

# Machine learning

# Mammographic Image Analysis

Team member: Chong Wei, Zheng Ni, Dan Kojis

# PART 1

## Motivation

---

Write something here  
Maybe some famous people' s  
word about mammography and  
machine learning

# Mammograms

- X ray images for breast cancer detection
- 1 in 8 U.S. women will get breast cancer
- 33 million mammograms annually in U.S.



# Motivation

- 1 in 5 screenings will not detect breast cancer if it is present (false-negative)
- About 10-20% of mammograms without breast cancer will be incorrectly diagnosed (false-positive)
- More effective diagnosis could reduce:
  - over diagnosis
  - overtreatment
  - fatality rates
  - healthcare costs
  - financial and emotional burden put on patients



# PART

# 2

## Data set introduction

---

Mammographic Image Analysis  
Society (MIAS) database

2

The screenshot shows a Jupyter Notebook interface for a dataset titled "MIAS Mammography: Looking for breast cancer". The dataset was created by Kevin Mader and updated a year ago (Version 3). The page includes tabs for Data, Overview (which is active), Kernels (8), Discussion (1), and Activity. A large preview image of mammogram scans is displayed, along with a sidebar showing 36 voters and a share button. Below the tabs, there are sections for Tags (with "health" selected) and a Description. A tooltip for the "health" tag is shown in Chinese, listing "电子战核(8)" (Electronic Warhead), "讨论(1)" (Discussion), and "活度" (Vitality). A link to "查看生词释义" (View glossary) is also present. The Content section contains a paragraph about the dataset being images and labels/annotations for mammography scans, with a note that more information can be found at [MIAS](#). The "Preview" kernel shows how Info.txt and PGM files can be parsed correctly.

Dataset

# MIAS Mammography

Looking for breast cancer

Kevin Mader · updated a year ago (Version 3)

36 voters  
share

Data Overview Kernels (8) Discussion (1) Activity

Download (205 MB) New Kernel

译文

Tags health 电子战核(8)讨论(1)活度

查看生词释义 >

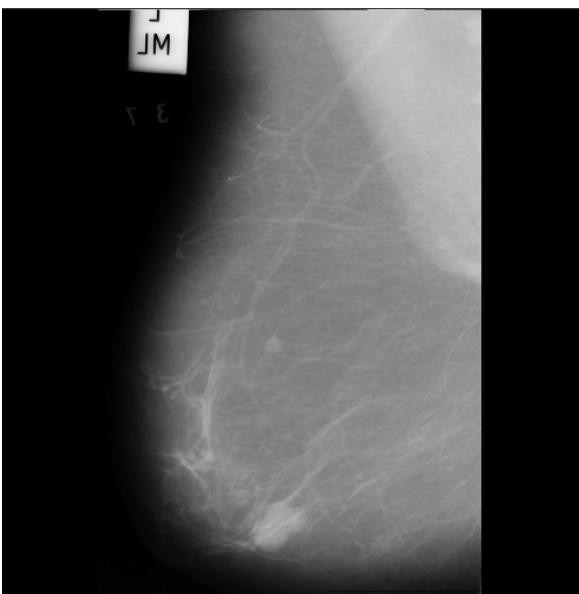
Description 基本图像2（具体步骤分析所

## Content

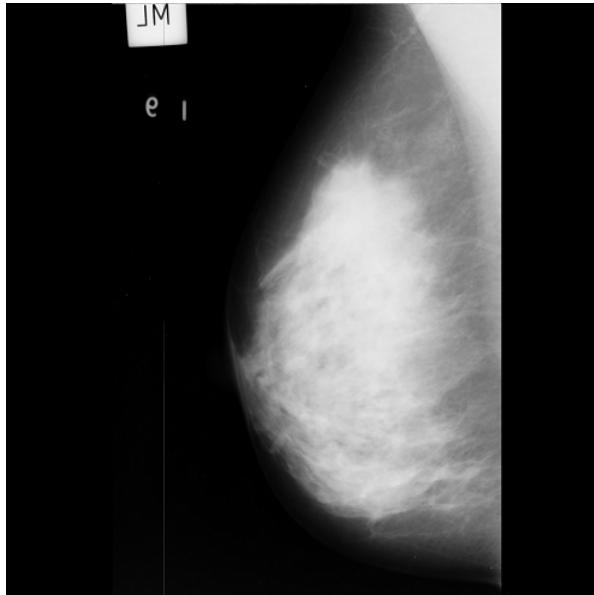
The data is images and labels / annotations for mammography scans. More about the database can be found at [MIAS](#). The 'Preview' kernel shows how the Info.txt and PGM files can be parsed correctly.

## 2

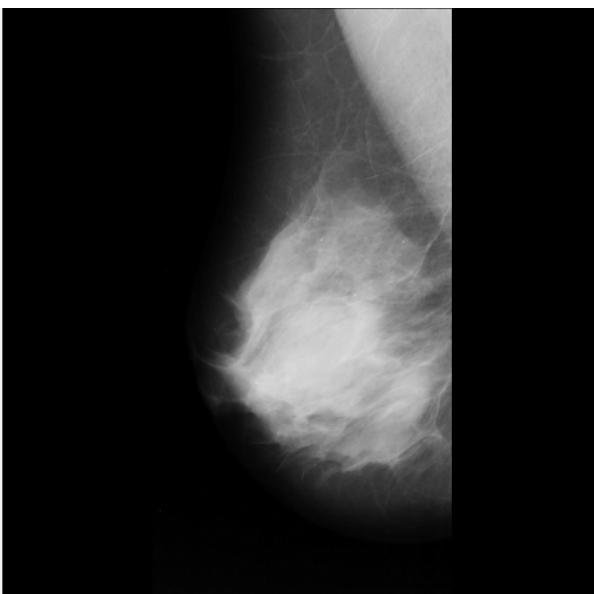
REFNUM	BG	CLASS	SEVERITY	X	Y	RADIUS
mdb001	G	CIRC	B	535	425	197
mdb002	G	CIRC	B	522	280	69
mdb003	D	NORM				
mdb004	D	NORM				
mdb005	F	CIRC	B	477	133	30
mdb005	F	CIRC	B	500	168	26
mdb006	F	NORM				
mdb007	G	NORM				
mdb008	G	NORM				
mdb009	F	NORM				
mdb010	F	CIRC	B	525	425	33
mdb011	F	NORM				
mdb012	F	CIRC	B	471	458	40
mdb013	G	MISC	B	667	365	31
mdb014	G	NORM				
mdb015	G	CIRC	B	595	864	68



