

# Homework 2

*Frank Chou, Milo Opdahl, Tejaswi Pukkalla*

*March 12, 2019*

## Question 2: A Hospital Audit

Hospital Audits are important to determine the effectiveness of hospital operations from an objective standpoint. In this particular case, the goal is to determine the performance of radiologists using a statistical audit of their recent patient interactions - a crucial link between modern data-science and hospital operations. Two overall questions are posited:

1. First question: are some radiologists more clinically conservative than others in recalling patients, holding patient risk factors equal?
2. Second question: when the radiologists at this hospital interpret a mammogram to make a decision on whether to recall the patient, does the data suggest that they should be weighing some clinical risk factors more heavily than they currently are?

At the core of each question is reducing the number of false negatives - where a radiologist recommends a patient to conduct further tests and thereby allows a patient to begin immediately; and false positives - where a radiologist recommends further tests but ultimately turns out that there was no cancer. By introducing a statistical model, the goal is to augment the predictive capabilities of radiologist and offer a better standard of care for patients.

This audit is structured in four parts: first is a brief summary of the data and how it is structured, second is a demonstration and presentation of answering question one, third is a similar approach for question two, fourth is a conclusion of the audit's findings and recommendations for improvement of future radiologist performance or audit effectiveness.

### Part One: Brief Summary of Data

##	radiologist	cancer	recall	age
##	radiologist13:198	Min. :0.00000	Min. :0.0000	age4049 :287
##	radiologist34:197	1st Qu.:0.00000	1st Qu.:0.0000	age5059 :284
##	radiologist66:198	Median :0.00000	Median :0.0000	age6069 :199
##	radiologist89:197	Mean :0.03749	Mean :0.1499	age70plus:217
##	radiologist95:197	3rd Qu.:0.00000	3rd Qu.:0.0000	
##		Max. :1.00000	Max. :1.0000	
##	history	symptoms	menopause	density
##	Min. :0.0000	Min. :0.00000	postmenoHT :321	density1: 89
##	1st Qu.:0.0000	1st Qu.:0.00000	postmenoNoHT :360	density2:332
##	Median :0.0000	Median :0.00000	postmenounknown: 35	density3:460
##	Mean :0.1763	Mean :0.04863	premeno :271	density4:106
##	3rd Qu.:0.0000	3rd Qu.:0.00000		
##	Max. :1.0000	Max. :1.00000		

The data of mammograms used in this audit were selected from a Hospital in Seattle, Washington. At this hospital, five radiologists were selected at random for the audit - where about 200 mammograms were randomly selected from the hospital for each. For a total of 987 mammograms covering 7 parameters:

- age: 40-49\*, 50-59, 60-69, 70 and older
- family history of breast cancer: 0=No\*, 1=Yes
- history of breast biopsy/surgery: 0=No\*, 1=Yes

- breast cancer symptoms: 0=No\*, 1=Yes
- menopause/hormone-therapy status: Pre-menopausal, Post-menopausal & no hormone replacement therapy (HT), Post-menopausal & HT\*, Post-menopausal & unknown HT
- previous mammogram: 0=No\*, 1=Yes
- breast density classification: 1=Almost entirely fatty, 2=Scattered fibroglandular tissue\*, 3=Heterogeneously dense, 4=Extremely dense

Of these factors, two are of special interest: [recall] and [cancer]. In the abstract [recall] can be explained as the following: upon seeing the medical history of a patient, they can either recommend either one of two actions: recall for further screening or not. It is presumed that radiologists utilize all of the information available before they make a decision. This implies that there is a inherent correlative factor between recall and patient history. On the other hand [cancer] is whether or not a patient, whether through the recall screening process, or through another pathway of discovery - develops cancer within a 12 month window after seeing the radiologist.

## Part Two: Clinical Conservativism

Without knowing how patients are assigned to radiologists, it is presumed that the relationship is random at best, and preferential at worst. With a random assignment, we can presume that each radiologist chosen for the audit would have seen, on average, the same makeup of patients that would necessitate a mammogram. A random assignment would entail a random drawing of cancer patients from the overall total cancer patient pool from the population. If preferential - meaning that a patient approaches a radiologist and requests care and upon the approval of the radiologist, we see an issue of sampling error within the audit data; as there is a bias introduced between patient selection and radiologist. Radiologist may either self-select for more difficult cases or easier based on preference and patients self-select based on their estimate of the reputation of the radiologist within the medical community.

