

UDRC Summer School: Machine Learning: Deep Neural Networks III

Dr Timothy Hospedales

Machine Intelligence Group <u>http://mig.inf.ed.ac.uk</u> University of Edinburgh

UDRC Co-Investigator

Samsung AI Centre, Cambridge

Dr. Timothy Hospedales

- Background
 - BA Computer Science, Cambridge 2002
 - Lecturer/Senior Lecturer, QMUL, 2012-16
 - Reader (Assoc. Prof), Edinburgh, 2016-
 - Head of Machine Intelligence Group
 - Funded by EPSRC, DSTL, EU Horizon 2020
 - Alan Turing Institute Fellow, 2018-
 - Principal Scientist, Samsung Al Research Center, 2019-
- Research Area:
 - Deep Learning. Meta-Learning. Vision. Vision & Language.
- Track record:
 - Over 50 papers in Tier 1 venues of AI, ML, Vision.
 - => CVPR, ICCV, ECCV, ICLR, AAAI, IJCAI, ICML. T-PAMI, IJCV.
 - Four best paper prizes.
 - Three patents



Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Part I: Sparse Data Deep Learning

Deep Learning Success





DeepMind Human-level control through deep reinforcement learning ?? Letter Deep Q-Learning 5

Sketch-a-Net that beats humans



Superhuman Pictionary





Mechanism of Deep Learning Era Success?

- Gather and annotate bigger datasets.
- Train bigger models.



Success Story of Deep Learning Era

- Gather and annotate bigger datasets.
- Train bigger models.

No saturation so far! Performance grows with log of train data.



[Sun et al, Revisiting Unreasonable Effectiveness of Data in Deep Learning Era, ICCV'17]

Do we need another paradigm?

Do we need another paradigm?

- Humans have one shot learning
 - Learn 5 objects per day for first 18 years.
- Annotating the long tail of object categories?
 - Emerging categories
- Annotating data in the long tail of domains
 - Underwater, Radar, Sonar, LIDAR, Medical, Satellite, et
- Defense applications training data, or rare

Why is Few-Shot Learning Hard? Overfitting

• Underfitting vs Overfitting. Linear regression example.



EG: 20th order polynomial.

EG: 2nd order polynomial.

EG: 0th order polynomial.

Why is Few-Shot Learning Hard? Overfitting

• Underfitting vs Overfitting. Linear regression example.



Question: How to diagnose over- vs under-fitting?

Overfitting?

- Low train error
- High test error

Underfitting?

- High train error
- High test error



Why is Few-Shot Learning Hard? Overfitting

• Underfitting vs Overfitting. Linear regression example.



EG: 20th order polynomial.

Classic Solution?

- Try several model complexities
- Evaluate validation set performance of each
- Pick the model with best validation performance

Issue for deep learning?

- Too many complexity parameters in DL (depth, width, non-linearity, etc)
- Few-Shot: Would pick a simple model that doesn't provide deep learning level performance.



EG: 0th order polynomial.

Common Techniques for Sparse Data Applications

Reducing Data Dependence

Common techniques for overfitting reduction: Regularisation

- Weight decay / L2 regularization
- Early Stopping
- Dropout
- Transfer Learning

Reducing Data Dependence

Weight decay. L2 regularisation.

• Other things being equal, prefer weights near zero.



Reducing Data Dependence

Transfer Learning

- Transfer Learning: "The application of skills, knowledge, and/or attitudes that were learned in one situation to another learning situation" (Perkins, 1992).
- Note:
 - Transfer Learning ≠ "Fine-Tuning"
 - Fine Tuning⊂ Transfer Learning



Transfer Learning: Fine-Tuning

• Initialise target parameters using source, and continue training



Transfer Learning: Fine-Tuning

- Application:
 - Across Task:

 $p(Y_s|X_s) \neq p(Y_t|X_t)$

Cross Task

- Change of label-space. $Y_s \neq Y_t$
- Across Domain
 - Change of data statistics. $p(X_s) \neq p(X_t)$





Transfer Learning: Fine-Tuning

- Why does fine-tuning work?
 - Neural network optimisation is non-convex.
 - With small learning rate, target task parameters do not change much.
 - Transfer initialization effectively regularizes target weights towards source weights, rather than towards zero.
 - Assume source task is relevant to target task.
 - => Better chance of good minima.





Transfer Learning: Fine-Tuning

- Very similar to explicit source \rightarrow target regularisation.
 - Used in many classic learning methods(*).

$$\min_{\boldsymbol{w}_t} \mathcal{L}(D_t; \boldsymbol{w}_t) + \lambda \|\boldsymbol{w}_t - \boldsymbol{w}_s\|_2^2$$

Compare:

$$\min_{\boldsymbol{w}} \bar{\mathcal{L}} = \mathcal{L}(D; \boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_{2}^{2}$$

$$\min_{\boldsymbol{w}} \bar{\mathcal{L}} = \mathcal{L}(D; \boldsymbol{w}) + \lambda \|\boldsymbol{w} - \boldsymbol{0}\|_{2}^{2}$$



[Improving SVM Accuracy by Training on Auxiliary Data Sources, ICML'04] [Cross-domain video concept detection using adaptive svms, ACM MM'07]