Fatty

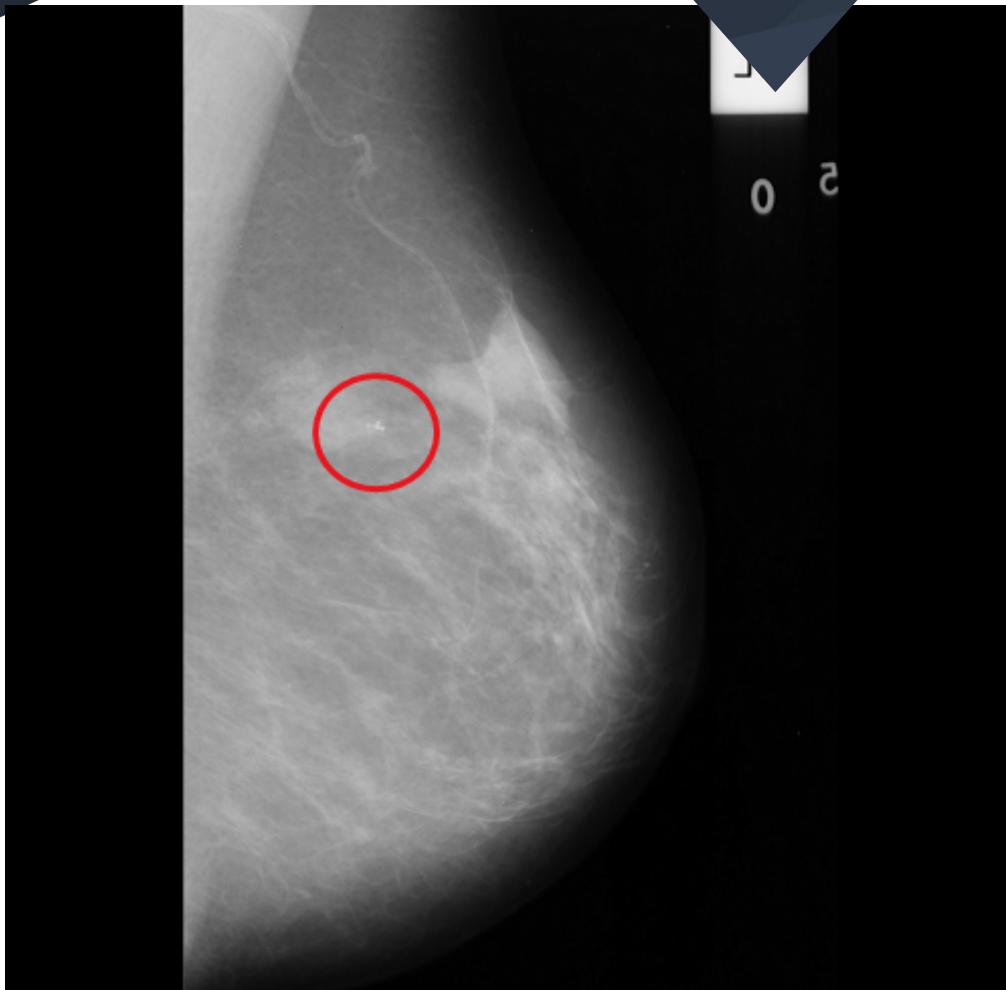


Dense-glandular



Fatty-glandular

2

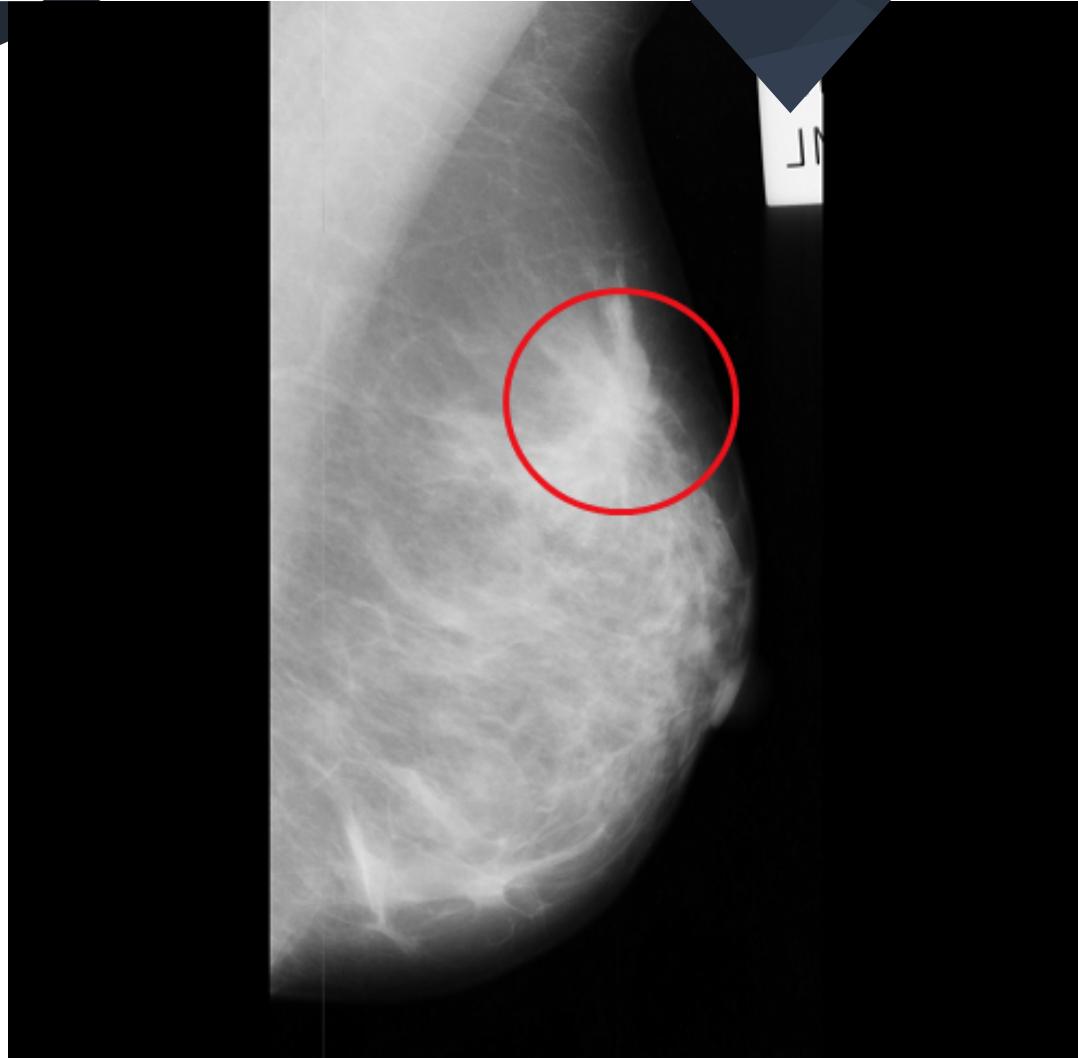


**Calcification**

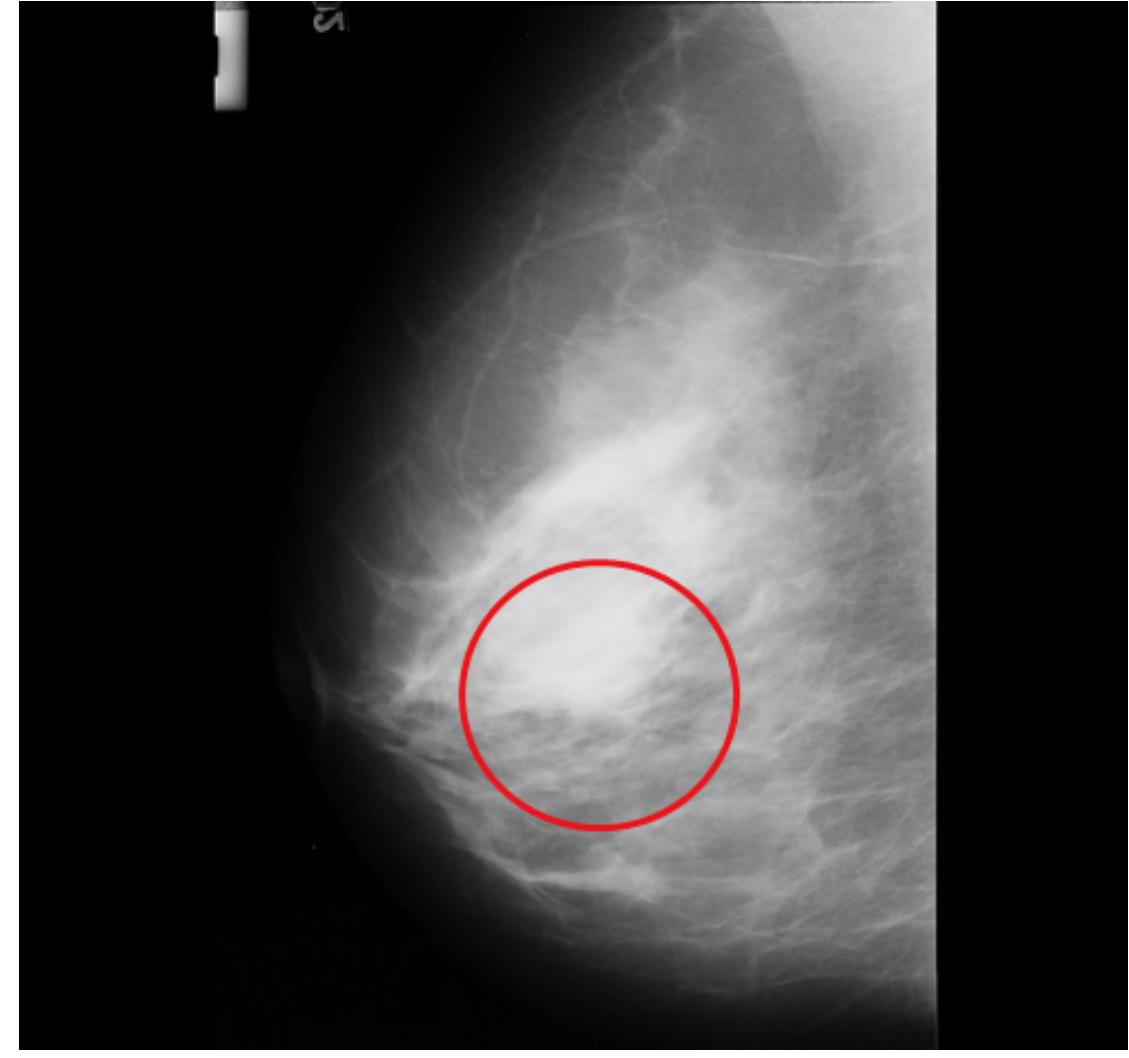


**Well-defined/circumscribed  
masses**

2

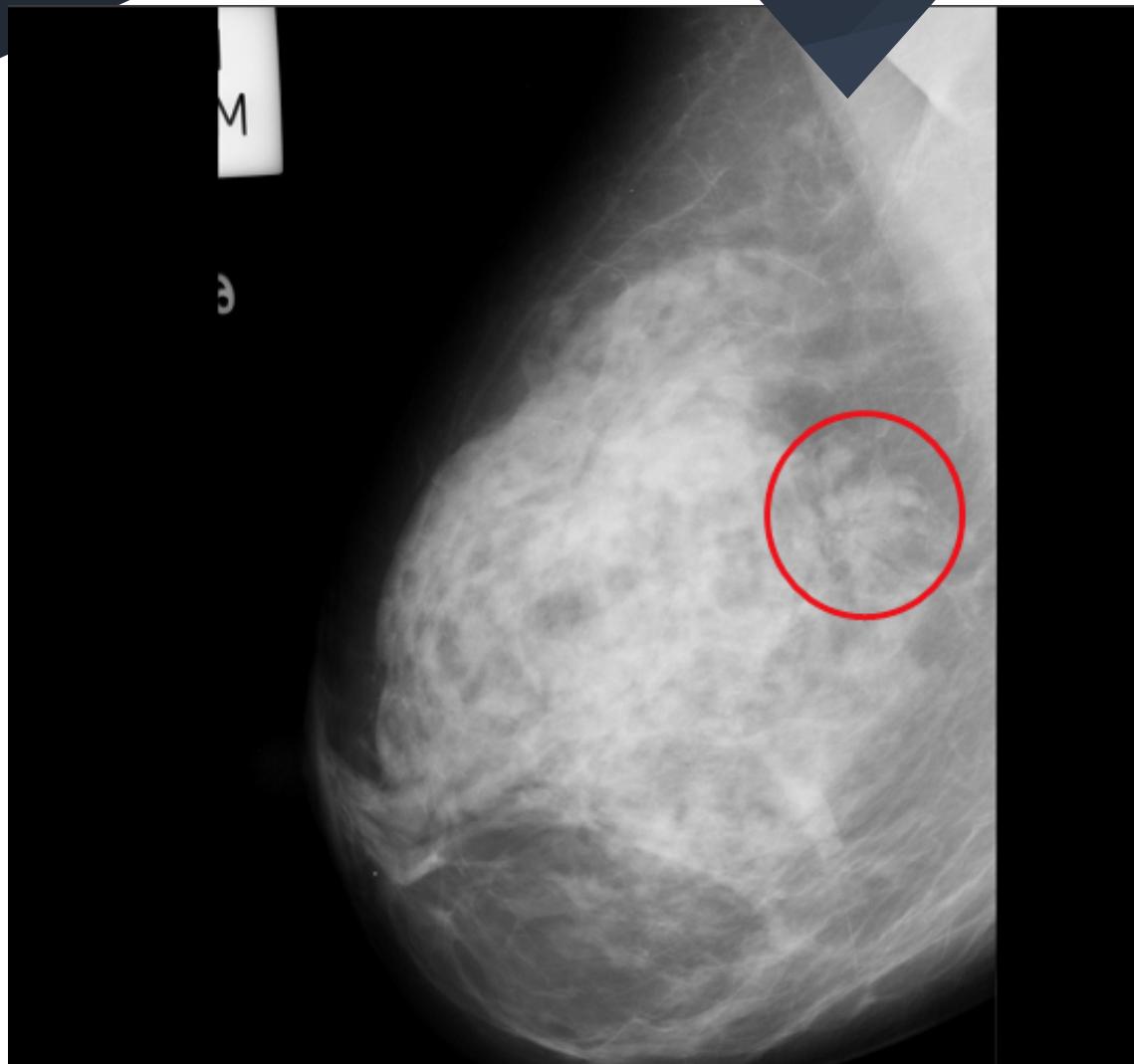


**Spiculated masses**

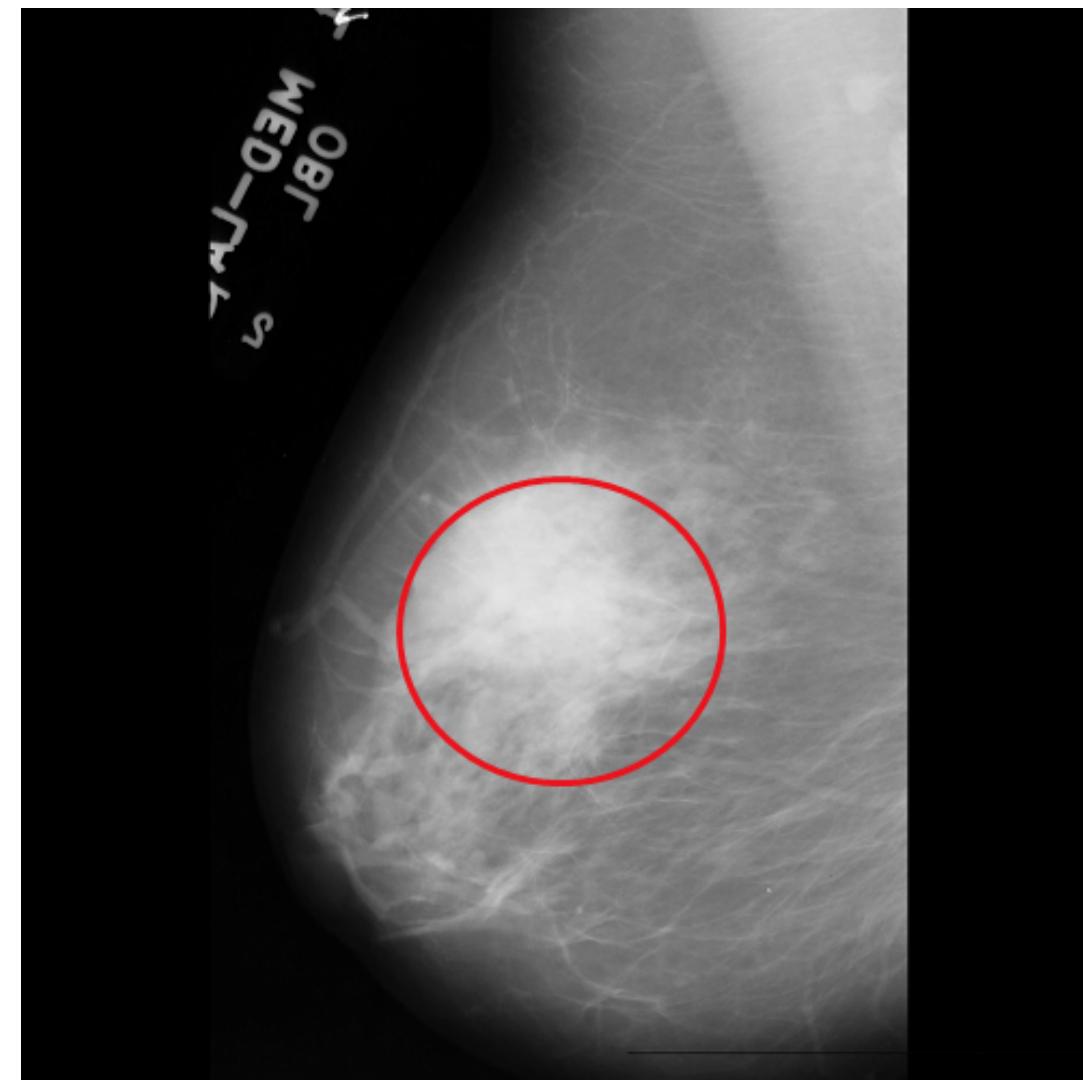


