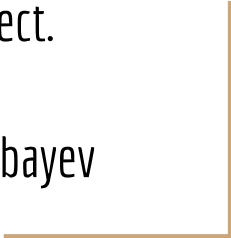




PaperScraper

MSCI-446 Term Project.
Andrew Andrade,
Andy Toulis, Baha Nurlybayev



Research Question

How has Artificial Intelligence and Predictive Analytics been used in the Petroleum Industry?

Problem

Data mining approach to aid multi-document classification for literature review.

Treating Uncertainties in Reservoir Performance Prediction with Neural Networks



Average from 0 ratings

Authors	Johann Peter Lechner (OMV A.G.) Georg Zangl (Schlumberger)
DOI	http://dx.doi.org/10.2118/94357-MS
Document ID	SPE-94357-MS
Publisher	Society of Petroleum Engineers
Source	SPE Europec/EAGE Annual Conference, 13-16 June, Madrid, Spain
Publication Date	2005
Document Type	Conference Paper
Language	English
ISBN	978-1-55563-943-3
Copyright	2005. Society of Petroleum Engineers
Disciplines	6.7.4 Probabilistic Methods , 6.5.5 Evaluation of Uncertainties
Downloads	2 in the last 30 days 476 since 2007 Show less detail

SPE Member Price: **USD 8.50**

SPE Non-Member Price: **USD 25.00**

[Add to cart](#) 

[View rights & permissions](#) [Export citation](#)

Abstract

In development projects reservoir parameters are only known within certain ranges - a fact that allows various realisations of the subsurface. Because of the computational time involved, not all of the possible parameter combinations can be covered by simulation models to obtain a probability distribution of possible outcomes. Creating a response surface that is based on a reduced number of simulation runs becomes necessary. Such a response surface can be utilized to approximate results for numerous different variations of input parameters. In contrast to the widely used methodology of fitting a polynomial model to the results of a limited number of simulation runs, an approach, where reservoir response is captured by an Artificial Neural Network (ANN) has been investigated.

Data

Features (Paper metadata, in text form):

- Abstracts
- Year
- Keywords
- Disciplines
- other meta data

Classes:

- Disciplines



Data

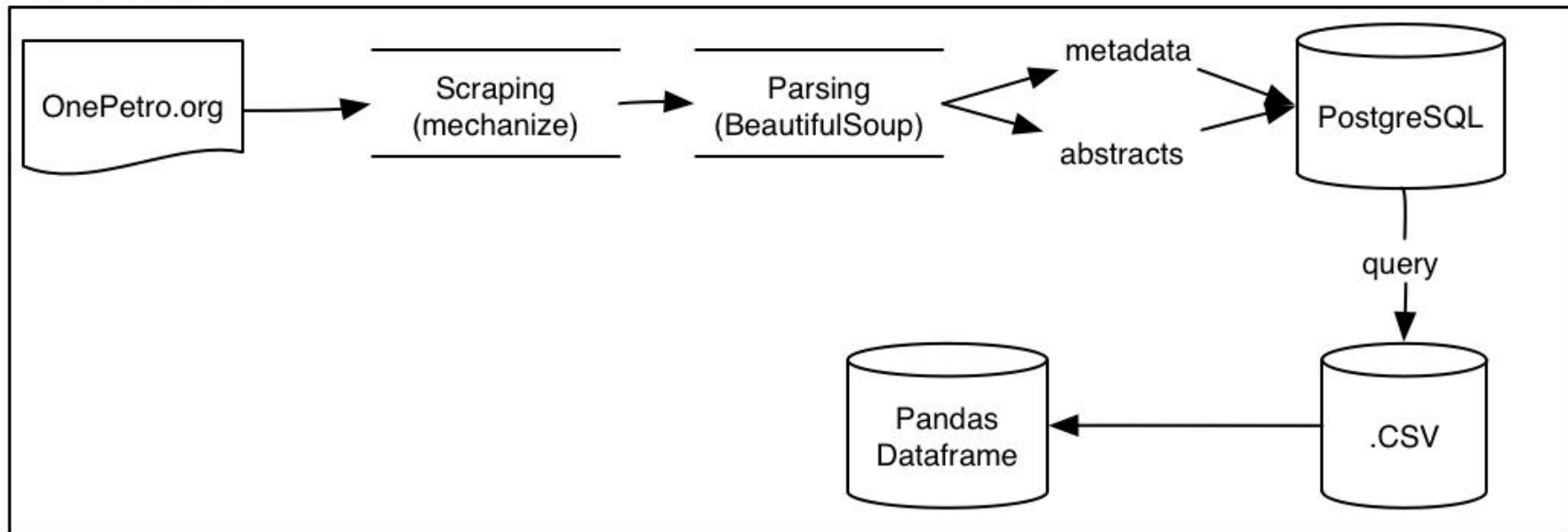
- 180 000+ papers on OnePetro
- 117 000 scraped, parsed*
- 22 250 labelled data
- Real world data is dirty (nightmarish)

