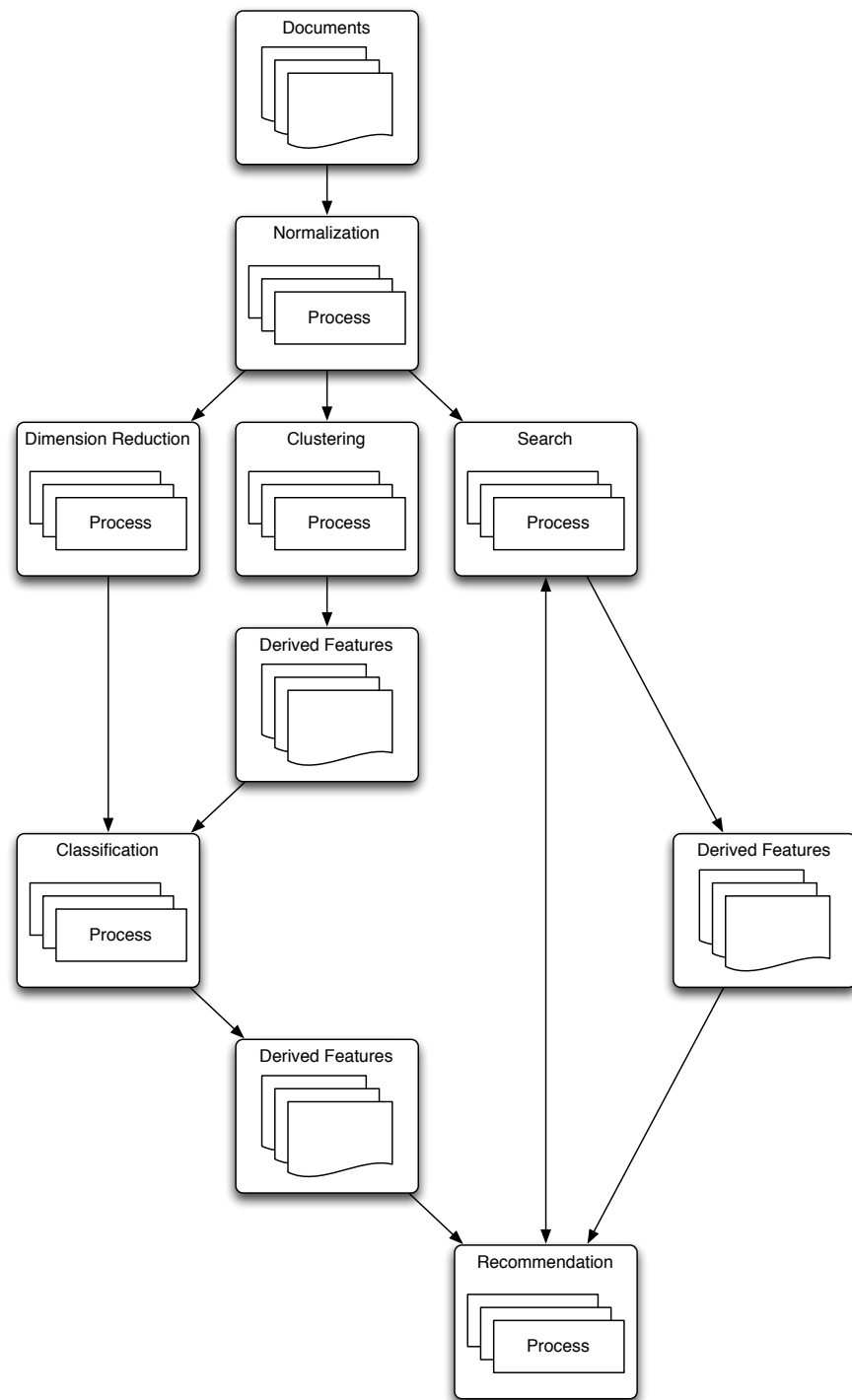


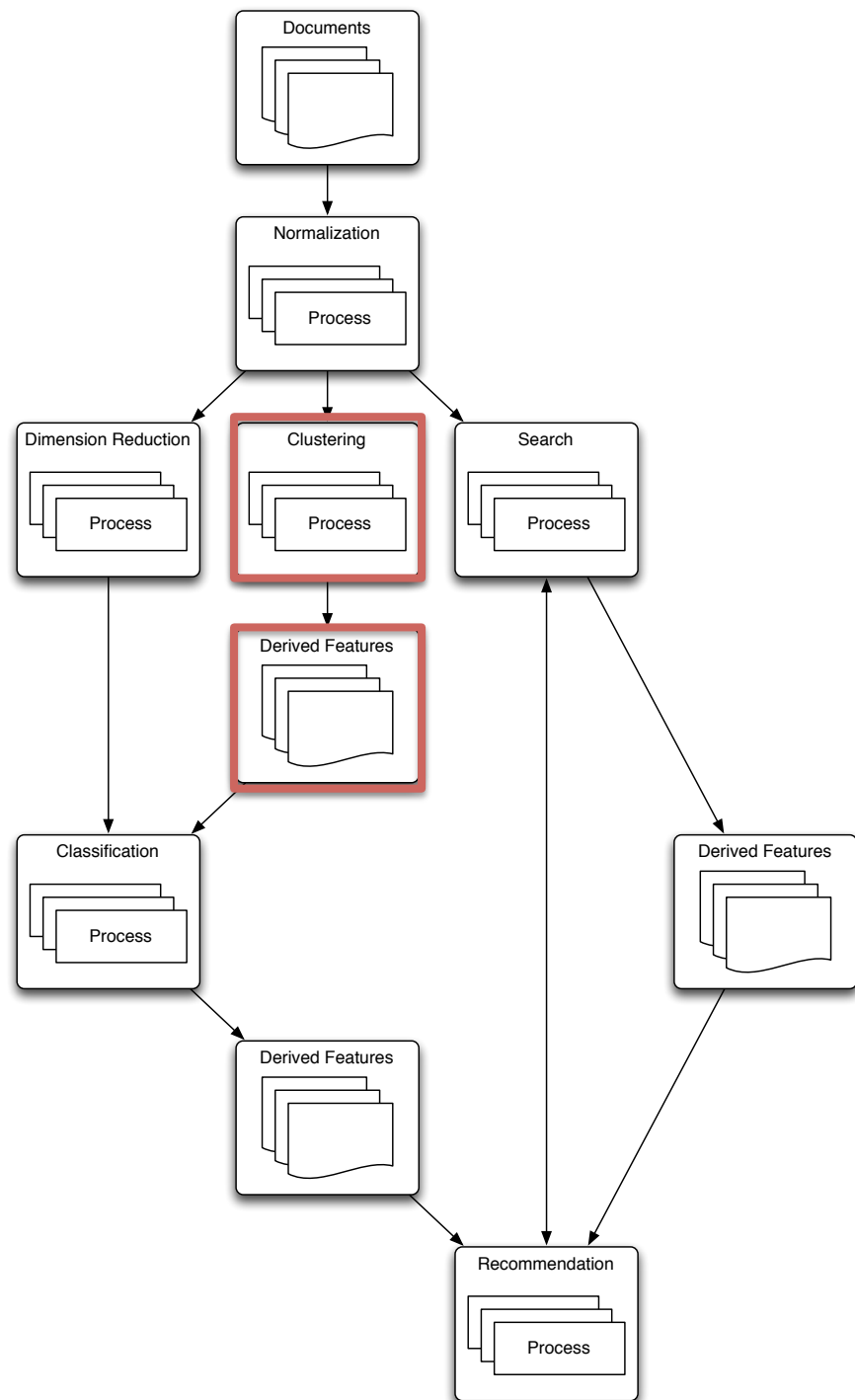
Mahout Workshop, Section 4

Allen Day, PhD