**Other, ill-defined masses**

2



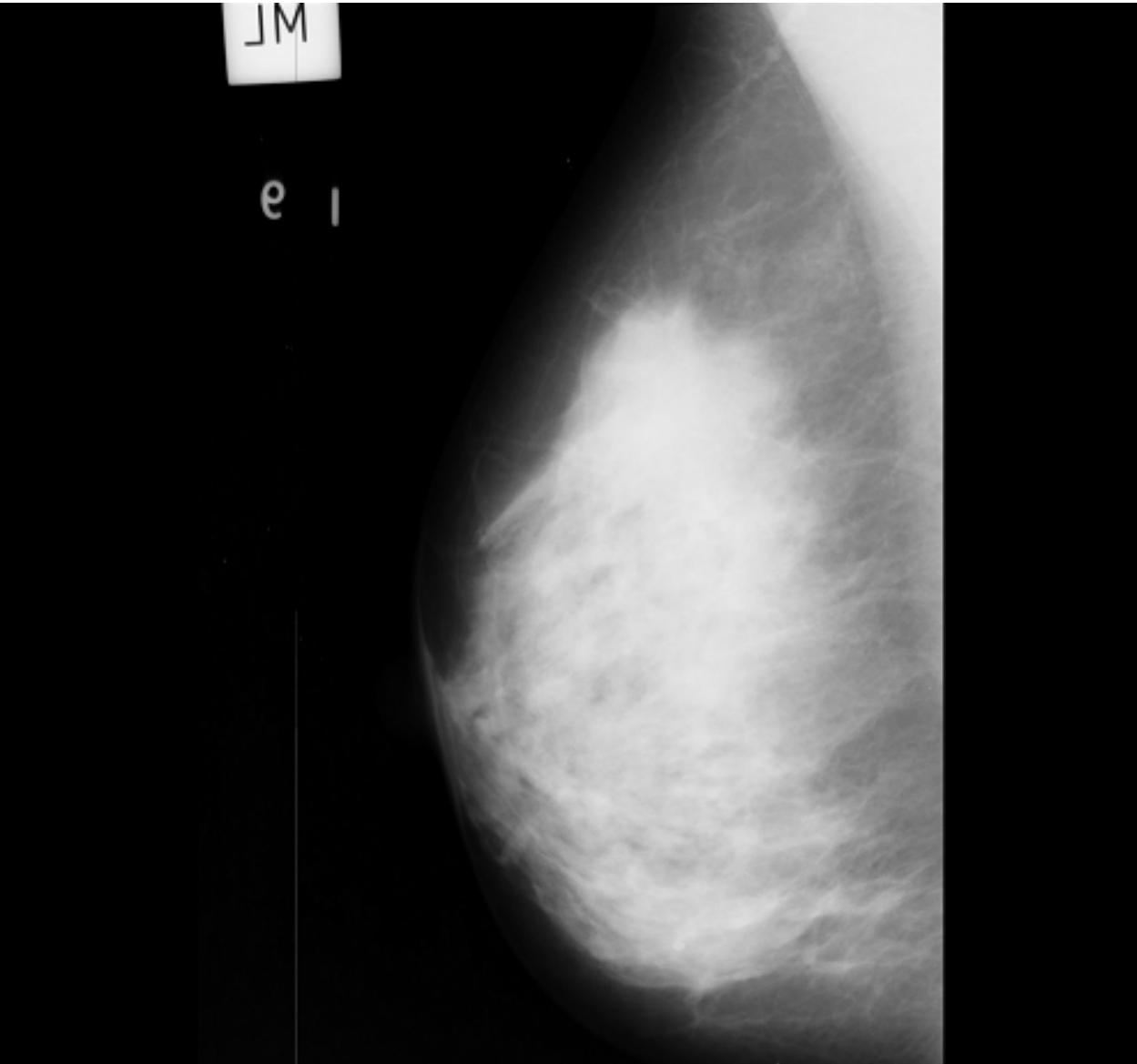
Architectural distortion



Asymmetry

2

Normal



REFNUM	BG	CLASS	SEVERITY	X	Y	RADIUS
mdb001	G	CIRC	B	535	425	197
mdb002	G	CIRC	B	522	280	69
mdb003	D	NORM				
mdb004	D	NORM				
mdb005	F	CIRC	B	477	133	30
mdb005	F	CIRC	B	500	168	26