*Scraper was originally written by Jonathon “Jay” Estrella as an intern at PetroPredict, later modified to include hidden metadata

Dirty Data

- UNICODE (special characters)
- Missing values
- non-sense
- mislabels
- inconsistencies

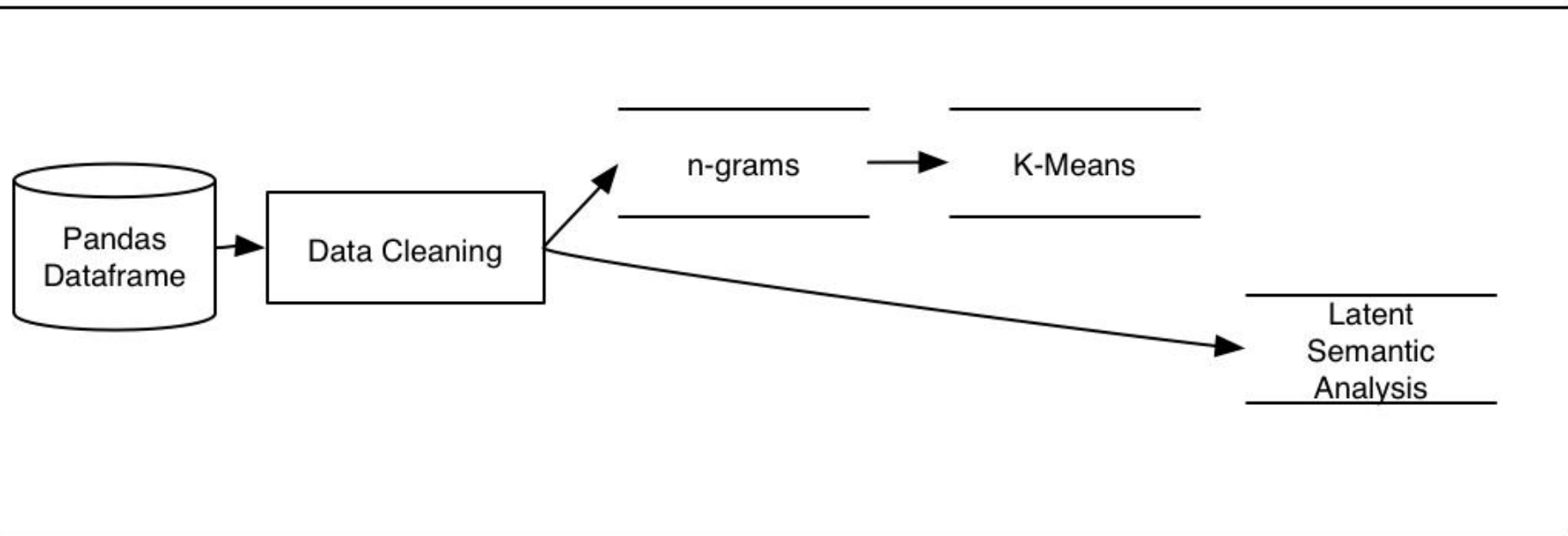
Data Collection



Approach

1. Data Collection and Storage
 - a. Scraping
 - b. Pandas (because we can)
2. **EDA & Clustering**

EDA



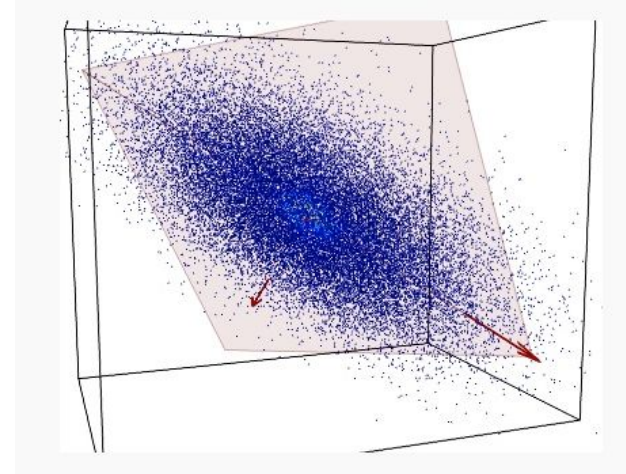
Latent Semantic Analysis

Feature: Abstracts (TF-IDF vector)

Dimensionality Reduction: Principal Component Analysis

Metric: Cosine Similarity

2 to 12 Clusters



<http://earlywarn.blogspot.ca/2010/10/extracting-signal-from-drought-noise.html>

Example: 12 Clusters

Cluster 1: steam crude polymer viscosity flooding produced sand surfactant

Cluster 2: strength waves structures results ice numerical soil tests structure failure study test loading

