

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Unsupervised Learning in PET Radiomics

### Permalink

<https://escholarship.org/uc/item/9n51s9w7>

### Authors

Liu, G  
Huang, S-Y  
Franc, B  
[et al.](#)

### Publication Date

2017-10-01

### DOI

10.1109/nssmic.2017.8532959

Peer reviewed



# HHS Public Access

Author manuscript

*IEEE Nucl Sci Symp Conf Rec (1997)*. Author manuscript; available in PMC 2019 January 08.

Published in final edited form as:

*IEEE Nucl Sci Symp Conf Rec (1997)*. 2017 October ; 2017: . doi:10.1109/NSSMIC.2017.8532959.

## Unsupervised Learning in PET Radiomics

**G. Liu,**

School of Computing, Florida Institute of Technology, Melbourne, FL

**S.-y. Huang,**

Radiology Department, University of California San Francisco

**B. Franc,**

Radiology Department, University of California San Francisco

**Y. Seo [Senior Member, IEEE],** and

Radiology Department, University of California San Francisco

**D. Mitra [Senior Member, IEEE]**

School of Computing, Florida Institute of Technology, Melbourne, FL

### Abstract

In this study, we investigated large scale radiomics on 116 breast cancer patients. We are particularly interested in unsupervised learning to bicluster patients and features in order to associate such biclusters with the disease characteristics. The results show that radiomics features with wavelet features have a better biclustering ability. And 172 radiomics features have shown a better classification capability.

### Keywords

Radiomics; Workflow; Unsupervised Clustering; PET; Breast Cancer

## I. Introduction

In this project, we have developed a pipeline for large scale radiomics analyses. We are particularly interested in unsupervised learning to bicluster patients and features in order to associate such biclusters with the disease characteristics. This is performed in order to understand usefulness of different features' sets against their capability to predict disease outcomes. We have run a pilot study with 116 breast cancer patients' PET datasets and present our findings here.

## II. Methods

One hundred sixteen breast cancer patient datasets were obtained for our study from patients' electronic health record and picture archiving and communication system (PACS). All of these patients underwent FDG-PET and DCE-MRI. The FDG-PET studies covered the whole-body while DCE-MRI covered the breast. Results presented here relates to PET studies only. The acquisition protocol for FDG-PET followed the established procedure guideline at UCSF. All data are anonymized for this study.

We studied nine disease outcome types: Progesterone Receptor (PR) status (108 total patients with available PR status) with values negative (44), positive (64); Estrogen Receptor (ER) status (108 total patients with available ER status) with values negative (37), positive (71); Human Epidermal growth factor Receptor 2 (HER2) status (107 total patients with available HER2 status) with values negative (77), positive (30); Triple negative status (108 total patients with available triple negative status) with value negative (87), positive (21); Recurrence (96 total patients with available recurrence numbers) with values 0 (78), 1 (18); Vital status (80 total patients with available vital status) with value dead (14), alive (66); Over all stage (93 total patients with available data) with values 0 (29), 1~2 (53), 3~4 (11); Tumor grade (103 total patients with available data) with values 0 (15), 1~2 (56), 3~4 (32); and Histology (114 total patients with available data) with values InvasiveLobular (6), InvasiveDuctal (95), Adenocarcinoma (5), MixedInvasiveDuctalandLobular (4) and Ductal carcinoma in situ DCIS (4), numbers in parentheses indicate number of patients in that group.

A major objective of the project is to build a script code pipeline for large scale data analyses task with as much generic capability and as less human intervention as possible. By “generic capability” we mean the pipeline should be adaptable easily for different tasks. Even though such a data analyses workflow uses mostly already written and established software, human intervention is almost always involved. However, for any large scale radiomics data analyses over thousands or, even millions of patients’ data, as contemplated by the community now [1, 2] is impossible with such human-in-the-loop semi-automated workflow architectures. Below we show a sample workflow for automation that we have nearly achieved.