mdb006	F	NORM				
mdb007	G	NORM				
mdb008	G	NORM				
mdb009	F	NORM				
mdb010	F	CIRC	B	525	425	33
mdb011	F	NORM				
mdb012	F	CIRC	B	471	458	40
mdb013	G	MISC	B	667	365	31
mdb014	G	NORM				
mdb015	G	CIRC	B	595	864	68

根据问题背景和所提供的信息，本文致力于解决以下两个问题：

1. 假设地球是球体的，卫星或飞船在地球上空300公里运行，所有测控站都与卫星或飞船的运行轨道共面，至少应该建立多少个测控站才能对其进行全程跟踪测控？

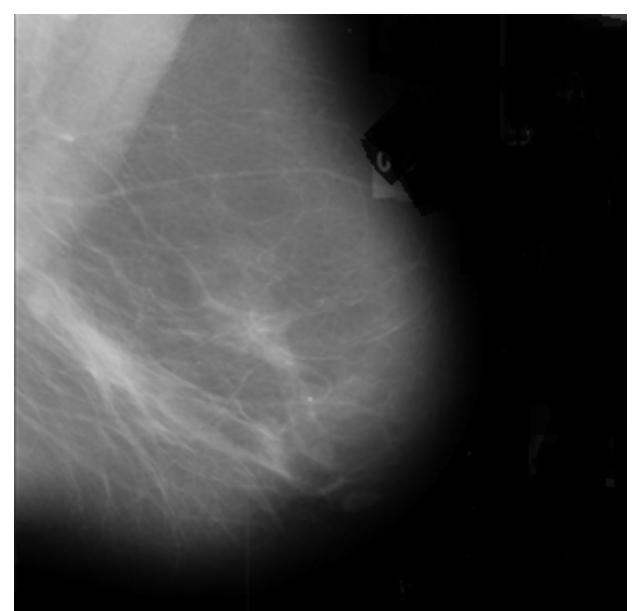
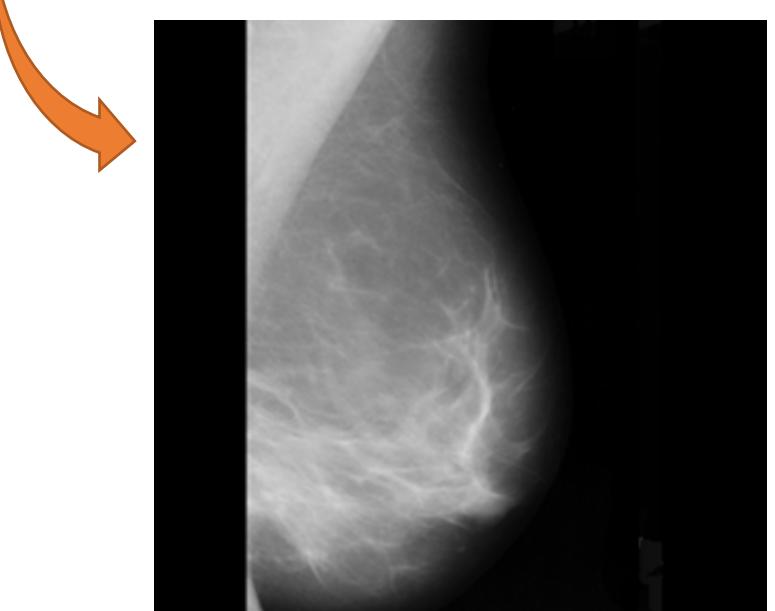
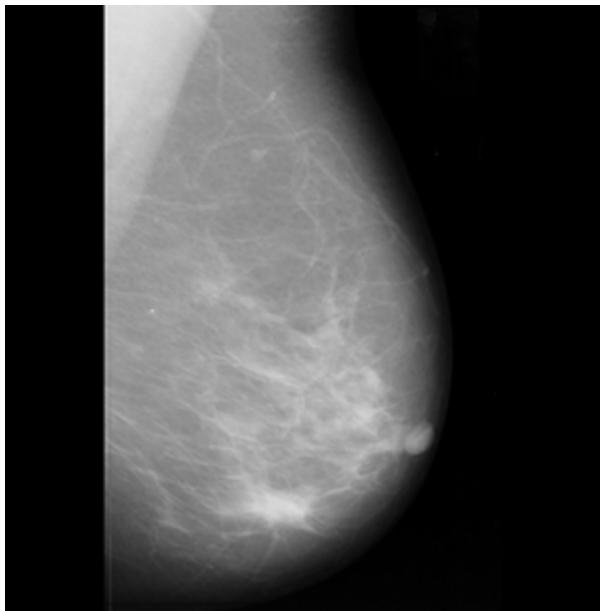
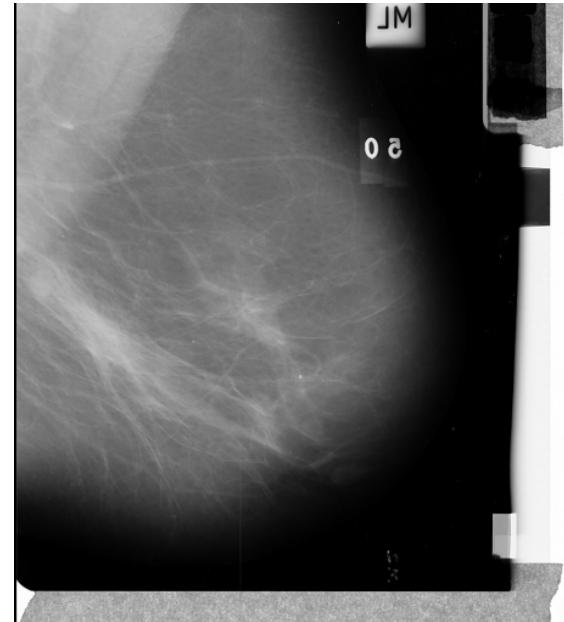
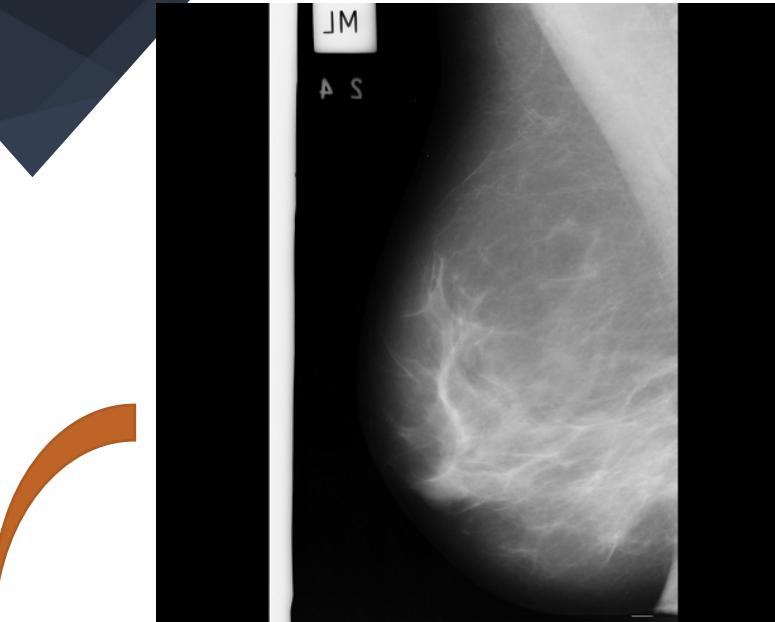
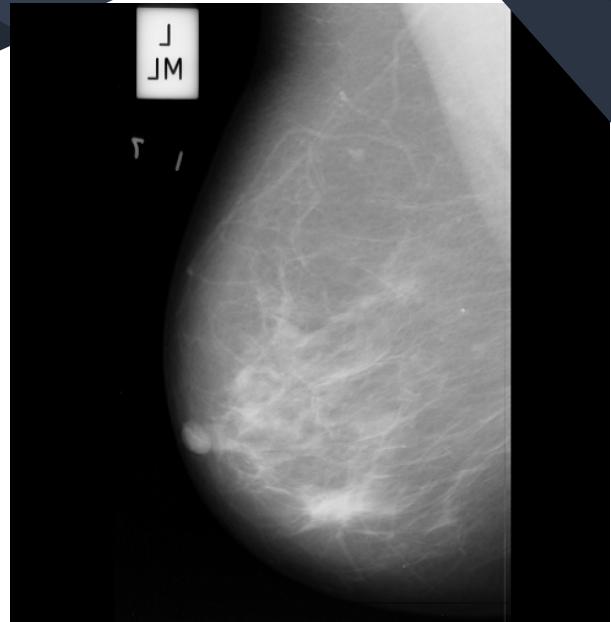
2. 假设地球是球体的，卫星或飞船在地球上空椭圆形轨道上运行，近地点为200公里，远地点为342公里。所有测控站都与卫星或飞船的运行轨道共面，至少应该建立多少个测控站才能对其进行全程跟踪测控？