 \mathbf{X}

Transfer Learning: Fine-Tuning

- Assumption:
 - Source task relevant to target.
- Practical Considerations. Questions:
 - How to control relevance degree?
 - Which layers to transfer?
 - What learning rate to use?
- Typical:
 - Relearn top while freeze bottom.
 - Then LR tune all w/ LR proportional depth



Question: Always good practice? $\min_{\boldsymbol{w}} \mathcal{L}(D_t; \boldsymbol{w}_t) + \boldsymbol{\lambda} \| \boldsymbol{w}_t - \boldsymbol{w}_s \|_2^2$

[How transferable are features in deep neural networks?, NIPS'14]

Transfer Learning: Fine-Tuning

- Assumption:
 - Source task relevant to target.
- Practical Considerations. Questions:
 - How to control relevance degree?
 - Which layers to transfer?
 - What learning rate to use?
- Typical:
 - Relearn top while freeze bottom.
 - Then LR tune all w/ LR proportional depth

Only for cross-task.

Assumption:

• Relevance proportional to depth.

 $\min_{\boldsymbol{w}} \mathcal{L}(D_t; \boldsymbol{w}_t) + \boldsymbol{\lambda} \| \boldsymbol{w}_t - \boldsymbol{w}_s \|_2^2$



[How transferable are features in deep neural networks?, NIPS'14]

Transfer Learning: Fine-Tuning

- Challenges:
 - What if the inputs are heterogeneous
 - (EG: Task 1: RGB, Task 2: Depth)?
 - How to know if a given source is relevant?
 - How to prevent negative transfer if source is irrelevant?
 - How to select the relevant source among many?

Transfer Learning: Fine-Tuning

- Issues:
 - What if the inputs are heterogeneous?
 - EG: Perceptual arithmetic network.





[EG: HOUDINI: Lifelong Learning as Program Synthesis, NIPS'18]

Transfer Learning: Fine-Tuning

- Issues:
 - How to select the relevant source among many?
- For linear models....
 - Optimising target model *w*, assuming a set of potentially relevant sources {*w_k*}:

 $\min_{\boldsymbol{w}}[\|\boldsymbol{y} - \boldsymbol{w}\boldsymbol{X}\| + \min_{k}\|\boldsymbol{w} - \boldsymbol{w}_{k}\|]$

• not so straightforward with deep models

E.g., Evgeniou et al, JMLR, 2005 E.g., Kang et al, ICML, 2011

Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Metric Learning

Nearest Neighbor Classifier

- Classic Method in ML
 - Given training data $\{x_i, y_i\}_{i=1}^N$
- 1-NN:
 - Classify a new point \boldsymbol{x}^* . Euclidean nearest neighbor:

•
$$y^* = y_{i^*}, i^* = \operatorname{argmin} \| x^* - x_i \|$$

- K-NN: For 1-hot **y**.
 - Classify a new point \boldsymbol{x}^* as:
 - $p(\mathbf{y}^*|\mathbf{x}^*) = \frac{1}{K} \sum_{\mathbf{x}_j \in NN(\mathbf{x}^*)} \mathbf{y}_j$



Can we measure similarity by a better distance than Euclidean?

Find a similarity measure where different categories are far, and same categories are near.

Classic Metric Learning

• Define a Mahalanobis Metric for NN classification:

$$d_M(x_i, x_j) = (x_i - x_j)^T M(x_i - x_j) = (x_i - x_j)^T L^T L(x_i - x_j) = (Lx_i - Lx_j)^T (Lx_i - Lx_j)$$

• Learn the Metric that separates the training data:

$$\min_{M} L(M) = \sum_{i,j \in Pos} d_M(x_i, x_j) - \sum_{i,j \in Neg} d_M(x_i, x_j) \qquad \bigoplus_{i,j \in Neg} d_M(x_i, x_j)$$

• Reuse the metric to solve 1-shot learning by NN matching:



[EG: Hirzer, ECCV'12; Mignon, CVPR'12]

Deep Metric Learning



• Linear => Non-linear feature transform

$$d_{M}(x_{i}, x_{j}) = (x_{i} - x_{j})^{T} M(x_{i} - x_{j}) = (Lx_{i} - Lx_{j})^{T} (Lx_{i} - Lx_{j})$$
$$d_{W}(x_{i}, x_{j}) = (f_{W}(x_{i}) - f_{W}(x_{j}))^{T} (f_{W}(x_{i}) - f_{W}(x_{j}))$$

• Common Losses:

$$\min_{W} L(W) = \sum_{i,j\in Pos} d_W(x_i, x_j) - \sum_{i,j\in Neg} d_W(x_i, x_j)$$

Minimize average positive distance, Maximize average negative distance

Every negative distance more than every positive distance

positive pair distances = 0 negative pair distances = 1

$$\begin{split} \min_{W} L(W) &= \sum_{i \in D} \sum_{j \in Pos(i), k \in Neg(i),} \left| d_{W}^{Pos}(x_{i}, x_{j}) + \Delta - d_{W}^{Neg}(x_{i}, x_{k}) \right|^{+} \\ \min_{W} L(W) &= \sum_{i, j \in Pos} \left\| d_{W}(x_{i}, x_{j}) - 0 \right\| + \sum_{i, j \in Neg} \left\| d_{W}(x_{i}, x_{j}) - 1 \right\| \end{split}$$

Deep Metric Learning: Example: Triplet Ranking for Similarity Search



[CVPR'14, Learning Fine-grained Image Similarity with Deep Ranking]

Deep Metric Learning: Challenges?

$$d_W(x_i, x_j) = (f_W(x_i) - f_W(x_j))^T (f_W(x_i) - f_W(x_j))$$

- Challenges? Training: Scalibility.
 - Num triplets: $O(N^3)$. L(W) =
 - Sampling useful ones

 $\left|d_W^+(x_i,x_j) + \Delta - d_W^-(x_i,x_k)\right|^+$

• => Hard negative mining, learning sampling policy.