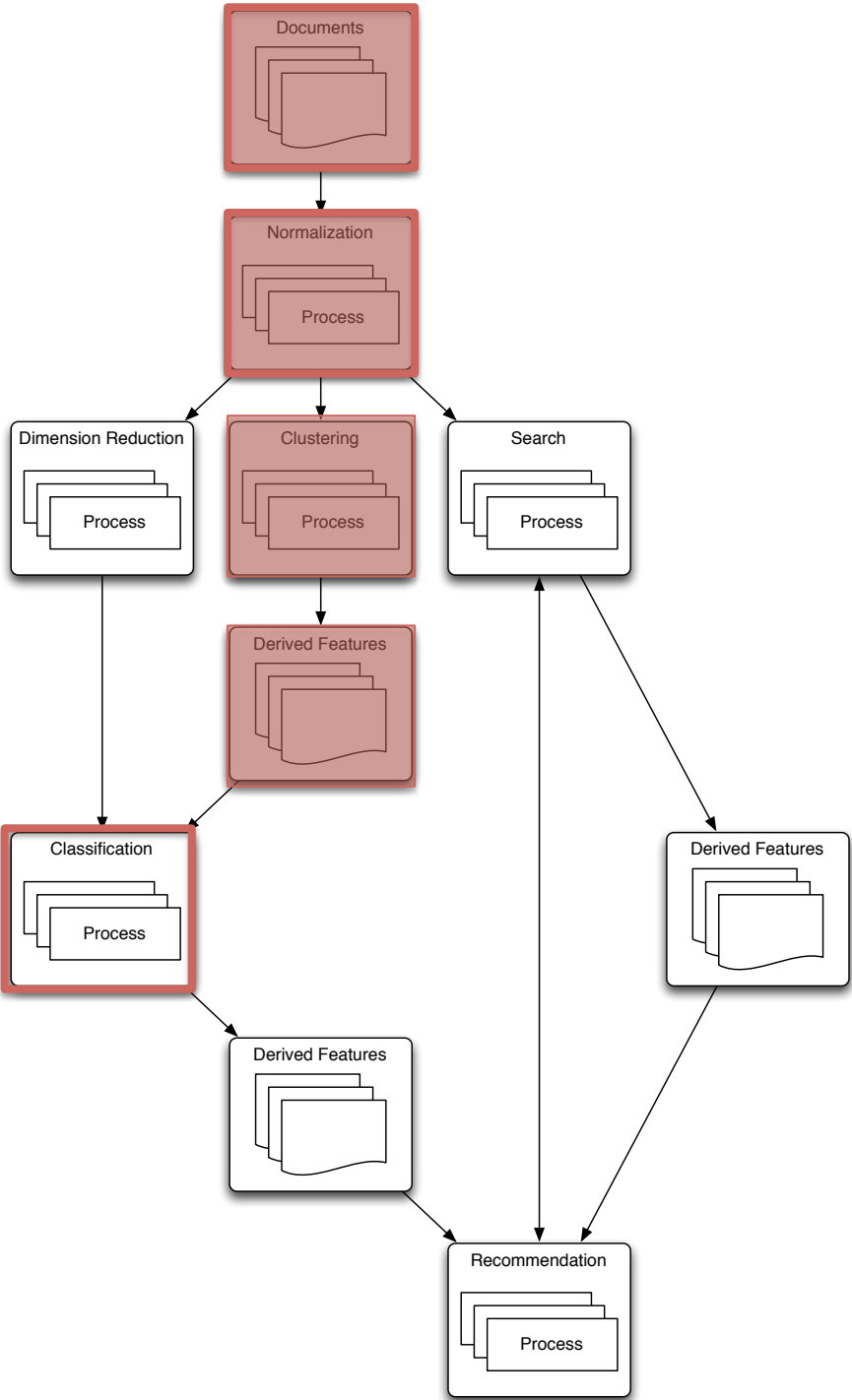
MapR Technologies



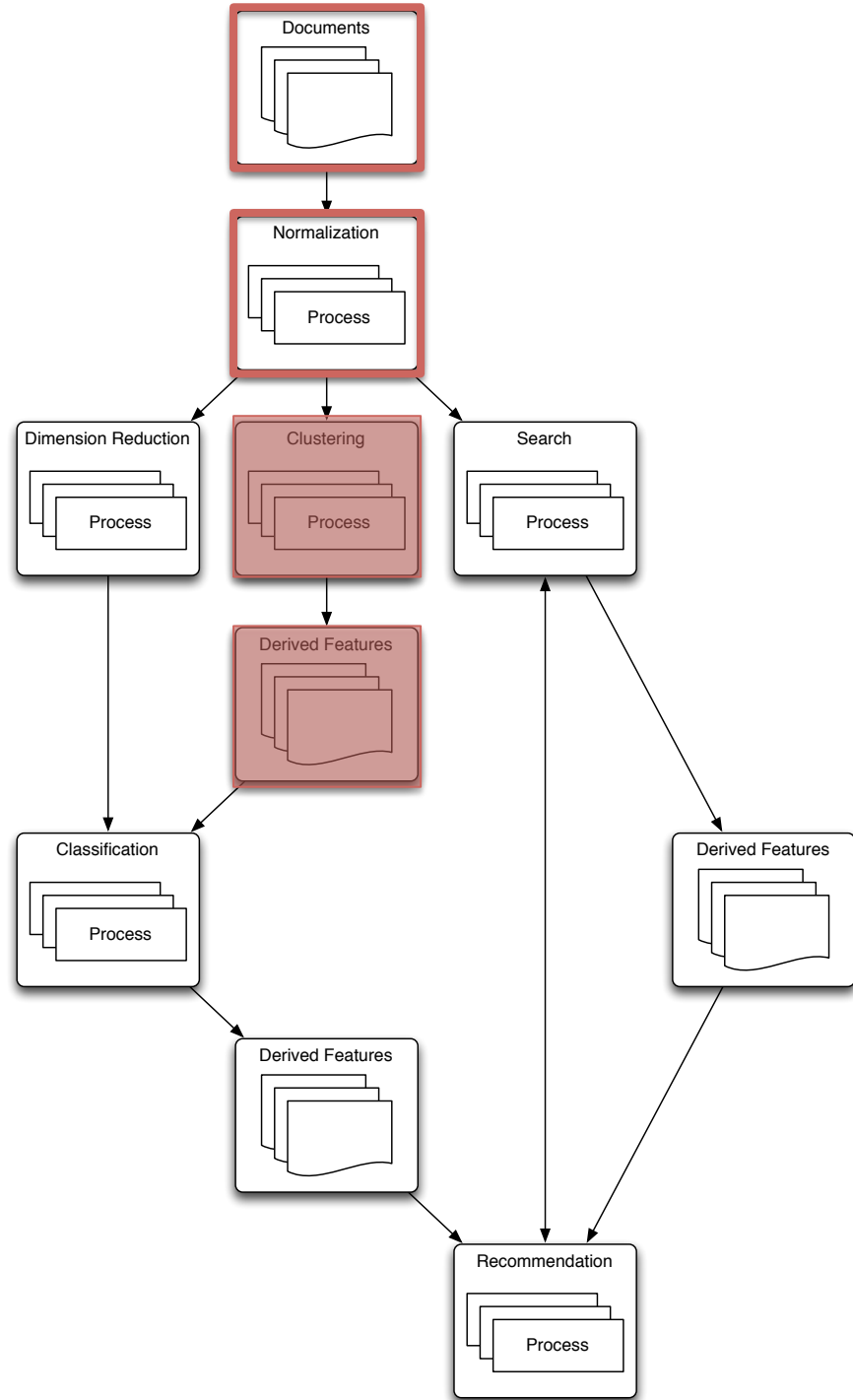
Section 3 (this one)