We have extracted 709 radiomics features from 116 PET whole body images using their corresponding breast cancer masks. For this purpose we have used the PyRadiomics package [3]. Then we clustered patients for first 92 features (F92) and all 709 features (F709) [4]. The first 92 features include first order texture features [5, 6] shape features, GLCM features, GLRLM features and GLSZM features [4]. F709 includes all 92 features over 8 coiflet wavelet features, beside the original 92 features. The patients are then clustered [7] according to both F92 and F709 features’ sets.

Before biclustering patients and radiomics features matrix we normalize the feature values with the corresponding Z-scores across all patients for each feature, since different features’ values span over different ranges and are not comparable with each other. Furthermore, our biclustering algorithm [7] normalizes values across features as well. We use different k values for eliciting k by k number of biclusters of patients with our clustering algorithm in this study, where k is classes number in each outcome. For example, for binary classes outcome, triple negative status, we use 2 by 2 biclusters to cluster the patients radiomics feature matrix. The heatmap before and after 2 by 2 and 3 by 3 biclustering shows in figure 2. And we then compare how many disease outcome-type values are identified within each patient cluster. The absolute patients number of each bicluster in each outcome-type for F709 feature set shows in figure 3.

### III. Result

Biclusters on F92 do not show patient number difference in each bicluster. Therefore, we only show the normalized patient number for F709 in this report. This result indicates that F709 have better biclustering ability than F92.

For F709 feature set, all outcome types are significantly grouped under one of the biclusters. For two by two biclusters, cluster 2 shows better clustering ability than cluster 1 for all outcome types. In marjan histology analysis, all patients under adniocarcinoma, mixed invasive ductal and lobular and DCIS fall in cluster 2. The radiomics features numbers of two biclusters are 172 and 537. And the first 172 radiomics features show better classification capability on two patient biclusters in the heatmap.

For three by three biclusters, patients in 1~2 overall stage more likely fall in cluster 2. In tumor grade analysis, cluster 3 has a better distinguish ability on grade 3 than other grades. The radiomics features numbers of three biclusters are 191, 172 and 346. And the same 172 radiomics features show better classification capability on cluster 3 in the heatmap.

### IV. Discussion

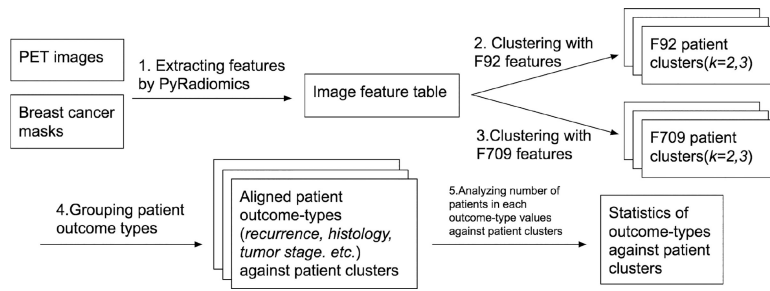
Our future work will be focused on (1) to understand the above mentioned radiomics type identified by F709 biclustering, and (2) to identify relevant features that provide better classification capability of known patient outcome-types. The combinatorics of such subsets of features is huge, e.g., for F709. Only a handful of such subsets are likely to provide us with good classification capability. On the other hand, many of the features are actually enhancing noise toward this objective. Identifying optimal feature set that provides the best biclustering will be a significant contribution to the field. We may need a supervised learning mode to identify the best set of features. Also, our results relate to only PET data now. We want to analyze MRI features using our pipeline as our immediate next stage. Eventually, the objective of radiomics is to develop reliable classifier that would be able to predict patient outcomes given a patient's radiology image. Our work is a step in that direction as well as in finding unknown radiomics types.