Cluster 3: safety spe management industry health risk companies business conference prepared training environmental

Cluster 4: drilling mud bit hole drill drilled wellbore casing rig fluids cuttings

Cluster 5: subsea offshore systems equipment project platform completion cement pipeline installation tubing

Cluster 6 : porosity models simulation properties log core saturation pore

Cluster 7: seismic inversion velocity migration imaging noise image processing acquisition survey source

Cluster 8: natural hydrate liquid condensate energy reserves methane coal

Cluster 9: preview atce commercially otc technologies readily limited engineers houston currently discussion

12 Clusters Cont.

Fracking:

Cluster 10: fracture fracturing fractures hydraulic proppant fractured stimulation treatment treatments

French!:

Cluster 11: eacute des les der egrave die agrave und dans pour une par von est ouml que sur sont qui mit

Cluster 12: corrosion steel alloy alloys resistance steels stainless metal coating materials hydrogen cracking

Approach

1. Data Collection and Storage

- a. Scraping
- b. SQL database
- c. Pandas (because we can)

2. Exploratory Data Analysis

- a. Data Visualization
- b. Determine papers using AI techniques
- c. K-means clustering Topic modeling

3. Document Classification

MultiLabel Classification (6 Disciplines of SPE)

Class 0: Drilling and Completions

Class 1: Health, Safety, Security, Environment and Social Responsibility

Class 2: Management and Information

Class 3: Project Facilities and Construction

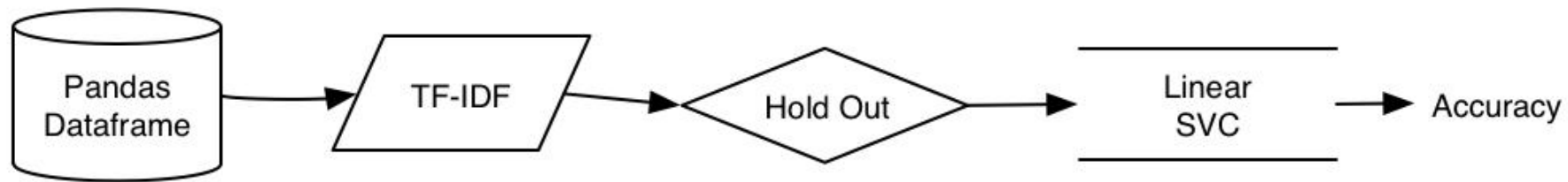
Class 4: Production and Operations

Class 5: Reservoir Description and Dynamics

Approach

1. Data Collection and Storage
 - a. Scraping
 - b. SQL database
 - c. Pandas (because we can)
2. Exploratory Data Analysis
 - a. Data Visualization
 - b. Determine papers using AI techniques
 - c. K-means clustering
3. Multi Label Document Classification
 - a. **Linear SVM**

Data Classification 1



Very High Accuracy!

71% - 86%

```
Class0accuracy:
0.813251201923
Class1accuracy:
0.842247596154
Class2accuracy:
0.861628605769
Class3accuracy:
0.815504807692
Class4accuracy:
0.714242788462
Class5accuracy:
0.792518028846
```

Management and Information

CM



0 921

0 5735



86% Accuracy!

Only 14% are about
Management

Labeling everything
as not Management!

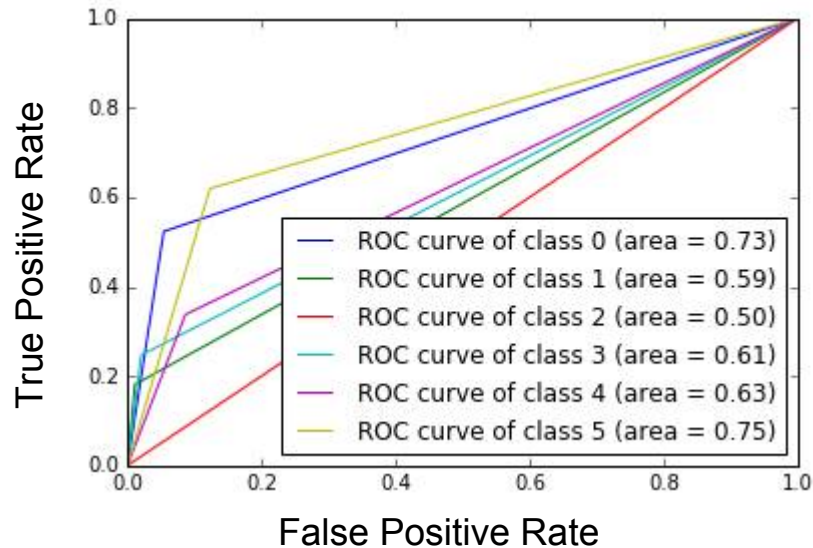


0 921

0 5735



Receiver Operating Characteristic and AUC



Approach

1. Data Collection and Storage

- a. Scraping
- b. SQL database
- c. Pandas (because we can)