Section 4 (next one)

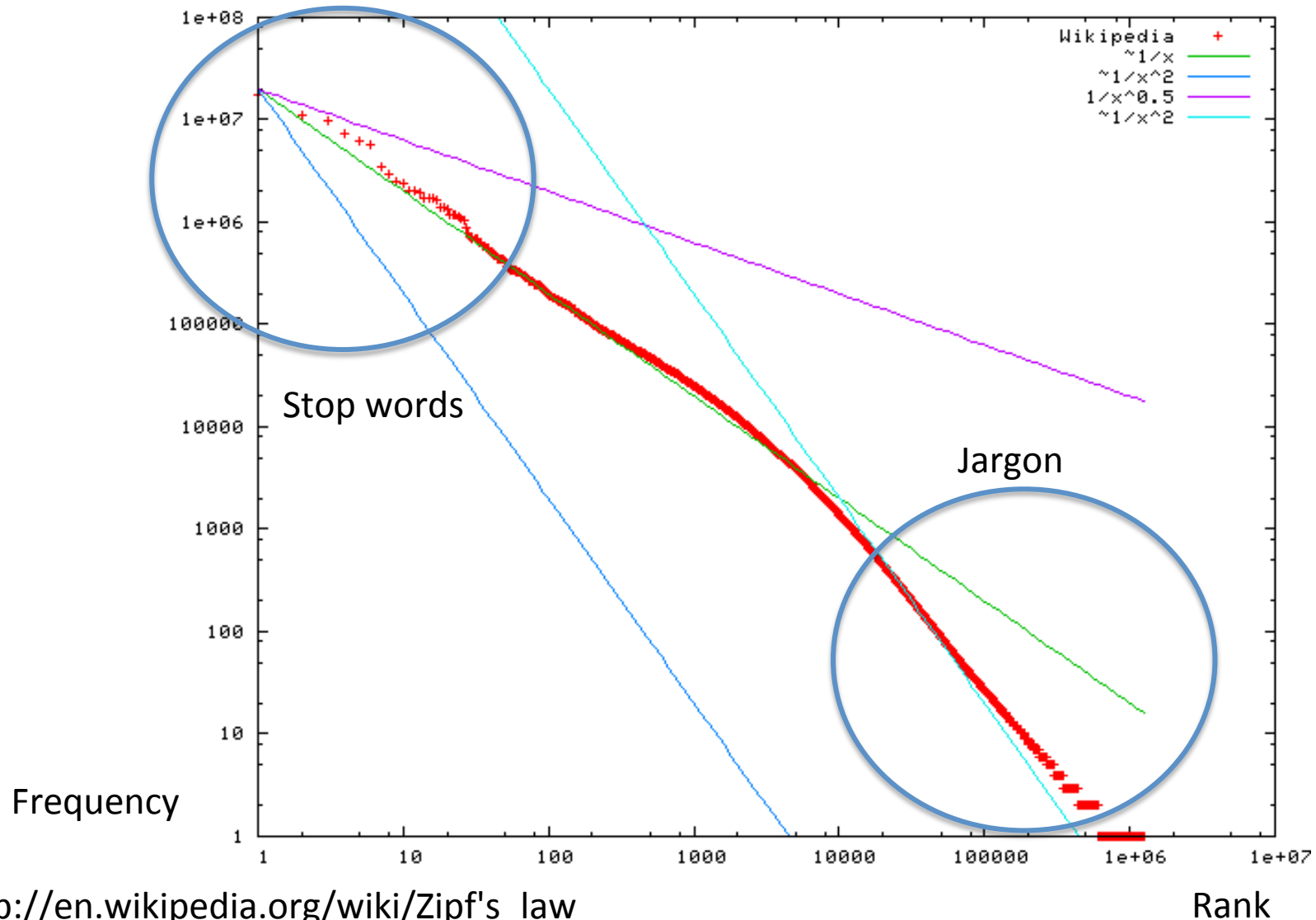


Section 4 (next one)



TF-IDF

Observation: Unequal Word Distributions



Motivation:

“Undo” Effect of Unequal Distribution

	doc1	doc2	doc3	TF (max)	DF(bool)
car	27	4	24	27	3
auto	3	33	0	33	2
insurance	0	33	29	33	2
best	14	0	17	17	2

	doc1	doc2	doc3	TF (max)	DF (bool)
car	$27/27 * \log(3/3)$	$4/27 * \log(3/3)$	$24/27 * \log(3/3)$	27	3
auto	$3/33 * \log(3/2)$	$33/33 * \log(3/2)$	$0/33 * \log(3/2)$	33	2
insurance	$0/33 * \log(3/2)$	$33/33 * \log(3/2)$	$29/33 * \log(3/2)$	33	2
best	$14/17 * \log(3/2)$	$0/17 * \log(3/2)$	$17/17 * \log(3/2)$	17	2

Motivation:

“Undo” Effect of Unequal Distribution

	doc1	doc2	doc3	TF (max)	DF (bool)
car	$27/27 * \log(3/3)$	$4/27 * \log(3/3)$	$24/27 * \log(3/3)$	27	3
auto	$3/33 * \log(2/3)$	$33/33 * \log(2/3)$	$0/33 * \log(2/3)$	33	2
insurance	$0/33 * \log(2/3)$	$33/33 * \log(2/3)$	$29/33 * \log(2/3)$	33	2
best	$14/17 * \log(2/3)$	$0/17 * \log(2/3)$	$17/17 * \log(2/3)$	17	2

	doc1	doc2	doc3	TF (max)	DF (bool)
car	0	0	0	27	3
auto	0.04	0.40	0	33	2
insurance	0	0.40	0.35	33	2
best	0.14	0	0.40	17	2

Note that I've used TF_{\max} and DF_{bool} here. Other tf*idf variants are also valid.

