

HIERACHICAL ACTIVE LEARNING BIOINFORMATICS APPLICATION

by

James D. Duin

A THESIS

Presented to the Faculty of
The Graduate College at the University of Nebraska
In Partial Fulfilment of Requirements
For the Degree of Master of Science

Major: Computer Science

Under the Supervision of Professor Stephen Scott

Lincoln, Nebraska

February, 2017

HIERACHICAL ACTIVE LEARNING BIOINFORMATICS APPLICATION

James D. Duin, M.S.

University of Nebraska, 2017

Adviser: Stephen Scott

I made the classifiers to predict the BioDataset, I made the plots for the Active Passive curves, I made the plots for the Fixed Fine Ratio experiments with various costs.

DEDICATION

This thesis is dedicated to my parents Paul and Vicki Duin and fiancée Anna Spady.

ACKNOWLEDGMENTS

I would like to thank my advisor Dr. Stephen Scott for guidance in selecting this research topic, Yugi Mo for his work in developing the HAL methodology, and Dr. Douglas Downey for his work on this topic. I would like to thank Juan Cui, Jiang Shu, Kevin Chiang for assistance accessing and understanding the protein dataset that is the subject of the paper.

Contents

Contents	v
1 Introduction	1
1.1 Machine Learning	1
1.2 Hierarchical Bioinformatics Data Set	1
1.3 Coarse Grained vs Fine Grained Trade Off	2
2 Background and Related Work	4
2.1 Active Learning	4
2.2 Other Papers cited by Yugi	4
2.3 Hierarchical Active Learning	5
2.4 Application to Dispatch Dataset	5
3 Bio HAL Application	7
3.1 Training and Testing Coarse Grain and Fine Grain Classifiers	7
3.2 Passive SVM Rbf kernel vs Logistic Reg	13
3.3 Active vs Passive curves	16
3.3.1 Plots for Logistic Regression Active vs Passive curves	17
3.3.2 Plots for SVM Active vs Passive curves	20
3.4 Plots for FFR experiments	24
4 Conclusions and Future Work	34
A Tuning the fine grained classes	35
Bibliography	36

Chapter 1

Introduction

1.1 Machine Learning

Machine Learning algorithms are defined as computer programs that learn from experience E with respect to some class of tasks T and performance measure P , if their performance at tasks in T , as measured by P , improves with experience E [4]. In the context of this paper, the machine learning algorithm that is used is a support vector machine (SVM) implementation by libSvm [3]. The performance measure is the classification of protein instances according to a label, for example, originated in the mitochondria or not. The experience is the number of training instances in the training set. The dataset is initially partitioned into training and test sets. The algorithm improves its classifier structure based on the features of each instance by iterating through all instances in the training set. When training has been completed, the classifier is then tested on the test set and the number of instances that the classifier correctly or incorrectly labels determines its accuracy, precision and recall scores. In our protein dataset there are 20,098 proteins with 449 features each relating to their structure. In our experiments an SVM classifier is applied to the dataset with the goal of achieving high precision scores for the label mitochondrion, that is, if the protein originates in the mitochondria or not.

1.2 Hierarchical Bioinformatics Data Set

The protein dataset is labeled according to where it originates in the cell. At the root is mitochondrion, then there is the sub level labels for if its native to the mitochondria or if it has a separate target compartment specification. The complete tree along with the number of instances belonging to the each label is included

in Figure ??.

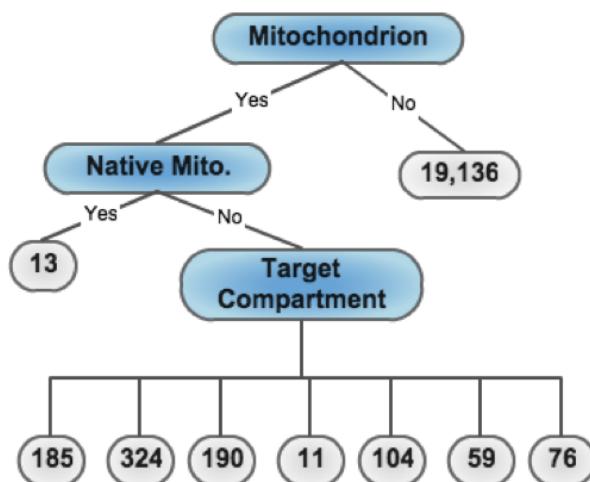


Figure 1.1: The protein dataset is labeled according to where it originates in the cell. At the root is mitochondrion, then there is the sub level labels for if its native to the mitochondria or if it has a separate target compartment specification. The complete tree along with the number of instances belonging to the each label is included in ??.

1.3 Coarse Grained vs Fine Grained Trade Off

The classifier that does not take advantage of any of the fine grained labels works off of the root labels for each instance and does not train separate classifiers for the fine grained labels. This classifier is referred to as the coarse grained classifier. The classifier that does use fine grained labels, and trains a separate classifier for each label, then combines them to generate a root level label is referred to as the fine grained classifier. It can be demonstrated through a dummy example that for certain datasets, a fine grained approach to the root level classifier can achieve higher levels of precision for the same level of recall. Such a dataset is shown in ??. The classifiers for this dataset can be thought of as a function of axis parallel rectangular boxes. For the course grained to have high recall and return all of the positive circle instances, it must encompass the entire dataset and incidentally return all of the negative diamond instances as positive also. A fine grained approach is preferable for the dummy dataset pictured. It is the intention of this study to demonstrate that the fine grained classification approach for a root level classifier will achieve higher levels of precision for the same level of recall when applied to the protein dataset.

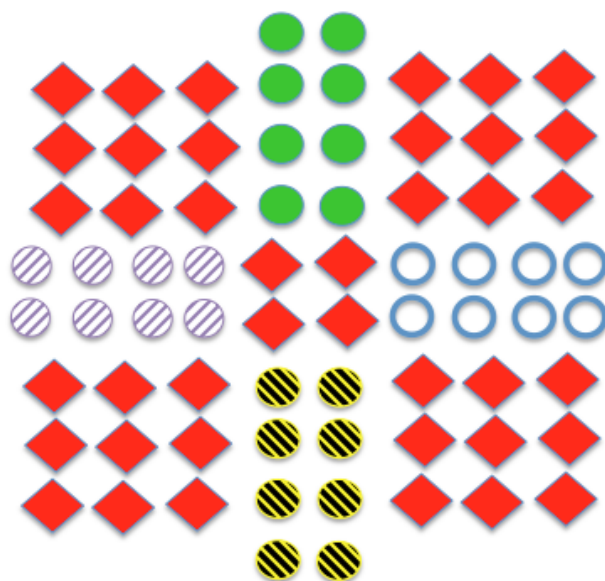


Figure 1.2: Demonstration of a dataset that would benefit from multiple fine grained learners for each circle type. In order for the coarse grain learner to have high recall, precision must be scarified and a large amount of false positives returned. By combining fine grained classifiers the same level of recall can be achieved with a higher level of precision because none of the false positive diamonds will be returned

Chapter 2

Background and Related Work

2.1 Active Learning

Active Learning relates to the coarse grained vs fine grained tradeoff because it is reasonable to assume that fine grained labels may not be as readily available as coarse grained labels, and thus have a higher cost. An active learning approach is used to determine how many fine grained labels to purchase in order to minimize the total cost to train the algorithm and maximize the precision and recall scores. The following equations for precision, recall, and a weighted F score are shown below in equations 1 through 3.

Precisioneqn

Recalleqn

F0.5eqn

The goal in an active learning approach is to maximize the F measure where α equals 0.5 [2]. The F-0.5 measure gives more weight to precision, as opposed to recall, so it gives incentive to purchase enough fine grained labels to increase the F-0.5 measure. The coarse grained labels will cost less than fine grained labels, but the increase in the F-0.5 measure justifies the increase in cost up to a certain point. The F-0.5 measure is used in the results section of this paper. ... That's why we use PRauc in the results.

2.2 Other Papers cited by Yugi

Describe some of the other papers that Yugi cited.

2.3 Hierarchical Active Learning

The Hierarchical Active Learning algorithm (HAL) is shown diagrammatically in Figure 3. Multiple fine grained classifiers are trained at each level of the Hierarchy of the dataset. Queries to the oracle are performed purchase the most cost effective labels to add to the training sets of the classifiers. The active learning cycle continues until a cost budget has been reached. The benefit of an active learning approach is to maximize the F-0.5 measure for a given cost budget. It was the goal of this study to apply the HAL algorithm to the protein dataset, however this is not achieved at this time. An existing application of HAL is briefly discussed in the following section.



Figure 2.1: Diagram of HAL approach

2.4 Application to Dispatch Dataset

HAL was applied to a Dispatch dataset. This dataset contains 375,026 manually labeled hierarchical names across 1,384 newspaper articles [2]. This is a clear example where fine grained labels have a higher cost since it is easier for a person to manually determine if the article pertains to an organization or not, rather than if it pertains to a railroad or a zoo, which would be sub labels of the organization root. The first analysis step was to determine that the F-0.5 measure is increased by using fine grained classifiers. The results are shown in Figure 4. The highest F-0.5 measure for a given iteration of purchases of training instances is obtained by using the active learning approach with all fine-grained labels. The passive learning curves were generated by selecting batches of instances randomly rather than querying the oracle for a specific

label type that offers that most gain in classifier accuracy. The active learning curves did take advantage of querying for specific labels in order to maximize gain in classifier accuracy. ... add other data sets.

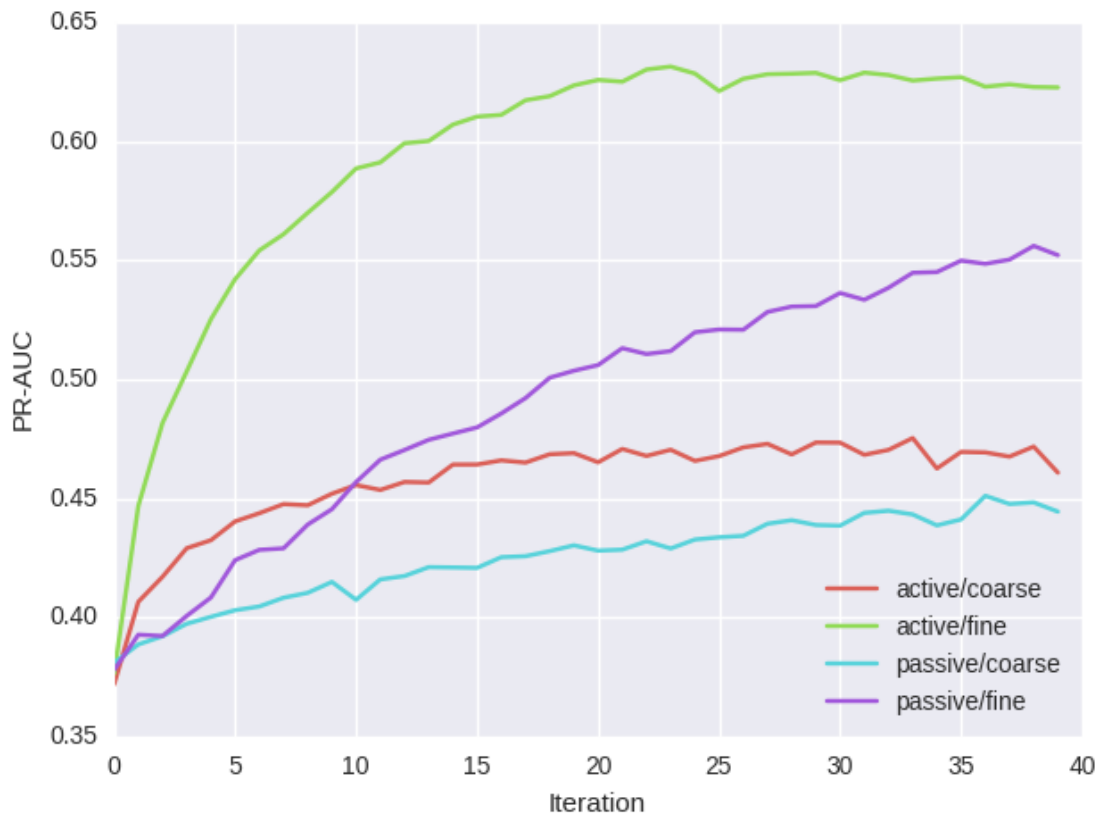


Figure 2.2: Application of HAL demonstrating the benefit of Actively selecting the type of labels to purchase for instances rather than randomly selecting labels to purchase, as in the Passive curves.

