### Discriminative Representations for Heterogeneous Images and Multimodal Data

Heather Couture

**Dissertation Defense** 

Nov. 19, 2018



Heterogeneous images

Imaging + genomics

# **Cancer Heterogeneity**



Nature Reviews | Cancer

Manusyk, "Intra-tumour heterogeneity: a looking glass for cancer?", 2012

### **Tissue Microarray**



# Motivations

### Prognosis

Assess tumor aggressiveness:



Intermediate grade highly variable amongst pathologists

Rakha, 2010

# Motivations

### Prognosis

Assess tumor aggressiveness:



Intermediate grade highly variable amongst pathologists

#### Rakha, 2010

### **Personalized treatment**

Target tumors based on molecular analysis:



### **Thesis Statement**

Learned representations for histology images of tissue can capture both intra- and inter-tumor heterogeneity, enabling discriminative models for tumor properties. Combining these image features with data from other modalities such as genomics in a taskdriven model can provide insight into the shared tumor properties and further improve predictions. These computational techniques using discriminative features can provide a lower cost and more repeatable alternative to molecular methods and insight into tumor **heterogeneity**.

# **Contributions to Breast Cancer Research**

- 1) Methods to capture **biologically-relevant** features by operating on the H&E stain intensities extracted from histology images
- A low cost and repeatable method for predicting histological, molecular, and genomic properties of tumors from H&E histology
- 3) A mechanism to find predicted **tumor heterogeneity** from H&E histology

Contributions to Computer Science will be discussed at the end

# Outline



2) Handling heterogeneous images

3) Combining imaging & genomics

# Outline



2) Handling heterogeneous images

3) Combining imaging & genomics

# **Problem Definition: Classification**

Infer tumor class from histology



high

Grade:

status:

## **Problem Definition: Classification**

### Infer tumor class from histology



### Stain Normalization & Extraction

### Original image



Hematoxylin: blue Eosin: pink Hematoxylin: red Eosin: green Residual: blue

Stain channels

### Background

### **Hand-crafted features**





nuclei segmentation Appearance & shape





region growing



best fit ellipses

Focused on cell by cell morphology

Difficult to adapt to new data sets

Delaunay triangulation

eosin

hematoxylin

### Background

### **Hand-crafted features**





nuclei segmentation Appearance & shape

best fit

ellipses





region co growing ł

convex hematoxylin hulls



Delaunay

triangulation

eosin

Focused on cell by cell morphology

# Difficult to adapt to new data sets

### **Learned features**

Dictionary learning Deep learning



training images



Adapted to given images

patches

features

# Method 1: Task-driven Dictionary Learning



# Method 1: Task-driven Dictionary Learning



### Method 2: Transferred Deep Features



(generalizable and discriminative)

### Method 1 & 2: Classification



### Feature & Classifier Comparison

AUC

Method	Log. reg.	Linear SVM	RBF SVM	DWD
Hand-crafted features	0.789 (0.032)	0.778 (0.027)	0.573 (0.040)	0.728 (0.022)
Dictionary learning (nuclei)	0.812 (0.020)	0.794 (0.017)	0.661 (0.045)	0.755 (0.035)
Dictionary learning (dense)	0.845 (0.020)	<b>0.855</b> (0.024)	0.631 (0.035)	0.799 (0.023)
Deep transfer learning	0.832 (0.029)	0.825 (0.032)	0.716 (0.039)	0.811 (0.034)

