



Language Technologies Institute



Multimodal Machine Learning

Lecture 7.1: Multimodal alignment

* Original version co-developed with Tadas Baltrusaitis

Lecture objectives

- Multimodal alignment
 - Implicit
 - Explicit
- Explicit signal alignment
 - Dynamic Time Warping
 - Canonical Time Warping
- Attention models in deep learning (implicit and explicit alignment)
 - Soft attention
 - Hard attention
 - Spatial Transformer Networks





Administrative Stuff



Language Technologies Institute

Upcoming Schedule

- First project assignment:
 - Proposal presentation (10/1 and 10/3)
 - First project report (Sunday 10/6)
- Midterm project assignment
 - Midterm presentations (11/5 and 11/7)
 - Midterm report (Sunday 11/10)
- Final project assignment
 - Final presentation (12/3 & 12/5)
 - Final report (Sunday 12/8)



Multi-modal alignment



Language Technologies Institute

Multimodal-alignment

- Multimodal alignment finding relationships and correspondences between two or more modalities
- Two types
 - **Explicit** alignment is the task in itself
 - Implicit / Latent alignment helps when solving a different task (for example "Attention" models)
- Examples ?
 - Images with captions
 - Recipe steps with a how-to video
 - Phrases/words of translated sentences





Explicit multimodal-alignment

- Explicit alignment goal is to find correspondences between modalities
 - Aligning speech signal to a transcript
 - Aligning two out-of sync sequences
 - Co-referring expressions





Implicit multimodal-alignment

- Implicit alignment uses internal latent alignment of modalities in order to better solve various problems
 - Machine Translation
 - Cross-modal retrieval
 - Image & Video Captioning
 - Visual Question Answering



What season does this appear to be? GT: fall Our Model: fall



What is soaring in the sky? GT: kite Our Model: kite





Explicit alignment



Language Technologies Institute



Temporal sequence alignment



Applications:

- Re-aligning asynchronous data

- Finding similar data across modalities (we can estimate the aligned cost)

- Event reconstruction from multiple sources



Let's start unimodal – Dynamic Time Warping

 We have two unaligned temporal unimodal signals

•
$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n_X}] \in \mathbb{R}^{d \times n_X}$$

•
$$\mathbf{Y} = \left[\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{n_y} \right] \in \mathbb{R}^{d \times n_y}$$

Find set of indices to minimize the alignment difference:

$$L(\mathbf{p}_{t}^{x}, \mathbf{p}_{t}^{y}) = \sum_{t=1}^{l} \left\| x_{\mathbf{p}_{t}^{x}} - y_{\mathbf{p}_{t}^{y}} \right\|_{2}^{2}$$

- Where p^{x} and p^{y} are index vectors of same length
- Dynamic Time Warping is designed to find these index vectors





Dynamic Time Warping continued

Lowest cost path in a cost matrix

- Restrictions?
 - Monotonicity no going back in time
 - Continuity no gaps
 - Boundary conditions start and end at the same points
 - Warping window don't get too far from diagonal
 - Slope constraint do not insert or skip too much





Dynamic Time Warping continued

Lowest cost path in a cost matrix

 Solved using dynamic programming while respecting the restrictions





DTW alternative formulation

$$L(p^{x}, p^{y}) = \sum_{t=1}^{l} ||x_{p_{t}^{x}} - y_{p_{t}^{y}}||_{2}^{2}$$
Replication doesn't change the objective
$$\int_{1}^{4} \int_{1}^{4} \int_{1}^{1} \int_{1}^{1}$$

Alternative objective:

$$L(\boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \left\| \boldsymbol{X}\boldsymbol{W}_{\boldsymbol{x}} - \boldsymbol{Y}\boldsymbol{W}_{\boldsymbol{y}} \right\|_{F}^{2}$$

Frobenius norm $||A||_F^2 = \sum_i \sum_j |a_{i,j}|^2$

X, Y — original signals (same #rows, possibly different #columns)

 W_{χ} , W_{γ} - alignment matrices



DTW – Some Limitations

Computationally complex



- Sensitive to outliers
- Unimodal!