 $i \in D, j \in Pos(i), k \in Neg(i)$

- Annotating them.
- Challenges? Testing: Scalability.
 - Binary embedding for fast lookup/hashing $f_W(x) = \text{sign}(Wh + b)$
 - => Gradient descent on discrete gradient operation

[Wang, PAMI'18, A Survey on Learning to Hash]



Metric Learning for Few-Shot

Few-Shot Learning: Setup

- Meta-train on Auxiliary Set.
 - EG: Vehicles. Disjoint labels: $\mathcal{Y}^{aux} \neq \mathcal{Y}^{sup}$
- Train on Support Set $D^{sup} = \{x_i, y_i\}_{i=1}^1$
 - EG: 1-example per animal $y^{sup} = y^{tst}$
- Test on Query Set. $D^{tst} = \{x_i, y_i\}_{i=1}^N$
 - EG: New animal images to classify.



 $D^{aux} = \{X^{aux}, Y^{aux}\} = \{x_i, y_i\}_{i=1}^{N^{aux}}$



Few-Shot Learning: Prototypical Nets (NIPS'17)

• Calculate a "Prototype" per class:

$$\boldsymbol{c}_{k} = \frac{1}{|S_{k}|} \sum_{(\boldsymbol{x}_{i}, \boldsymbol{y}_{i}) \in S_{k}} f_{\theta}(\boldsymbol{x}_{i})$$

• Classify with

$$p(y = k | x) \propto \exp(-\|f_{\theta}(\boldsymbol{x}_i) - \boldsymbol{c}_k\|^2)$$



[Snell et al, NIPS'17, Prototypical Networks for Few Shot Learning]

Few-Shot Deep Learning: Relation Net

- Relation Net: Similar to Prototypical, but:
 - Architecture: Concatenate branch feature maps and convolve more.
 - Loss: Regress to 0/1.



[Sung, CVPR'18, "Learning to Compare: Relation Networks for Few Shot Learning"]
Few-Shot Deep Learning: Architectures Deep Metric Learning Neg Exp Euclidean fθ CNN Distance to NN [1,0,0,0,0] **Prototypical Nets** fθ CNN $\blacktriangleright \{\boldsymbol{c}_k\}_{k=1}^K$ fθ CNN Neg Exp Euclidean Distance to $\exp(-\|f_{\theta}(\mathbf{x}_{i}) - \mathbf{c}_{k}\|^{2}) - \mathbf{c}_{k}\|^{2}$ fθ CNN [1,0,0,0,0] Prototype **Relation Net** fθ CNN Regress $\operatorname{CNN} g_{\phi}$ [0, 0, 0, 1, 0]Concatenate CNN f_{θ}

Few-Shot Deep Learning: Relation Net

• Relation Net: Episodic Training



Episodic Training. Each mini-batch is a K-way N-shot "episode" => Trains features to support few-shot recognition

[Sung, CVPR'18, "Learning to Compare: Relation Networks for Few Shot Learning"]

Relation Net: Analysis

- PrototypeNet/RelationNet: Pros & Cons?
 - Pros:
 - Train+test cardinality don't have to match.
 - Training is super fast if K-way/N-shot is small.
 - Cons:
 - Testing is slow if K-way is large.
 - (RelationNet) Can't cache features for fast binary/approx NN lookup
 - Neutral
 - Can do old+new classes together. (But cost + accuracy....)



[Sung, CVPR'18, "Learning to Compare: Relation Networks for Few Shot Learning"]

Relation Net

- Summary:
 - State of the art on miniImageNet and tieredImageNet
 - Simple to implement and fast
 - Also supports zero-shot learning



Model (SENet)	5-way 1-shot Acc					
MAML	55.9%					
Prototypical	51.6%					
RelationNet	57.4%					
RelationNet 2.0	62.8%					

[Sung, CVPR'18, "Learning to Compare: Relation Networks for Few Shot Learning"] [Zhang, arXiv'18, "Relation Columns for Few Shot Learning"]

Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Tensor-Based Models

Multi-task Learning

Deep Multi-Task Learning

- Typically:
 - Share feature extraction layers.
 - Multiple predictions & losses.





[TPAMI 2017, HyperFace: A Deep Multi-task Learning Framework for Face Detection...; CVPR'17, UberNet: Training a `Universal' Convolutional Neural Network]



Deep Multi-Task Learning

- Typically:
 - Share feature extraction layers.
 - Multiple predictions & losses.



- Why/When to use MTL?
 - Care about performance on all tasks (else use TL)
 - Share compute & memory of feature extraction.

[TPAMI 2017, HyperFace: A Deep Multi-task Learning Framework for Face Detection... ; CVPR'17, UberNet: Training a `Universal' Convolutional Neural Network]

Architecture for Classification



Tensor-Based Models

Transfer Learning

Multi-Source Knowledge Transfer Scenarios















Hit



Throw

Multi-Source Knowledge Transfer Scenarios

Source

Target





Tensor Composition vs Decomposition

U

Х

- You may have heard of matrix factorisation.
 - We also have tensor factorisation .