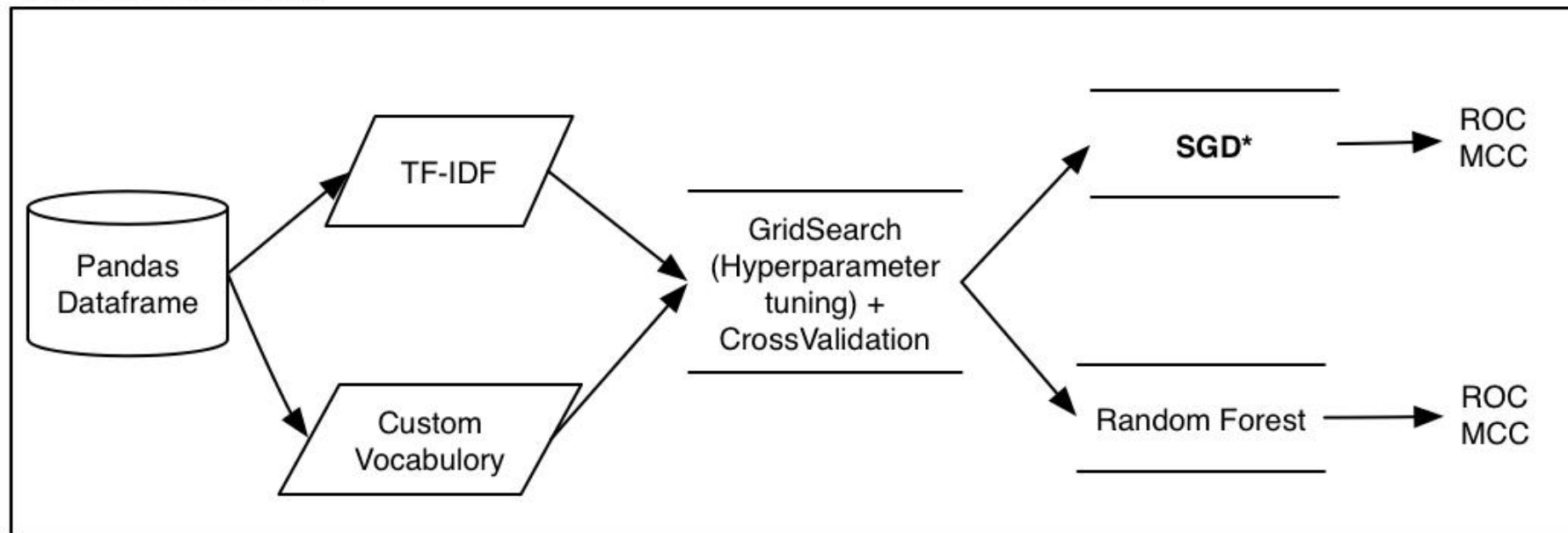
2. Exploratory Data Analysis

- a. Data Visualization
- b. Determine papers using AI techniques
- c. K-means clustering