# PART Data pre-processing

---

3

3



# 3

```
from numpy import *

'''通过方差的百分比来计算将数据降到多少维是比较合适的，  
函数传入的参数是特征值和百分比percentage，返回需要降到的维度数num'''  
def eigValPct(eigVals,percentage):  
  
#This cell is simply a function i write to replace the PCA in sklearn, ignore the chinese words.  
  
sortArray=sort(eigVals) #使用numpy中的sort()对特征值按照从小到大排序  
sortArray=sortArray[-1::-1] #特征值从大到小排序  
arraySum=sum(sortArray) #数据全部的方差arraySum  
tempSum=0  
num=0  
for i in sortArray:  
    tempSum+=i  
    num+=1  
    if tempSum>=arraySum*percentage:  
        return num  
  
'''pca函数有两个参数，其中dataMat是已经转换成矩阵matrix形式的数据集，列表示特征；  
其中的percentage表示取前多少个特征需要达到的方差占比，默认为0.9'''  
def pca(dataMat,percentage=0.9):  
    meanVals=mean(dataMat, axis=0) #对每一列求平均值，因为协方差的计算中需要减去均值  
    meanRemoved=dataMat-meanVals  
    covMat=cov(meanRemoved, rowvar=0) #cov()计算方差  
    eigVals,eigVects=linalg.eig(mat(covMat)) #利用numpy中寻找特征值和特征向量的模块linalg中的eig()方法  
    k=4 #要达到方差的百分比percentage，需要前k个向量  
    eigValInd=argsort(eigVals) #对特征值eigVals从小到大排序  
    eigValInd=eigValInd[-(k+1):-1] #从排好序的特征值，从后往前取k个，这样就实现了特征值的从大到小排列  
    redEigVects=eigVects[:,eigValInd] #返回排序后特征值对应的特征向量redEigVects (主成分)  
    lowDDDataMat=meanRemoved*redEigVects #将原始数据投影到主成分上得到新的低维数据lowDDDataMat  
    reconMat=(lowDDDataMat*redEigVects.T)+meanVals #得到重构数据reconMat  
    return lowDDDataMat, reconMat
```

# PART First experiment

---

4

Using knn, random forest and cnn

Accuracy 64.65%					Accuracy 63.64%				
	precision	recall	f1-score	support		precision	recall	f1-score	support
ARCH	0.00	0.00	0.00	6	ARCH	0.50	0.17	0.25	6
ASYM	0.00	0.00	0.00	4	ASYM	0.00	0.00	0.00	4
CALC	1.00	0.33	0.50	9	CALC	0.50	0.11	0.18	9
CIRC	0.00	0.00	0.00	8	CIRC	0.50	0.12	0.20	8
MISC	0.00	0.00	0.00	4	MISC	1.00	0.25	0.40	4
NORM	0.64	0.98	0.78	62	NORM	0.65	0.95	0.77	62
SPIC	0.00	0.00	0.00	6	SPIC	0.00	0.00	0.00	6
micro avg	0.65	0.65	0.65	99	micro avg	0.64	0.64	0.64	99
macro avg	0.23	0.19	0.18	99	macro avg	0.45	0.23	0.26	99
weighted avg	0.49	0.65	0.53	99	weighted avg	0.56	0.64	0.55	99