Canonical Correlation Analysis reminder

maximize:
$$tr(U^T \Sigma_{XY} V)$$

subject to: $U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$, $u_{(i)}^T \Sigma_{XY} v_{(i)} = 0$ for $i \neq j$

Linear projections maximizing correlation



Orthogonal projections

Unit variance of the projection vectors





Language Technologies Institute

Canonical Correlation Analysis reminder

- When data is normalized it is actually equivalent to smallest RMSE reconstruction
- CCA loss can also be re-written as:

 $L(\boldsymbol{U},\boldsymbol{V}) = \|\boldsymbol{U}^T\boldsymbol{X} - \boldsymbol{V}^T\boldsymbol{Y}\|_F^2$

subject to:
$$U^T \Sigma_{YY} U = V^T \Sigma_{YY} V = I$$
, $u_{(j)}^T \Sigma_{XY} v_{(i)} = 0$





Canonical Time Warping

Dynamic Time Warping + Canonical Correlation Analysis
 = Canonical Time Warping

$$L(\boldsymbol{U}, \boldsymbol{V}, \boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \left\| \boldsymbol{U}^{T} \boldsymbol{X} \boldsymbol{W}_{\boldsymbol{x}} - \boldsymbol{V}^{T} \boldsymbol{Y} \boldsymbol{W}_{\boldsymbol{y}} \right\|_{F}^{2}$$

- Allows to align multi-modal or multi-view (same modality but from a different point of view)
- W_x , W_y temporal alignment
- U, V cross-modal (spatial) alignment

[Canonical Time Warping for Alignment of Human Behavior, Zhou and De la Tore, 2009]





Canonical Time Warping

$$L(\boldsymbol{U}, \boldsymbol{V}, \boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \left\| \boldsymbol{U}^{T} \boldsymbol{X} \boldsymbol{W}_{\boldsymbol{x}} - \boldsymbol{V}^{T} \boldsymbol{Y} \boldsymbol{W}_{\boldsymbol{y}} \right\|_{F}^{2}$$

Optimized by Coordinate-descent – fix one set of parameters, optimize another

Generalized Eigen-decomposition



[Canonical Time Warping for Alignment of Human Behavior, Zhou and De la Tore, 2009, NIPS]





Generalized Time warping

 Generalize to multiple sequences all of different modality

$$L(\boldsymbol{U}_{i}, \boldsymbol{W}_{i}) = \sum_{i=1}^{T} \sum_{j=1}^{T} \left\| \mathbf{U}_{i}^{T} \mathbf{X}_{i} \mathbf{W}_{i} - \mathbf{U}_{j}^{T} \mathbf{X}_{j} \mathbf{W}_{j} \right\|_{F}^{2}$$

- *W_i* set of temporal alignments
- *U_i* set of cross-modal (spatial) alignments



(1) Time warping(2) Spatial embedding



[Generalized Canonical Time Warping, Zhou and De la Tore, 2016, TPAMI]



Alignment examples (unimodal)

CMU Motion Capture Subject 1: 199 frames Subject 2: 217 frames Subject 3: 222 frames



Weizmann

Subject 1: 40 frames Subject 2: 44 frames

Subject 3: 43 frames





Alignment examples (multimodal)







Canonical time warping - limitations

- Linear transform between modalities
- How to address this?





$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{W}_{\boldsymbol{x}}, \boldsymbol{W}_{\boldsymbol{y}}) = \left\| f_{\boldsymbol{\theta}_1}(\mathbf{X}) \mathbf{W}_{\mathbf{x}} - f_{\boldsymbol{\theta}_1}(\mathbf{Y}) \mathbf{W}_{\mathbf{y}} \right\|_F^2$$

Could be seen as generalization of DCCA and GTW



[Deep Canonical Time Warping, Trigeorgis et al., 2016, CVPR]



$$L(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{W}_x, \boldsymbol{W}_y) = \left\| f_{\boldsymbol{\theta}_1}(\mathbf{X}) \mathbf{W}_{\mathbf{X}} - f_{\boldsymbol{\theta}_1}(\mathbf{Y}) \mathbf{W}_{\mathbf{y}} \right\|_F^2$$