Motivation:

“Undo” Effect of Unequal Distribution

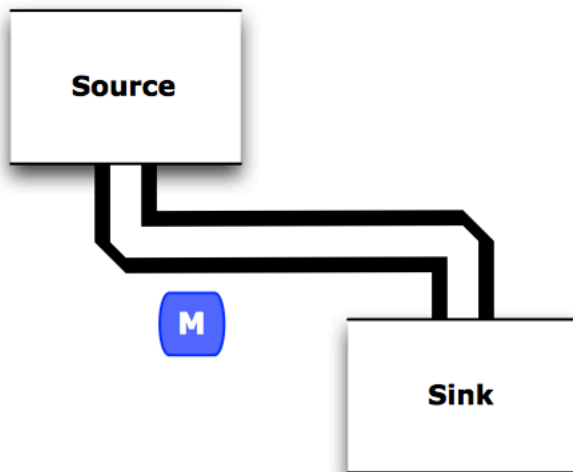
	doc1	doc2	doc3	TF (sum)	DF (bool)
car	$27/55 * \log(3/3)$	$4/55 * \log(3/3)$	$24/55 * \log(3/3)$	55	3
auto	$3/36 * \log(3/2)$	$33/36 * \log(3/2)$	$0/36 * \log(3/2)$	36	2
insurance	$0/62 * \log(3/2)$	$33/62 * \log(3/2)$	$29/62 * \log(3/2)$	62	2
best	$14/31 * \log(3/2)$	$0/31 * \log(3/2)$	$17/31 * \log(3/2)$	31	2

	doc1	doc2	doc3	TF (sum)	DF (bool)
car	0	0	0	27	3
auto	0.03	0.37	0	33	2
insurance	0	0.21	0.19	33	2
best	0.18	0	0.22	17	2

Note that I've used TF_{\max} and DF_{bool} here. Other tf*idf variants are also valid.

TF-IDF IN CASCADING (ON HADOOP)

1: copy



1 mapper
0 reducers
10 lines code

```
public class
    Main
    {
        public static void
        main( String[] args )
        {
            String inPath = args[ 0 ];
            String outPath = args[ 1 ];

            Properties props = new Properties();
            AppProps.setApplicationJarClass( props, Main.class );
            HadoopFlowConnector flowConnector = new HadoopFlowConnector( props
        );

            // create the source tap
            Tap inTap = new Hfs( new TextDelimited( true, "\t" ), inPath );

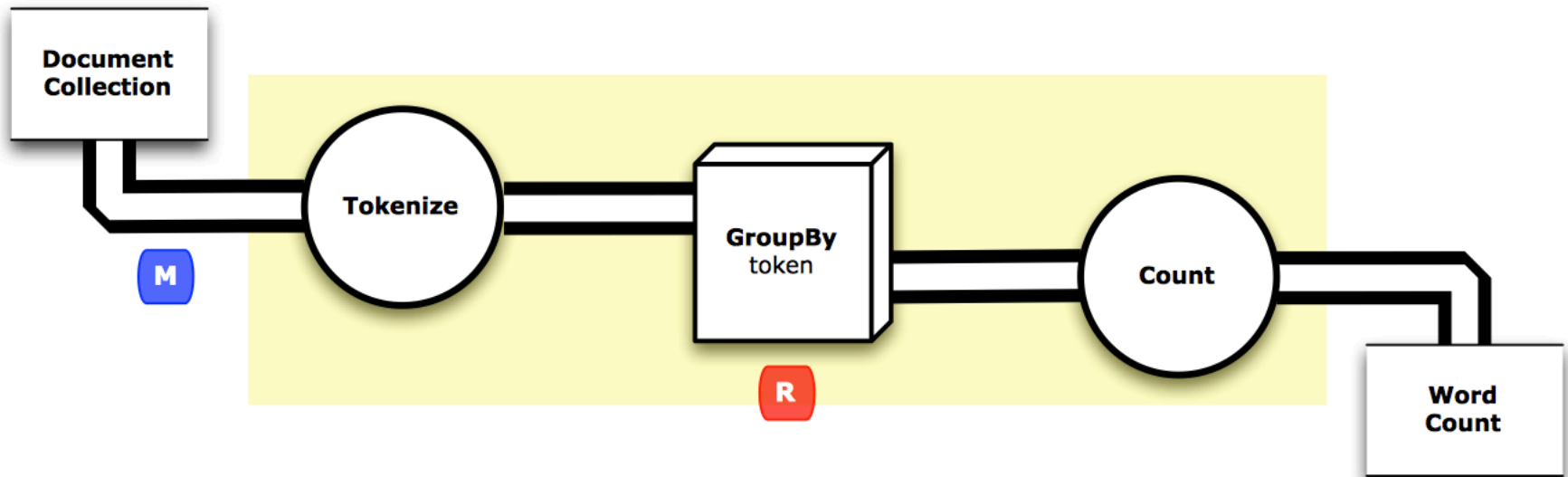
            // create the sink tap
            Tap outTap = new Hfs( new TextDelimited( true, "\t" ), outPath );

            // specify a pipe to connect the taps
            Pipe copyPipe = new Pipe( "copy" );

            // connect the taps, pipes, etc., into a flow
            FlowDef flowDef = FlowDef.flowDef().setName( "copy" )
                .addSource( copyPipe, inTap )
                .addTailSink( copyPipe, outTap );

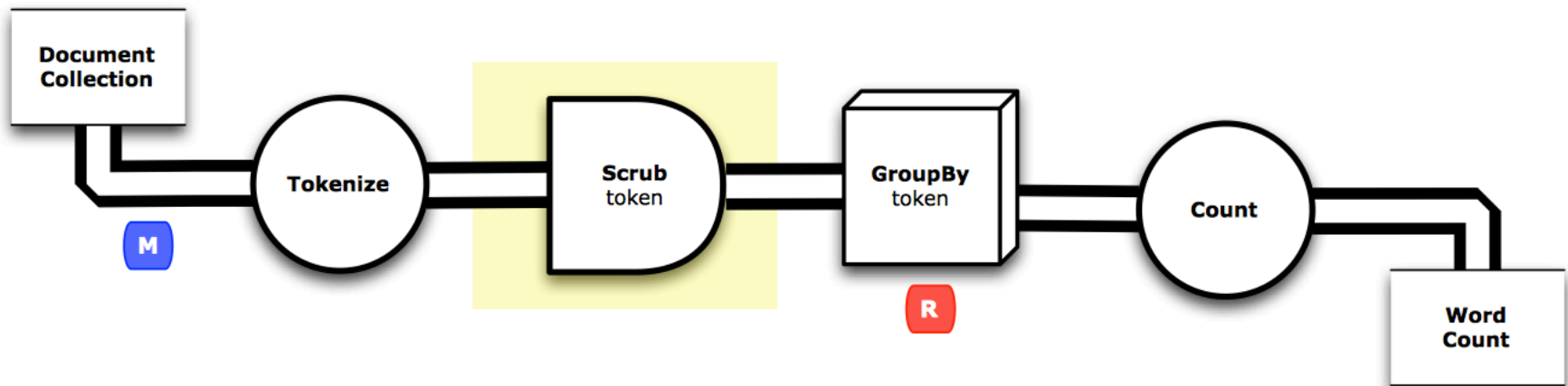
            // run the flow
            flowConnector.connect( flowDef ).complete();
        }
    }
```