3. Multi Label Document Classification

- a. Linear Support Vector Machine
- b. **Grid search + SGD**

Data Classification 2

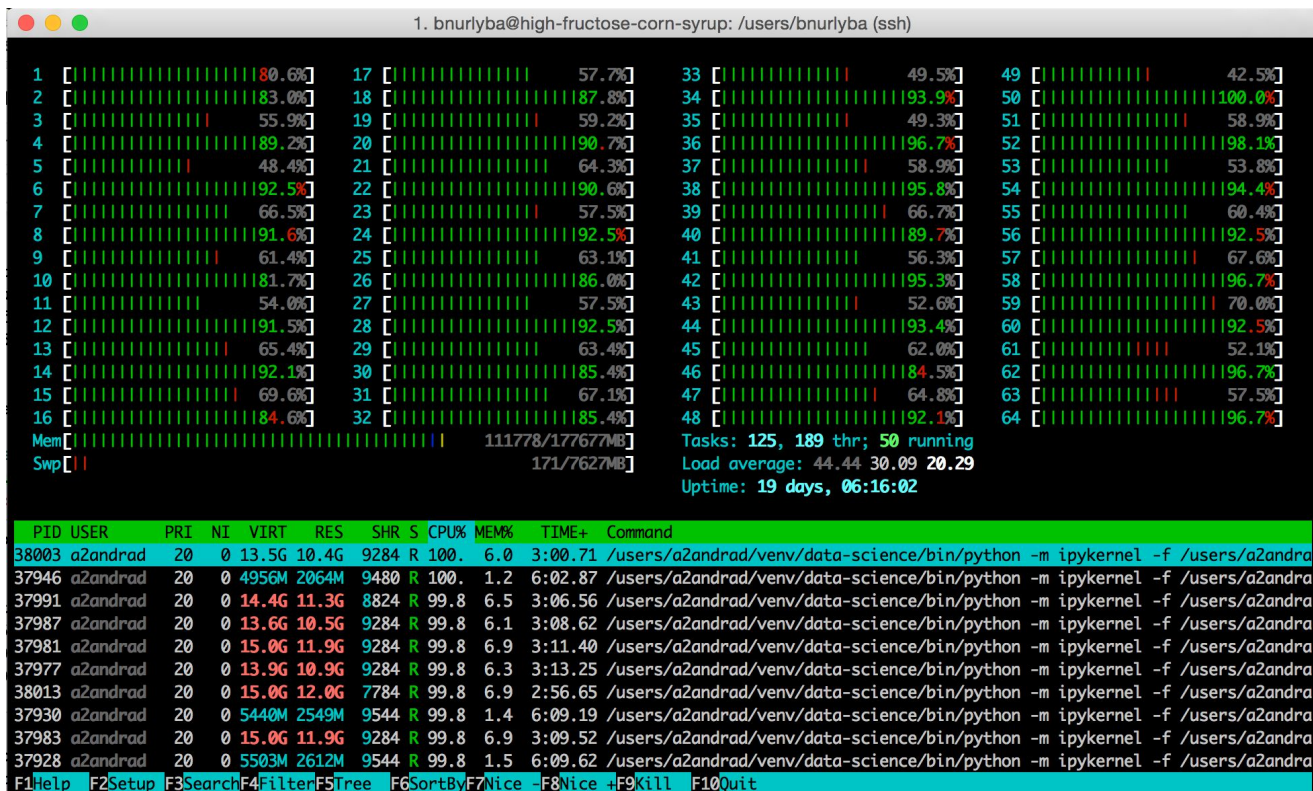


Pipelines and Hyperparameters

```
pipeline = Pipeline([
    ('vect', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
    ('clf', SGDClassifier(loss='log', n_iter=10, alpha=0.0001)),
])

parameters = {
    'vect__min_df': (0.1, 0.2),
    'vect__max_df': (0.8, 0.9),
    #'vect__max_features': (None, 5000, 10000, 50000),
    'vect__ngram_range': ((1, 1), (1, 2), (1, 3)),
    #'tfidf__use_idf': (True, False),
    #'tfidf__norm': ('l1', 'l2'),
    'clf__penalty': ('l2', 'elasticnet')
}
```

GridSearch



GridSearch



GridSearch

Hello Andrew,

This morning (Wednesday, November 30) your python processes running on HFCS consumed all available RAM and CPU. Your processes have been killed by the Systems Committee as they were in violation of our Machine Usage Agreement (http://csclub.uwaterloo.ca/services/machine_usage), specifically the section titled "User Responsibilities".

If you restart the processes, please ensure that they do not consume all available resources on the machine.