- The projections are orthogonal (like in DCCA)
- Optimization is again iterative:
 - Solve for alignment (W_x, W_y) with fixed projections (θ_1, θ_2)
 - Eigen decomposition
 - Solve for projections (θ_1, θ_2) with fixed alignment (W_x, W_y)
 - Gradient descent
 - Repeat till convergence

[Deep Canonical Time Warping, Trigeorgis et al., 2016, CVPR]



Implicit alignment



Language Technologies Institute

Implicit alignment

- We looked how to explicitly align temporal data
- Could use that as an internal (hidden) step in our models?
- Can we instead encourage the model to align data when solving a different problem?
- Yes!
 - Graphical models
 - Neural attention models (focus of today's lecture)



Attention models



Language Technologies Institute

Attention in humans

- Foveal vision we only see in "high resolution" in 2 degrees of vision
- We focus our attention selectively to certain words (for example our names)
- We attend to relevant speech in a noisy room







Attention models in deep learning

- Many examples of attention models in recent years!
- Why:
 - Allows for implicit data alignment
 - Good results empirically
 - In some cases faster (don't need to focus on all the image)
 - Better Interpretability





Types of Attention Models

- Recent attention models can be roughly split into three major categories
 - 1. Soft attention
 - Acts like a gate function. Deterministic inference.
 - 2. Transform network
 - Warp the input to better align with canonical view
 - 3. Hard attention
 - Includes stochastic processes. Related to reinforcement learning.





Soft attention



Language Technologies Institute



Machine Translation

• Given a sentence in one language translate it to another



 Not exactly multimodal task – but a good start! Each language can be seen almost as a modality.





Machine Translation with RNNs

- A quick reminder about encoder decoder frameworks
- First we encode the sentence
- Then we decode it in a different language

Context / embedding / sentence representation

Dog

on



Encoder

chien

sur

la

plage

le

Carnegie Mellon University

beach

Decode

the

Machine Translation with RNNs

- What is the problem with this?
- What happens when the sentences are very long?
- We expect the encoders hidden state to capture everything in a sentence, a very complex state in a single vector, such as







Decoder – attention model

 Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states





Decoder – attention model

 Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states





Decoder – attention model

 Before encoder would just take the final hidden state, now we actually care about the intermediate hidden states





How do we encode attention

- Before:
 - $p(y_i|y_1, ..., y_{i-1}, x) = g(y_{i-1}, s_i, z)$, where $z = h_T$, and s_i - the current state of the decoder
- Now:

•
$$p(y_i|y_1, ..., y_{i-1}, x) = g(y_{i-1}, s_i, z_i)$$

- Have an attention "gate"
 - A different context z_i used at each time step!

•
$$\mathbf{z}_i = \sum_{j=i}^{T_x} \alpha_{ij} \mathbf{h}_j$$

 α_{ij} - the (scalar) attention for word j at generation step i



MT with attention

So how do we determine α_{ij} ,

•
$$\alpha_{i,j} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_{\chi}} \exp(e_{ik})}$$
 - softmax, making sure they sum to 1

Where:

•
$$e_{ij} = \boldsymbol{v}^T \, \sigma \big(W \boldsymbol{s_{i-1}} + U \boldsymbol{h_j} \big)$$

a feedforward network that can tell us given the current state of decoder how important the current encoding is now

v, W, U- learnable weights

$$z_i = \sum_{j=i}^{T_{\mathcal{X}}} \alpha_{ij} h_j$$

expectation of the context (a fancy way to say it's a weighted average)



MT with attention

Basically we are using a neural network to tell us where a neural network should be looking!

- We can use with RNN, LSTM or GRU
- Encoder being used is the same structure as before
 - Can use uni-directional
 - Can use bi-directional
- Model can be trained using our regular back-propagation through time, all of the modules are differentiable



Does it work?