2: word count



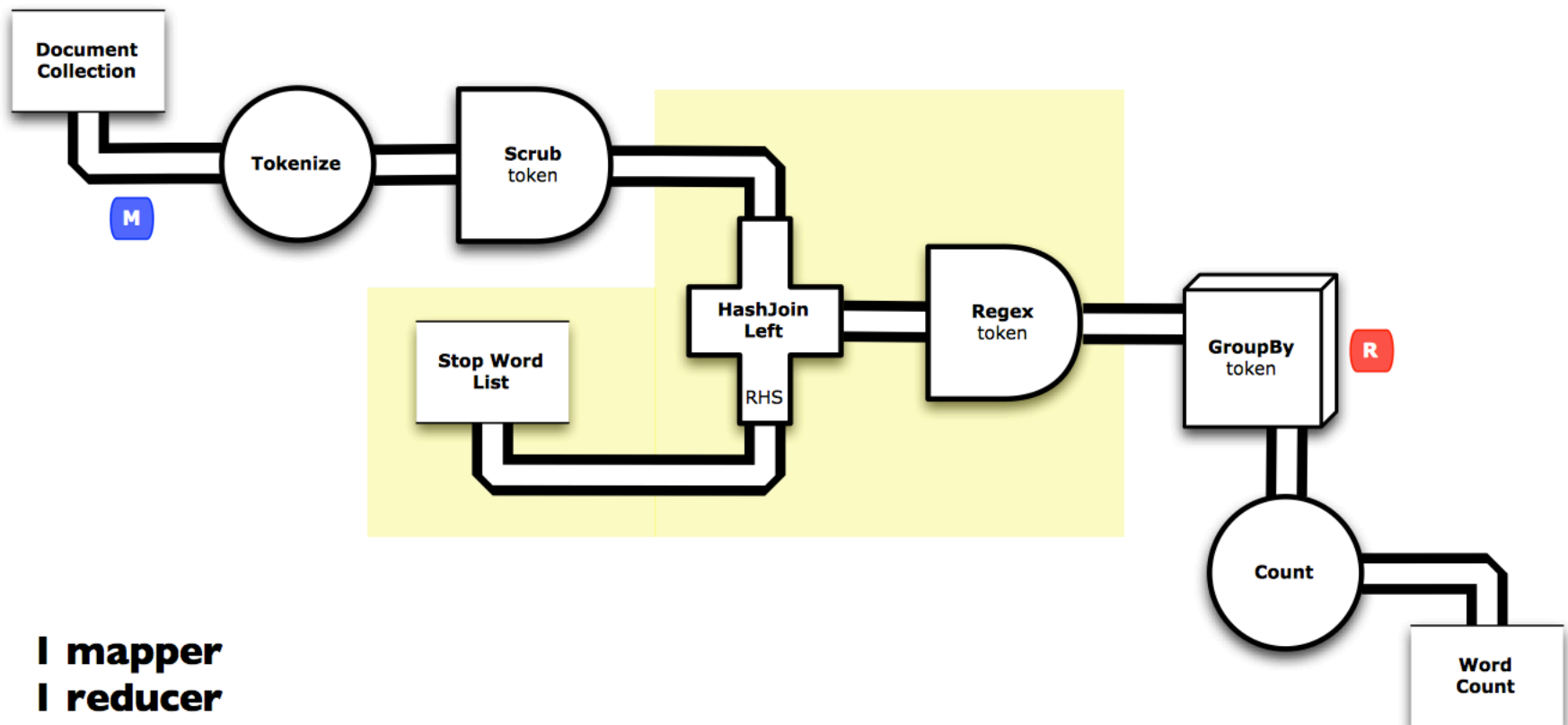
1 mapper
1 reducer
18 lines code

3: wc + scrub



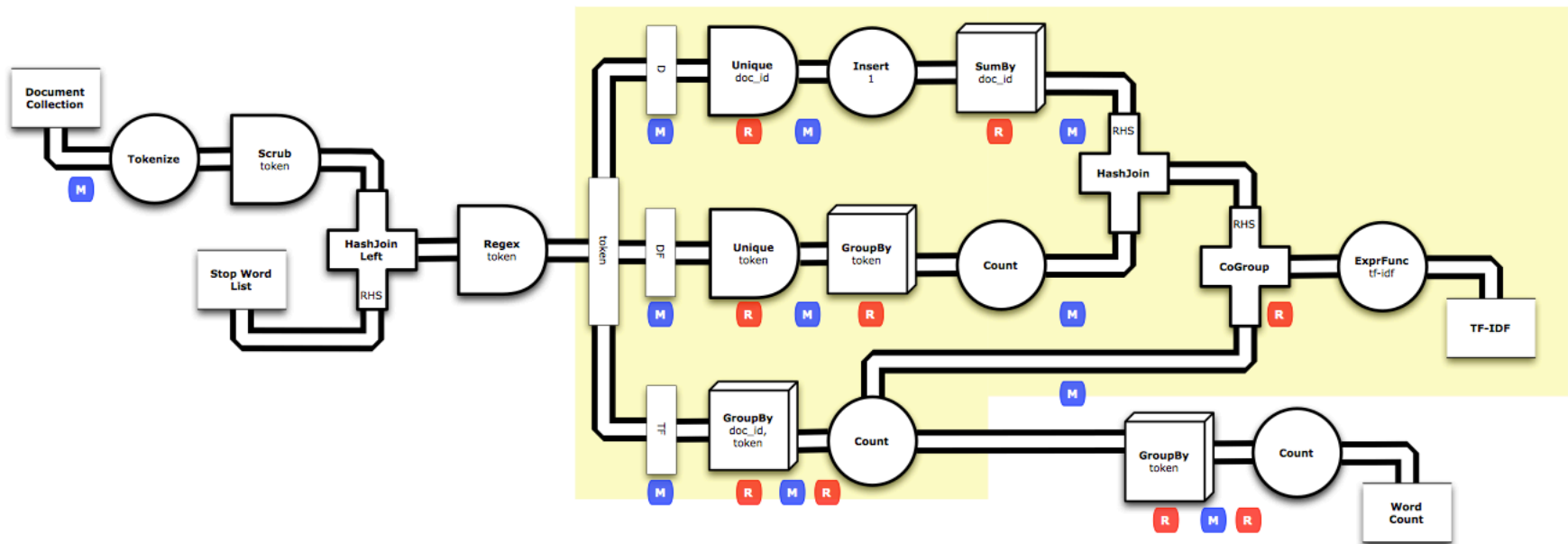
1 mapper
1 reducer
22+10 lines code

4: wc + scrub + stop words



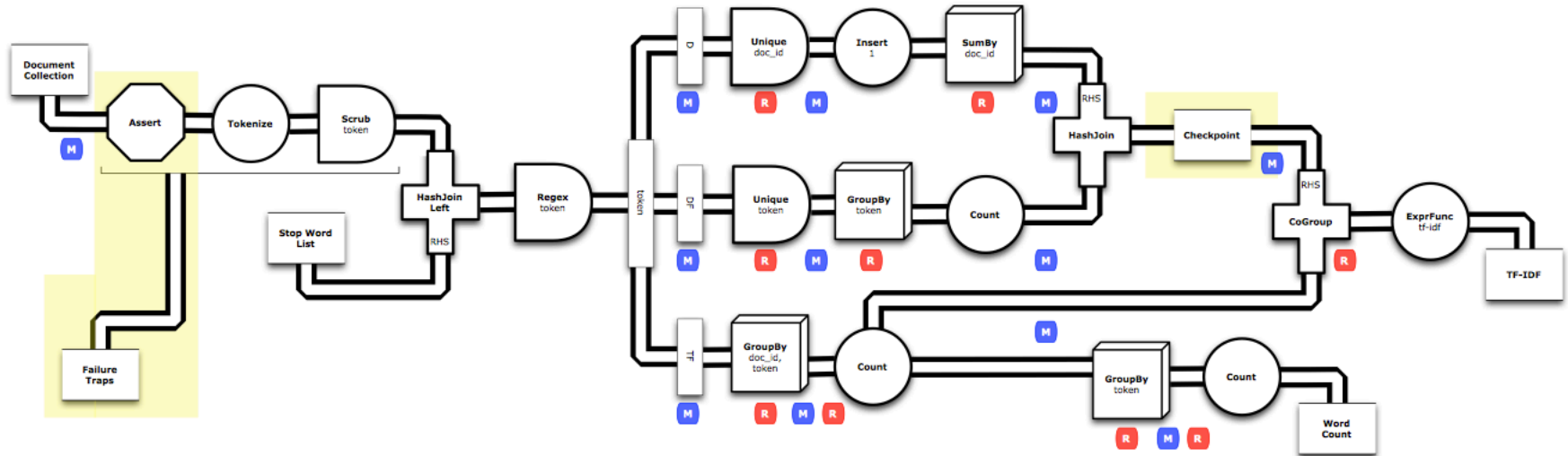
1 mapper
1 reducer
28+10 lines code

5: tf-idf



11 mappers
9 reducers
65+10 lines code

6: tf-idf + tdd



12 mappers
9 reducers
76+14 lines code

CLASSIFICATION INTRO

What is it?

- D -dimensional observations $X (x_1, \dots, x_n)$
- 2+ known classes K on some subset of X
- Find function $f(x)$ to assign each X to a class K
 - Make it perform well
 - Precision vs Recall. Resource and time complexity

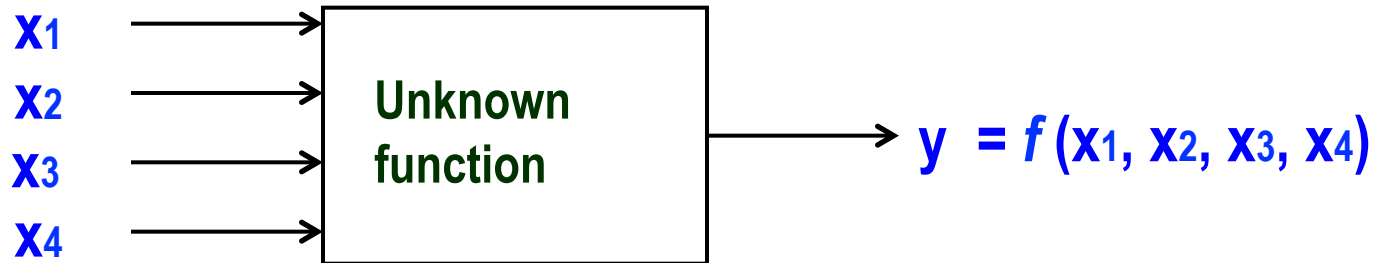
What is it?

- D -dimensional observations $X (x_1, \dots, x_n)$
- 2+ known classes K on some subset of X
- Find function $f(x)$ to assign each X to a class K
 - Make it perform well
 - Precision vs Recall. Resource and time complexity