#### **SPECS data set:**

43 Basal, 42 Luminal A 2 cores/patient 5-fold cross-validation

#### **Dictionary learning:**

Dictionary size 256 Patch size 17x17 Nuclei-centered patches vs. dense patches

#### **Deep transfer learning:**

AlexNet, conv4

### Best results with dictionary learning

Deep transfer learning is promising Fine-tuning could improve further

# **Dictionary & Deep Learning Comparison**



Input image	AUC			Accuracy		
	Basal	ER	Grade	Basal	ER	Grade
Dictionary						
Original RGB	0.810 (0.008)	0.843 (0.009)	0.905 (0.007)	0.790 (0.012)	0.797 (0.009)	0.828 (0.010)
Normalized RGB	0.817 (0.010)	0.850 (0.009)	0.911 (0.010)	0.780 (0.010)	0.790 (0.012)	0.821 (0.013)
Stain channels	<b>0.822</b> (0.014)	<b>0.860</b> (0.008)	<b>0.927</b> (0.009)	<b>0.795</b> (0.012)	<b>0.805</b> (0.009)	<b>0.848</b> (0.010)
Deep transfer						
Original RGB	0.784 (0.017)	0.819 (0.009)	0.906 (0.014)	0.775 (0.013)	0.767 (0.006)	0.819 (0.013)
Normalized RGB	<b>0.807</b> (0.013)	<b>0.857</b> (0.011)	0.928 (0.005)	<b>0.784</b> (0.014)	0.792 (0.010)	<b>0.848</b> (0.006)
Stain channels	0.785 (0.015)	0.851 (0.008)	<b>0.933</b> (0.009)	0.778 (0.015)	<b>0.798</b> (0.009)	0.842 (0.015)

**CBCS data set**: 1713 patient samples, 4 cores/patient Linear SVM, 5-fold cross-validation

Stain normalization improves results

Deep transfer learning works on non-RGB images

Dictionary learning slightly better than deep transfer learning

# Task-driven Dictionary Learning

### **Classification Accuracy**

	Unsupervised	Task-driven
Patch-level	0.507	0.520
Patient-level Mean of patch probabilities Sum of log of patch probabilities Linear SVM on histogram of features	0.646 0.729 0.698	0.642 0.664 0.713

#### SPECS data set:

43 Basal, 42 Luminal A 2 cores/patient 5-fold cross-validation

Dictionary size 128 Patch size 9x9 Task-driven extension successful for patch-level accuracy

Less clear for patient-level accuracy

Patient-level labels are weak when applied to small patches

# Outline



### **Multiple Instance Learning**



### Method Overview



# Method Overview



# Multiple Instance Terminology

### **Standard assumption:**

- Negative bag: all instances negative
- Positive bag: one or more instances positive





Good for diagnosis

Classes not treated symmetrically

# Multiple Instance Terminology

### **Standard assumption:**

- Negative bag: all instances negative
- Positive bag: one or more instances positive





Good for diagnosis

Classes not treated symmetrically

### Alternative assumption: majority vote



Remove assumption: learn how to aggregate probabilities





# Method 1: Quantile Aggregation



# Method 1: Quantile Aggregation



# Method 1: Quantile Aggregation



### Method 1: Quantile Aggregation Prediction



### Method 2: Iterative MI with Majority Vote



1) Initialize instance labels



1) Initialize instance labels



1) Initialize instance labels




### Method 2: Iterative MI with Majority Vote Learn Instance Labels



## Method Overview



### Method 3: Fine-tune CNN



## Method Comparison

#### **Classification Accuracy**

Method	Basal vs. Non-Basal	ER Status	Grade 1 vs. 3
AlexNet			
Baseline: Majority vote	0.776	0.772	0.853
Method 1: Quantile aggregation	0.799	0.815	0.876
Method 2: Iterative MI with majority vote	0.788	0.807	0.870
Method 3: Fine-tune CNN	0.831	0.841	0.954
VGG16			
Baseline: Majority vote	0.807	0.823	0.897
Method 1: Quantile aggregation	0.824	0.853	0.908
Method 2: Iterative MI with majority vote	0.812	0.846	0.905
Method 3: Fine-tune CNN	0.833	0.879	0.973

**CBCS data set:** 1713 patient samples, 4 cores/patient 5x random split: <sup>1</sup>/<sub>2</sub> train, <sup>1</sup>/<sub>2</sub> test

