

# GemIdent: An Interactive Statistical Image Segmentation System

---

Adam Kapelner\* & Susan Holmes\*\*

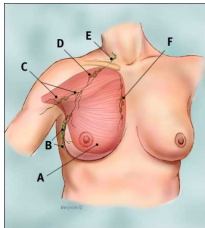
\*Department of Statistics  
The Wharton School, University of Pennsylvania

\*\*Department of Statistics  
Stanford University

---

# Breast Cancer and Lymphnodes

The breast drains lymph fluid to surrounding lymph nodes.

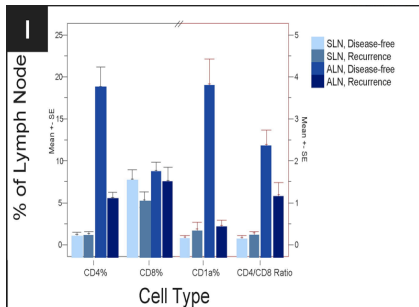


In breast cancer patients, cancer can metastasize in the lymph nodes.

This is a major prognostic indicator — most patients do not do well after lymphnode metastasis.

# Criminals in Police Headquarters?

How is it possible that cancer can set up shop in lymph nodes, the immune cell centers?

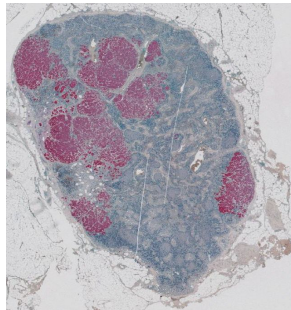


How is it possible that some patients can weather the storm and some cannot? Kohrt et al, 2005 showed that certain populations of immune cells is correlated with survival.

But what is happening biologically? Maybe the answer lies in how the cancer and the immune cells are *spatially distributed*.

## A typical lymphnode

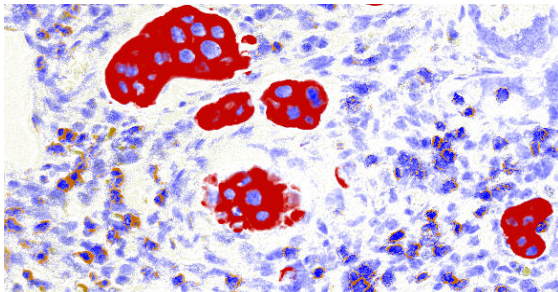
Here is a typical lymph node measuring about 1 cm in diameter. It was surgically removed and sliced into sections 3  $\mu\text{m}$  thick. Cell populations of interest have been stained with different colors using immunohistochemistry.



In the above image, cancer is **red**, and CD4 T-cells (a type of immune cell) are **brown**, and all nuclei are “counterstained” **blue**.

## Under the microscope

The entire node has been imaged using automated microscopy at 200X creating a few thousand subimages, for example:



If the goal is to find spatial distributions of cells, we must find the location of every cell nucleus in the lymphnode and determine the type of cell, *i.e.* its “phenotype”.

# The basic classification problem

Finding cells in images can be thought of as a classification problem for *each and every* pixel  $i, j$ :

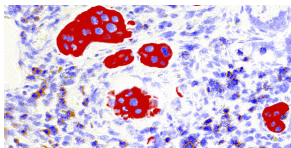
$$\hat{y}_{ij} = f(\mathbf{x}_{ij}), \quad \hat{y}_{ij} \in \{0, 1, 2, \dots, P\}$$

Where:

- $\mathbf{x}_{ij}$  is a vector of features / covariates
- “ $f$ ” is a classification machine
- $\hat{y}_{ij}$  represents the machine’s best guess for the pixel’s phenotype.  
“0” represents not belonging to any phenotype.

# Covariates?

Which features are important *i.e.* how do you build each  $\mathbf{x}_{ij}$ ?