The next analysis step is to apply a given ratio of fine grained vs coarse grained labels to purchase at each batch request to the oracle. The results of varying the percentage of fine grained labels purchased are shown in Figure 5. The figure shows that even a small amount of fine grained labels purchased, that is, 20perc provides a significant increase in the F-0.5 measure for a given iteration. ... pr auc results. Add explanation of Bandit approach and results.

Chapter 3

Bio HAL Application

3.1 Training and Testing Coarse Grain and Fine Grain Classifiers

The first step is to examine the dataset.

Classes	All	All	Folds	0	1	2	3	4	5	6	7	8
0	19136	1	2010	1914	1	19	32	19	1	11	6	7
1	13	2	2010	1914	1	19	32	19	1	11	6	7
2	185	3	2010	1914	1	19	32	19	1	11	5	8
3	324	4	2010	1914	1	19	32	19	1	10	6	8
4	190	5	2010	1914	1	18	33	19	1	10	6	8
5	11	6	2010	1914	1	18	33	19	1	10	6	8
6	104	7	2010	1913	2	18	33	19	1	10	6	8
7	59	8	2010	1913	2	18	33	19	1	10	6	8
8	76	9	2009	1913	2	18	32	19	2	10	6	7
Total	20098	10	2009	1913	1	19	32	19	1	11	6	7
Shape	450	Total	20098	19136	13	185	324	190	11	104	59	76

(a) Classes

(b) Folds

Table 3.1: This is what the dataset looks like there are 20098 instances total with 450 features each. This is what the folds of the dataset look like.

Next the dataset is partitioned into 10 folds, each fold contains a representative proportion of each of the classes, the instances are added to each partition at random. The total partitioning looks like Table ??

Then 9 of the folds are compressed into the test set and the fold held out is the test set the totals of each class in the train and test set for fold 1 is shown in Table 3.4.

train	all	0	1	2	3	4	5	6	7	8
Total	18088	17222	12	166	292	171	10	93	53	69
test	all	0	1	2	3	4	5	6	7	8
Total	2010	1914	1	19	32	19	1	11	6	7

Table 3.2: This is what the train and test set look like.

coarse just used marked all positives as 1 and ran a binary classifier. fine grained classifiers trained 8 separate classifiers for the 8 fine grained classes, so for fine grain classifier for 1, all other fine grained classes are marked as 0 in addition to all the coarse instances being marked as 0. Also the Pr auc and Roc auc are the primary metrics for determining the performance of the classifier. The next step is to determine what classifier can be applied to 'learn' the classes of this dataset. Because the experiment will involve running multiple rounds with increasing the instances to be trained on iteratively I tested each classifier against the full dataset and then a reduced dataset with one fifth of the negative instances.

I tried using SVM. Throughout this project I used the python library sci-kit learn [1]. The support vector machine implemented by this library has the following default parameters. SVC C=1.0, kernel='rbf', degree=3, gamma='auto', coef0=0.0, shrinking=True, probability=False, tol=0.001, cache-size=200, class-weight=None, verbose=False, max-iter=-1, decision-function-shape=None, random-state=None.

Classes	All	All	Folds	0	1	2	3	4	5	6	7	8
0	3827	1	479	383	1	19	32	19	1	11	6	7
1	13	2	479	383	1	19	32	19	1	11	6	7
2	185	3	479	383	1	19	32	19	1	11	6	7
3	324	4	479	383	1	19	32	19	1	11	5	8
4	190	5	479	383	1	19	32	19	1	10	6	8
5	11	6	479	383	1	18	33	19	1	10	6	8
6	104	7	479	383	1	18	33	19	1	10	6	8
7	59	8	479	382	2	18	33	19	1	10	6	8
8	76	9	479	382	2	18	33	19	1	10	6	8
Total	4789	10	478	382	2	18	32	19	2	10	6	7
Shape	450	Total	4789	3827	13	185	324	190	11	104	59	76

(a) Classes Subset

(b) Folds Subset

Table 3.3: This is for the partitions subset.

train	all	0	1	2	3	4	5	6	7	8
Total	4310	3444	12	166	292	171	10	93	53	69
test	all	0	1	2	3	4	5	6	7	8
Total	479	383	1	19	32	19	1	11	6	7

Table 3.4: This is what the train and test set look like for the subset.

coarse-pr	fine-pr	coarse-roc	fine-roc	coarse-acc	fine-acc	coarse-f1	fine-f1
0.807	0.796	0.779	0.768	0.816	0.802	0.214	0.021
0.848	0.822	0.828	0.790	0.825	0.804	0.263	0.041
0.846	0.821	0.810	0.765	0.818	0.802	0.243	0.021
0.860	0.832	0.826	0.775	0.831	0.802	0.319	0.021
0.859	0.829	0.828	0.783	0.833	0.804	0.298	0.041
0.796	0.763	0.748	0.715	0.816	0.806	0.214	0.061
0.838	0.825	0.797	0.792	0.818	0.800	0.243	0.020
0.836	0.816	0.803	0.770	0.823	0.800	0.309	0.020
0.863	0.845	0.833	0.805	0.829	0.797	0.305	0.000
0.844	0.806	0.806	0.758	0.836	0.807	0.339	0.061
avg 0.840	avg 0.815	avg 0.806	avg 0.772	avg 0.825	avg 0.802	avg 0.275	avg 0.031

Table 3.5: SVMDef result stats

coarse-tn	fine-tn	coarse-fp	fine-fp	coarse-fn	fine-fn	coarse-tp	fine-tp
379	383	4	0	84	95	12	1
380	383	3	0	81	94	15	2
378	383	5	0	82	95	14	1
379	383	4	0	77	95	19	1
382	383	1	0	79	94	17	2
379	383	4	0	84	93	12	3
378	382	5	1	82	95	14	1
375	382	7	0	78	96	19	1
379	382	3	0	79	97	18	0
379	382	3	0	75	92	20	3
avg 378.8	avg 382.6	avg 3.9	avg 0.1	avg 80.1	avg 94.6	avg 16.0	avg 1.5

Table 3.6: SVMDef conf matrix

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.840	0.806	0.825	0.275	(378.8 / 80.1)	(3.9 / 16.0)
fine	0.815	0.772	0.802	0.031	(382.6 / 94.6)	(0.1 / 1.5)

Table 3.7: SVMDef

I tried different scaling methods (min max scaler, std scaler), I settled on std scaler.

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.881	0.855	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)
fine	0.840	0.810	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)

Table 3.8: SVMMinMax

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.801	0.791	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)
fine	0.636	0.615	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)

Table 3.9: SVMNorm

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.912	0.882	0.881	0.631	(372.7 / 47.1)	(10.0 / 49.0)
fine	0.879	0.848	0.809	0.094	(382.7 / 91.3)	(0.0 / 4.8)

Table 3.10: SVMStandard

I tried different feature select measures, I settled on 75 perc feature select.

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.907	0.875	0.877	0.623	(370.7 / 47.1)	(12.0 / 49.0)
fine	0.854	0.823	0.806	0.068	(382.7 / 92.7)	(0.0 / 3.4)

Table 3.11: SVMsel25

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.913	0.885	0.879	0.632	(371.3 / 46.4)	(11.4 / 49.7)
fine	0.874	0.842	0.810	0.097	(382.7 / 91.2)	(0.0 / 4.9)

Table 3.12: SVMsel50

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.913	0.883	0.878	0.622	(372.1 / 47.9)	(10.6 / 48.2)
fine	0.880	0.848	0.809	0.089	(382.7 / 91.6)	(0.0 / 4.5)

Table 3.13: SVMsel75

I tried different different C costs, kernels, decision function shape, gamma, tolerance settled on `classif = svm.SVC(C=1.0, kernel='rbf', decisionfunctionshape='ovo', gamma=0.0025, tol=0.00001)`. The default gamma setting is 0.002967 or $(1/\text{num-features})$ or $(1/337)$. Default cost is actually 1.0, and the default class weight is balanced wiced weights each class by the number of instances it has in the train set. Tolerance default is 0.001.

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.913	0.883	0.878	0.621	(372.1 / 48.0)	(10.6 / 48.1)
fine	0.880	0.847	0.809	0.089	(382.7 / 91.6)	(0.0 / 4.5)

Table 3.14: SVM-ovr

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.913	0.883	0.878	0.622	(372.1 / 47.9)	(10.6 / 48.2)
fine	0.880	0.848	0.809	0.089	(382.7 / 91.6)	(0.0 / 4.5)

Table 3.15: SVM-ovo

I left the class weight as balanced for this part, the results did not show much advantage for using the fine grained classifier.

next I tried using a Logistic regression classifier.

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.887	0.862	0.867	0.615	(364.1 / 45.0)	(18.6 / 51.1)
fine	0.854	0.837	0.833	0.395	(372.8 / 69.9)	(9.9 / 26.2)

Table 3.16: LogRegDef

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.864	0.849	0.846	0.583	(353.8 / 44.8)	(28.7 / 51.3)
fine	0.833	0.816	0.831	0.471	(362.0 / 60.2)	(20.5 / 36.0)

Table 3.17: LogRegStandard

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.790	0.761	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)
fine	0.767	0.735	0.799	0.000	(382.7 / 96.1)	(0.0 / 0.0)

Table 3.18: LogRegNorm

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.891	0.867	0.864	0.581	(368.6 / 50.9)	(14.1 / 45.2)
fine	0.888	0.862	0.812	0.130	(382.1 / 89.3)	(0.6 / 6.8)

Table 3.19: LogRegMinMax

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.872	0.848	0.849	0.497	(370.8 / 60.3)	(11.9 / 35.8)
fine	0.869	0.845	0.804	0.052	(382.2 / 93.5)	(0.5 / 2.6)

Table 3.20: LogRegSel25

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.875	0.849	0.849	0.497	(370.8 / 60.3)	(11.9 / 35.8)
fine	0.872	0.846	0.803	0.050	(382.2 / 93.6)	(0.5 / 2.5)

Table 3.21: LogRegSel50

title	pr	roc	acc	f1	conf (tn/fn)	conf (fp/tp)
coarse	0.871	0.847	0.848	0.493	(370.6 / 60.6)	(12.1 / 35.5)
fine	0.869	0.845	0.803	0.048	(382.0 / 93.7)	(0.7 / 2.4)

Table 3.22: LogRegSel75

I tried different scaling methods, min max scaler, std scaler, I settled on min max scaler.

I tried different feature select measures, decided to use all of the features.

I tried different C costs, tolerances, and class weights. I settled on C=0.1, tol = 0.00001, and weight equal to the balanced, adjusted via a scaling line.