MI techniques for training always beneficial

Fine-tuning CNN gives largest improvement

### Method 3: Fine-tune CNN MI Learning



### Method 3: Fine-tune CNN MI Learning



### Method 1: Quantile Aggregation Heterogeneity



### Method 3: Fine-tune CNN Heterogeneity



### Method 3: Fine-tune CNN Heterogeneity





## **Future Work**

- Validation of heterogeneity
- Outcome prediction
- Model visualization and interpretation
- Other cancer and disease types

## Outline



### **Tissue Microarray**



### **Discriminative Common Space**



Extract shared components of data

Improve discriminability

### Standard Solution: Canonical Correlation Analysis



#### **Challenges:**

Features may not be discriminative

Not robust to small training set size or high dimensional, low sample size (HDLSS) data

#### My solution:

Add task-driven component to a deep variant of CCA

## **Background: Canonical Correlation Analysis**



## **Background: Canonical Correlation Analysis**



Solved with SVD

### **Background: Deep Canonical Correlation Analysis**



### Deep CCA: Correlation vs. Accuracy



### Deep CCA Challenges







Compute linear CCA projections from  $A_1$  and  $A_2$  after DNN optimization

### Task-driven Deep CCA

Goal: add task-driven objective (e.g., for classification)



### Task-driven Deep CCA

Goal: add task-driven objective (e.g., for classification)



Linear CCA:

- 1) Maximize sum correlation
- 2) Such that projections are orthogonal

### Solution 1: Eigendecomposition

### DDCCA-ED



## Solution 2: Whitening



## Solution 2: Whitening



### Solution 3: Soft Decorrelation

### **DDCCA-SD**



Maximize sum correlation
Such that projections are orthogonal

Encourage orthogonality using regularization

aka DeCov (Cogswell, 2016) or Soft Decorrelation Loss (Chang, 2018)

$$\mathscr{L}_{Decorr}(\Sigma) = \sum_{i=1}^{d_o} \sum_{j \neq i}^{d_o} |\Sigma_{i,j}|$$
$$\Sigma_1 = \frac{1}{n-1} A_1 A_1^T \qquad \Sigma_2 = \frac{1}{n-1} A_2 A_2^T$$

Penalize off-diagonal elements of covariance matrix

### Solution 4: No Explicit Decorrelation

### **DDCCA-ND**



Maximize sum correlation
Such that projections are orthogonal

Rely on task objective to decorrelate as needed

## **Benefits of All DDCCA Models**



## **Experiments: MNIST Split**

#### **Cross-modal classification**

![](_page_65_Figure_2.jpeg)

Robust to small training set size and HDLSS data

### MNIST Split: Visualization with t-SNE

![](_page_66_Figure_1.jpeg)

Digits are better clustered with task-driven method

8

## Carolina Breast Cancer Study: Cross-modal

#### Train PAM50, Test Image

# Cross-modal classification accuracy

1003 patients

#### **Image features:**

VGG16 – output of  $4^{th}$  set of conv layers + mean pool  $\rightarrow$  512 D

Method	Basal	ER	Grade
CCA	0.732 (0.010)	0.637 (0.008)	0.741 (0.005)
RCCA	0.815 (0.008)	0.811 (0.003)	0.877 (0.010)
PLS-SVD	0.650 (0.016)	0.656 (0.003)	0.797 (0.010)
DCCA	0.787 (0.010)	0.785 (0.011)	0.867 (0.012)
SoftCCA	0.780 (0.010)	0.769 (0.014)	0.848 (0.015)
DDCCA-ED	0.802 (0.015)	0.803 (0.029)	0.852 (0.011)
DDCCA-W	<b>0.820</b> (0.008)	<b>0.828</b> (0.006)	<b>0.917</b> (0.019)
DDCCA-SD	0.796 (0.004)	0.811 (0.004)	0.874 (0.019)
DDCCA-ND	0.766 (0.013)	0.805 (0.007)	0.878 (0.011)