My naive thinking was:

- I stain (color) information is obviously important
- II rotational invariance
- III the presence of stains at certain distances from the pixel of interest is key

## Color separation via Mahalanobis distance

Images are in RGB format which means each pixel has an 8-bit red value, 8-bit green value, and 8-bit blue value.

The image would be more useful if decomposed into  $\{0, 1, \dots, S\}$  relevant stains (where 0 represents the background, usually white). How to do this?

After getting a few examples of the color, use the Mahalanobis distance function:

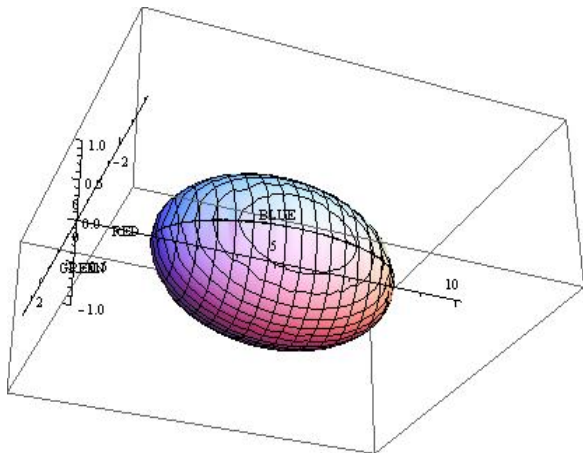
$$d(\mathbf{t}_{ij}) = \sqrt{(\mathbf{t}_{ij} - \boldsymbol{\mu}_{\text{blue}})^T \hat{\boldsymbol{\Sigma}}_{\text{blue}}^{-1} (\mathbf{t}_{ij} - \boldsymbol{\mu}_{\text{blue}})}$$

where  $\mathbf{t}_{ij}$  is the  $3 \times 1$  RGB-color vector at the pixel of interest,  $\boldsymbol{\mu}_{\text{blue}}$  is the mean RGB-color vector for the “blue” stain, and  $\hat{\boldsymbol{\Sigma}}_{\text{blue}}$  is the  $3 \times 3$  sample covariance matrix of all the examples of “blue” pixels.



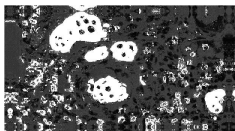
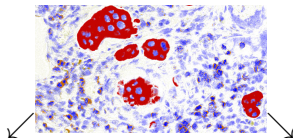
## Color separation via Mahalanobis distance

The distance function “ $d$ ” assigns large values to pixels with colors far away from the “blue” stain, and small values to pixels with colors close to the “blue:”



## Color Decomposition

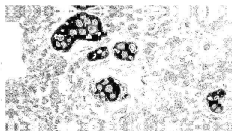
Thereby, our image can be decomposed into score matrices for each stain (as well as the background) which looks sort of like this:



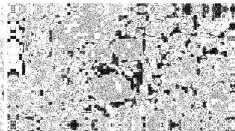
Stain: Blue



Stain: Brown



Stain: Red



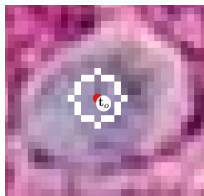
Stain: Background

We denote each of these score matrices  $F_s$ .

## Ring Scores

For each pixel, we know that stain information away from the pixel itself is important, and we know that the images are rotationally invariant.

A natural choice would be to create “ring scores” which reflect “how much” of a certain stain exists a certain radius away from the pixel:



For each stain, rings  $\{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_R\}$  (where  $R$  is the maximum radius) are considered important.