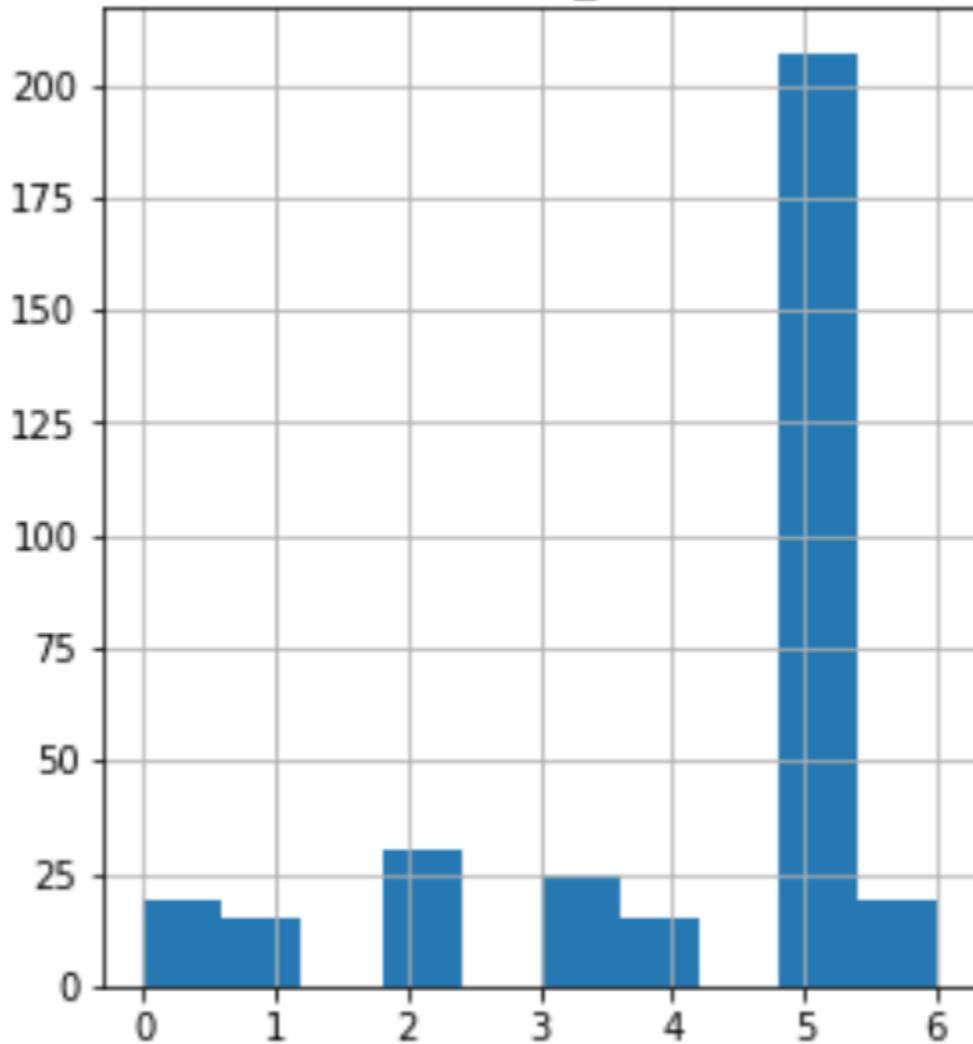
Knn, n = 9

Random forest, n\_estimators = 20

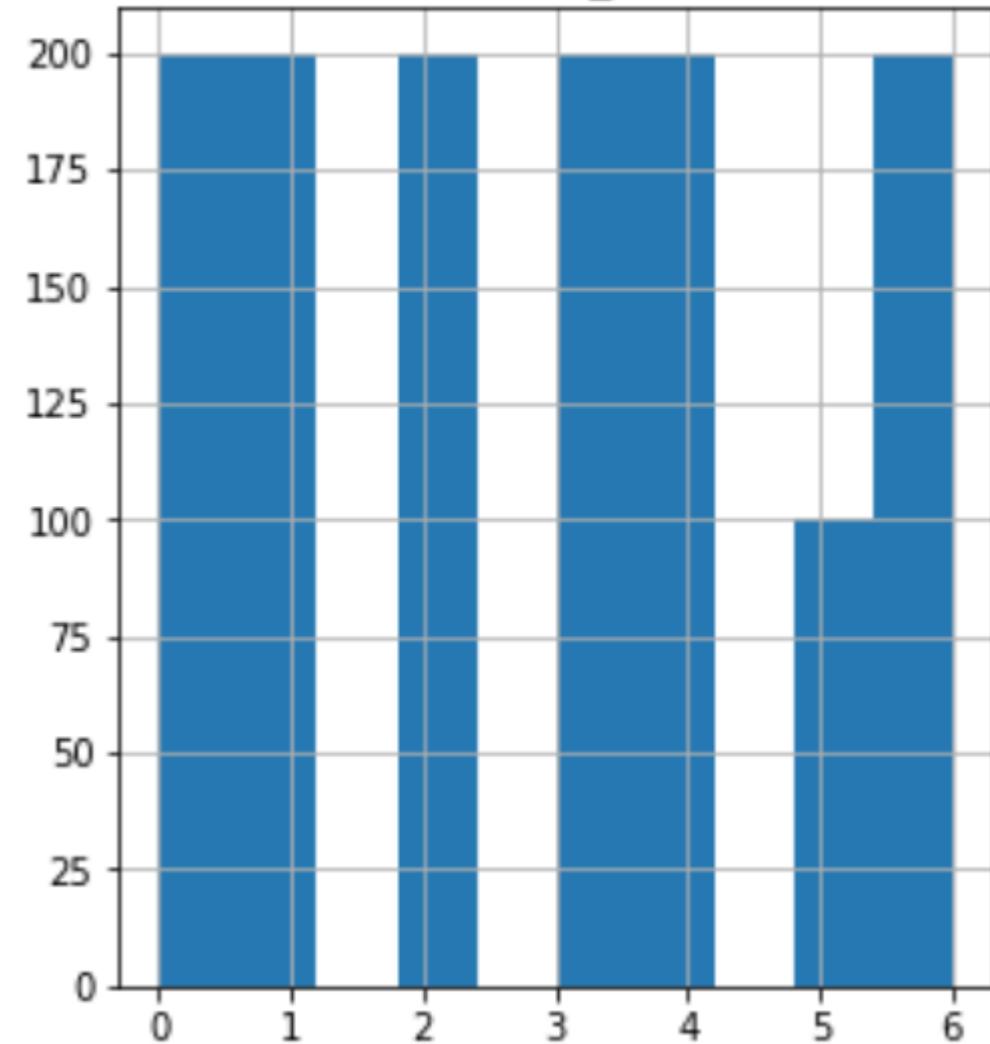


4

CLASS\_ID



CLASS\_ID



	precision	recall	f1-score	support
ARCH	0.20	0.40	0.27	5
ASYM	0.20	0.33	0.25	3
CALC	0.25	0.88	0.39	8
CIRC	0.33	0.50	0.40	6
MISC	0.15	0.50	0.24	4
NORM	1.00	0.02	0.04	52
SPIC	0.24	0.80	0.36	5
avg / total	0.72	0.24	0.15	83