Applications

- Security – Does an observation fall outside of normal behavior **boundaries**?
- Marketing – Can a user's gender be **inferred** by browsing behavior? Search terms?
- Search – Does the user's query pattern **imply** she wants documents of a particular class?
- Health – Emergency Room triage. Does this patient need to be seen immediately?

Learning is impossible, unless...



Given:

Training examples $(x, f(x))$
of **unknown** function f

Find:

A good **approximation** to f

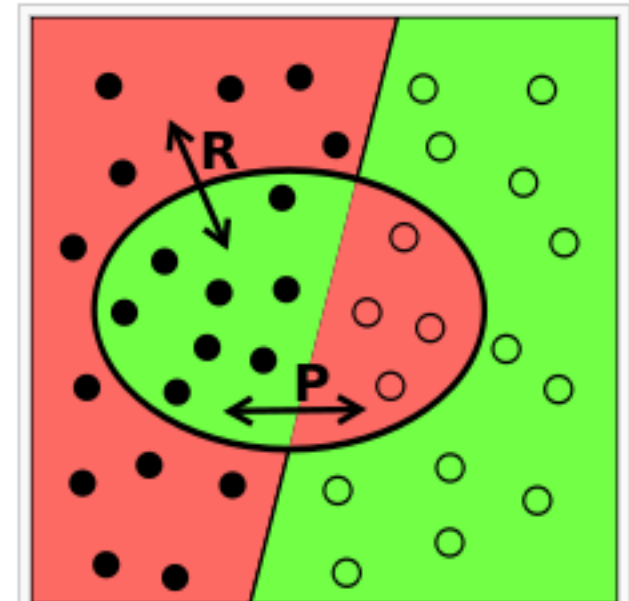
Example	x_1	x_2	x_3	x_4	y
1	0	0	1	0	0
2	0	1	0	0	0
3	0	0	1	1	1
4	1	0	0	1	1
5	0	1	1	0	0
6	1	1	0	0	0
7	0	1	0	1	0

Precision vs. Recall aka Sensitivity vs. Specificity

predicted class (expectation)	actual class (observation)	
	tp (true positive) Correct result	fp (false positive) Unexpected result
	fn (false negative) Missing result	tn (true negative) Correct absence of result