### V. Conclusion

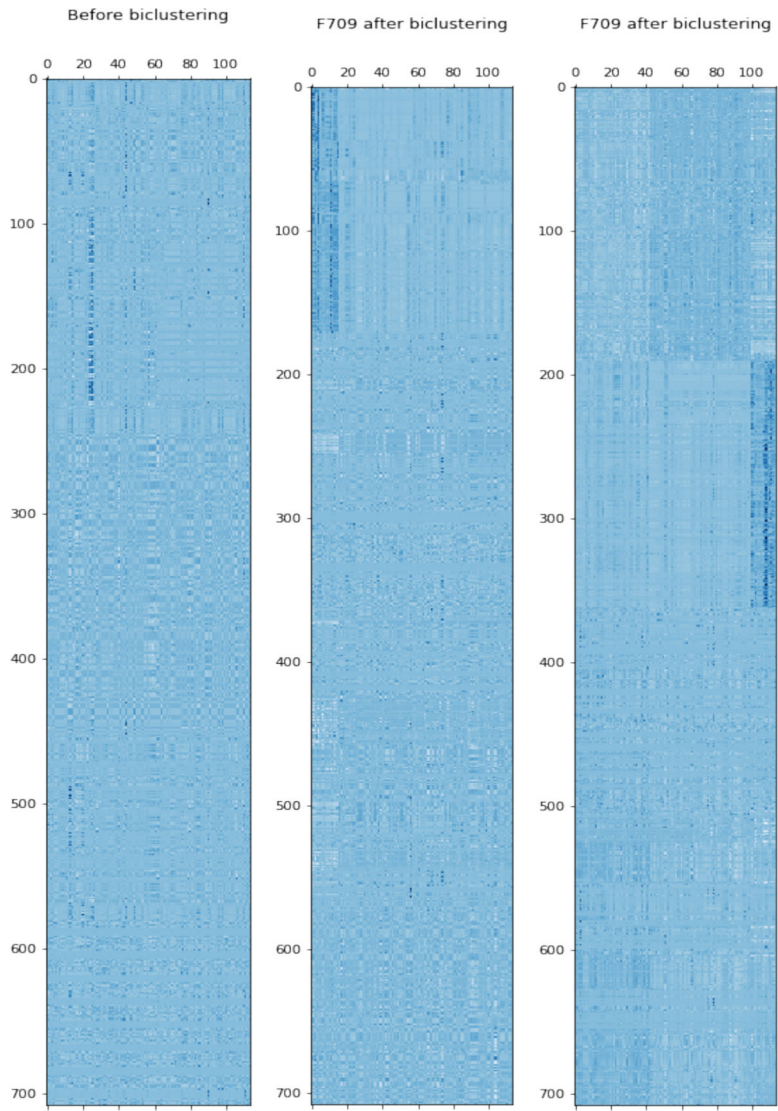
In this work, we have used two different sets of image features for biclustering 116 breast cancer patients and aligned the clusters against nine different types of disease outcomes: progesterone receptor (PR) status, estrogen receptor (ER) status, human epidermal growth factor receptor 2 (HER2) status, triple negative status, recurrence of tumors, vital status, tumor stages and histology results. We observe some predictability of these outcome types from our biclustering results and possibly an unknown radiomics characteristic for which we have no known disease outcome type. It is possible that we are observing some new topology in the data set that were not expressed with available disease outcome types [8].