Regardless of assignment, the primary method of which we rank the clinical conservatism is to create a model that is trained on each of the radiologists' and then test the model on data from both the radiologist and other patients not seen by the radiologist in question. The goals behind this approach are twofold: one is to recreate a evaluation profile of the radiologist through a linear model of determining whether or not a patient should be recalled, two to determine whether or not a patient who is recalled or not develops cancer within a 12 month time frame.

The table below depicts the Root Mean Squared Error (RMSE) of each radiologist's model tested on a small sub-sample of the radiologist's test data and other radiologists' testing data.

##	lm1	lm1.w	lm2	lm2.w
## radiologist13	0.365743440	0.363424033	0.43912457	0.41722775
## radiologist34	0.285635286	0.388601868	0.34351880	0.43011474
## radiologist66	0.392460315	0.369352289	0.50307706	0.43786460
## radiologist89	0.408167944	0.400619339	0.47049508	0.45182382
## radiologist95	0.349184207	0.382951933	0.42744528	0.51249380
## SuperRad	0.361069858	0.354034509	0.37239727	0.34937231
## Rad13.compare	0.004673582	0.009389524	0.06672730	0.06785544
## Rad34.compare	-0.075434573	0.034567359	-0.02887846	0.08074243
## Rad66.compare	0.031390457	0.015317780	0.13067979	0.08849229
## Rad89.compare	0.047098086	0.046584830	0.09809781	0.10245151
## Rad95.compare	-0.011885651	0.028917424	0.05504802	0.16312149

Example:

- **radiologist13:** we have a the same linear model, `lm1 = glm(recall ~ .-cancer, data=brca_train, maxit = maxit)`, trained to 20% of radiologist13's sample data as well as the whole mammogram data - excluding radiologist13's.

- **SuperRad**: is a model trained on a 20% random sample of the whole data set and tested on the remainder of the whole data set. This pseudo-radiologist serves as the benchmark for comparing radiologists to an artificial standard if one radiologist had access and saw all of the patients from the data set.
- **Rad13.compare**: is determined by subtracting the model RMSE result of **\*radiologist13** by **SuperRad**. A positive value means that a model trained on **radiologist13's** training data did worse once it was tested on out of sample testing data and vice-versa.

##		lm1	lm1.w	lm2	lm2.w
##	radiologist13	0.365743440	0.363424033	0.43912457	0.41722775
##	radiologist34	0.285635286	0.388601868	0.34351880	0.43011474
##	radiologist66	0.392460315	0.369352289	0.50307706	0.43786460
##	radiologist89	0.408167944	0.400619339	0.47049508	0.45182382
##	radiologist95	0.349184207	0.382951933	0.42744528	0.51249380
##	SuperRad	0.361069858	0.354034509	0.37239727	0.34937231
##	Rad13.compare	0.004673582	0.009389524	0.06672730	0.06785544
##	Rad34.compare	-0.075434573	0.034567359	-0.02887846	0.08074243
##	Rad66.compare	0.031390457	0.015317780	0.13067979	0.08849229
##	Rad89.compare	0.047098086	0.046584830	0.09809781	0.10245151
##	Rad95.compare	-0.011885651	0.028917424	0.05504802	0.16312149