## Feature computation and the feature vector

The score for stain  $s$  at radius  $r$  can be computed as follows by just adding up the values:

$$\ell_{s,r}(\mathbf{t}_o) = \sum_{\mathbf{t} \in \mathbf{c}_r} F_s(\mathbf{t}_o + \mathbf{t})$$

Now the feature vector for pixel  $i, j$  can be created by concatenating all the ring scores for each stain at each radius:

$$\mathbf{x}_{ij} = [\ell_{1,1}, \ell_{1,2}, \dots, \ell_{1,R}, \ell_{2,1}, \ell_{2,2}, \dots, \ell_{2,R}, \dots, \ell_{S,1}, \ell_{S,2}, \dots, \ell_{S,R}]$$

Note: somewhat naive — no interactions, no transformations

# Supervised Statistical Learning

What would make sense would be to give a few examples of each phenotype, then create a machine that would make future predictions:

$$\underbrace{\begin{bmatrix} \text{— } \mathbf{x}_{ij} \text{ —} \\ \text{— } \mathbf{x}_{ij} \text{ —} \\ \text{— } \mathbf{x}_{ij} \text{ —} \\ \text{— } \mathbf{x}_{ij} \text{ —} \\ \vdots \\ \text{— } \mathbf{x}_{ij} \text{ —} \\ \text{— } \mathbf{x}_{ij} \text{ —} \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ \vdots \\ P \\ P \end{bmatrix}}_{\text{"Training Data"}} \xrightarrow{\text{Machine Creation}} \boxed{f}$$

What machines can you create that perform  $\hat{y}_{ij} = f(\mathbf{x}_{ij})$ ?

# Regression & Machine Learning

**Regression** *e.g.* OLS (a la Stat 102)

- Effect Tests
- p-values
- Confidence Intervals
- Prediction

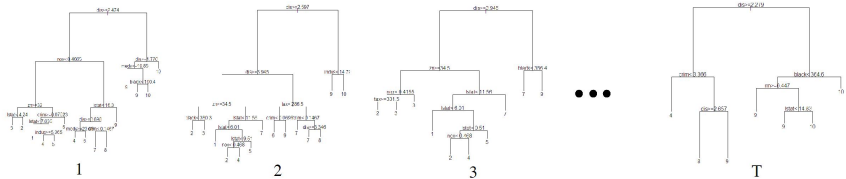
## Machine Learning

- ~~Effect Tests~~
- ~~p-values~~
- ~~Confidence Intervals~~
- Prediction

As long as we only care about prediction. . .

# Random Forests

For pixel classification, we chose the “machine” with the lowest error rate. At the time of conception (2006), Random Forests beat out the contenders for our data.



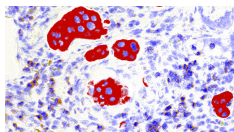
**Random Forests Algorithm - Decision trees (i.e. CARTs) with a few twists:**

- 1 Bootstrap training data
- 2 Select a subset of features to evaluate at each node
- 3 Repeat 1-2 to build  $T$  trees

To predict, let trees vote democratically

## Classification

With 5 training points for cancer, and 50 NON points, we build a random forest, then proceed to classify *each and every* pixel of the example image:

 $f \rightarrow$ 

These blobs are then post-processed using a simple erosion algorithm to yield centroids *i.e.* coordinates of each cell.



# Introduction to the software

This algorithm, coupled with a graphical user interface, was developed in Java into a program called GemIdent and open-sourced under GPL.

