

# On Deep Multi-View Representation Learning (2015 ICML)

Weiran Wang, Raman Arora, Karen Livescu, Jeff Bilmes

Discussed by: Yizhe Zhang

April 15th, 2016

# Outline

1 Introduction

2 Method

3 Experiments

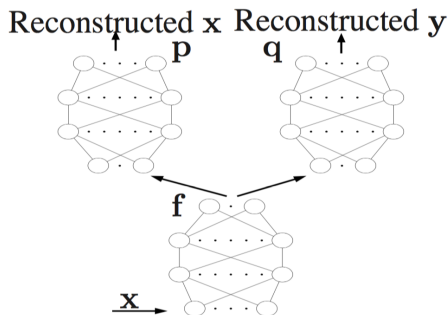
# Multi-view representation learning

- In many applications, we have access to multiple views of data:
  - Audio + video (Kidron et al., 2005; Chaudhuri et al., 2009)
  - Audio + articulation (Arora & Livescu, 2013)
  - Images + text (Hardoon et al., 2004; Socher & Li, 2010; Hodosh et al., 2013)
  - Parallel text in two language (Vinokourov et al., 2003; Haghighi et al., 2008)
- The task is to leverage multiple-view information to learn a better representation (than single view).
- At training time, one attempt to learn the *latent representations* from paired two-view training set.
- At test time, only primary view is available.

# DNN-based multiview feature learning

- Suppose we have access to paired observations from two views, denoted  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ , where  $N$  is the sample size.
- Split autoencoders (Ngiam et al. (2011)) minimize the sum of reconstruction errors for the two views

$$\min_{\mathbf{W}_f, \mathbf{W}_p, \mathbf{W}_q} \frac{1}{N} \sum_{i=1}^N (\|\mathbf{x}_i - \mathbf{p}(\mathbf{f}(\mathbf{x}_i))\|^2 + \|\mathbf{y}_i - \mathbf{q}(\mathbf{f}(\mathbf{x}_i))\|^2) \quad (1)$$



(a) SplitAE

# Canonical correlation analysis (CCA)

- Given two data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T$ ,  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_M]^T$ , canonical-correlation analysis (CCA) first seeks vectors  $\mathbf{u}_1 \in \mathbb{R}^N$  and  $\mathbf{v}_1 \in \mathbb{R}^M$  to maximize the correlation  $\rho = \text{corr}(\mathbf{u}_1^T \mathbf{X}, \mathbf{v}_1^T \mathbf{Y})$ .
- For identifiability issue, one may add constraint on  $\mathbf{u}_1^T \mathbf{X} \mathbf{X}^T \mathbf{u}_1$  and  $\mathbf{v}_1^T \mathbf{Y} \mathbf{Y}^T \mathbf{v}_1$
- Then one seeks vectors minimizing the same correlation, subject to the constraint that  $(\mathbf{u}_2^T \mathbf{X}, \mathbf{v}_2^T \mathbf{Y})$  are uncorrelated with the  $(\mathbf{u}_1^T \mathbf{X}, \mathbf{v}_1^T \mathbf{Y})$
- This procedure may be continued up to  $\min(N, M)$  times.

# Deep canonical correlation analysis (DCCA)

- Andrew et al. (2013) propose a DNN extension of CCA termed deep CCA.
- In DCCA, the canonical correlation of the extract features for each view is maximized

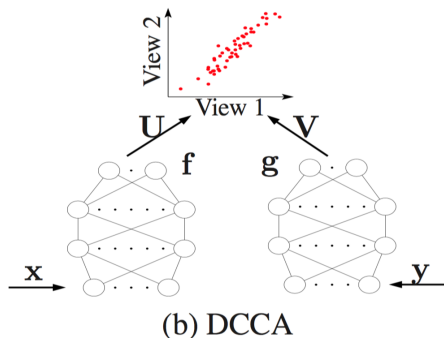


Figure: DCCA

# Deep canonical correlation analysis (DCCA) Cont'd

- Specifically, the optimization function can be written as <sup>1</sup>

$$\min_{\mathbf{W}_f, \mathbf{W}_g, \mathbf{U}, \mathbf{V}} \frac{1}{N} \text{tr}(\mathbf{U}^T \mathbf{f}(\mathbf{X}) \mathbf{g}(\mathbf{Y}) \mathbf{V}^T) \quad (2)$$

$$\text{s.t. } \mathbf{U}^T \left( \frac{1}{N} \mathbf{f}(\mathbf{X}) \mathbf{f}(\mathbf{X})^T + r_x \mathbf{I} \right) \mathbf{U} = \mathbf{I}, \text{ (whitening)} \quad (3)$$

$$\mathbf{V}^T \left( \frac{1}{N} \mathbf{g}(\mathbf{Y}) \mathbf{g}(\mathbf{Y})^T + r_y \mathbf{I} \right) \mathbf{V} = \mathbf{I}, \text{ (whitening)} \quad (4)$$

$$\mathbf{u}_i^T \mathbf{f}(\mathbf{X}) \mathbf{g}(\mathbf{Y})^T \mathbf{v}_j = 0, \text{ for } i \neq j \text{ (orthogonal)} \quad (5)$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ ,  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N]$ ,  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_L]$  and  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]$ .  $L$  denote the feature dimension.  $r_x, r_y$  are regularization parameters (De Bie & De Moor, 2003).

