Square Error (NRMSE)

$$NRMSE(s_p, s_g) = \frac{1}{n * D} \sum_{i=1}^n \|s_{pi} - s_{gi}\|_2,$$

s_p, s_g predicted, ground truth shape. D inter-ocular distance of s_g

Cumulative Distribution Function: CDF_{NRMSE}

$$CDF_x = \frac{CARD(NRMSE \leq x)}{N}$$

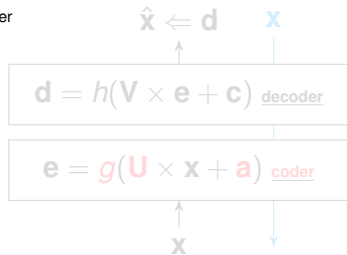
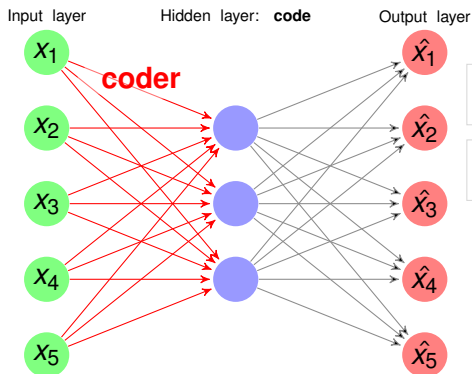
$CARD(.)$ cardinal of a set. N number of images.

e.g. $CDF_{0.1} = 0.4$ means that 40% of images have an NRMSE error less or equal than 0.1

Area Under the CDF Curve (AUC) ****new****: more numerical precision

- Plot a CDF_{NRMSE} curve by varying NRMSE in $[0, 0.5]$.
- Calculate the area under this curve.

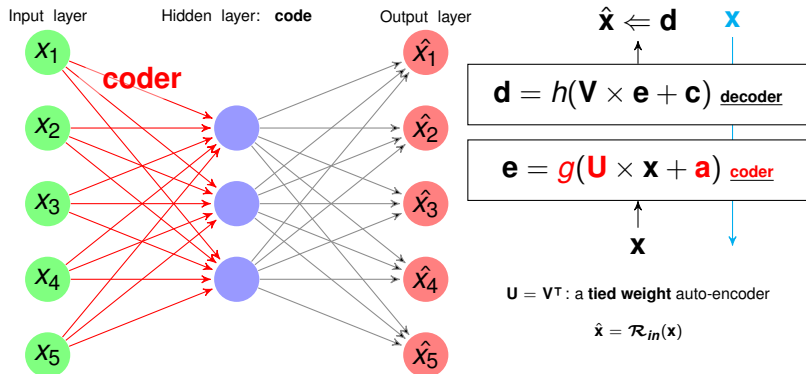
Input layer pre-training using auto-encoders (1)



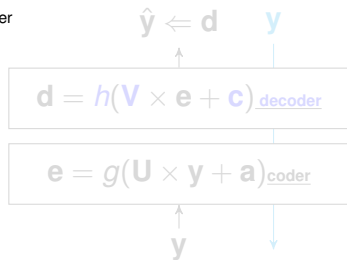
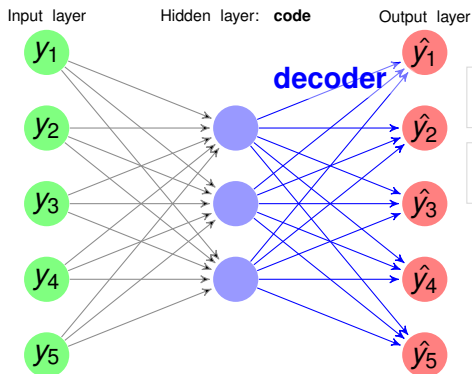
$\mathbf{U} = \mathbf{V}^T$: a tied weight auto-encoder

$$\hat{x} = \mathcal{R}_{in}(x)$$

Input layer pre-training using auto-encoders (1)



Output layer pre-training using auto-encoders (2)



$\mathbf{U} = \mathbf{V}^T$: a tied weight auto-encoder

$$\hat{y} = \mathcal{R}_{out}(y)$$

Output layer pre-training using auto-encoders (2)