I also further tuned the fine grained classifiers starting from the initial scaling from the line an then multiplying that by a ratio. I got this vector of ratios [0.87, 0.4, 0.78, 0.65, 3.48, 0.78, 1.74, 0.87]

3.2 Passive SVM Rbf kernel vs Logistic Reg

This shows the advantage to fine grained labels to justify the experiment.

coarse-pr	fine-pr	coarse-roc	fine-roc	coarse-acc	fine-acc	coarse-f1	fine-f1
0.898	0.901	0.905	0.896	0.767	0.945	0.259	0.474
0.870	0.869	0.847	0.846	0.803	0.944	0.272	0.456
0.897	0.907	0.895	0.901	0.792	0.947	0.287	0.500
0.864	0.866	0.852	0.848	0.778	0.943	0.256	0.430
0.855	0.865	0.859	0.859	0.795	0.947	0.269	0.451
0.867	0.869	0.874	0.865	0.785	0.939	0.263	0.417
0.871	0.887	0.873	0.881	0.784	0.940	0.269	0.442
0.835	0.845	0.843	0.842	0.794	0.940	0.258	0.388
0.870	0.878	0.869	0.871	0.784	0.939	0.262	0.417
0.873	0.873	0.891	0.890	0.786	0.933	0.279	0.368
avg 0.870	avg 0.876	avg 0.871	avg 0.870	avg 0.787	avg 0.942	avg 0.268	avg 0.434

Table 3.23: Here are the results for the logistic regression passive 10 folds.

coarse-tn	fine-tn	coarse-fp	fine-fp	coarse-fn	fine-fn	coarse-tp	fine-tp
1460	1849	454	65	14	46	82	50
1540	1851	374	63	22	49	74	47
1508	1851	406	63	12	43	84	53
1486	1853	428	61	19	53	77	43
1521	1859	393	55	20	52	76	44
1501	1843	413	71	19	52	77	44
1496	1841	417	72	17	49	80	48
1524	1852	389	61	25	59	72	38
1498	1841	415	72	18	51	77	44
1497	1836	416	77	13	57	83	39
avg 1503.1	avg 1847.6	avg 410.5	avg 66.0	avg 17.9	avg 51.1	avg 78.2	avg 45.0

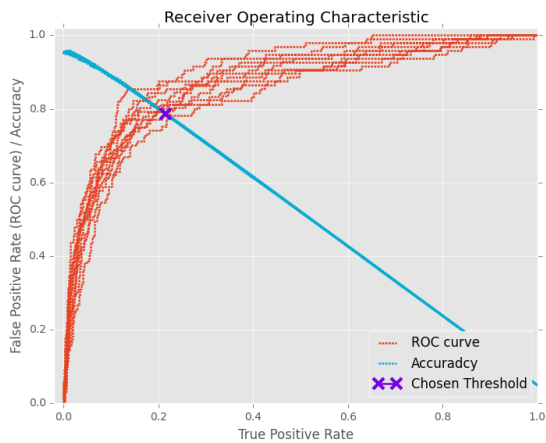
Table 3.24: Here are the results for the logistic regression confusion matrices. The main source of the advantage for fine is from the decreased amount of false negatives.

coarse-pr	fine-pr	coarse-roc	fine-roc	coarse-acc	fine-acc	coarse-f1	fine-f1
0.904	0.912	0.883	0.891	0.940	0.957	0.508	0.491
0.885	0.889	0.863	0.870	0.941	0.957	0.494	0.456
0.886	0.887	0.851	0.858	0.944	0.959	0.477	0.461
0.893	0.884	0.874	0.858	0.935	0.956	0.472	0.467
0.879	0.874	0.856	0.856	0.931	0.953	0.420	0.390
0.882	0.878	0.859	0.858	0.934	0.956	0.436	0.418
0.898	0.892	0.882	0.869	0.942	0.959	0.473	0.458
0.889	0.893	0.872	0.869	0.938	0.958	0.461	0.472
0.913	0.913	0.901	0.907	0.942	0.960	0.491	0.491
0.898	0.899	0.892	0.882	0.938	0.959	0.446	0.465
avg 0.893	avg 0.892	avg 0.873	avg 0.872	avg 0.939	avg 0.957	avg 0.468	avg 0.457

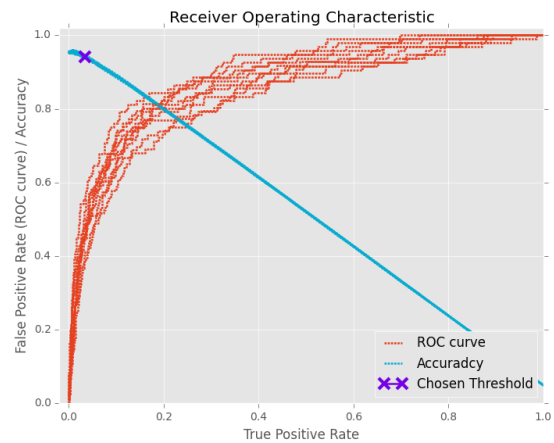
Table 3.25: Here are the results for the SVM passive 10 folds.

coarse-tn	fine-tn	coarse-fp	fine-fp	coarse-fn	fine-fn	coarse-tp	fine-tp
1828	1881	86	33	34	54	62	42
1833	1888	81	26	38	60	58	36
1847	1893	67	21	45	61	51	35
1822	1882	92	32	38	57	58	39
1822	1886	92	28	46	66	50	30
1827	1889	87	25	45	64	51	32
1842	1892	71	21	45	62	52	35
1833	1887	80	26	44	59	53	38
1836	1888	77	25	39	56	56	39
1835	1890	78	23	46	60	50	36
avg 1832.5	avg 1887.6	avg 81.1	avg 26.0	avg 42.0	avg 59.9	avg 54.1	avg 36.2

Table 3.26: Here are the results for the SVM confusion matrices. Here the fine returns less false negatives than the coarse, but not as many true positives compared to coarse.

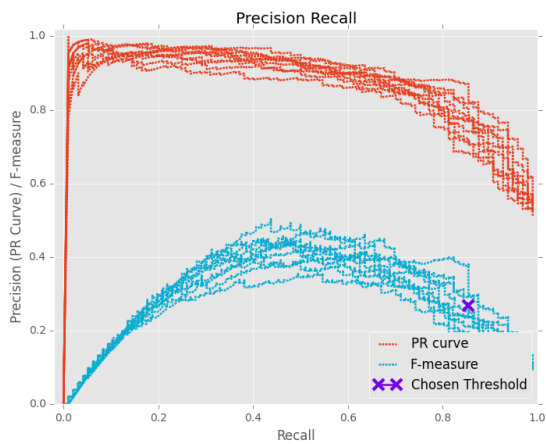


(a) Log Reg ROC Curves - coarse

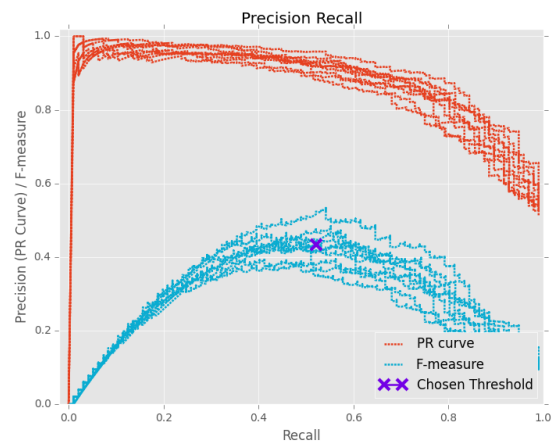


(b) Log Reg ROC Curves - fine

Figure 3.1: Fine has a higher accuracy than coarse at the default threshold.

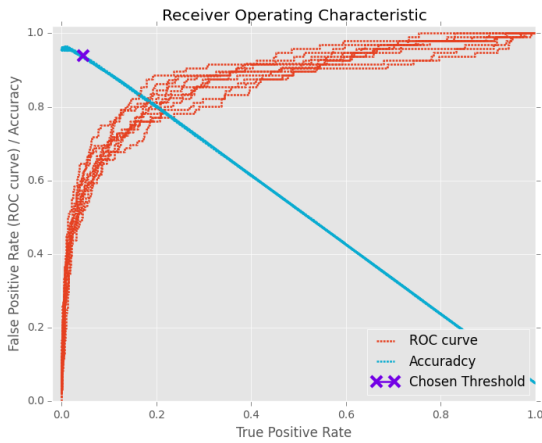


(a) Log Reg Pr Curves - coarse

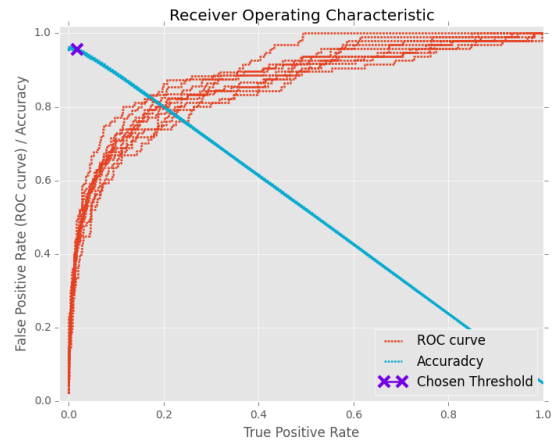


(b) Log Reg Pr Curves - fine

Figure 3.2: The fine threshold occurs at a point on the pr-curve associated with a higher f-measure than the coarse curves.

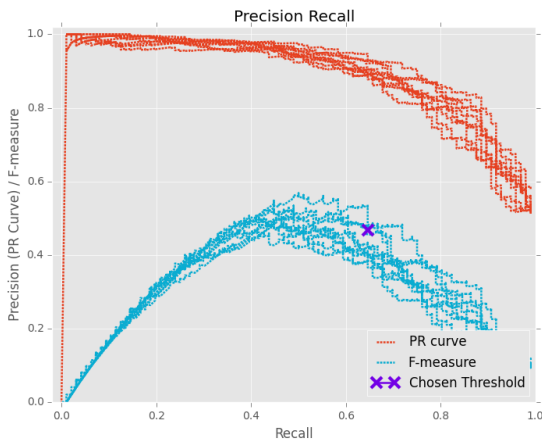


(a) SVM ROC Curves - coarse

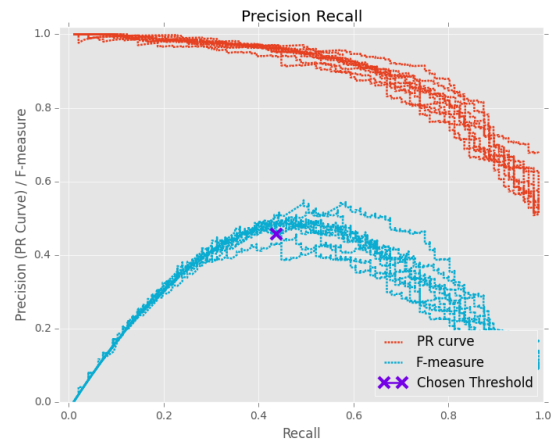


(b) SVM ROC Curves - fine

Figure 3.3: SVM results are similar between coarse and fine.



(a) SVM Pr Curves - coarse



(b) SVM Pr Curves - fine

Figure 3.4: SVM results for PR-curves and F-measure have coarse and fine picking different parts of the curves for their respective thresholds, coarse f1 avg is slightly higher at 0.468 compared to 0.457 for fine.

3.3 Active vs Passive curves

The following plots were obtained with a round batch size of 100 and a starter set of 1040 instances out of the total 2098 instances. The plots are the average of 10 folds, for each fold a test set of 2010 containing representatives of each class was held out, out of the remaining 18088, the starter set was selected which again contained representatives of each class. Coarse and fine classifiers share the same starter set. During each round coarse and fine classifiers are trained on their corresponding sets, metrics are outputted on the

held out test set, then confidence estimates are ran on the remaining eligible instances. Eligible instances are kept in separate sets for coarse and fine, 100 of the most uncertain instances are removed from each eligible set and added to its corresponding coarse or fine set to be trained on for the next round.