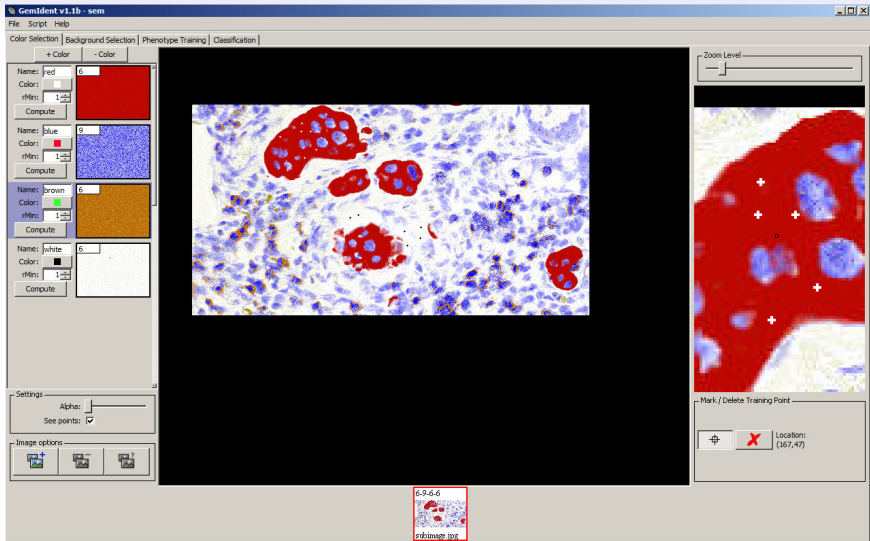
The software has the following main features organized into panels:

- Color Selection
- Phenotype Training
- Classification
- Data Analysis (not shown)



[www.gemident.com](http://www.gemident.com)

# Color Selection



# Phenotype Training

**GemIdent v1.1b - sem**

File Script Help

Color Selection | Background Selection | **Phenotype Training** | Classification

+ Phen - Phen view points

Name: NON 50  
Color:   
rMin: 1  
rMax: 1

Name: cancer 10  
Color:   
rMin: 1  
rMax: 9

Pix/Cent: ☒ ☒

Name: T\_cells 20  
Color:   
rMin: 1  
rMax: 8

Pix/Cent: ☒ ☒

Name: backg 21  
Color:   
rMin: 1  
rMax: 8

Div 1/2: ☒ ☒

Settings

Visualize: ☐

See points: ☒

Confusion Blinker: ☐

Confusion Slider:

See Type 1 Errors: ☐

Image options

50-10-20-21  
gemident  
subimage.jpg

Zoom Level:

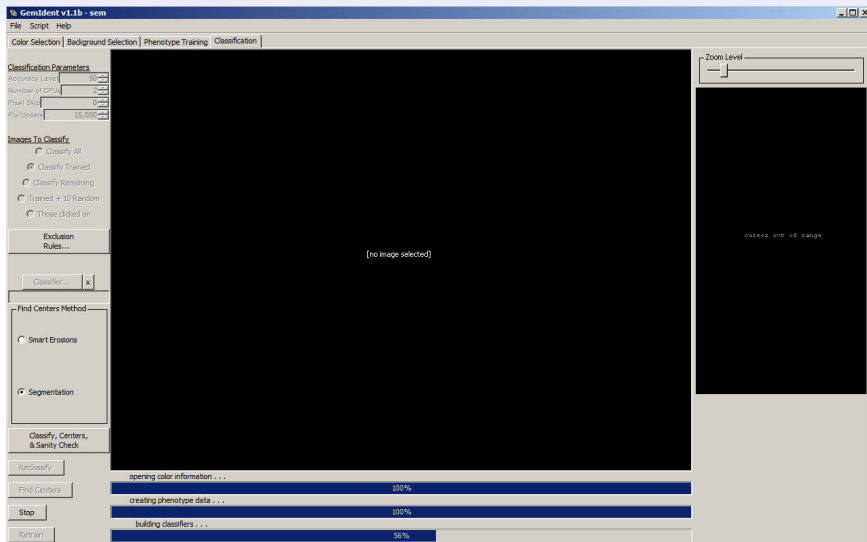
Mark / Delete Training Point

Location: (67, 178)