GridSearch

tmux -2 attach -t 3

```
1 [|||||||||||||||||||||||||||||||||||||||||100.0%] 17 [|||||||||||||||||||||||||||||||||||||0.0%] 33 [|||||||||||||||||||||||||||||||||||||0.0%]
2 [|||||||||||||||||||||||||||||||||||||0.0%] 18 [|||||||||||||||||||||||||||||||||||||0.0%] 34 [|||||||||||||||||||||||||||||||||||||0.0%]
3 [|||||||||||||||||||||||||||||||||||||0.0%] 19 [|||||||||||||||||||||||||||||||||||||0.0%] 35 [|||||||||||||||||||||||||||||||||||||0.0%]
4 [|||||||||||||||||||||||||||||||||||||100.0%] 20 [|||||||||||||||||||||||||||||||||||||100.0%] 36 [|||||||||||||||||||||||||||||||||||||0.0%]
5 [|||||||||||||||||||||||||||||||||||||0.7%] 21 [|||||||||||||||||||||||||||||||||||||100.0%] 37 [|||||||||||||||||||||||||||||||||||||0.0%]
6 [|||||||||||||||||||||||||||||||||||||99.3%] 22 [|||||||||||||||||||||||||||||||||||||0.0%] 38 [|||||||||||||||||||||||||||||||||||||0.0%]
7 [|||||||||||||||||||||||||||||||||||||4.7%] 23 [|||||||||||||||||||||||||||||||||||||100.0%] 39 [|||||||||||||||||||||||||||||||||||||0.0%]
8 [|||||||||||||||||||||||||||||||||||||1.3%] 24 [|||||||||||||||||||||||||||||||||||||0.0%] 40 [|||||||||||||||||||||||||||||||||||||0.0%]
9 [|||||||||||||||||||||||||||||||||||||66.9%] 25 [|||||||||||||||||||||||||||||||||||||72.0%] 41 [|||||||||||||||||||||||||||||||||||||0.0%]
10 [|||||||||||||||||||||||||||||||||||||34.9%] 26 [|||||||||||||||||||||||||||||||||||||0.0%] 42 [|||||||||||||||||||||||||||||||||||||0.0%]
11 [|||||||||||||||||||||||||||||||||||||31.1%] 27 [|||||||||||||||||||||||||||||||||||||100.0%] 43 [|||||||||||||||||||||||||||||||||||||0.0%]
12 [|||||||||||||||||||||||||||||||||||||16.1%] 28 [|||||||||||||||||||||||||||||||||||||0.0%] 44 [|||||||||||||||||||||||||||||||||||||0.0%]
13 [|||||||||||||||||||||||||||||||||||||82.8%] 29 [|||||||||||||||||||||||||||||||||||||100.0%] 45 [|||||||||||||||||||||||||||||||||||||0.0%]
14 [|||||||||||||||||||||||||||||||||||||17.4%] 30 [|||||||||||||||||||||||||||||||||||||0.0%] 46 [|||||||||||||||||||||||||||||||||||||0.0%]
15 [|||||||||||||||||||||||||||||||||||||69.5%] 31 [|||||||||||||||||||||||||||||||||||||70.9%] 47 [|||||||||||||||||||||||||||||||||||||0.0%]
16 [|||||||||||||||||||||||||||||||||||||0.0%] 32 [|||||||||||||||||||||||||||||||||||||2.7%] 48 [|||||||||||||||||||||||||||||||||||||0.0%]
Mem[|||||||||||||||||||||||||||||||||||||42136/177677MB] Tasks: 169, 852 thr; 28 running
Swp[|||||||||||||||||||||||||||||||||||||5183/7627MB] Load average: 7.72 4.20 2.91
Uptime: 22 days, 07:54:45
```

PID	USER	PRI	NI	VIRT	RES	SHR	S	CPU%	MEM%	TIME+	Command
35008	a2andrad	30	10	6055M	2533M	33916	R	2156	1.4	31:39.53	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
16870	a2andrad	30	10	13.4G	10.1G	34168	R	98.8	5.8	3h16:57	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38456	a2andrad	30	10	6055M	2532M	33916	S	82.5	1.4	0:10.27	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38452	a2andrad	30	10	6055M	2532M	33916	S	81.9	1.4	0:20.01	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38439	a2andrad	30	10	6055M	2532M	33916	S	81.9	1.4	0:38.89	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38451	a2andrad	30	10	6055M	2532M	33916	S	81.9	1.4	0:22.03	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38433	a2andrad	30	10	6055M	2532M	33916	S	81.2	1.4	0:41.81	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38442	a2andrad	30	10	6055M	2532M	33916	S	81.2	1.4	0:36.16	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38450	a2andrad	30	10	6055M	2532M	33916	S	81.2	1.4	0:23.74	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh
38435	a2andrad	30	10	6055M	2532M	33916	S	80.6	1.4	0:40.89	/users/a2andrad/venv/data-science/bin/python -m ipykernel -f /users/a2andrad/.local/sh

Approach

1. Data Collection and Storage

- a. Scraping
- b. SQL database
- c. Pandas (because we can)

2. Exploratory Data Analysis

- a. Data Visualization
- b. Determine papers using AI techniques
- c. K-means clustering

3. Multi Label Document Classification

- a. Linear Support Vector Machine
- b. **Grid search + SGD+ Dimensionality Reduction**

Reduce Dimensionality

Disciplines have text!