3.3.1 Plots for Logistic Regression Active vs Passive curves

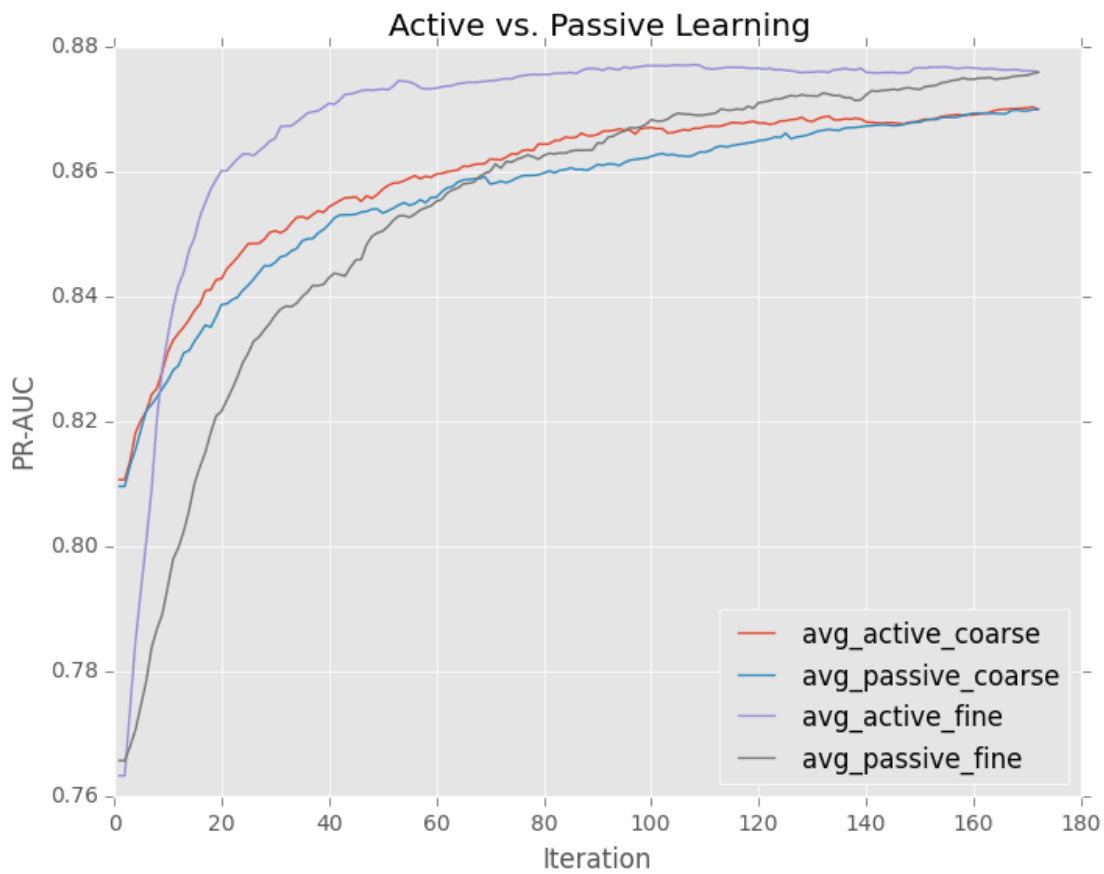


Figure 3.5: The PR AUC curves for rounds with the Logistic Regression classifier conforms to expectations, with active-fine having the highest performance. Active-coarse outperforms passive-coarse. Passive-fine doesn't outperform the coarse classifiers until rnd 100.

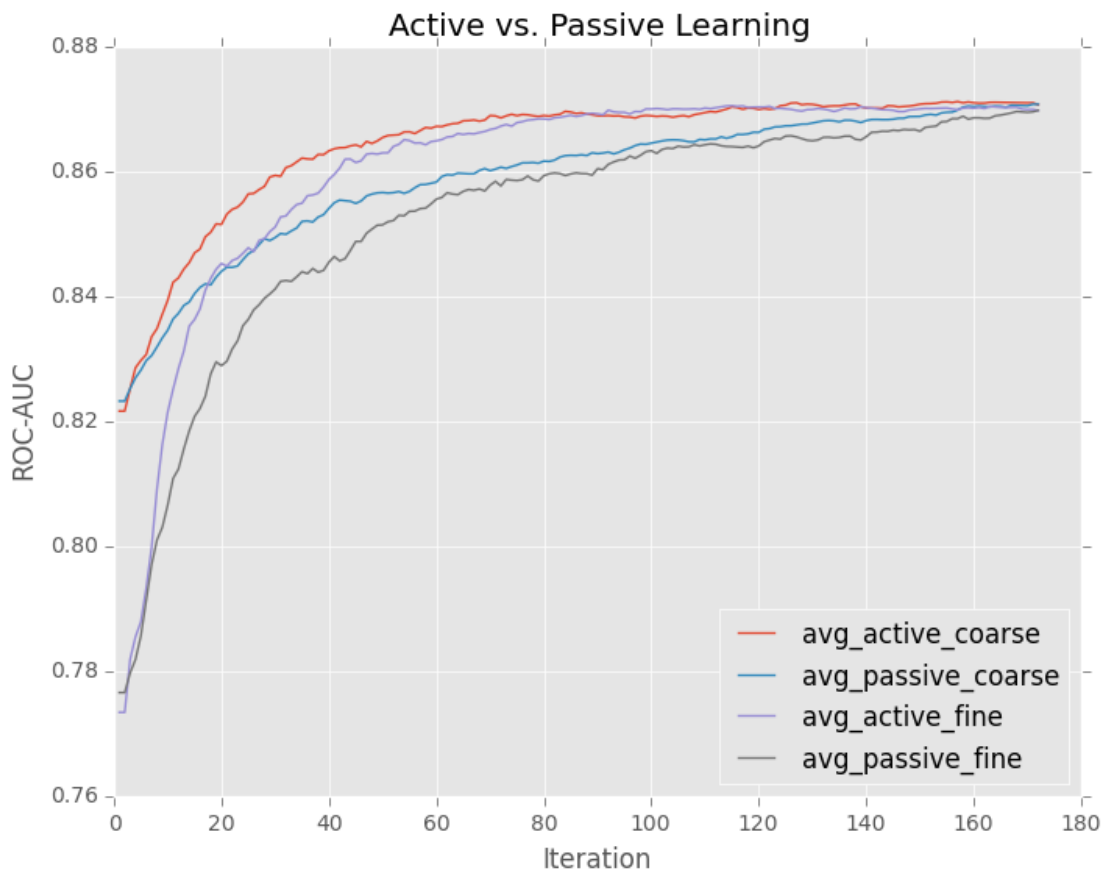


Figure 3.6: The ROC AUC curves for rounds with the Logistic Regression classifier. The active curves beat out the passive curves for both coarse and fine. Coarse roc starts with an advantage over fine as in the PR curves. Both converge to the same rate after roc auc level after 80.

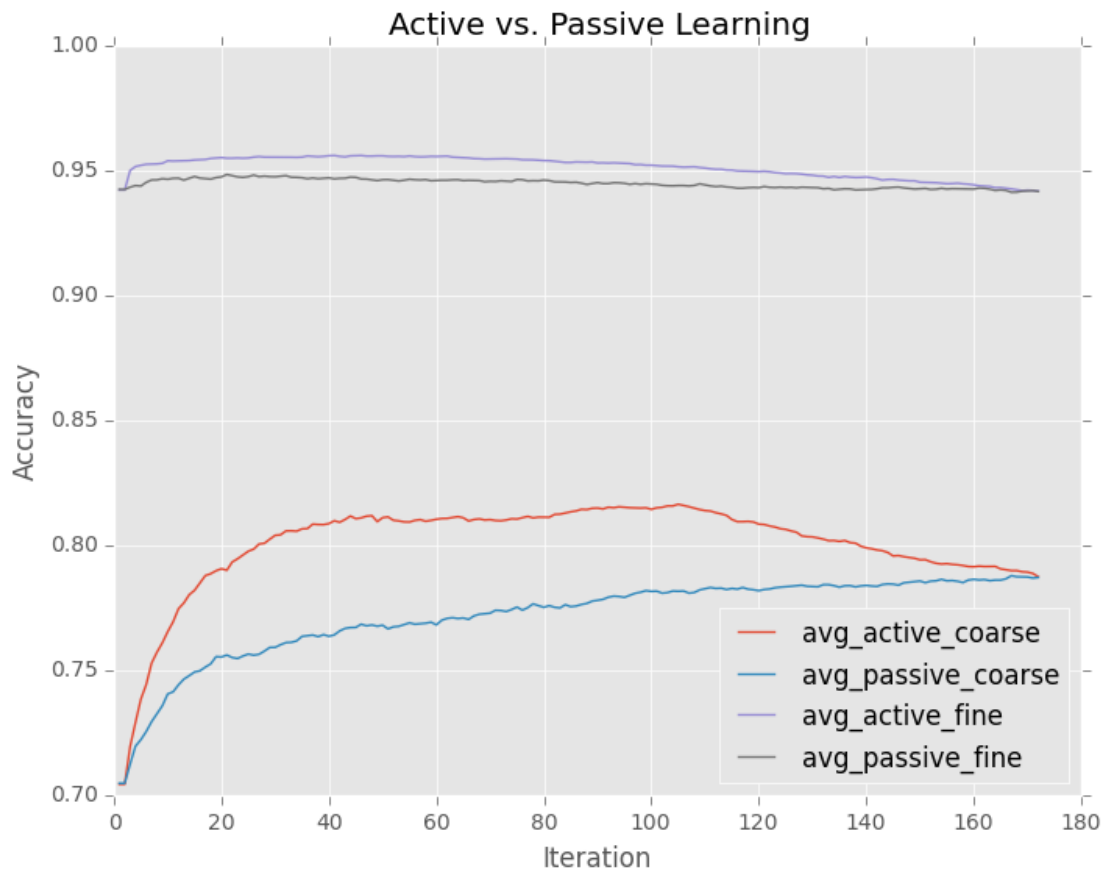


Figure 3.7: The accuracy of the fine classifiers stays at roughly the same rate throughout the rounds, this is due to an effective weighting scheme for the fine grained classifiers. The active coarse accuracy drops towards the end due to an increase in false positives as more negative instances are added in the later rounds.