#### **Train Image, Test PAM50**

Method	Basal	ER	Grade		
CCA	0.943 (0.005)	0.869 (0.006)	0.828 (0.011)		
RCCA	0.978 (0.003)	0.905 (0.008)	0.836 (0.009)		
PLS-SVD	0.912 (0.018)	0.865 (0.005)	0.794 (0.012)		
DCCA	0.976 (0.005)	0.874 (0.010)	0.854 (0.015)		
SoftCCA	0.978 (0.004)	0.897 (0.010)	0.843 (0.007)		
DDCCA-ED	0.971 (0.009)	0.889 (0.010)	0.833 (0.010)		
DDCCA-W	<b>0.983</b> (0.006)	<b>0.908</b> (0.009)	0.845 (0.005)		
DDCCA-SD	0.978 (0.005)	<b>0.908</b> (0.009)	<b>0.848</b> (0.004)		
DDCCA-ND	0.975 (0.004)	0.896 (0.005)	0.817 (0.005)		

#### **Genomic features:**

PAM50 – expression for 50 genes

**Train/validation/test**: Random split <sup>1</sup>/<sub>2</sub>, <sup>1</sup>/<sub>4</sub>, <sup>1</sup>/<sub>4</sub>

## Carolina Breast Cancer Study: Regularization

Training:

 $OO\cdots OOB_{2}$ 

 $\begin{bmatrix} 00 \cdots 00 \end{bmatrix} A_2$ 

Whiten

00...00

(training only)

Corr

 $\cdot 00$ 

00

OO

Whiten

 $B_1 \bigcirc \bigcirc$ 

 $A_1 \bigcirc$ 

OO

![](_page_68_Picture_2.jpeg)

Testing:

![](_page_68_Picture_4.jpeg)

	Method	Training data	Basal	ER	Grade
	Linear SVM	Image only	0.785 (0.004)	0.838 (0.003)	0.897 (0.006)
	DNN	Image only	0.796 (0.007)	<b>0.852</b> (0.008)	0.907 (0.009)
	DDCCA-W	Image+PAM50	0.827 (0.006)	0.839 (0.007)	<b>0.911</b> (0.012)
	DDCCA-SD	Image+PAM50	0.820 (0.010)	0.826 (0.009)	0.859 (0.021)
	DDCCA-W	Image+GE163	<b>0.840</b> (0.010)	0.838 (0.009)	0.910 (0.017)
	DDCCA-SD	Image+GE163	0.812 (0.011)	0.815 (0.020)	0.891 (0.020)
DNN f		-			

#### **Genomic features:**

PAM50 – expression for 50 genes GE163 – larger set of 163 genes

Large improvement for some tasks

## Future Work

- Shared and individual representations
- Fine-tune CNN as one modality
- Data augmentation

## **Contributions to Computer Science**

- Discriminative representations for histology images using dictionary learning or deep transfer learning
- 2) Multiple instance learning methods for handling large, heterogeneous images with an SVM on any type of feature set or with a CNN for end-to-end training
- 3) A set of **multimodal** methods to find a shared space that is also **discriminative**
- 4) Techniques for **deep learning** on problems traditionally viewed as **"small data"**

### Acknowledgements

Advisor: Marc Niethammer

**Other committee members**: Steve Marron, Chuck Perou, Steve Pizer, Alex Berg

**Collaborators**: Melissa Troester, Lindsay Williams, Joseph Geradts, David Eberhard, Nancy Thomas, Susan Wei, Jayson Miedema

**Funding**: Royster Society, Lineberger Comprehensive Cancer Center, CISMM (NIH)
### **Questions?**



2) Handling heterogeneous images 3) Combining imaging & genomics

### Extra slides

# Publications

- **ISBI 2015**: Hierarchical task-driven feature learning for tumor histology
- **MICCAI 2018**: Multiple instance learning for heterogeneous images: training a CNN for histopathology
- **npj Breast Cancer 2018**: Image analysis with deep learning to predict breast cancer grade, ER status, histologic subtype, and intrinsic subtype