PCA, SVD, etc

Х

- $\min_{U,V} \|X UV\|^2$
 - $Dim(X) = N \times D$ $Dim(U) = N \times R$ $Dim(V) = R \times D$



Tensor-based Transfer: Background

1. Find the low-rank subspace spanned by the source tasks.



Tensor-based Transfer: Background

1. Find the low-rank subspace spanned by the source tasks.







Original Tensor W is size: $T \times D \times H$. Approximate it as product of low-rank factors: $W \approx S \times_1 U^T \times_2 U^D \times_3 U^H$

dim(S): $R \times R \times R$ (Agnostic) dim(U^T): $R \times T$ (Task specific) dim(U^D): $R \times D$ (Input specific) dim(U^H): $R \times H$ (Hidden unit specific)

Compress params from TDH to $R^3 + RT + RD + RH$

Note: Tensor Composition vs Decomposition

- You may have heard of matrix/tensor decomposition.
 - Note that we are doing matrix/tensor composition.



Tensor-based Transfer: Multi-Task Source Training

1. Find the low-rank subspace spanned by the source tasks.





$$\min_{..W.} \sum_{i} \sum_{k} \mathcal{L}(y_i^k, f_{\phi_k}(g_W(x_i)))$$
$$\min_{..S, U^T, U^D, U^H, \dots} \sum_{i} \sum_{k} \mathcal{L}(y_i^k, f_{\phi_k}(g_W(x_i)))$$



Original Tensor W is size: $T \times D \times H$. Approximate it as product of low-rank factors: $W \approx S \times_1 U^T \times_2 U^D \times_3 U^H$

Where: R < T, D, Hdim(S): $R \times R \times R$ (Agnostic) dim(U^T): $R \times T$ (Task specific) dim(U^D): $R \times D$ (Input specific) dim(U^H): $R \times H$ (Hidden unit specific)

Compress params from TDH to $R^3 + RT + RD + RH$

Tensor-based Transfer: Idea

- 1. Find the low-rank subspace spanned by the source tasks.
- 2. Train only the task-specific parameters for new task.



Tensor-based Transfer: Context



[Yang & Hospedales, ICLR'15; Yang & Hospedales ICLR'17]

Tensor-based Transfer:

Indul

Input

F

F

Ĩ

Question: What is the difference between these two strategies?



Tensor-Based Transfer Learning: Some Examples

Robot Control

• Consider robot with an arm.

- Model: $\boldsymbol{a}_{t} = f_{\boldsymbol{w}}\left(\boldsymbol{x}_{t}\right)$
 - *x_t*: Angle of each joint, time, (ball position, target position)
 - *a_t*: Force to apply at each joint at time *t*.
 - f_w : Neural network controller
- Train with reinforcement learning
 - $R(x_t) = \begin{cases} +1 \text{ if ball hits target} \\ -0.1 \text{ otherwise} \end{cases}$
 - Train: $\mathbf{w} = \operatorname{argmax}_{\mathbf{w}} \sum_{t} R(\mathbf{x}_{t})$





- Consider robot with an arm.
 - Model: $\boldsymbol{a}_{t} = f_{\boldsymbol{w}}\left(\boldsymbol{x}_{t}\right)$
 - x_t : Angle of each joint, time, (ball position, target position)
 - *a_t*: Force to apply at each joint at time *t*.
 - f_w : Neural network controller
 - Train with reinforcement learning
 - $R(x_t) = \begin{cases} +1 \text{ if ball hits target} \\ -0.1 \text{ otherwise} \end{cases}$
 - Train: $\mathbf{w} = \operatorname{argmax}_{\mathbf{w}} \sum_{t} R(\mathbf{x}_{t})$

- Issue:
 - Requires lots of data: $\left\{ \{a_t, x_t, R(x_t)\}_{t=1}^T \right\}_{n=1}^N$
 - Each is trial slow and expensive.
 - Redo for each new task.



• Given some known tasks, use knowledge transfer to learn new task quicker.



Hit



Throw

• Given some known tasks, use knowledge transfer to learn new task quicker.



Throw

Cast

Ball-In-Cup

Tensor Transfer for Robot Control

• In this case vectorize entire neural network....



[Yang et al, ICLR'17. Zhao et al, Tensor Based Knowledge Transfer Across Skill Categories for Robot Control, IJCAI, 2017]

Tensor Transfer for Robot Control

- Dramatic Improvement in Learning Efficiency for Target Task:
 - Now possible to learn BIC with RL.



[Yang et al, ICLR'17. Zhao et al, Tensor Based Knowledge Transfer Across Skill Categories for Robot Control, IJCAI, 2017]

Tensor-Based Transfer Learning: Some Examples

Multiple Visual Domains

Visual Decathlon Challenge

- A small sub-challenge of Artificial General Intelligence (AGI):
 - A general visual feature extractor.
 - Current best results: train one deep feature per dataset.
 - But for general visual intelligence, prefer a single deep network.
- Visual Decathlon Challenge [Rebuffi, NIPS'17, Learning multiple visual domains with residual adapters]
 - 10 datasets: Aircraft, Handwriting, Flowers, Pedestrians, Traffic Signs...
 - Goal is to learn them incrementally, without forgetting, and with minimal parameter growth.