## VI. References

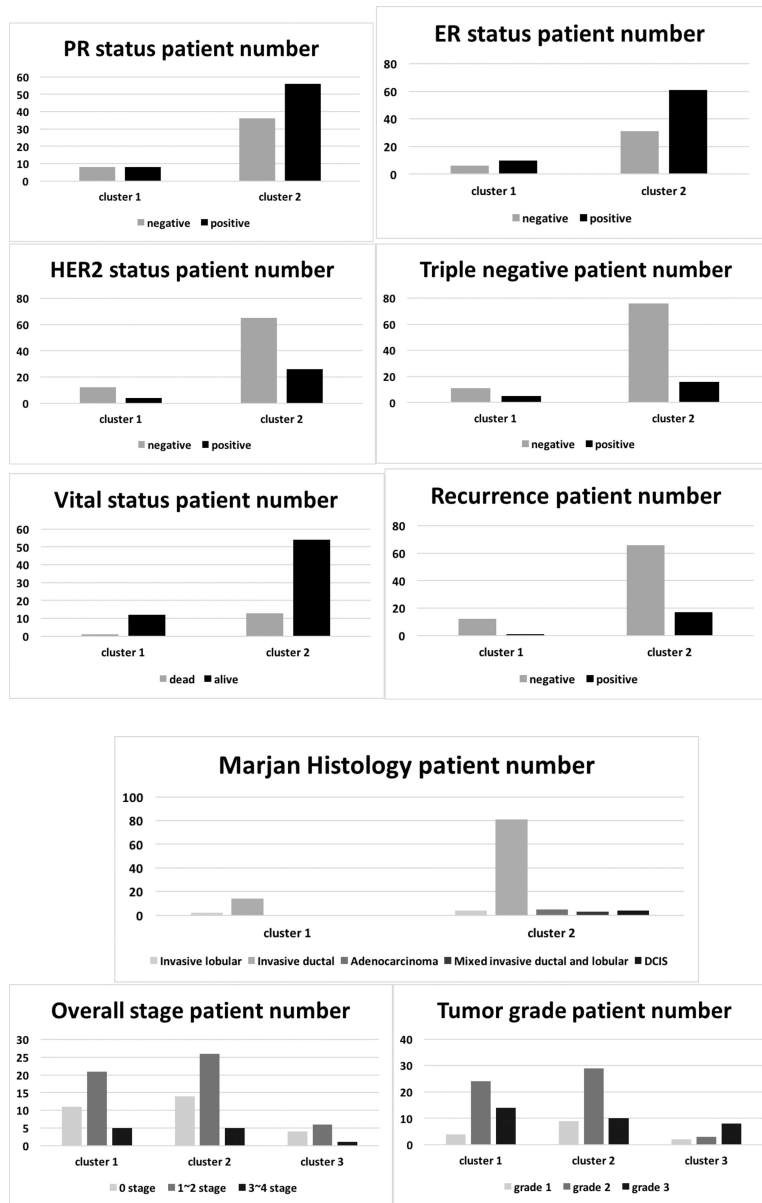
- [1]. Aerts HJWL, Velaquez ER, et al. “Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach,” *Nature Communications*/ 5:4006, 2014 [6]
- [2]. Cook GJR, Siddique M, et al. “Radiomics in PET: principles and applications,” *Clin. Transl. Imaging* 2:269–276, 2014.
- [3]. van Griethuysen JJM, Fedorov A, et al. “Computational Radiomics System to Decode the Radiographic Phenotype”; Submitted 2017.
- [4]. Zwanenburg A, Leger S, Vallières M and Löck S “Image biomarker standardisation initiative, feature definitions, version 1.2,” arXiv:1612.07003v2, 1 2017.
- [5]. Chicklore S, Goh V, et al. “Quantifying tumour heterogeneity in 18F-FDG PET/CT imaging by texture analysis.” *Eur J Nucl Med Mol Imaging* 40:133–140, 2011.
- [6]. Hatt M, Tixier F, et al. “Characterization of PET/CT images using texture analysis: the past, the present... any future?” *Eur J Nucl Med Mol Imaging*, 44(1):151–165, 2016. [PubMed: 27271051]
- [7]. Kluger Y, Basri R, et al. “Spectral biclustering of microarray data: coclustering genes and conditions.” *Genome research* 13.4: 703–716, 2003.
- [8]. Lum PY, Singh G “Extracting insights from the shape of complex data using topology,” *SCIENTIFIC REPORTS* 3: 1236, 2013. [PubMed: 23393618]



**Fig. 1.**  
Data processing pipeline.



**Fig 2.** Heatmaps before and after 2 by 2 and 3 by 3 biclustering



**Fig 3.** Patient number of each bicluster in each outcome-type for F709 feature set