# Method 1: Task-driven Dictionary Learning



## **Task-Driven Dictionary Learning**

Simplification: classify individual patches

Train classifier by minimizing expected loss Optimize over classifier and dictionary (Mairal, 2012) 1 1 5 4 3 7 5 3 5 3 5 5 9 0 6 3 5 2 0 0 applied to small images e.g., MNIST

separating

hyperplane

(classifier)

$$f(y, x, \alpha) = E[loss(y, w, \alpha(x, D))] + \frac{v}{2} ||w||_{2}^{2}$$
  
loss labels sparse encodings  
e.g., logistic classifier regularization  
$$log(1 + e^{-yw^{T}\alpha(x, D)})$$

### **Task-Driven Dictionary Learning**

$$f(y, x, \alpha) = E[loss(\underline{y}, \underline{w}, \alpha(\underline{x}, \underline{D}))] + \frac{\underline{v}}{2} ||w||_{2}^{2}$$
  
labels / / dictionary  
classifier patches  
sparse encodings

#### **Solution**: stochastic gradient descent

Initialize dictionary D with unsupervised learning

Initialize classifier w with logistic regression on set of patch encodings Repeat until convergence:

Select set of labels y and patches x from training set

Compute sparse encodings  $\alpha$ 

Update classifier w using  $\nabla_{w} f(y,x,\alpha)$ 

Update dictionary D using  $\nabla_{D} f(y,x,\alpha)$ 

# Method 1: Pre-trained CNN Sample Weighting

#### Sample weighting for unbalanced data



#### Sample weighting by grade



### Method 1: Quantile Aggregation Classification Accuracy

	Sensitivity (%)	Specificity (%)	Accuracy (%)	Карра
Grade low-int vs high			82	0.64
ER status	84	72	84	0.64
Basal vs non-Basal	78	73	77	0.47
ROR-PT low-med vs high risk	79	74	76	0.47
Histologic subtype ductal vs lobular	71	96	94	0.66

CBCS data set: 859 patient samples, 4 cores/patient Random division into 2/3 training, 1/3 test

### Method 1: Quantile Aggregation Inter-rater Agreement

#### Grade:

	Accuracy	Карра
Image analysis	82%	0.64
Agreement between two pathologists	89%	0.78
Reported by other groups		0.6-0.7, 0.5

#### **ER Status Kappa:**

Image analysis	0.64
Centralized pathology and SEER classifications	0.70
Different IHC antibodies	0.6-0.8
Medical records and staining of tissues	0.62

#### **PAM50 Subtype Accuracy:**

Image analysis	77%
IHC vs. RNA-based for Basal	90%
IHC vs. RNA-based for Luminal A	77%

### Method 3: Fine-tune CNN Classification Accuracy (Multi-class)

	Mean	Quantile
Histologic subtype ductal vs lobular	0.931 (0.004)	0.952 (0.003)
ER status	0.833 (0.008)	0.841 (0.006)
Grade	0.680 (0.003)	0.676 (0.006)
ROR-PT	0.595 (0.003)	0.582 (0.008)
Intrinsic subtype	0.548 (0.006)	0.544 (0.003)

# **Background: Canonical Correlation Analysis**

**Given**: input data  $X_1 \in \Re^{d_1 \times n}$ ,  $X_2 \in \Re^{d_2 \times n}$  (mean centered) Compute covariance matrices:

$$\Sigma_1 = \frac{1}{n-1} X_1 X_1^T \qquad \Sigma_2 = \frac{1}{n-1} X_2 X_2^T \qquad \Sigma_{12} = \frac{1}{n-1} X_1 X_2^T$$