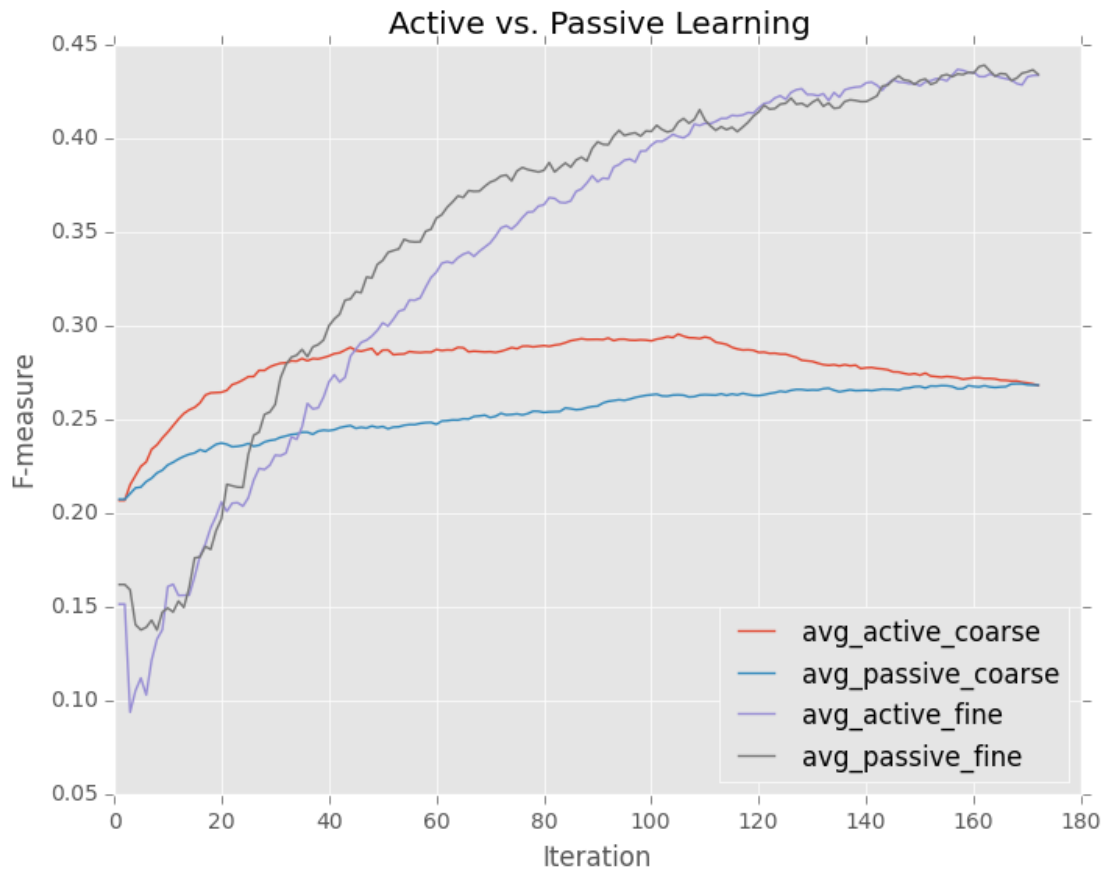


Figure 3.8: The F-measure of the the fine classifiers increases throughout the rounds as more true positives are predicted. The active coarse again decreases at later rounds due to increased false positives.

3.3.2 Plots for SVM Active vs Passive curves

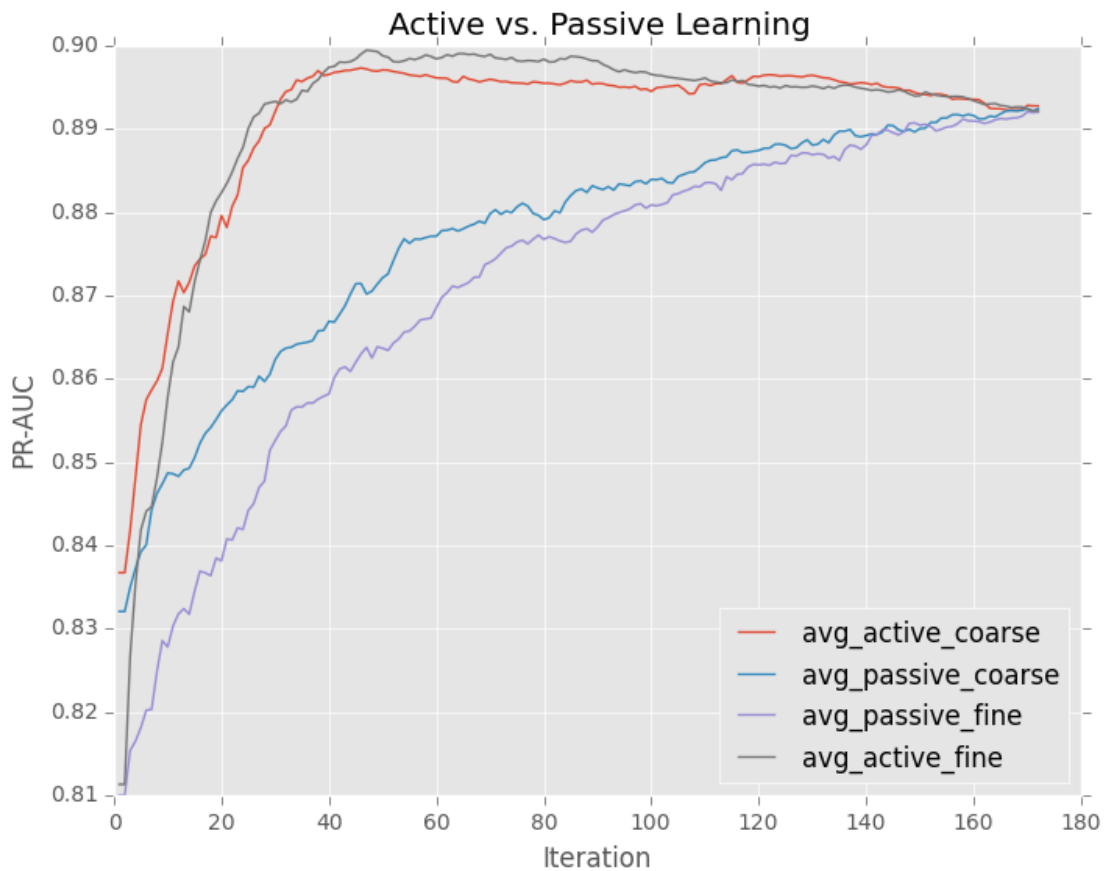


Figure 3.9: The PR AUC curves for rounds with SVM show little advantage for fine. The results are slightly different than the ones shown on 2/14 due to fixing a bug with the code that wasn't performing the preprocessing scaling for the SVM case at the same stage as it was being done for the logistic regression classifier.

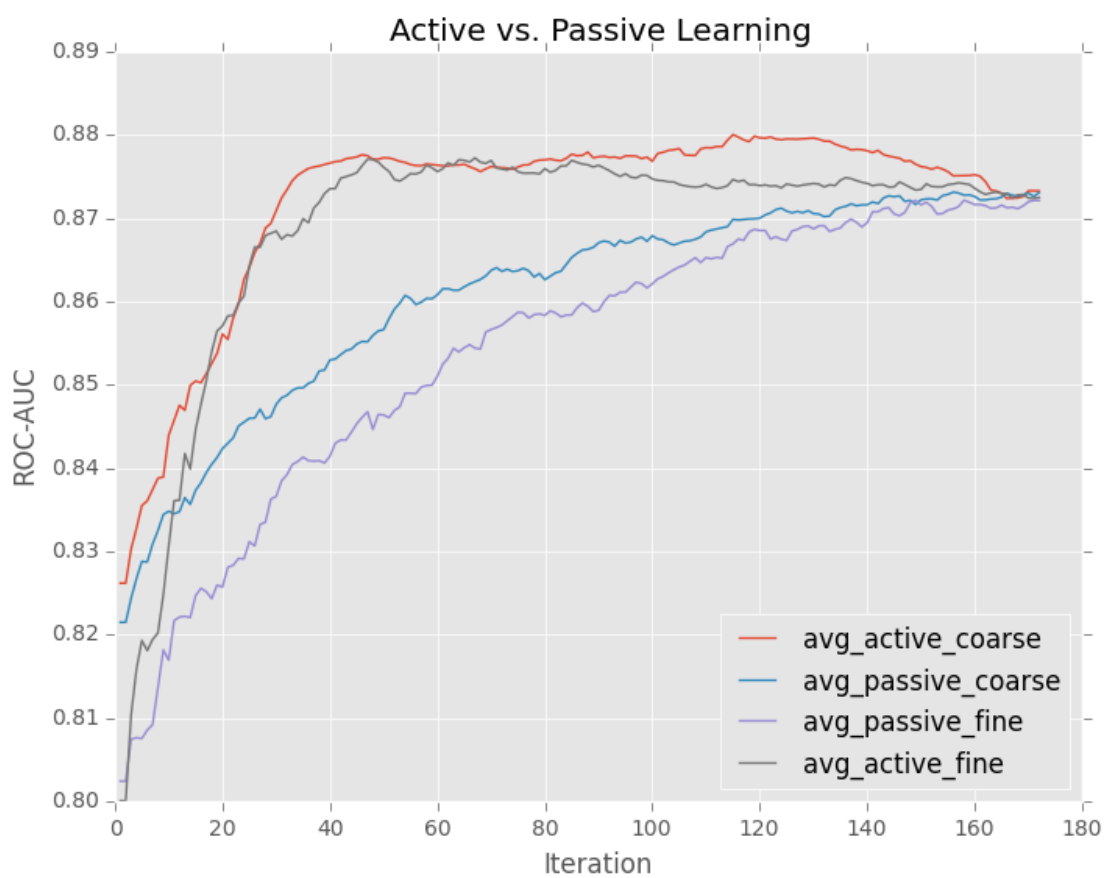


Figure 3.10: The ROC curves show more of an advantage for coarse classifiers.

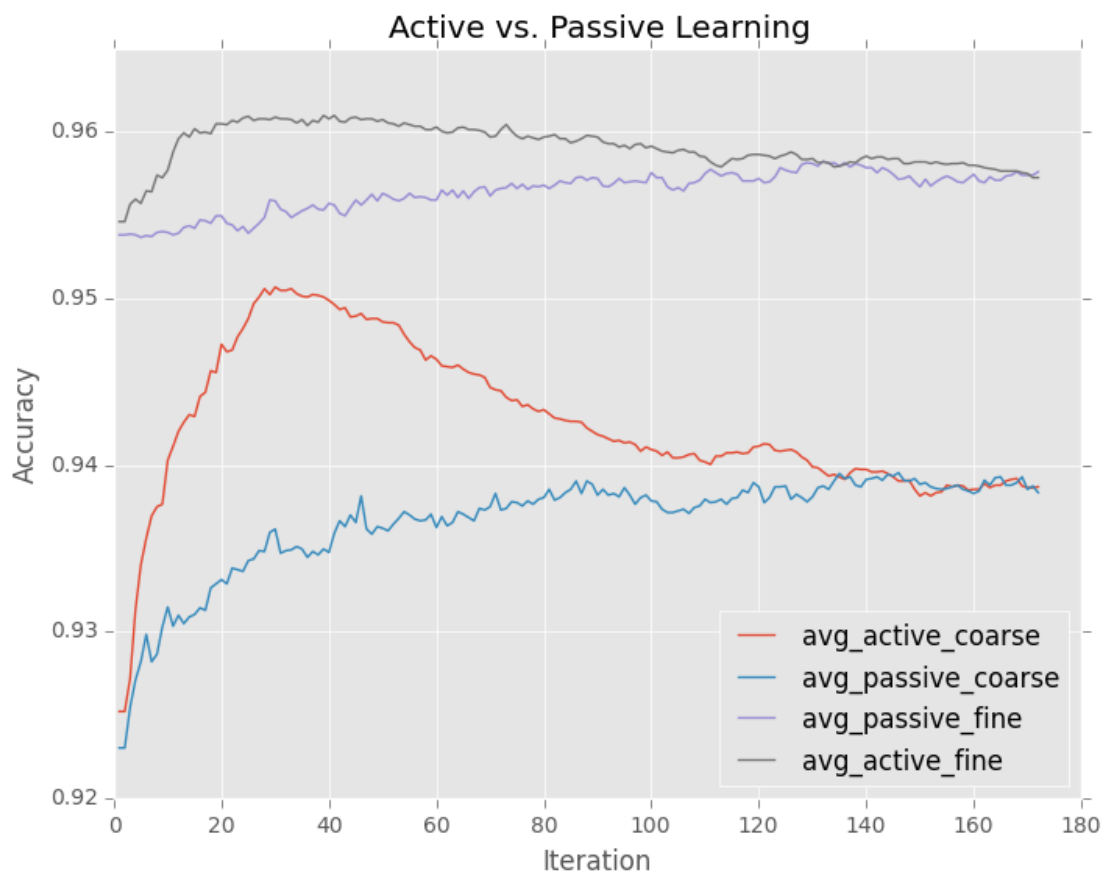


Figure 3.11: The accuracy for the coarse decreases sharply due to coarse predicting steadily more false positives, behaving similar to the Log Reg case. Fine accuracy is higher due to predicting less false positives than coarse. Fine also predicts less true positives, compare apx. 37 to apx. 60 t.p. for coarse at round 60.