---

<sup>1</sup>optimization the DCCA requires full data by whitening constraints. The authors claims a sufficiently large minibatch will be enough for estimating the covariances, thus SGD is still valid

# Outline

1 Introduction

2 Method

3 Experiments



# Deep canonically correlated autoencoders (DCCA)

- They propose a model that optimizes the combination of *canonical correlation* and the *reconstruction errors* of the autoencoders.
- **Interpretation:** “The DCCA objective offers a trade-off between the information captured in the (input, feature) mapping within each view, and the information in the (feature, feature) relationship”

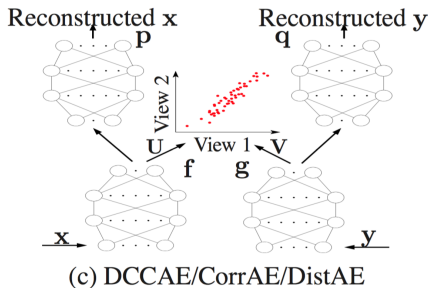


Figure: DCCA

# Deep canonically correlated autoencoders (DCCA) Cont'd

- Specifically, the optimization function can be written as

$$\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_p, \mathbf{W}_q, \mathbf{U}, \mathbf{V} - \frac{1}{N} \text{tr}(\mathbf{U}^T \mathbf{f}(\mathbf{X}) \mathbf{g}(\mathbf{Y})^T \mathbf{V}) \quad (6)$$

$$+ \frac{\lambda}{N} \sum_{i=1}^N (\|\mathbf{x}_i - \mathbf{p}(\mathbf{f}(\mathbf{x}_i))\|^2 + \|\mathbf{y}_i - \mathbf{q}(\mathbf{g}(\mathbf{y}_i))\|^2) \quad (7)$$

$$\text{s.t. } \mathbf{U}^T \left( \frac{1}{N} \mathbf{f}(\mathbf{X}) \mathbf{f}(\mathbf{X})^T + r_x \mathbf{I} \right) \mathbf{U} = \mathbf{I}, \quad (8)$$

$$\mathbf{V}^T \left( \frac{1}{N} \mathbf{g}(\mathbf{Y}) \mathbf{g}(\mathbf{Y})^T + r_y \mathbf{I} \right) \mathbf{V} = \mathbf{I}, \quad (9)$$

$$\mathbf{u}_i^T \mathbf{f}(\mathbf{X}) \mathbf{g}(\mathbf{Y})^T \mathbf{v}_j = 0, \text{ for } i \neq j \quad (10)$$

# Correlated autoencoders (CorrAE)

- They further relax the uncorrelated feature constraint of DCCAE (i.e.  $\mathbf{u}_i^T \mathbf{f}(\mathbf{X})$ ,  $\mathbf{v}_j^T \mathbf{g}(\mathbf{Y})$  etc. can be correlated), which they called correlated autoencoders (CorrAE)

$$\min_{\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_p, \mathbf{W}_q, \mathbf{U}, \mathbf{V}} -\frac{1}{N} \text{tr}(\mathbf{U}^T \mathbf{f}(\mathbf{X}) \mathbf{g}(\mathbf{Y})^T \mathbf{V}) \quad (11)$$

$$+ \frac{\lambda}{N} \sum_{i=1}^N (\|\mathbf{x}_i - \mathbf{p}(\mathbf{f}(\mathbf{x}_i))\|^2 + \|\mathbf{y}_i - \mathbf{q}(\mathbf{g}(\mathbf{y}_i))\|^2) \quad (12)$$

$$\text{s.t. } \mathbf{u}_i^T \mathbf{f}(\mathbf{X}) \mathbf{f}(\mathbf{X})^T \mathbf{u}_i = \mathbf{v}_i^T \mathbf{g}(\mathbf{Y}) \mathbf{g}(\mathbf{Y})^T \mathbf{v}_i = N, 1 \leq i \leq L. \quad (13)$$

- They demonstrate in experiments that this relaxation results in a large performance gap. Therefore, the constraint 10 is necessary.

# Minimum-distance autoencoders (DistAE)

- The whitening constraints complicate the optimization of CCA-based objectives. Thus they also consider another discrepancy objective to substitute CCA objective.

$$\min_{\mathbf{W}_f, \mathbf{W}_g, \mathbf{W}_p, \mathbf{W}_q} \frac{\|\mathbf{f}(\mathbf{x}_i) - \mathbf{g}(\mathbf{y}_i)\|^2}{\|\mathbf{f}(\mathbf{x}_i)\|^2 + \|\mathbf{g}(\mathbf{y}_i)\|^2} \quad (14)$$

$$+ \frac{\lambda}{N} \sum_{i=1}^N (\|\mathbf{x}_i - \mathbf{p}(\mathbf{f}(\mathbf{x}_i))\|^2 + \|\mathbf{y}_i - \mathbf{q}(\mathbf{g}(\mathbf{y}_i))\|^2) \quad (15)$$

- This objective is unconstrained and can be factorized by each training sample, so normal SGD applies using small minibatches.