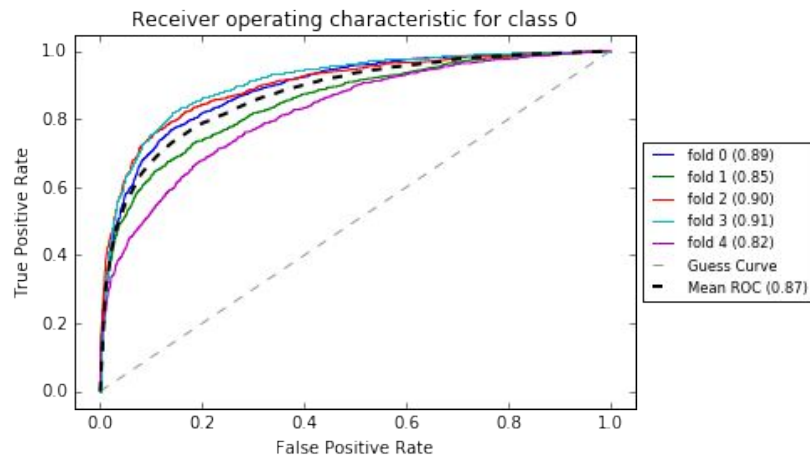
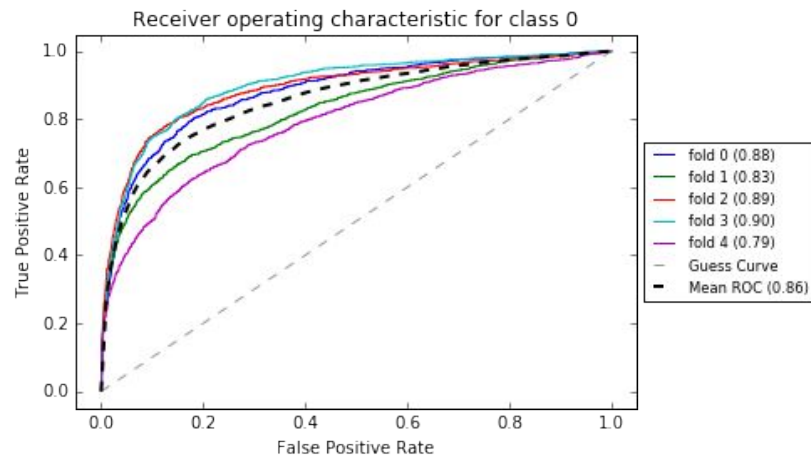
Build dict from the discipline text

Use document term vector from discipline

Full Abstract

vs.

Custom Dictionary



Approach

1. Data Collection and Storage

- a. Scraping
- b. SQL database
- c. Pandas (because we can)

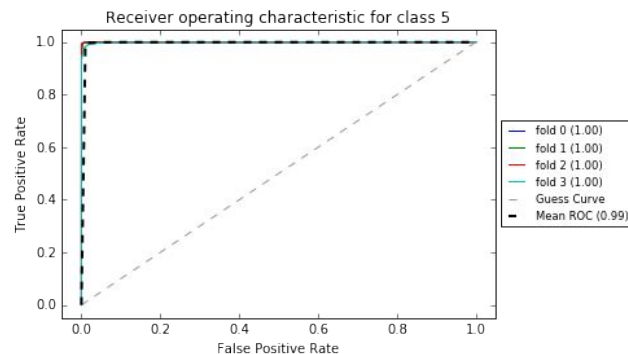
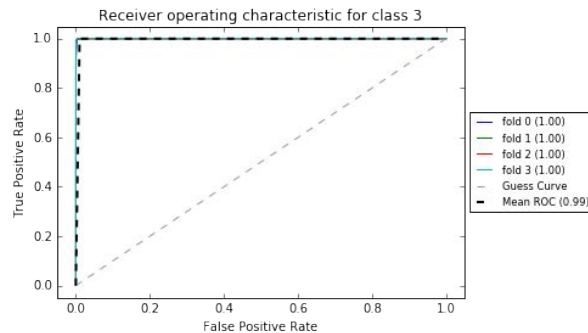
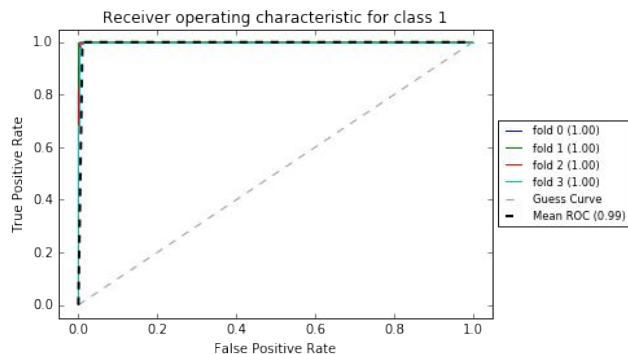
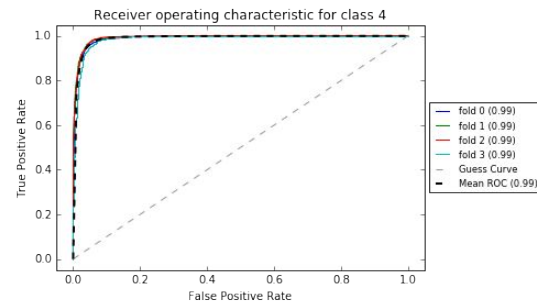
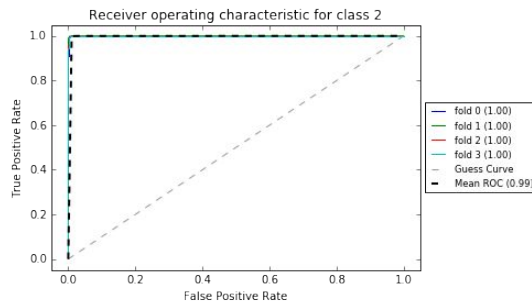
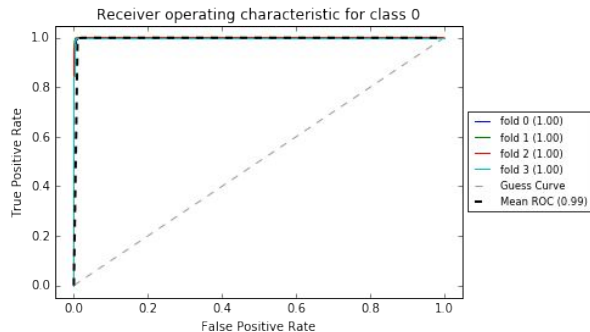
2. Exploratory Data Analysis