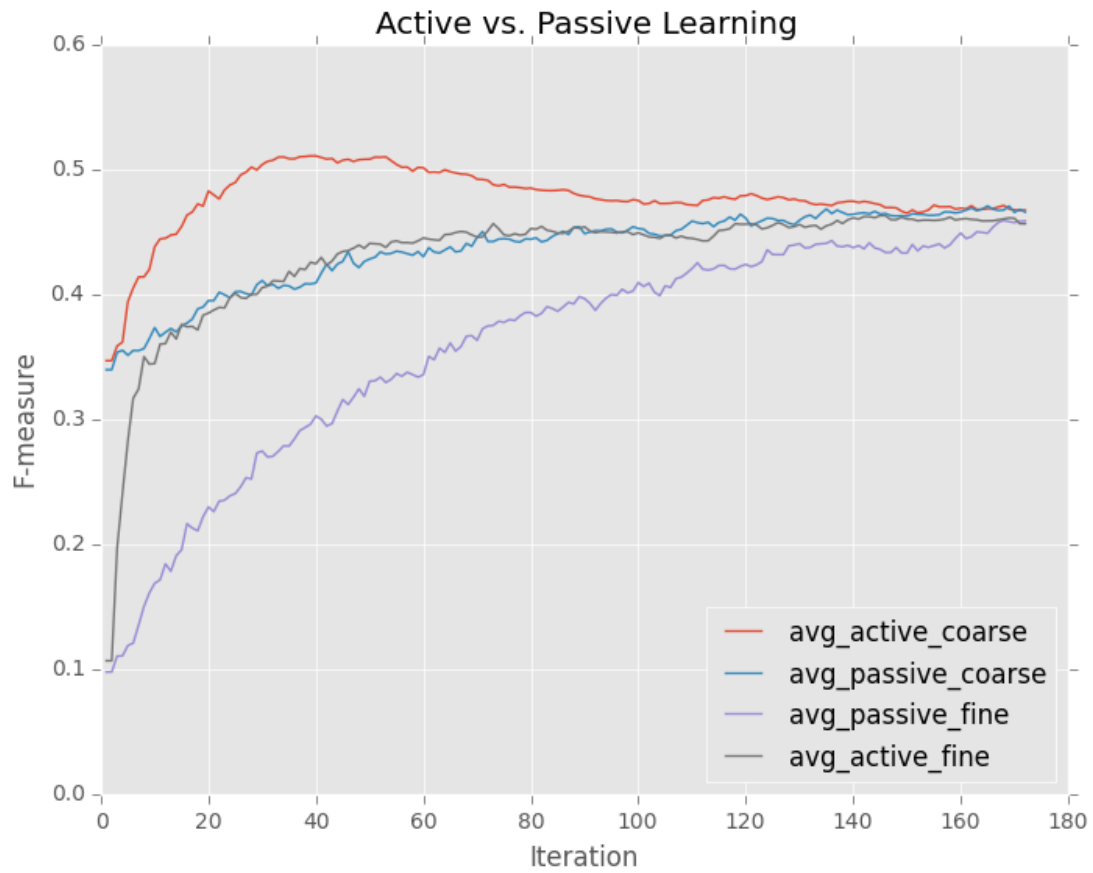


Figure 3.12: The F-measure favors coarse, and trends to the same level for both coarse and fine.

3.4 Plots for FFR experiments

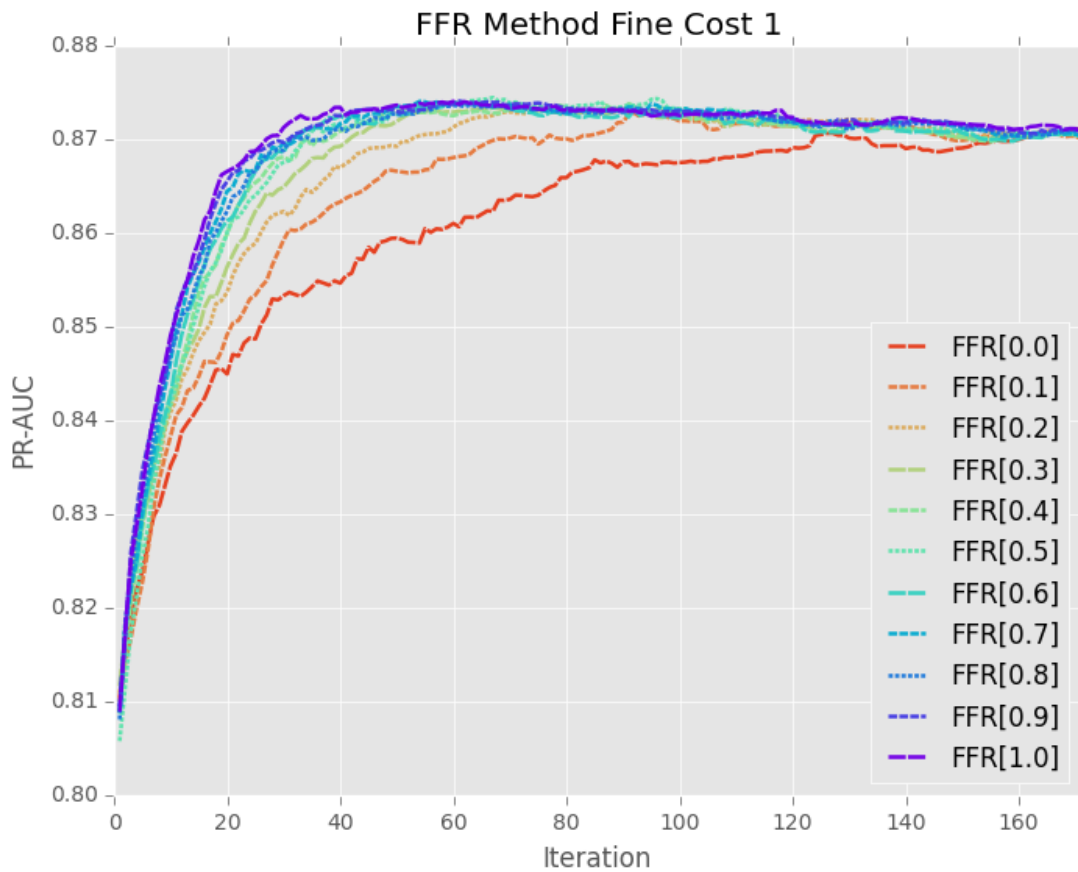


Figure 3.13: The strategy is changed from purchasing a set number of instances per round to having a set budget per round and spending a portion of that budget on fine and coarse grained labels. For this curve the fine and coarse grain labels both have a cost of 1. The purple 1.0 curve shows that if only fine grained labels are purchased, the highest performing PR-AUC can be obtained. The results are an average of 10 folds.

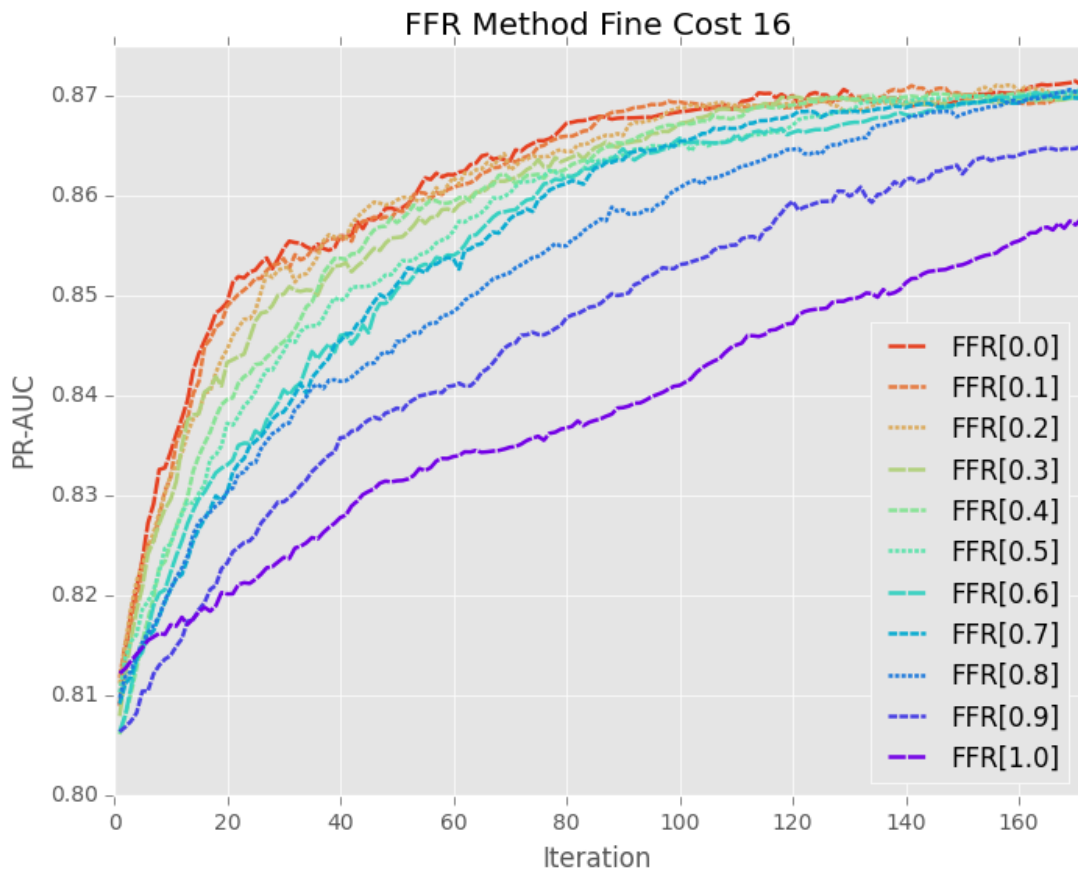


Figure 3.14: The fine cost is increased to 16. The budget for each iteration is 100, and for the case of the 0.5 curve, 50 instances are bought for coarse and 3.125 instances are bought for fine. The remainder 0.125 is then turned into a 0.125 chance for any round to purchase an extra fine label. The round size for the FFR 1.0 curve is very small, with only 7 labels purchased per iteration. The cost is too high for the fine label advantage to offset the decreased number of instances purchased.