Visual Decathlon: Tensor Solution

- Train Deep networks for the first K domains
- Extract and fix the core tensor
 - (no forgetting)
- For each new domain, train only task specific factor
 - (small growth pert task)



Tensor-based Transfer: Results

• Visual Decathlon [Bulat et al, arXiv'19]



		Dataset											
Model	#param	ImNet	Airc.	C100	DPed	DTD	GTSR	Flwr	OGlt	SVHN	UCF	mean	Scor
#images	-	1.3M	7k	50k	30k	4k	40k	2k	26k	70k	9k	-	-
Rebuffi et al. [29]	2×	59.23	63.73	81.31	93.30	57.02	97.47	83.43	89.82	96.17	50.28	77.17	2643
Rosenfeld et al. [32]	$2 \times$	57.74	64.11	80.07	91.29	56.54	98.46	86.05	89.67	96.77	49.38	77.01	2851
Mallaya et al. [22]	$1.28 \times$	57.69	65.29	79.87	96.99	57.45	97.27	79.09	87.63	97.24	47.48	76.60	2838
Series Adap. [30]	$2 \times$	60.32	61.87	81.22	93.88	57.13	99.27	81.67	89.62	96.57	50.12	77.17	3159
Parallel Adap. [30]	$2 \times$	60.32	64.21	81.91	94.73	58.83	99.38	84.68	89.21	96.54	50.94	78.07	3412
Parallel SVD [29]	1.5×	60.32	66.04	81.86	94.23	57.82	99.24	85.74	89.25	96.62	52.50	78.36	3398
Ours	$1.35 \times$	61.48	67.36	80.84	93.22	59.10	99.64	88.99	88.91	96.95	47.90	78.43	3585

Tensor-Based Transfer Learning: Some Examples

Fact Induction

Knowledgebase Completion

- Induce Missing Facts:
 - Completion model induces missing tuples:
 - <Cambridge, locatedIn, England>
 - Task := Relation
 - \rightarrow Independent tasks.
- Typical Solution: Train Embeddings θ:
- vec_θ("Cambs.")+vec_θ("locatedIn")≅vec_θ("England")
- vec_θ("Cambs.")*matrix_θ("locatedIn")≅vec_θ("England")
- vec_θ("Cambs.")*matrix_θ("locatedIn")*vec_θ("England") ≅ 1

Regression Problem. One per relation.



Knowledgebase Completion

• Induce Missing Facts:

Typical Solution: Train Embeddings θ :

- vec_θ("Cambs.")*matrix_θ("locatedIn")≅vec_θ("England")
- vec_θ("Cambs.")*matrix_θ("locatedIn")*vec_θ("England") ≅ 1



- CoreTensor $x_1 \operatorname{vec}_{\theta}(\text{``locatedIn''}) x_2 \operatorname{vec}_{\theta}(\text{``Cambs.''}) x_3 \operatorname{vec}_{\theta}(\text{``England''}) \cong 1$
- CoreTensor x₁ vec_θ("locatedIn") x₂ vec_θ("Cambs.") ≅ vec_θ("England")

Generates a DxD matrix

 \Rightarrow View as a weight generation net.

To learn one relation: Before: D² parameters. **Now**: R<D parameters.



Share across regressions (and entities)

Knowledgebase Completion

- Induce Missing Facts:
 - Completion model induces missing tuples:
 - <Cambridge, locatedIn, England>
 - Task := Relation
 - \rightarrow Independent tasks.



91

			WN	18 R R			FB15k-237					
	Linear	MRR	Hits@10	Hits@3	Hits@1	MRR	Hits@10	Hits@3	Hits@1			
DistMult (Yang et al., 2015)	yes	.430	.490	.440	.390	.241	.419	.263	.155			
ComplEx (Trouillon et al., 2016)	yes	.440	.510	.460	.410	.247	.428	.275	.158			
Neural LP (Yang et al., 2017)	no	—	—		—	.250	.408	—	-			
R-GCN (Schlichtkrull et al., 2018)	no	—	-	_	—	.248	.417	.264	.151			
MINERVA (Das et al., 2018)	no	—	-	-	—	—	.456	—	—			
ConvE (Dettmers et al., 2018)	no	.430	.520	.440	.400	.325	.501	.356	.237			
HypER (Balažević et al., 2018)	no	.465	.522	.477	.436	.341	.520	.376	.252			
M-Walk (Shen et al., 2018)	no	.437	—	.445	.414	—	_	—	a a			
TuckER (ours)	yes	.470	.526	.482	.443	.358	.544	.394	.266			

Share across regressions (and entities)

Tensor Knowledge Transfer

Summary:

- Pros:
 - Fast
 - Sample-efficient
 - Parameter-efficient



Task Agnostic

Task Specific

Input

- Relatedness is not crucial if enough ranks.
- Cons:
 - Need many source tasks (cf: Few-Shot vs "Fine-tuning")
 - Tasks need same "shape"
Tensor vs Metric Based Few-Shot

- Tensor-based
 - Suitable for any task (classification, regression, detection)
 - Weak assumption on known source task relevance (vs fine-tune).
 - Requires (-) and exploits (+) multiple sources-.

- Metric-based
 - Suitable for classification.
 - Suitable to extend an existing label-space.
 - No requirement (+), and limited exploitation (-) of multi-source.