MT with attention recap

- Get good translation results (especially for long sentences)
- Also get a (soft) alignment of sentences in different languages
 - Extra interpretability of method functioning
- How do we move to multimodal?





Visual captioning with soft attention



throwing(0.33),









frisbee(0.37)



is(0.37)

in(0.21)









[Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Xu et al., 2015]





Recap RNN for Captioning



Why might we not want to focus on the final layer?





Looking at more fine grained features







Soft attention

- Allows for latent data alignment
- Allows us to get an idea of what the network "sees"
- Can be optimized using back propagation
- Good at paper naming!
 - Show, Attend and Tell (extension of Show and Tell)
 - Listen, Attend and Walk
 - Listen, Attend and Spell
 - Ask, Attend and Answer







Language Technologies Institute

Some limitations of grid based attention

Can we fixate on small parts of image but still have easy end-to-end training?



A woman is throwing a frisbee in a park.

- A dog is standing on a hardwood floor.



A <u>stop</u> sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of <u>people</u> sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.









Language Technologies Institute





Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{1,1} & \theta_{1,2} & \theta_{1,3} \\ \theta_{2,1} & \theta_{2,2} & \theta_{2,3} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$





Idea: Function mapping pixel coordinates (x^t, y^t) of output to pixel coordinates (x^s, y^s) of input



Network "attends" to input by predicting θ



Repeat for all pixels in *output* to get a **sampling grid**







Differentiable "attention / transformation" module



Insert spatial transformers into a classification network and it learns to attend and transform the input







Examples on real world data

Results on traffic sign recognition





Code available http://torch.ch/blog/2015/09/07/spatial_transformers.html





Recap on Spatial Transformer Networks

- Differentiable so we can just use back-prop for training end-to-end
- Can be used with complex transformations to focus on an image
 - Affine and Piece-Wise Affine, Perspective, This Plate Splines
- We can use it instead of grid based soft and hard attention for multimodal tasks







Glimpse Network (Hard Attention)



Language Technologies Institute

Hard attention

- Soft attention requires computing a representation for the whole image or sentence
- Hard attention on the other hand forces looking only at one part
- Main motivation was reduced computational cost rather than improved accuracy (although that happens a bit as well)
- Saccade followed by a glimpse how human visual system works

[Recurrent Models of Visual Attention, Mnih, 2014] [Multiple Object Recognition with Visual Attention, Ba, 2015]





Hard attention examples







Glimpse Sensor

Looking at a part of an image at different scales



- At a number of different scales combined to a single multichannel image (human retina like representation)
- Given a location l_t output an image summary at that location
 [Recurrent Models of Visual Attention, Mnih, 2014]



Glimpse network

• Combining the Glimpse and the location of the glimpse into a joint network



- The glimpse is followed by a feedforward network (CNN or a DNN)
- The exact formulation of how the location and appearance are combined varies, the important thing is combining what and where
- Differentiable with respect to glimpse parameters but not the location



Overall Architecture - Emission network

- Given an image a glimpse location *l_t*, and optionally an action *a_t*
- Action can be:
 - Some action in a dynamic system – press a button etc.
 - Classification of an object
 - Word output
- This is an RNN with two output gates and a slightly more complex input gate!





Recurrent model of Visual Attention (RAM)

- Sample locations of glimpses leading to updates in the network
- Use gradient descent to update the weights (the glimpse network weights are differentiable)
- The emission network is an RNN
- Not as simple as backprop but doable
- Turns out this is very similar and in some cases equivalent to reinforcement learning using the REINFORCE learning rule [Williams, 1992]





Multi-modal alignment recap



Language Technologies Institute

Multimodal-alignment recap

- Explicit alignment aligns two or more modalities (or views) as an actual task. The goal is to find correspondences between modalities
 - Dynamic Time Warping
 - Canonical Time Warping
 - Deep Canonical Time Warping
- Implicit alignment uses internal latent alignment of modalities in order to better solve various problems
 - Attention models
 - Soft attention
 - Spatial transformer networks
 - Hard attention