Boost on classified images

<none>  
Image 0/0

# Classification



# Phenotype Retraining

The screenshot displays the GemIdent v1.1b - sem software interface. The main window is titled "Phenotype Retraining" and contains a large image of a histological section with several clusters of cells. The interface includes a sidebar on the left with a list of phenotypes: "NGN" (50), "cancer" (10), "T\_cells" (20), and "backg" (21). Each phenotype has a color selection box and a small thumbnail image. The "cancer" phenotype is highlighted in red, and the "T\_cells" phenotype is highlighted in green. The "backg" phenotype is highlighted in blue. The "NGN" phenotype is highlighted in black. The "cancer" phenotype is also highlighted in red in the main image. The "T\_cells" phenotype is also highlighted in green in the main image. The "backg" phenotype is also highlighted in blue in the main image. The "NGN" phenotype is also highlighted in black in the main image. The main image shows a histological section with several clusters of cells. The "cancer" phenotype is highlighted in red, and the "T\_cells" phenotype is highlighted in green. The "backg" phenotype is highlighted in blue. The "NGN" phenotype is highlighted in black. The "cancer" phenotype is also highlighted in red in the main image. The "T\_cells" phenotype is also highlighted in green in the main image. The "backg" phenotype is also highlighted in blue in the main image. The "NGN" phenotype is also highlighted in black in the main image. The interface also includes a "Settings" section with options for "Visualize", "See points", "Confusion Blinker", and "Confusion Slider". There is also a "Image options" section with buttons for "Add", "Remove", and "Reset". A "Classification Overlay Adjustment for Retraining (subimage.jpg)" dialog box is open, showing sliders for "cancer identified pixels visibility", "cancer centroids visibility", "backg identified pixels visibility", "backg centroids visibility", "T\_cells identified pixels visibility", and "T\_cells centroids visibility". A "Mark / Delete Training Point" section is also visible, with a "Location" field showing (64, 175) and a "Boost on classified images" section with buttons for "Previous", "Next", and "None".

GemIdent v1.1b - sem

File Script Help

Color Selection Background Selection Phenotype Training Classification

+ Phen - Phen view points

Name: NGN 50  
Color:   
rMin: 1  
rMax: 1

Name: cancer 10  
Color:   
rMin: 1  
rMax: 9

Phx/Cent: ☒ ☒

Name: T\_cells 20  
Color:   
rMin: 1  
rMax: 8

Phx/Cent: ☒ ☒

Name: backg 21  
Color:   
rMin: 1  
rMax: 8

Phx/Cent: ☒ ☒

Settings

Visualize: ☐  
See points: ☒  
Confusion Blinker: ☐  
Confusion Slider:   
See Type 1 Errors: ☐

Image options

Classification Overlay Adjustment for Retraining (subimage.jpg)

cancer identified pixels visibility (blink: ☐)  
cancer centroids visibility (blink: ☐)

backg identified pixels visibility (blink: ☐)  
backg centroids visibility (blink: ☐)

T\_cells identified pixels visibility (blink: ☐)  
T\_cells centroids visibility (blink: ☐)

Mark / Delete Training Point

Location: (64, 175)

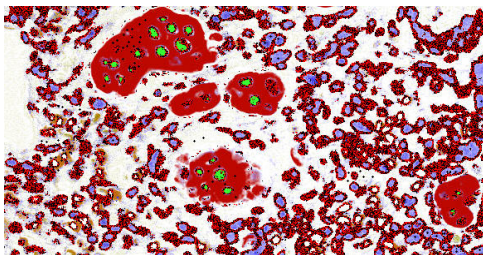
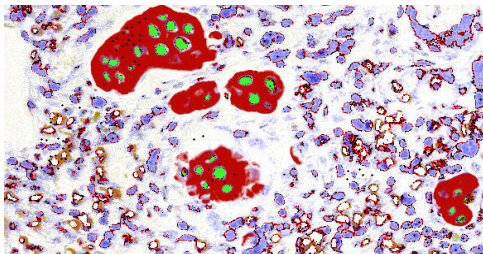
Boost on classified images

<none> Image 0/0

50-10-20-21  
subimage.jpg

## “Confusion” Measure

The ultimate in interactive boosting...when trees fight?