Tensors: Bonus Topics

Tensor-Based Network Compression

Tensor-Factorisation for Model Compression

- Aside: These tensor methods also widely used for model compression & acceleration.
 - Fully Connected Compression: $D^1 \times D^2 \rightarrow D^1 \times r + D^2 \times r$



Tensors: Bonus Topics

Tensor-Based Modality Fusion

How to fuse multi-modal data?

- Often: One prediction, Two signals: Audio, video. Still, Motion.
- Typical answer: $y = w^T [x_a, x_b] = w^T_a x_a + w^T_b x_b$
 - Shallow model: concatenate (implies linear combination)
 - Deep model: concatenate and pass through MLP. $y = V\sigma(W[x_a, x_b])$
- Also: Tensor-based data fusion
 - Tensor fusion layer: $y = W \times_1 x_a \times_2 x_b$
 - Question: What's the problem with this?
 - Too many O(N³) parameters:
 - Low-rank tensor fusion layer: $y = S \times_1(U x_a) \times_2(V x_b)$

Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Part II: Learning with Multiple Domains

What is the difference between transfer learning and domain adaptation?

- Confusion caused by vagueness of term "transfer learning"
 - If "TL" means <u>any</u> knowledge transfer: DA is a subset of TL.
 - IF "TL" means "supervised" or more specifically "fine-tune":
 - (Unsupervised) Domain adaptation is disjoint to TL.
- Problem: Domain-Shift vs Task Shift

Data:	$n(X_{\star}) \neq n(X_{\star})$	$p(Y_{a} X_{a}) \neq p(Y_{b} X_{b})$	
			Confusing:
$\{X_s, Y_s\} \to \{X_t, Y_t\}$	Supervised Domain Adaptation. (Incl: Fine-tuning, Tensor).	Supervised Task Adaptation. (Incl Fine-tuning. Tensor).	e.g., FT, Tensor, can apply to
$\{X_s, Y_s\} \to \{X_t\}$	Unsupervised DA.	Impossible. Unless ZSL.	both.
$\{X_S, Y_S\} \to \{\}$	Domain Generalization	Impossible. Unless ZSL.	

Learning with Multiple Domains

Solving Domain-Shift:

- Model still works, even when data statistics change:
 - Different domains are a nuisance variable:

Cross Domain Matching:

• Your goal is to match across domains.









Cross-Domain Matching

• Text->Image Retrieval



a white plate topped with a cut in half sandwich.



a city street filled with lots of traffic.

a train traveling down tracks next to a forest.



Image->Caption Retrieval



a white plate topped with a cut in half sandwich.



a city street filled with lots of traffic.

a train traveling down tracks next to a forest.





• Bilingual Dictionary Induction



[Mikolov, arXiv'13, Exploiting similarities among languages for machine translation]

• Sketch-based Image Retrieval: Online Shopping



[Yu et al, CVPR 2016, Sketch me that shoe]

• Person Re-Identification







• Description-Based Person Search

Description-based person search:

• E.g., "Suspect had blue jeans, brown jacket, and blonde hair"

Red-Shirt

irt

The man is wearing a white shirt and a pair of brown pants, and a black backpack.









Red-Shirt + Blue-Trousers





• Mind Reading



[Mitchell, Science'08, Predicting Human Brain Activity Associated with the Meanings of Nouns]

Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Regression and Energy Function Methods

Matching Approaches: Regression

Linear

• Train:

$$\min_{\boldsymbol{W}} \sum_{i} \|\boldsymbol{y}_{i} - \boldsymbol{W}\boldsymbol{x}_{i}\|$$

Deep • Train:



$$\min_{\boldsymbol{W}} \sum_{i} \|\boldsymbol{y}_{i} - f_{\boldsymbol{W}}(\boldsymbol{x}_{i})\|$$

• Test: Nearest Neighbour

• Test: Nearest Neighbour

$$\operatorname{argmin}_{\boldsymbol{y}_j \in Gallery} \| \boldsymbol{y}_j - \boldsymbol{W} \boldsymbol{x}^* \|$$

$$\operatorname{argmin}_{\boldsymbol{y}_j \in Gallery} \left\| \boldsymbol{y}_j - f_{\boldsymbol{W}}(\boldsymbol{x}^*) \right\|$$

Matching Approaches: Energy Functions

• Bilinear Train:

$$\min_{\boldsymbol{W}} \sum_{(i,j)\in Pos,(i,k)\in Neg} \left| -\boldsymbol{x}_i^T \boldsymbol{W} \boldsymbol{y}_i - \Delta + \boldsymbol{x}_i^T \boldsymbol{W} \boldsymbol{y}_k \right|^+ \bullet \text{Low-Rank}$$

$$\min_{\boldsymbol{U},\boldsymbol{V}} \sum_{(i,j)\in Pos,(i,k)\in Neg} |-\boldsymbol{x}_i \boldsymbol{U} \boldsymbol{V} \boldsymbol{y}_i - \Delta + \boldsymbol{x}_i \boldsymbol{U} \boldsymbol{V} \boldsymbol{y}_k|^+$$

• Bilinear Test:

 $\operatorname{argmax}_{\boldsymbol{y}_{j} \in Gallery} \| \boldsymbol{y}_{j}^{T} \boldsymbol{W} \boldsymbol{x}^{*} \|$