# Outline

1 Introduction

2 Method

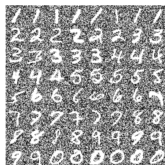
3 Experiments

# Experiments

- They compare the following methods in the multi-view learning setting on MNIST classification, speech recognition, and word pair semantic similarity.
  - **Baseline:** Original data without any transformation.
  - **DNN-based models:** SplitAE, CorrAE, DCCA, DCCAE, and DistAE.
  - **Linear CCA (CCA)**
  - **Kernel CCA approximations:** Two approximation methods: FKCCA (random Fourier features), NKCCA (Nyström approximation)

# Noisy MNIST digits

- MNIST dataset. Rescale to  $[0, 1]$ , with 60K/10K images for training/testing.
  - **View 1:** Add independent random uniform noise.
  - **View 2:** Rotate the images at angles uniformly drawn from  $[\pi/4, \pi/4]$ .
- “A good multi-view learning algorithm should be able to extract features that disregard the noise”.
- Criteria:
  - **ACC:** (Clustering accuracy) How well the spectral clustering of projected view 1 matches with Ground truth.
  - **NMI:** (Clustering accuracy) Normalized mutual information.
  - **Error:** Classification error of a linear SVM on the projections.



Method	ACC (%)	NMI (%)	Error (%)
Baseline	47.0	50.6	13.1
CCA ( $L = 10$ )	72.9	56.0	19.6
SplitAE ( $L = 10$ )	64.0	69.0	11.9
CorrAE ( $L = 10$ )	65.5	67.2	12.9
DistAE ( $L = 20$ )	53.5	60.2	16.0
FKCCA ( $L = 10$ )	94.7	87.3	5.1
NKCCA ( $L = 10$ )	95.1	88.3	4.5
DCCA ( $L = 10$ )	<b>97.0</b>	<b>92.0</b>	<b>2.9</b>
DCCAE ( $L = 10$ )	<b>97.5</b>	<b>93.4</b>	<b>2.2</b>

Figure: Performance comparison on the test set of noisy MNIST digits (view 1)

# Noisy MNIST digits

- They measure the class separation visually by t-SNE embedding.
- The DCCA and DCCAE manage to map digits of the same identity to similar locations while suppressing the rotational variation.
- Overall, DCCAE gives the cleanest embedding.
- The learned mixing weight hyperparameter  $\lambda$  is very small ( $10^{-3}$ ).

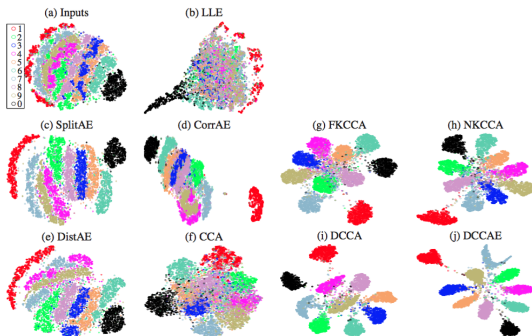


Figure: t-SNE embedding of the projected test set of noisy MNIST (view 1)



## Acoustic-articulatory data for speech recognition

- **Wisconsin X-Ray MicroBeam (XRMB) corpus.** Two views: *speech* and *articulatory measurements*.
- Split the XRMB data to 35/8/2/2 speakers for feature learning/recognizer training/tuning/testing.
- **Criteria:** Phone error rates (PERs) by using an HMM recognizer.

Method	Mean (std) PER (%)
Baseline	34.8 (4.5)
CCA	26.7 (5.0)
SplitAE	29.0 (4.7)
CorrAE	30.6 (4.8)
DistAE	33.2 (4.7)
FKCCA	26.0 (4.4)
NKCCA	26.6 (4.2)
DCCA	<b>24.8 (4.4)</b>
DCCAE	<b>24.5 (3.9)</b>

**Figure:** Mean and standard deviations of PERs over 6 folds obtained by each

# Multilingual data for word embeddings

- Learn a vectorial embedding representation of bigram (AN: adjective-noun, VN verb-noun). Two views: English and German.
- Tuning and test splits (of size 649/1,972) for each subset.
- **Criteria:** Spearman's correlation of the bigram representation between two languages.

Method	AN	VN	Avg.
Baseline	45.0	39.1	42.1
CCA	46.6	37.7	42.2
SplitAE	47.0	<b>45.0</b>	46.0
CorrAE	43.0	42.0	42.5
DistAE	43.6	39.4	41.5
FKCCA	46.4	42.9	44.7
NKCCA	44.3	39.5	41.9
DCCA	48.5	42.5	45.5
DCCAE	<b>49.1</b>	43.2	<b>46.2</b>

Figure: Spearman's correlation ( $\rho$ ) for bigram similarities.