Just by viewing the table, it can be clearly discerned under **lm1** that on average, radiologists 13, 66, and 89 had worse performance than the benchmark **SuperRad** when looking at the RadXX.compare values for each radiologist; while 34 and 95 had better performance. But when we examine the results of each radiologists' model tested on the global data set, we find that on average, all radiologists were worse off. However **lm1** is a linear regression involving non-interacting variables from the data set. If we were to examine **lm2 <- glm(recall ~ (.-cancer)^2, data=brca\_train, maxit = maxit)** where we interact every variable with itself and another we find different results. Radiologist 95's model performance flips and becomes worse with 95's within-sample data. But once tested on the global data set, all radiologists' models performed worse than the benchmark. The takeaway from this analysis demonstrates that human radiologists, on average, are not as effective in determining whether or not a patient should be recalled than a statistical model. Although this might increase the number of false positives and false negatives, the overall increase in cancer detection would allow immediate treatment for true positives who otherwise would have gone undiagnosed. As for whether or not this behavior can be determined to be clinically conservative, meaning that radiologist will opt to recall a patient even if the clinical factors do not signal a need to recall, the distinction is minimal at best and hard to determine as all of the radiologists selected in the audit perform marginally better or worse than the benchmark.

## Part Three: Weighing Different Clinical Risk Factors

We first approach this question by developing four linear models that attempts to predict cancer rates based on the parameters available in the data set.

- **lm3 <- glm(cancer ~ recall, data=brca\_train, maxit = maxit)**
- **lm4 <- glm(cancer ~ recall + history, data=brca\_train, maxit = maxit)**
- **lm5 <- glm(cancer ~ ., data=brca\_train, maxit = maxit)**
- **lm6 <- glm(cancer ~ (.)^2, data=brca\_train, maxit = maxit)**

Because the goal of this question is to determine whether or not radiologists are effectively utilizing all of a patient's clinical data to determine whether or not to recall a patient, we first examine **lm3** and **lm4**. Both are linear models designed to find the partial effect of whether or not a patient was recalled and if they developed cancer within the next 12 months. However the distinction is that **lm3** only has recall as its x variable while **lm4** has both recall and family history.

```
summary(lm3)
```

```
##
## Call:
## glm(formula = cancer ~ recall, data = brca_train, maxit = maxit)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.15385  -0.01783  -0.01783  -0.01783   0.98217
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.017831   0.007138   2.498  0.0127 *
## recall       0.136016   0.018547   7.334 5.57e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.03428528)
##
##      Null deviance: 28.861  on 789  degrees of freedom
## Residual deviance: 27.017  on 788  degrees of freedom
## AIC: -418.78
##
## Number of Fisher Scoring iterations: 2
```

```
summary(lm5)
```

```
##
## Call:
## glm(formula = cancer ~ . - recall, data = brca_train, maxit = maxit)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.13967  -0.05402  -0.03076  -0.01538   0.99279
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      -0.013158   0.037904  -0.347  0.72859
## radiologistradiologist34 -0.005238   0.021682  -0.242  0.80918
## radiologistradiologist66 -0.009295   0.021687  -0.429  0.66835
## radiologistradiologist89 -0.009161   0.022178  -0.413  0.67968
## radiologistradiologist95 -0.009552   0.021770  -0.439  0.66095
## ageage5059         0.032176   0.022905   1.405  0.16049
## ageage6069         0.025219   0.027285   0.924  0.35563
## ageage70plus       0.068629   0.027004   2.541  0.01123 *
## history            0.014234   0.018145   0.784  0.43302
## symptoms           0.003980   0.030258   0.132  0.89537
## menopausepostmenoNoHT  -0.014608   0.016983  -0.860  0.38995
## menopausepostmenounknown -0.006012   0.039544  -0.152  0.87920
## menopausepremeno     0.016505   0.024185   0.682  0.49516
## densitydensity2      0.013153   0.026395   0.498  0.61840
## densitydensity3      0.026385   0.026551   0.994  0.32065
## densitydensity4      0.089434   0.033242   2.690  0.00729 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.03639461)
```

```
##
## Null deviance: 28.861 on 789 degrees of freedom
## Residual deviance: 28.169 on 774 degrees of freedom
## AIC: -357.78
##
## Number of Fisher Scoring iterations: 2
```