Matching Approaches: Energy Functions

• Deep Train: • Deep Test

$$\min_{U,V} \sum_{(i,j)\in Pos,(i,k)\in Neg} \left| f_U(x_i)^T f_V(y_i) - f_U(x_i)^T f_V(y_j) - \Delta \right|^+ \qquad \max_{y_j\in Gal} f_U(x^*)^T f_V(y_j)$$

Connection to Metric Learning:

- "Siamese"
- Max energy/min distance

$$\min_{\substack{U,V\\ (i,j)\in Pos\\(i,k)\in Neg}} \sum_{\substack{(i,j)\in Pos\\(i,k)\in Neg\\(i,k)\in Neg}} \left| d^{-} \left(f_{U}\left(x_{i}\right), f_{V}\left(x_{k}\right) \right) - d^{+} \left(f_{U}\left(x_{i}\right), f_{V}\left(y_{j}\right) \right) - \Delta \right|^{+}$$

Discussion

- Multiple losses
 - Sometimes good to use classification, matching, and ranking.
- Siamese. Homo. vs Hetero. Branches.
 - Depends on modality heterogeneity.
- Non-CNN encoders:





Discussion

- Multiple losses
 - Sometimes good to use classification, matching, and ranking.
- Siamese. Homo. vs Hetero. Branches.
 - Depends on modality heterogeneity.
- Non-CNN encoders:





Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Cross-Domain Regression Examples

Sketch-to-Image: Sketch-based Image Retrieval



[CVPR'16 – Sketch me that shoe, ICCV'17, CVPR'18]

Image-to-Sketch: Learning to Draw



[Song, CVPR'18, Learning to Sketch with Shortcut Cycle Consistency]

Zero-Shot Learning

What is Zero-Shot Learning?

Live Demonstration!

- Audience Task: Recognise the Wampimuk.
 - Impossible?
- Solution: Semantic Transfer
 - Domain Ontology:
 - Wampimuk := small, horns, furry, cute
 - Wikipedia Page:

WIKIPEDIA Englis Españo The Free Encyclopedia La enciclopedia libre 4 853 000+ articles 1 172 000+ articulos Zero-Shot: Deutsch Русский Die freie Enzyklopädie Своболная энциклопели 1. Pattern recognition with no train)rices 'encyclopédie libre 2. Solved by semantic transference 1 614 000+ articles Italiano L'enciclopedia libera 814 000+ (81) 1 193 000+ voc Português Wolna encyklopedia A enciclopédia livre 1 106 000+ base 871 000+ artico Image: Ima wampimuk C English



From Supervised to Zero-Shot Pattern Recognition



From Supervised to Zero-Shot Pattern Recognition

Key is to embed categories as vectors

Task: Transferred from external source





Data-category vector map can generalize to new categories.

Where to get Category Vectors?

- "Supervised" sources:
 - Manual annotation of class properties
 - Vector encoding of a taxonomic class hierarchy
- "Unsupervised" sources: Existing unstructured data
 - Word (token) co-occurrence.
 - OR: word2vec: Representation of word prediction neural net.

Category Vectors

 \star

*

*

*

 V_2

• => Automatic+Free word vector for any "nameable" category.


Regression/Classification Approach

- A simple category vector: class-level attribute description.
 - Wampimuk := small, cute, furry, horns.
- Train:
 - Given some known class-category vectors v and data x:
 - Learn image-attribute classifiers/regressors v=f(x).
 - E.g., SVM/SVR. Deep Neural Nets.
- Test:
 - Specify vec v* for new class
 - Map new data f(x*)
 - NN matching of v* vs f(x*)
- Pros:
 - Easy and fast!
- Cons
 - Category separability





Why Does Zero-Shot Recognition Work?

- Very little theory.Intuition:
- - IF training category vectors and image vectors lie in the same relative positions on their respective manifolds....
 - Then a few examples can establish correspondence between the two spaces
 - And any new category can also be recognized.



ZSL Summary

- Family of ways to recognize (or regress, etc) with no training examples.
- Pros:
 - Avoid data bottleneck.
 - ZSL performance can be similar to K (1-5) shot performance
- Cons:
 - Need category embeddings
 - Performance is worse than many-shot learning.

From Supervised to Unsupervised Cross-Domain Matching

A Dictionary Induction Example

[Mukherjee & Hospedales, Learning Unsupervised Word Translations Without Adversaries, EMNLP 2018]

Intro: Bilingual Dictionary Induction

Goal: Inferring a bilingual lexicon from data

- Bilingual Lexicons:
 - Important for multi- and cross-lingual tasks: Multilingual word embedding, cross-lingual transfer learning, etc.
 - Not all pairs of languages have them.
 - Can we infer them from data?



Intro: Bilingual Dictionary Induction

Supervised Lexicon Induction

- Given some matching word pairs, deduce the full bilingual lexicon.
- Famously studied by Mikolov et al, 2013.
- Approach:
 - Monolingual word embeddings
 - Train a regressor on known pairs
 - Apply regressor to match unknown pairs





Intro: Bilingual Dictionary Induction

Unsupervised Lexicon Induction

- Related to statistical decipherment (Knight, Coling'06)
- How to deduce the full bilingual lexicon starting with no pairs?



Deep Distribution Matching: Matching

Our Approach: Statistical dependency:

• EG: Hilbert Schmid Independency Criteria (HSIC). Squared loss mutual information (SMI). Search for the representation and the pairing, that maximises statistical dependency.