$$\text{Precision} = \frac{tp}{tp + fp}$$

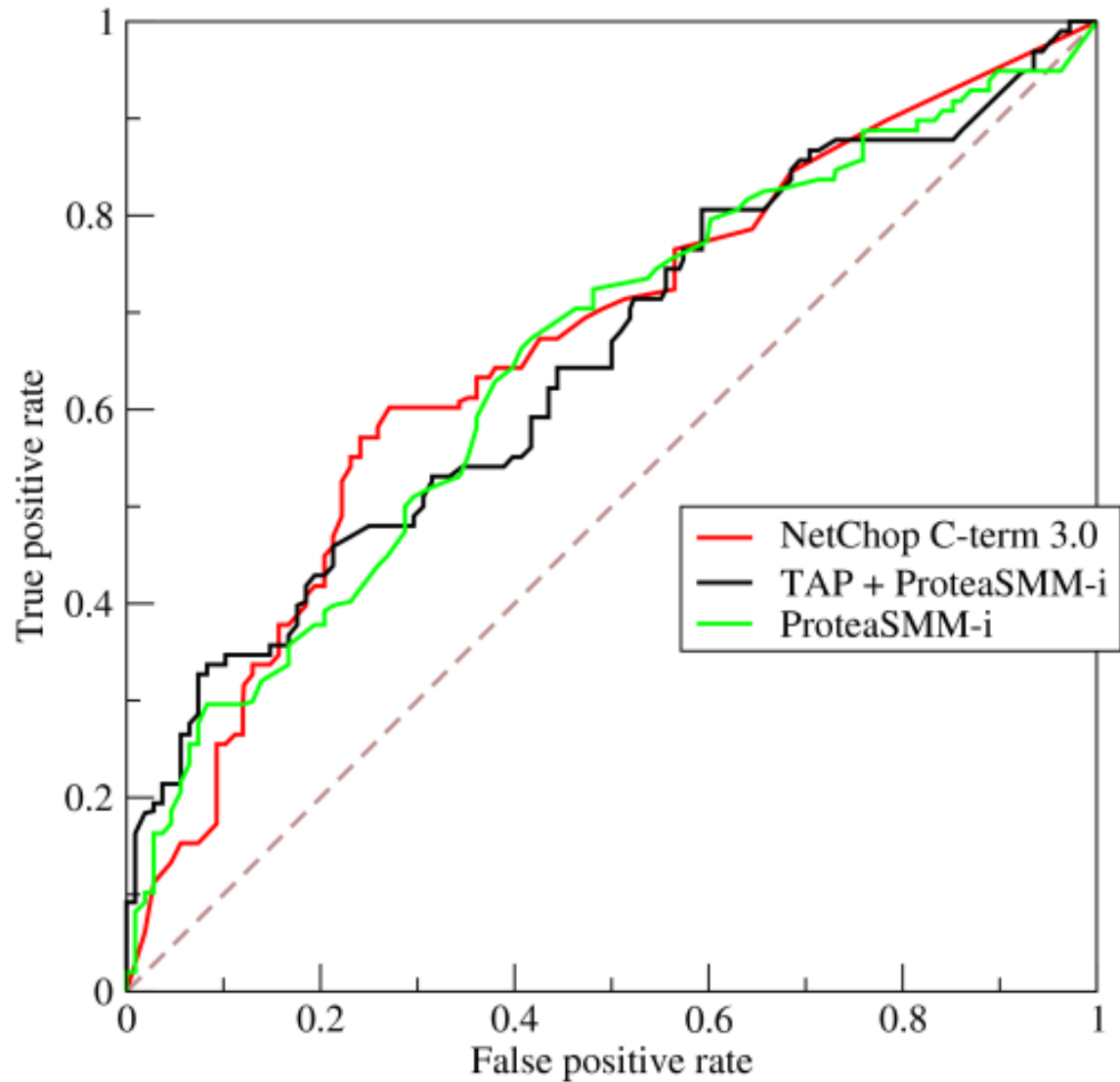
$$\text{Recall} = \frac{tp}{tp + fn}$$



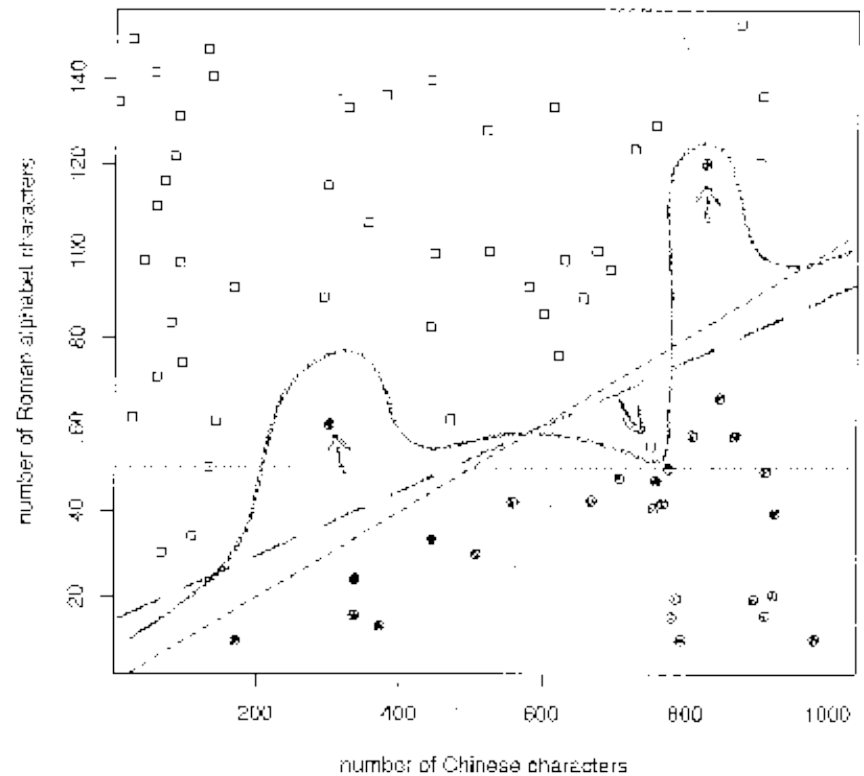
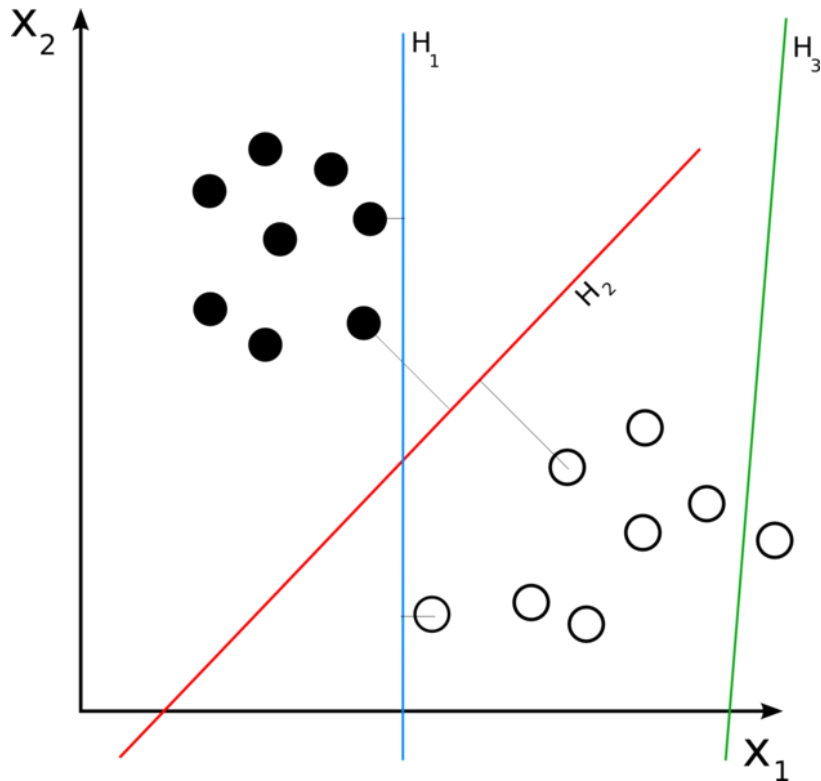
In this figure the relevant items are to the left of the straight line while the retrieved items are within the oval. The red regions represent errors. On the left these are the relevant items not retrieved (*false negatives*), while on the right they are the retrieved items that are not relevant (*false positives*).