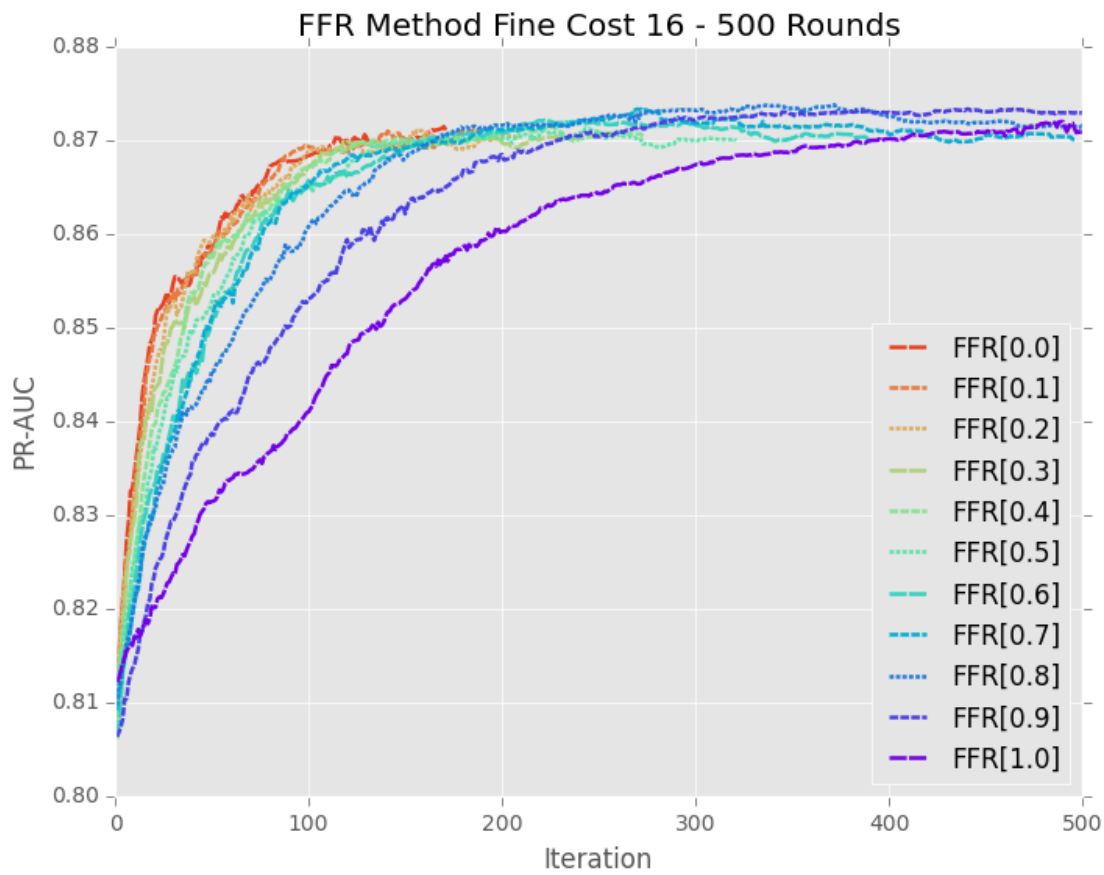


Figure 3.15: This shows the iterations continuing through round 500, the curves with the higher fine rates eventually settle to the same end point that the curves with the high rates of coarse labels purchased achieved at previous iterations.

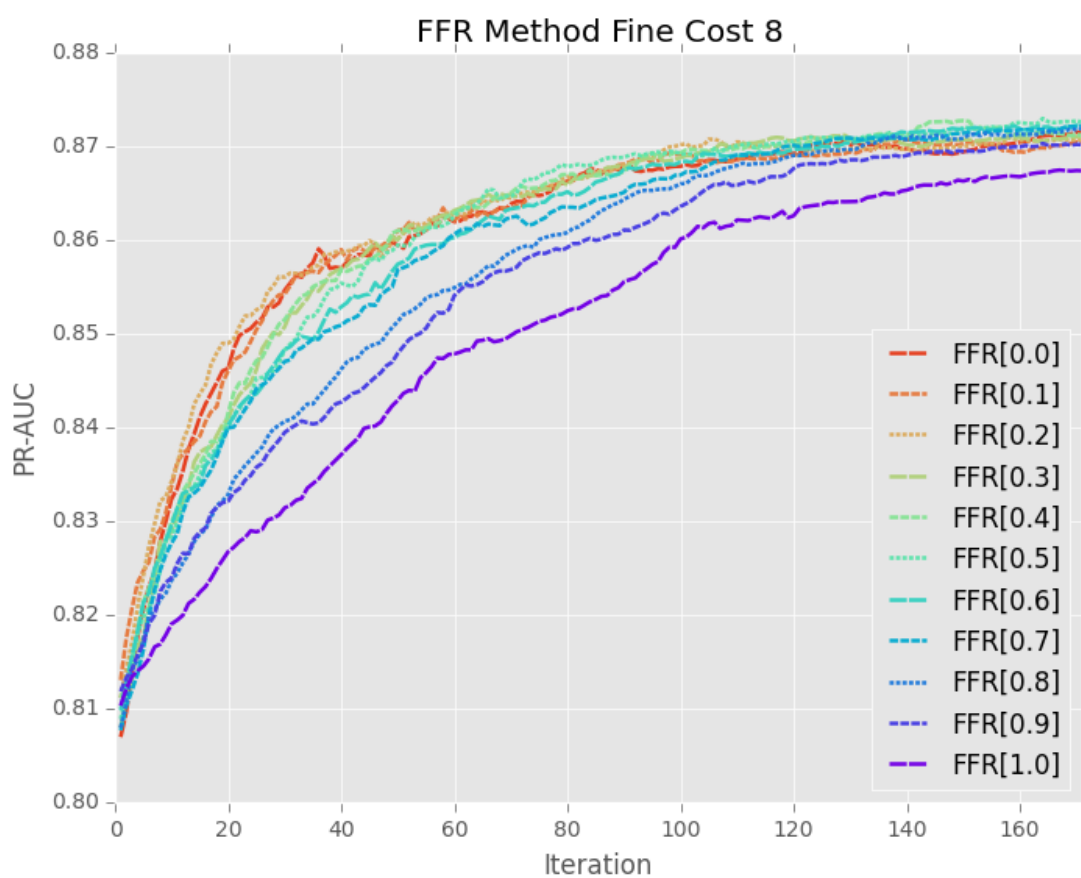


Figure 3.16: At fine cost 8 the FFR 0.0 rate is no longer the best option, 0.1 generally outperforms 0.0 slightly.

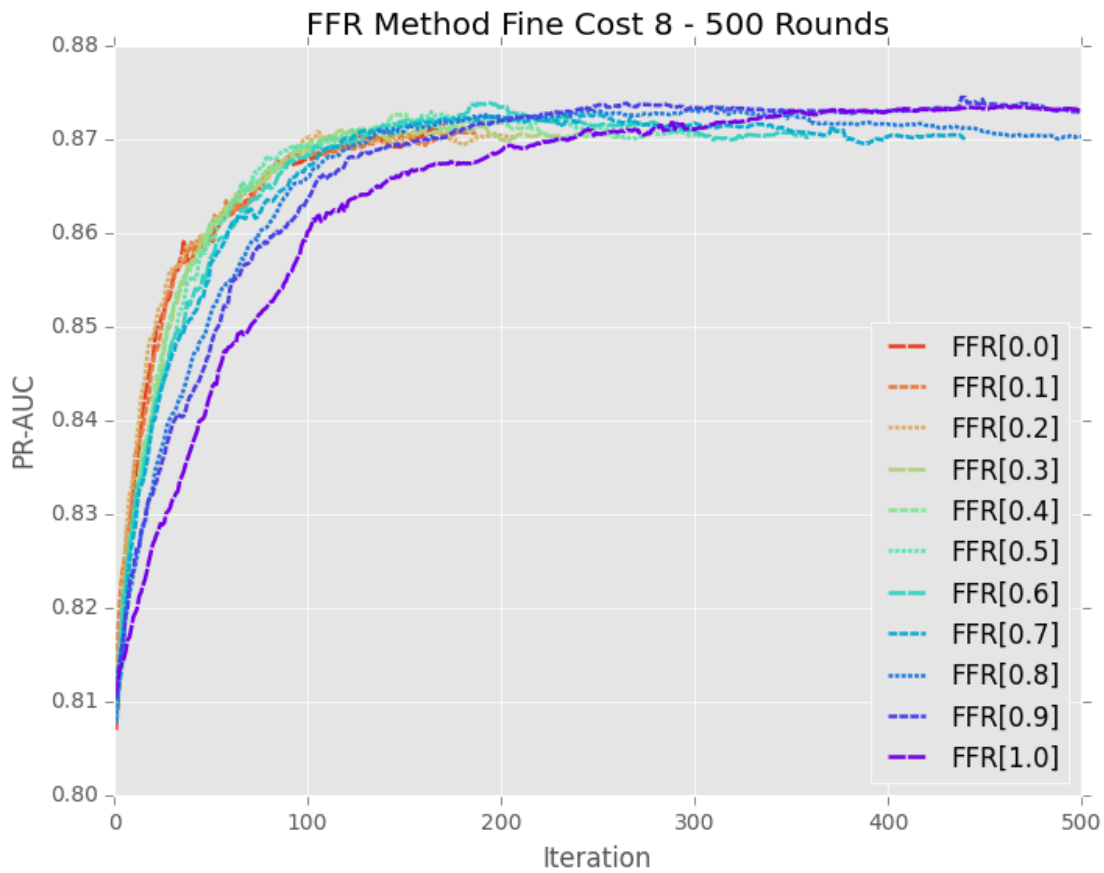


Figure 3.17: The extended picture of the FFR cost 8. The round size for FFR 1.0 is small, only 13 instances purchase per iteration.

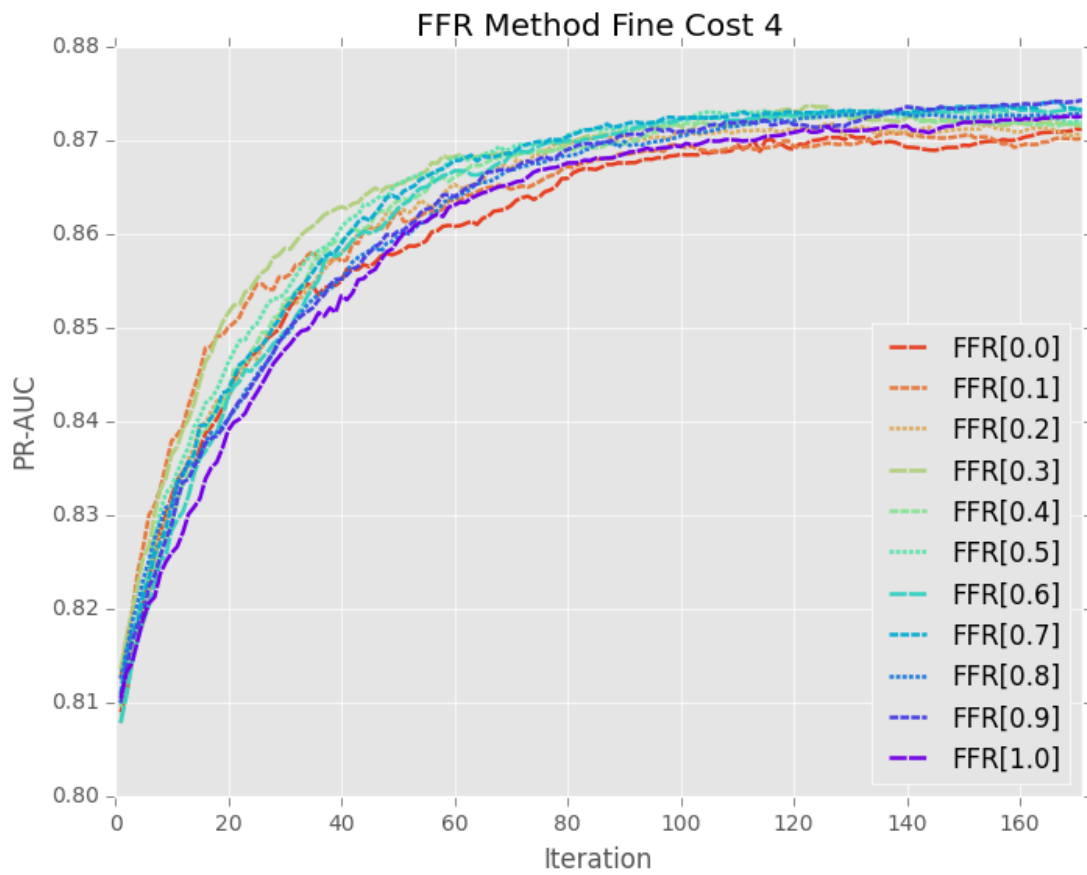


Figure 3.18: At fine cost 4, FFR 0.3 appears to be the highest performing rate.