# PART

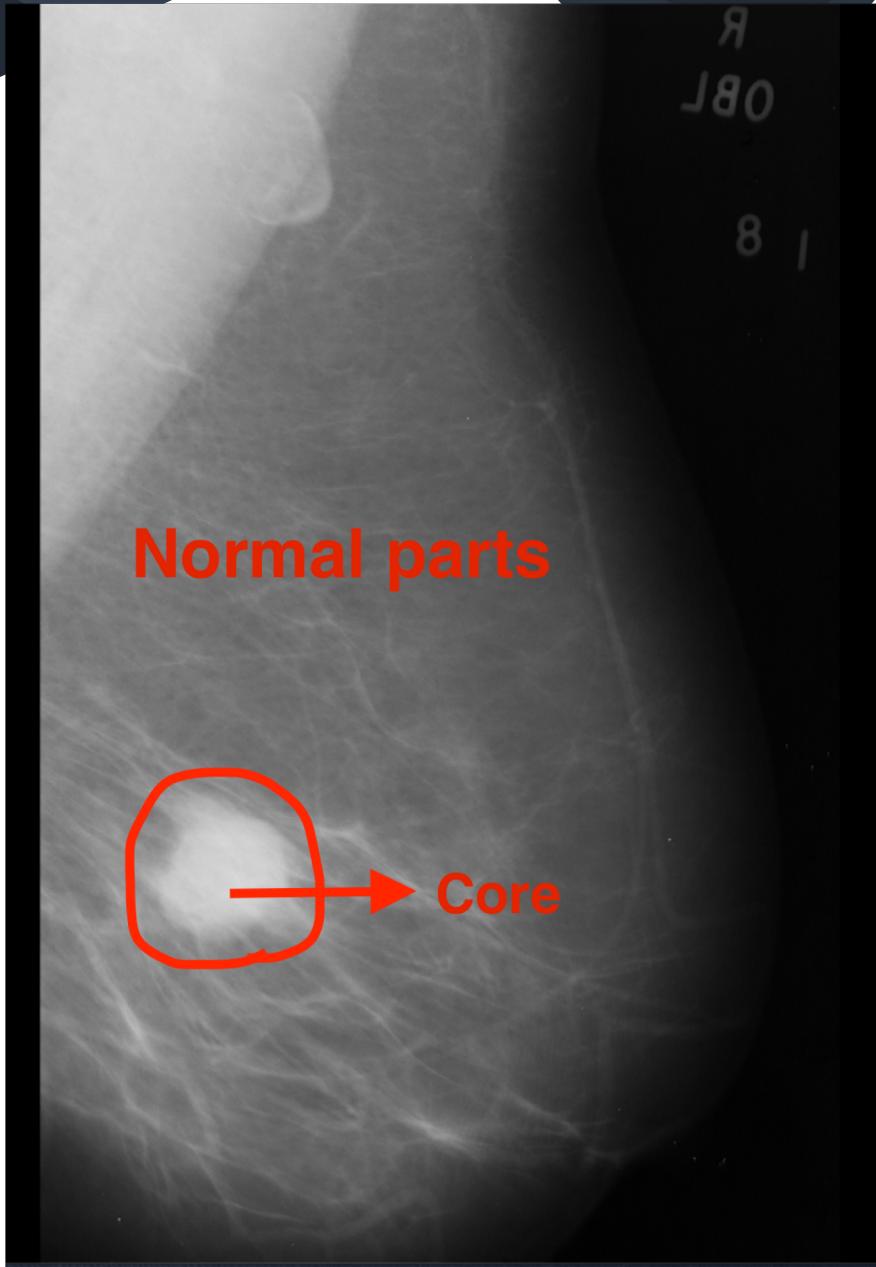
# 5

## Second experiment

---

Improved knn and random forest

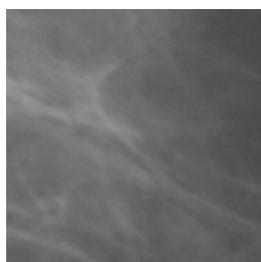
5



Core\_pic

1  
:  
2

Try to improve the accuracy of the classification of the core group for knn

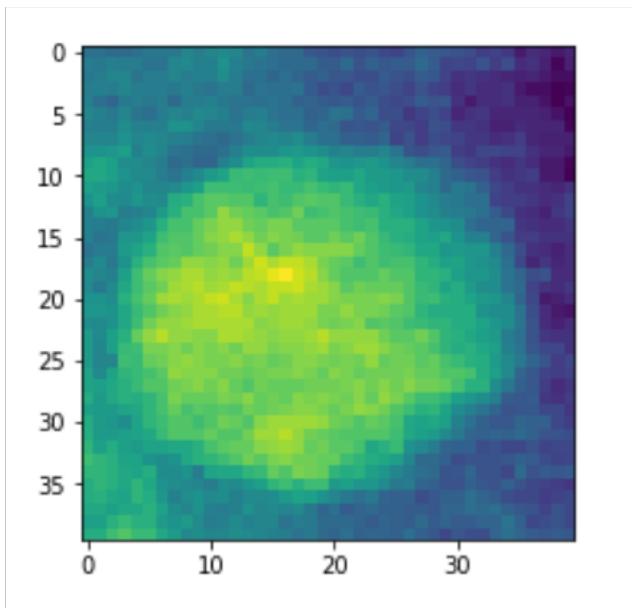


Normal parts

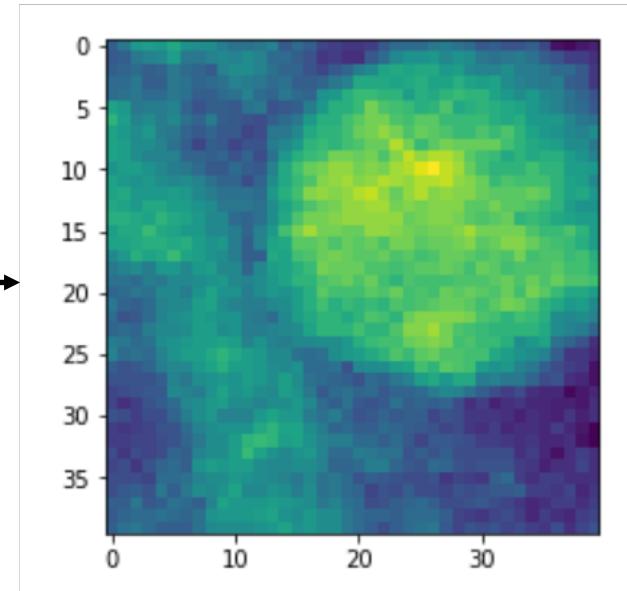


# 5

## Distribution balance



After Rotation and translation



# 5



A lot of pieces of  
squired pics

## 5

Accuracy 86.15%				
	precision	recall	f1-score	support
BALA	0.75	0.46	0.57	26
NORM	0.88	0.96	0.92	104
avg / total	0.85	0.86	0.85	130

Knn

Accuracy 84.62%				
	precision	recall	f1-score	support
BALA	0.64	0.54	0.58	26
NORM	0.89	0.92	0.91	104
avg / total	0.84	0.85	0.84	130

Random forest

# Machine learning

THANK YOU~~