**Goal**: find projections  $w_1$  and  $w_2$  such that the data are maximally correlated

$$\underset{w_{1,w_{2}}}{\operatorname{argmax}}\operatorname{corr}(w_{1}^{T}X_{1,}w_{2}^{T}X_{2}) = \underset{w_{1,w_{2}}}{\operatorname{argmax}} \frac{w_{1}^{T}\Sigma_{12}w_{2}}{\sqrt{w_{1}^{T}\Sigma_{1}w_{1}w_{2}^{T}\Sigma_{2}w_{2}}}$$

Alternatively, constrain projections to have unit variance:

maximize:  $w_1^T \Sigma_{12} w_2$  covariance subject to:  $w_1^T \Sigma_1 w_1 = w_2^T \Sigma_2 w_2 = 1$  unit variance

# **Background: Canonical Correlation Analysis**

Find **multiple pairs**  $(w_1^{(i)}, w_2^{(i)})$  such that  $(w_1^{(i)})^T \Sigma_1 w_1^{(j)} = (w_2^{(i)})^T \Sigma_2 w_2^{(j)} = 0$  for  $i \neq j$ orthogonal projections

Let 
$$W_1 = [W_1^{(1)}, ..., W_1^{(k)}]$$
 and  $W_2 = [W_2^{(1)}, ..., W_2^{(k)}]$   
maximize:  $tr(W_1^T \Sigma_{12} W_2)$   
subject to:  $W_1^T \Sigma_1 W_1 = W_2^T \Sigma_2 W_2 = I$ 

#### Solution:

Let  $T = \Sigma_1^{-1/2} \Sigma_{12} \Sigma_2^{-1/2}$ SVD  $T = U_1 \operatorname{diag}(\sigma) U_2^T$ 

Compute  $W_1$  and  $W_2$  from top k singular values of T:

$$W_1 = \Sigma_1^{-1/2} U_1^{(1:k)} \qquad W_2 = \Sigma_2^{-1/2} U_2^{(1:k)}$$

### Background: Deep CCA

Andrew et al., ICML, 2013

**Given**:  $A_1 = f_1(X_1, \Theta_1)$ ,  $A_2 = f_2(X_2, \Theta_2)$ ,  $A_1 \in \Re^{d_o xn}$ ,  $A_2 \in \Re^{d_o xn}$ Compute covariance matrices:



Compute linear CCA after DNN optimization complete

# Task-driven Deep CCA

**Problem**: maximizing sum correlation does not always result in improved cross-modal classification accuracy

Solution: add task-driven objective (e.g., for classification)

$$\mathcal{L}_{\text{task}}$$
 +  $\lambda$   $\mathcal{L}_{\text{CCA}}$ 

**CCA objective:**  

$$argmax_{W_1,W_2,\Theta_1,\Theta_2} \underbrace{\operatorname{tr}(W_1^T \Sigma_{12} W_2)}_{W_1,W_2,\Theta_1,\Theta_2} equivalent to argmin_{W_1,W_2,\Theta_1,\Theta_2} \underbrace{||W_1^T A_1 - W_2^T A_2||_F^2}_{\ell_2 \text{ distance}}$$
when  $\underbrace{(w_1^{(i)})^T \Sigma_1 w_1^{(i)} = (w_2^{(i)})^T \Sigma_2 w_2^{(i)} = 1}_{Unit \text{ variance}} \text{ for all i}$ 

$$unit \text{ variance}$$

Challenge: need to compute CCA projection in network

### **Batch Normalization**

**Input:** Values of x over a mini-batch:  $\mathcal{B} = \{x_{1...m}\};$ Parameters to be learned:  $\gamma$ ,  $\beta$ **Output:**  $\{y_i = BN_{\gamma,\beta}(x_i)\}$  $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i$  $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2$ // mini-batch mean // mini-batch variance  $\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize  $y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathbf{BN}_{\gamma,\beta}(x_i)$ // scale and shift

# **ZCA** Whitening



### **Batch Size**



### Visualization with t-SNE