1. Update: Association to maximize statistical dependency (given representation)





Spanish, p(y)

Roughly: Search for the sorting that makes the two within domain kernel matrices "look" similar.

$$\operatorname{tr}(\boldsymbol{K}\boldsymbol{L}) \implies \operatorname{tr}(\boldsymbol{K}\boldsymbol{\Pi}^{\top}\boldsymbol{L}\boldsymbol{\Pi})$$



Method: Deep Distribution Matching

Our Approach: Statistical dependency:

- EG: Hilbert Schmid Independency Criteria (HSIC). Squared loss mutual information (SMI). Search for the representation and the pairing, that maximises statistical dependency.
- Update: Association
 Linear Assignment Problem
- 2. Update: Representation Normal back-prop.
 - Good:
 - Every update decreases objective
 - No adversarial min-max.

$$\begin{split} & \underset{\boldsymbol{\Theta}_{x},\boldsymbol{\Theta}_{y},\boldsymbol{\Pi}}{\min} \underbrace{\Omega(\mathcal{D};\boldsymbol{\Theta}_{x},\boldsymbol{\Theta}_{y})}_{Regularizer} - \underbrace{\lambda D_{\boldsymbol{\Pi}}(\mathcal{D};\boldsymbol{\Theta}_{x},\boldsymbol{\Theta}_{y})}_{Dependency}, \\ & D_{\boldsymbol{\Pi}}(\mathcal{D};\boldsymbol{\Theta}_{x},\boldsymbol{\Theta}_{y}) = D_{\boldsymbol{\Pi}}(\{\boldsymbol{g}_{x}(\boldsymbol{x}_{i}),\boldsymbol{g}_{y}(\boldsymbol{y}_{\pi(i)})\}_{i=1}^{n}), \\ & \Omega(\mathcal{D};\boldsymbol{\Theta}_{x},\boldsymbol{\Theta}_{y}) = \sum_{i=1}^{n} \|\boldsymbol{x}_{i} - \boldsymbol{f}_{x}(\boldsymbol{g}_{x}(\boldsymbol{x}_{i}))\|_{2}^{2} \end{split}$$



Results

- Results Summary:
 - Promising results compared to recent adversarial methods.
 - Both much better than prior non-adversarial methods.
 - Simple convergence compared to adversarial methods.



Outline

- Part I: Sparse Data Deep Learning
 - Intro & common techniques.
 - Metric learning
 - Tensor-based
- Part II: Learning with multiple domains
 - Regression and Energy Functions for matching
 - Vision and Language Examples
- Links to Practical Examples & Exercises

Practical Exercises

- Examples and exercises on metric-based few-shot recognition.
 - Python/pytorch on google colab platform.
- Link to walkthrough (google docs):
 - <u>http://tiny.cc/ye2s8y</u>

References

- [Sun et al, Revisiting Unreasonable Effectiveness of Data in Deep Learning Era, ICCV'17]
- [Improving SVM Accuracy by Training on Auxiliary Data Sources, ICML'04]
- [Cross-domain video concept detection using adaptive svms, ACM MM'07]
- [How transferable are features in deep neural networks?, NIPS'14]
- [Taskonomy: Disentangling Task Transfer Learning, CVPR'18]
- [HOUDINI: Lifelong Learning as Program Synthesis, NIPS'18]
- [Evgeniou, Learning Multiple Tasks with Kernel Methods, JMLR'05]
- [Kang, Learning with Whom to Share in Multi-task Feature Learning, ICML'15]
- [Yang & Hospedales, ICLR'15, A unified perspective on multi-task and multi-domain learning]
- [Yang & Hospedales, ICLR'17, Deep Multi-task Representation Learning: A Tensor Factorisation Approach]
- [Rebuffi, NIPS'17, Learning multiple visual domains with residual adapters]
- [Bulat, arXiv'19, Incremental multi-domain learning with network latent tensor factorization]
- [Balažević , arXiv'19, TuckER: Tensor Factorization for Knowledge Graph Completion]
- [Lebedev, ICLR'15, Speeding-up Convolutional Neural Networks Using Fine-tuned CP-Decomposition]
- [Kossaifi, CVPR'19, T-Net: Parametrizing Fully Convolutional Nets with a Single High-Order Tensor]
- [Hu et al, Neural Tensor Fusion Networks, ICCV'17]
- [EG: Hirzer, ECCV'12; Mignon, CVPR'12]
- [Learning Fine-grained Image Similarity with Deep Ranking, CVPR'14]
- [Wang, PAMI'18, A Survey on Learning to Hash]
- [Snell et al, NIPS'17, Prototypical Networks for Few Shot Learning]
- [Fort, arXiv'17, Gaussian Prototypical Networks for Few-Shot Learning]
- [Sung, CVPR'18, "Learning to Compare: Relation Networks for Few Shot Learning"]
- [Zhang, arXiv'18, "Relation Columns for Few Shot Learning"]
- [Alemi, ICLR'17, Deep Variational Information Bottleneck]
- [Finn, ICML'17, Model Agnostic Meta Learning]
- [Nichol, arXiv'18, On First-Order Meta-Learning Algorithms]
- [Wang, PAMI'18, A Survey on Learning to Hash]
- [Mukherjee, EMNLP'16, Gaussian Visual-Linguistic Embedding for Zero-Shot Recognition]