By itself, we can see that the **recall** variable has a very significant (p-value close to 0) and large effect on whether or not a patient develops cancer. This makes sense because upon evaluating a patient, a radiologist will then determine whether or not the patient will be recalled and receive additional testing. Based on their experience and education, they will want to find the factors that most likely contributes to cancer. At the same time however, we also see significant (in terms of p-value and magnitude) effects from **age**, **menopause/hormone-therapy status**, and **breast density classification**. In light of these factors, a series of model efficacy tests were conducted to determine the effectiveness of different models.

```
##          lm3          lm3.w          lm4          lm4.w
## radiologist13 0.326750727 0.337005247 0.3332733946 0.336809528
## radiologist34 0.213122379 0.311894469 0.2255037862 0.319724743
## radiologist66 0.377001486 0.335069096 0.3791503679 0.337329619
## radiologist89 0.376251046 0.320967727 0.3945406249 0.360075362
## radiologist95 0.296198425 0.319017089 0.3011243104 0.331615728
## SuperRad      0.332031569 0.329502308 0.3330038087 0.330490614
## Rad13.compare -0.005280841 0.007502939 0.0002695859 0.006318914
## Rad34.compare -0.118909190 -0.017607838 -0.1075000225 -0.010765870
## Rad66.compare 0.044969917 0.005566788 0.0461465593 0.006839005
## Rad89.compare 0.044219477 -0.008534581 0.0615368162 0.029584748
## Rad95.compare -0.035833144 -0.010485219 -0.0318794982 0.001125114
##          lm5          lm5.w          lm6          lm6.w
## radiologist13 0.373740388 0.378109275 0.406027074 0.40226902
## radiologist34 0.279916211 0.397728621 0.311721104 0.41455655
## radiologist66 0.406538005 0.362819638 0.416143141 0.38860249
## radiologist89 0.432300265 0.396129977 0.474467845 0.45781915
## radiologist95 0.349830246 0.387427297 0.386435267 0.43368063
## SuperRad      0.377678967 0.374250986 0.380752134 0.37497029
## Rad13.compare -0.003938579 0.003858289 0.025274940 0.02729873
## Rad34.compare -0.097762756 0.023477635 -0.069031030 0.03958626
## Rad66.compare 0.028859038 -0.011431347 0.035391007 0.01363220
## Rad89.compare 0.054621298 0.021878992 0.093715711 0.08284886
## Rad95.compare -0.027848721 0.013176311 0.005683133 0.05871034
```

Looking across **SuperRad** we see that the RMSE of each model remains fairly consistent throughout the different implementation and test of each model - except when we exclude **recall** in models **lm5** and **lm6**. The exclusion of **recall** has a meaningful impact models' ability to guess the cancer rate for each patient. Given this puzzling outcome, the next step would be to examine **lm5.w** and **lm6.w** where we take models that exclude **recall** - after all, as recall determinations occur after a radiologist sees a patient and not before, we cannot use it to predict cancer; and see which radiologist model performs the best. Iteration terms seems to be resulting in higher RMSE in the predictive model than by itself. Given the summary results from earlier regarding the significance of some variables over others, it can be concluded that radiologists weigh **age**, **menopause/hormone-therapy status**, and **breast density classification** as indicators of cancer than other factors excluding recall.

## Part Four: Conclusion

Ultimately, it can be determined that human radiologists may appear to be more conservative than a statistical model, but the underlying analysis claims otherwise - the difference is small in nature and not of sufficient

significance to sacrifice patient care for a more effective diagnosing mechanism. The number of false positives and false negatives remain small in comparison when the model changes from one to another.

```
## [1] "lm3 Confusion Table"
```

```
##      yhat
## y      0   1
##   0 163  27
##   1   3   4
```

```
## [1] "lm4 Confusion Table"
```

```
##      yhat
## y      0   1
##   0 162  28
##   1   3   4
```

```
## [1] "lm5 Confusion Table"
```

```
##      yhat
## y      0   1
##   0 184   6
##   1   7   0
```

```
## [1] "lm6 Confusion Table"
```

```
##      yhat
## y      0   1
##   0 164  26
##   1   5   2
```

- Pair-wise guesses and actual cancer results.
- (0,0) means that a patient did not have cancer and was not recalled.
- (1,0) means that a patient had cancer but was not recalled.
- (0,1) means that a patient did not have cancer but was recalled.
- (1,1) means that a patient had cancer and was successfully recalled.