# The Output

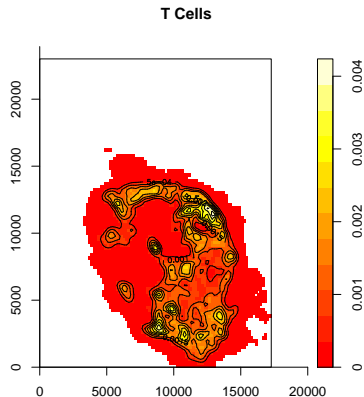
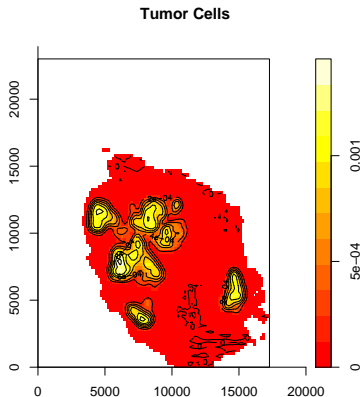
A big CSV file with every coordinate for every cell type across every subimage:

```
filename,locX,locY,globalX,globalY  
stage_0147,63,49,6416,10263  
stage_0147,92,122,6445,10336  
stage_0147,95,243,6448,10457  
stage_0147,109,243,6462,10457  
stage_0147,247,719,6600,10933  
stage_0147,276,757,6629,10971  
stage_0147,389,830,6742,11044  
stage_0147,397,394,6750,10608  
stage_0147,407,820,6760,11034
```

There are  $10^5$  to  $10^6$  cells per lymphnode.

Spatial Analysis done in R with `spatstat`, `spdep`, `DCluster` packages.

# Density plots in R of a typical lymphnode

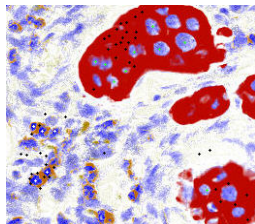




# Supervised learning

Supervised learning is . . . supervised.

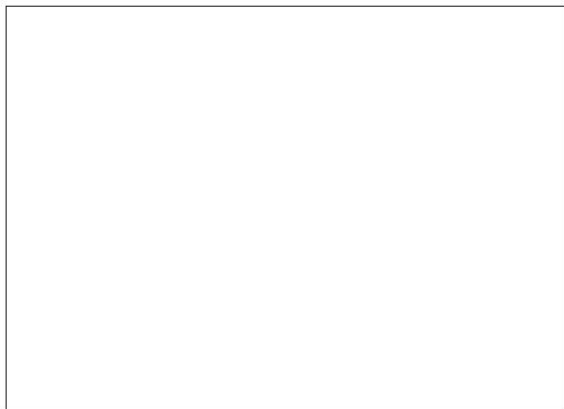
This means someone has to sit and click on examples of cells (as well as NON-phenotypes).



This is actually quite expensive. In fact, this is what users of GemIdent complain about the most.

Imagine we have 100's of lymphnodes requiring careful training. Could there be an easier way?

## Imagine a world...



where inside of a box, you can place any task, and have anyone in the world do it for a price you set...

# MTurk

The most popular crowdsourcing marketplace is MTurk where you can freely post tasks.

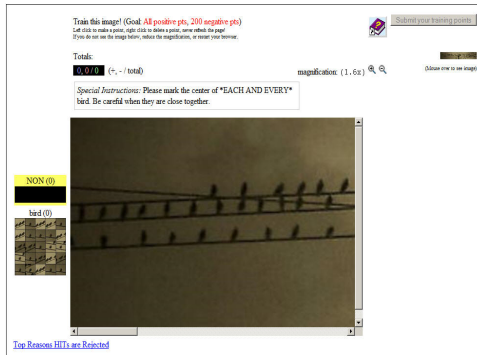
The screenshot shows the Amazon Mechanical Turk interface. At the top, there are navigation tabs: "Your Account", "HITS", and "Qualifications". A status bar indicates "89,569 HITS available now". Below the navigation, there's a search bar with filters: "All HITS", "HITS Available To You", and "HITS Assigned To You". The search criteria are set to "HITS" containing "that pay at least \$ 0.00 for which you are qualified". The results show "1-10 of 1627 Results". The list of HITs includes:

HIT Title	Requester	HIT Expiration Date	Time Allotted	Reward	HITS Available
Categorize the sentiment of a sentence	Amazon Requester, Inc.	Jan 23, 2011 (26 weeks 4 days)	30 minutes	\$0.01	5257
Choose the best category for this product	retaldata	May 16, 2010 (6 days 14 hours)	60 minutes	\$0.01	4233
Nebraska 3	Crowd Task	May 13, 2010 (1 day 11 hours)	20 hours	\$0.00	3529
Quick Survey: evaluate a short phrase	Amazon Requester, Inc.	May 14, 2010 (2 days 11 hours)	60 minutes	\$0.02	3415

Why not ask the world to click on cells for us for 10 cents an image?

# DistributeEyes

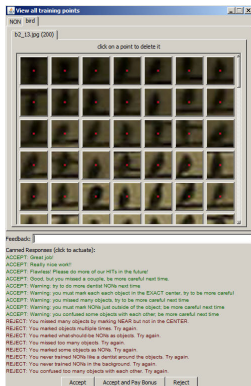
We created an add-on to GemIdent that takes images, and crowdsources the phenotype training step. Below is a DistributeEyes task that involves clicking on pigeons.



We can ensure the quality of the work by making workers watch an instructional video and passing a quiz.

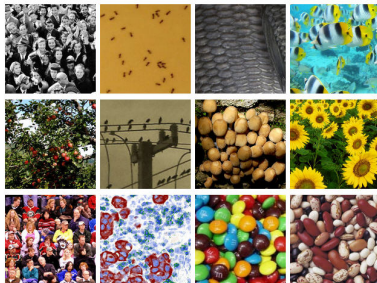
## Checking the training data

From within GemIdent, we can see the worker's points and choose to pay them or not. If they did a good job, we can even give them a bonus.



## Verifying the usefulness

We ran an experiment with 600 image-labeling tasks (of a variety of images) to check accuracy and find patterns in workers that do the best job.



# The hope

The results were very positive. Crowdsourcing can accelerate data collection.

We hope that this concept will be adopted by the medical research community as well as the statistical community because the price of supervised machine learning applications will drop.

# References



L. Breiman, "Random forests," Machine Learning, vol. 45, no. 1, pp. 5–32, 2001. [Online]. Available: [citeseer.ist.psu.edu/breiman01random.html](http://citeseer.ist.psu.edu/breiman01random.html)



Ihaka, R., and R. Gentleman (1996): "R: A Language for Data Analysis and Graphics," Journal of Computational and Graphical Statistics, 5(3), 299–314.  
[cran.r-project.org/](http://cran.r-project.org/)



S. Holmes ,A. Kapelner and P.P. Lee, An Interactive Java Statistical Image Segmentation System: GEMIDENT. Journal of Statistical Software, vol 30.



A. Kapelner, S. Holmes and P.P. Lee, "An Interactive Statistical Image Segmentation and Visualization System", medivis, pp. 81-86, International Conference on Medical Information Visualisation - BioMedical Visualisation (MediViz 2007), 2007.



H. E. Kohrt, N. Nouri, K. Nowels, D. Johnson, S. Holmes, and P. P. Lee, "Profile of immune cells in axillary lymph nodes predicts disease-free survival in breast cancer," PLoS Med, vol. 2(9), p. e284, 2005.



R. Tibshirani, T. Hastie and J. Friedman, The Elements of Statistical Learning. NY.: Springer, 2001.



# Acknowledgements

Peter P. Lee  
Holbrook Kohrt  
Francesca Setiadi  
Kyle Woodward

Funding from NIH/ NIGMS R01 and NSF-DMS.