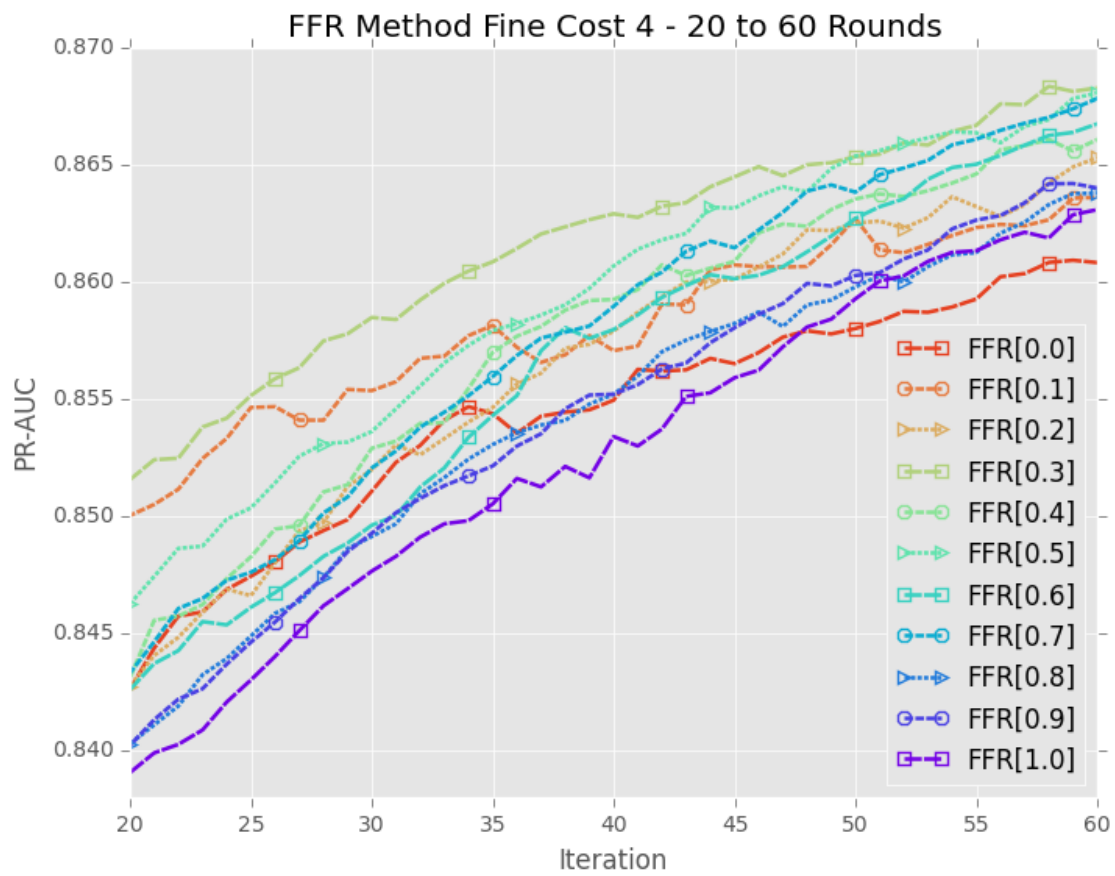


Figure 3.19: The fine cost 4 curves shown expanding the rounds 20-60.

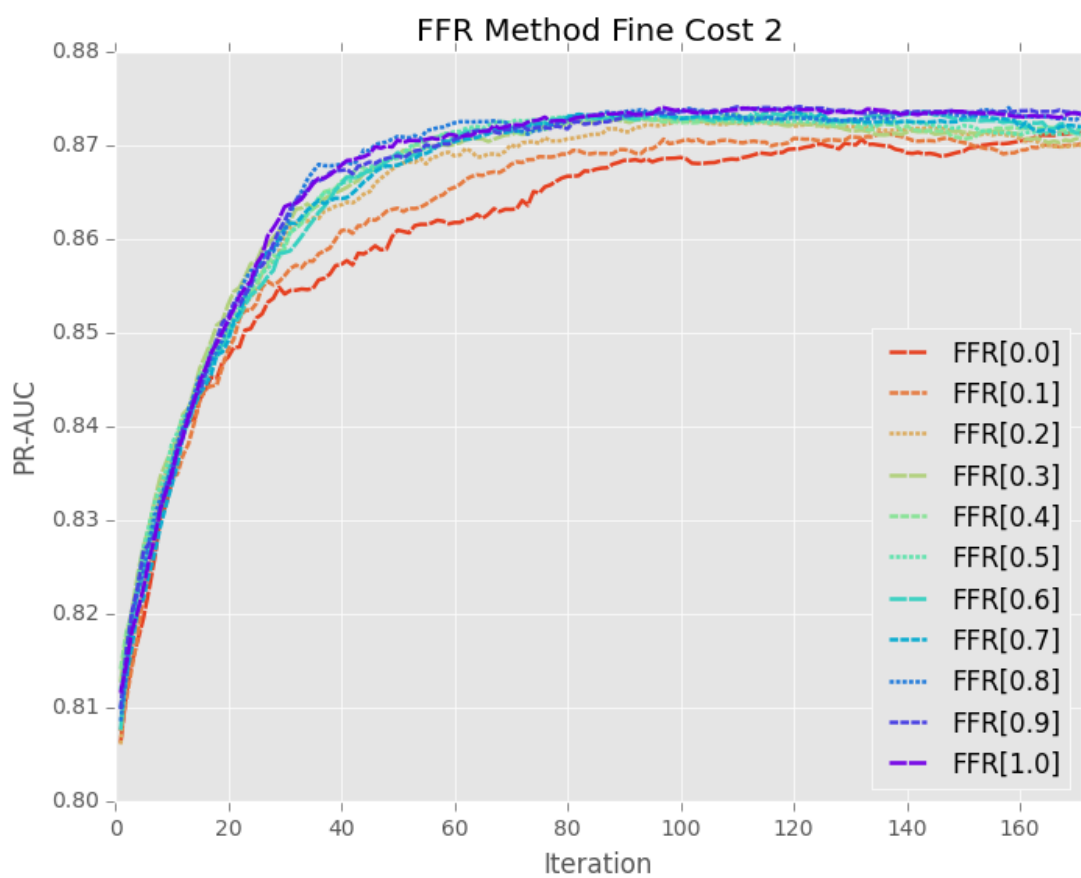


Figure 3.20: At fine cost 2, the preferred rate jumps up to 0.8, similar to the cost 1 results.

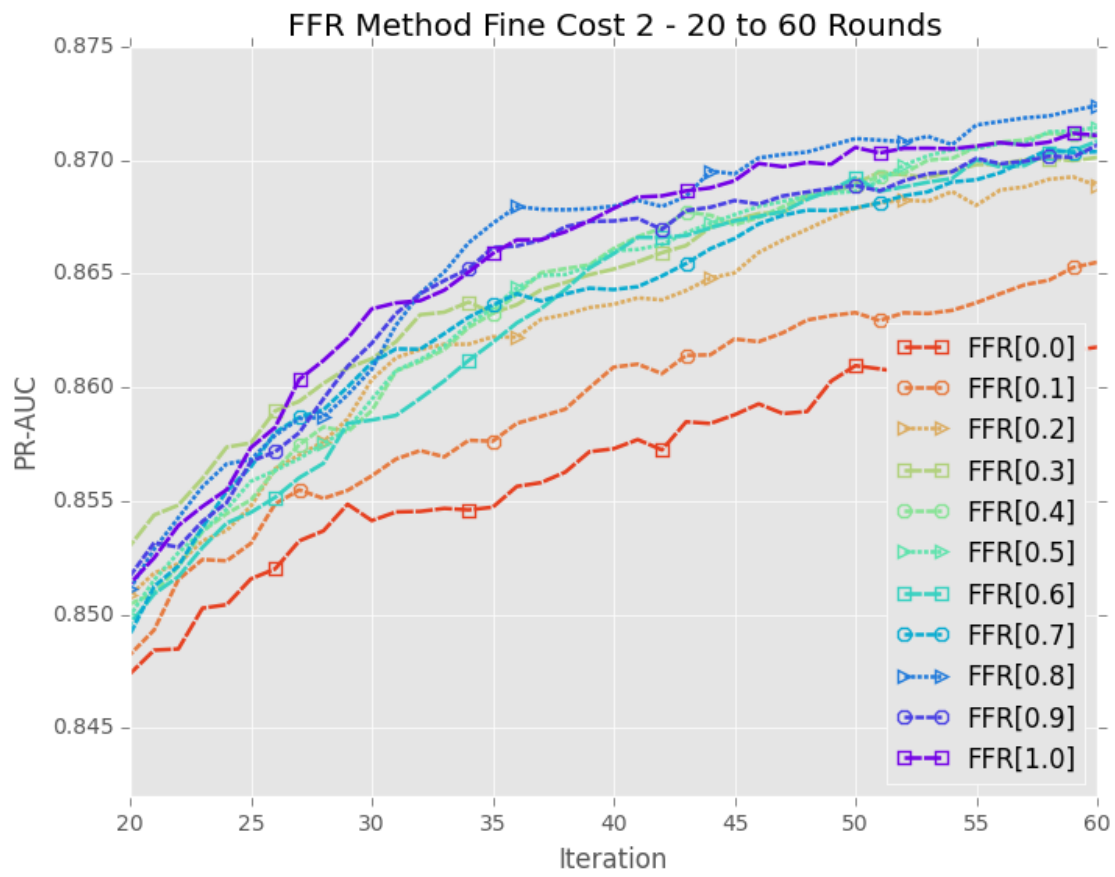


Figure 3.21: The fine cost 2 curves shown expanding rounds 20-60.

Chapter 4

Conclusions and Future Work

I should probably do the Bandit experiments.

Appendix A

Tuning the fine grained classes

All the data and results for tuning the fine grained classes.

Bibliography

- [1] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011. 3.1
- [2] A. Merialdo, "Improving Collaborative Filtering For New-Users By Smart Object Selection," *In Proceedings of International Conference on Media Features (ICMF)*, May 2001. [Online]. Available: <http://www.eurecom.fr/publication/670><https://www.eurecom.fr/fr/publication/670/download/mm-kohrar-010508.pdf>
- [3] W. Luo, A. Schwing, and R. Urtasun, "Latent structured active learning," in *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [4] S. Dasgupta and D. Hsu, "Hierarchical sampling for active learning," *Proceedings of the 25th international conference on Machine learning - ICML '08*, pp. 208–215, 2008. [Online]. Available: <http://portal.acm.org/citation.cfm?doid=1390156.1390183>
- [5] X. Ling and D. Weld, "Fine-grained entity recognition," *Proceedings of the 26th Conference on Artificial Intelligence*, 2012. [Online]. Available: <http://www.cs.washington.edu/ai/pubs/ling-aaai12.pdf><http://www.aaai.org/ocs/index.php/AAAI/AAAI12/paper/download/5152/5124>
- [6] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, and G. Varoquaux, "API design for machine learning software: experiences from the scikit-learn project," in *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, 2013, pp. 108–122.