- D -dimensional observations $X (x_1, \dots, x_n)$
- 2+ known classes K on some subset of X
- Find function $f(x)$ to assign each X to a class K
 - Make it perform well
 - Precision vs Recall. Resource and time complexity

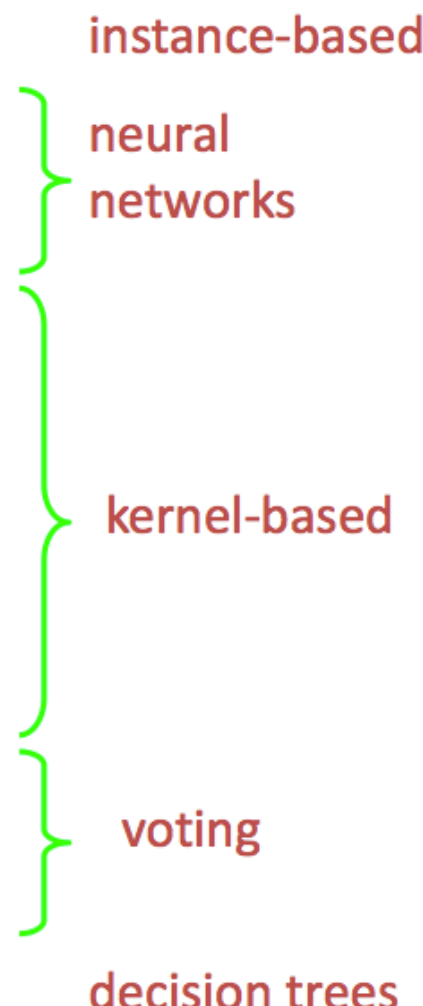
ROC Curve



Linear vs. Non-linear

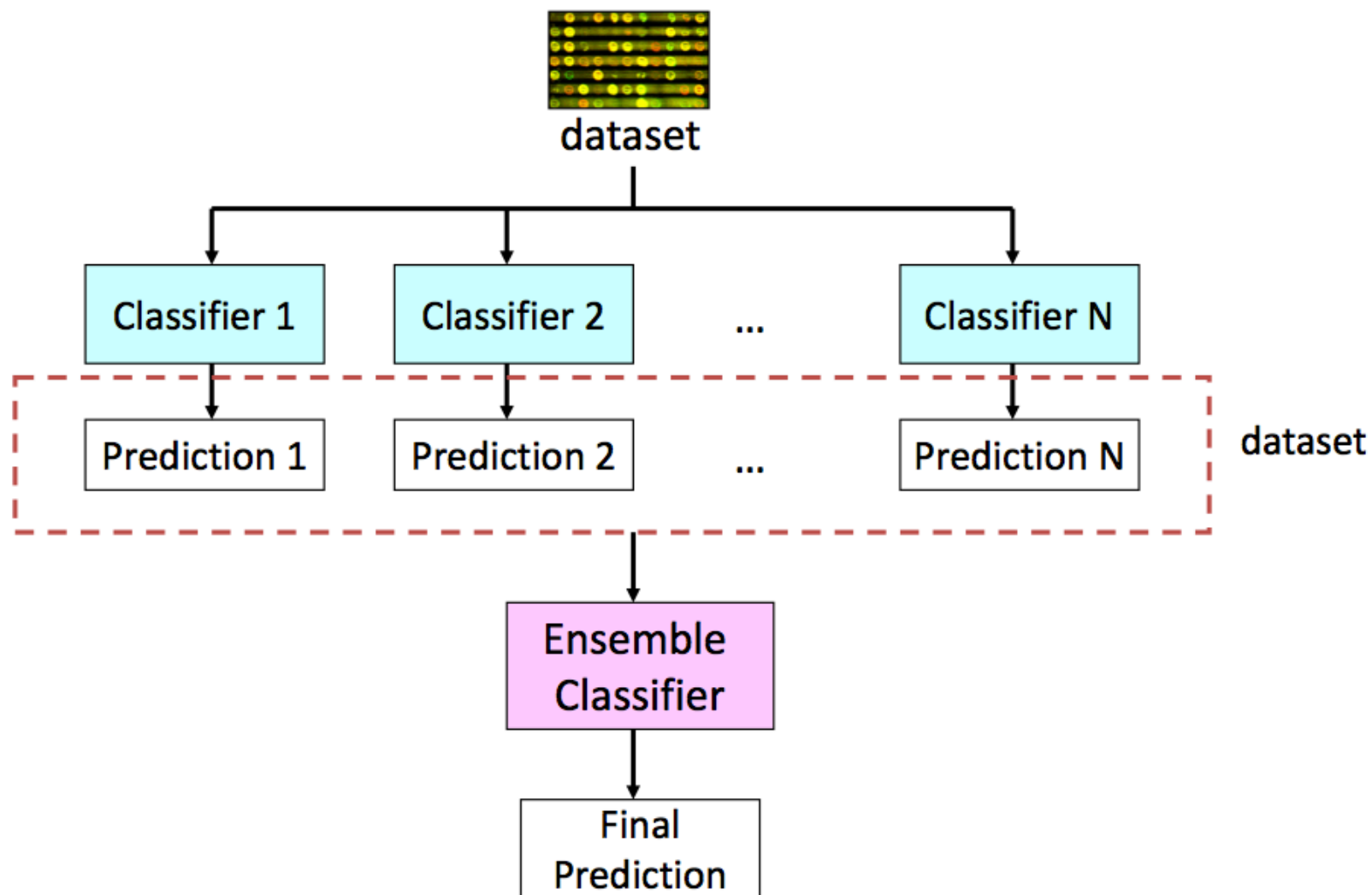


Classifiers

- K-Nearest Neighbors (**KNN**)
 - Backpropagation Neural Networks (**NN**)
 - Probabilistic Neural Networks (**PNN**)
 - Multi-Class SVM: One-Versus-Rest (**OVR**)
 - Multi-Class SVM: One-Versus-One (**OVO**)
 - Multi-Class SVM: **DAGSVM**
 - Multi-Class SVM by Weston & Watkins (**WW**)
 - Multi-Class SVM by Crammer & Singer (**CS**)
 - Weighted Voting: One-Versus-Rest
 - Weighted Voting: One-Versus-One
 - Decision Trees: CART
- 
- The diagram uses green curly braces to group the classifiers into five categories on the right side of the list:
- instance-based**: Groups K-Nearest Neighbors (KNN).
 - neural networks**: Groups Backpropagation Neural Networks (NN) and Probabilistic Neural Networks (PNN).
 - kernel-based**: Groups Multi-Class SVM: One-Versus-Rest (OVR), Multi-Class SVM: One-Versus-One (OVO), Multi-Class SVM: DAGSVM, Multi-Class SVM by Weston & Watkins (WW), and Multi-Class SVM by Crammer & Singer (CS).
 - voting**: Groups Weighted Voting: One-Versus-Rest and Weighted Voting: One-Versus-One.
 - decision trees**: Groups Decision Trees: CART.

- Naïve Bayesian Classifiers
- Support Vector Machines

Ensemble classifiers



NAÏVE BAYES

What is it?

- A classifier that:
- Can train on a small # of observations
- Assumes independence between features
- Simple to implement

Applications

- Typically $K=2$ -class classifiers
 - $K>2$ possible
- E.g.
 - Spam vs. Non-spam
 - Female vs. Male

Prior vs. Posterior Probabilities

- Description vs. Prediction

How does it work?

- Given the prior probability of a class k in K
- ...the prior probabilities of N object features
- ...and their associations with k
- Calculate the posterior probability that a given object belongs to class k
- Easier to work through an example
 - <http://bit.ly/10alkWY>

Making Best use of Data

- Train vs. Test
 - Cross-validation
 - ROC
- Ensembles

HIERARCHICAL CLUSTERING