- a. Data Visualization
- b. Determine papers using AI techniques
- c. K-means clustering

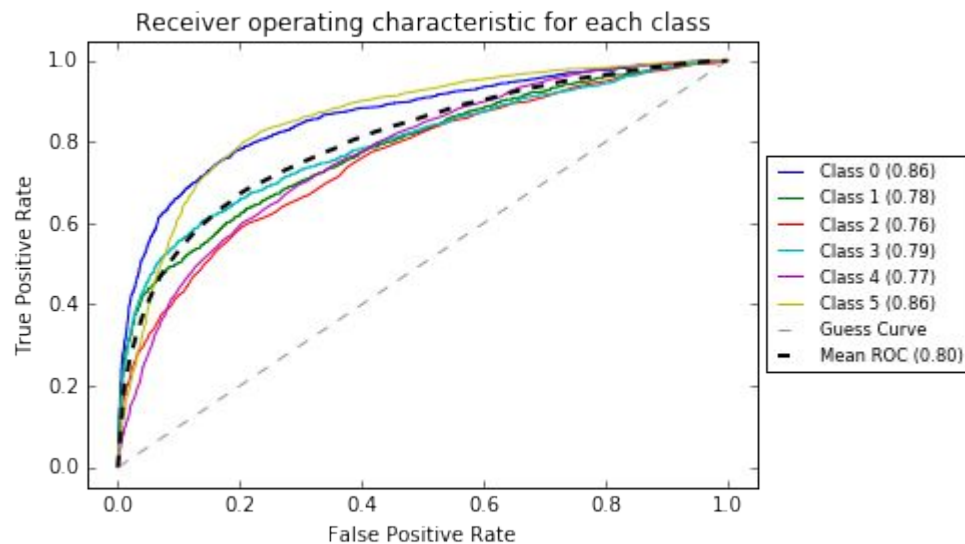
3. Multi Label Document Classification

- a. Linear Support Vector Machine
- b. Grid search + SGD + Dimensionality Reduction
- c. **Random Forest**

Random Forest! (severely overfitting)



30% hold out (model on 70% k folds)



Top 10 features: Random Forest

Drilling and Completions

wellbor
drill
string
bit
reservoir
cement
complet
format
mud
design

Health, Safety, Security, Environment and Social Responsibility

health
manag
environment
assess
emiss
hse
inject
model
train
impact

Management and Information

develop
integr
decis
oper
project
pressur
engin
time
safeti
drill

Top 10 features: Random Forest

Project Facilities and Construction

offshor
instal
facil
riser
design
moor
reservoir
float
use
format

Production and Operations

pump
drill
product
hydraul
fractur
reservoir
water
enhanc
use
flow

Reservoir Description and Dynamics

simul
interpret
reservoir
recoveri
log
inject
miscibl
model
drill

Approach

1. Data Collection and Storage

- a. Scraping
- b. SQL database
- c. Pandas (because we can)

2. Exploratory Data Analysis

- a. Data Visualization
- b. Determine papers using AI techniques
- c. K-means clustering

3. Multi Label Document Classification

- a. Linear Support Vector Machine
- b. Grid search + SGD + Dimensionality Reduction
- c. Random Forest
- d. **Final SGD Model**

Final SGD Model

```
final_pipeline = Pipeline([
    ('vect', CountVectorizer(ngram_range=(1,3), max_df = 0.4, min_df=0001)),
    ('tfidf', TfidfTransformer()),
    ('clf', SGDClassifier(loss='modified_huber', n_iter=10, penalty='elasticnet')),
])
```

Final SGD Model

MCC 70% --> 30% holdout

Class 0: 0.60 --> 0.62

Class 1: 0.46 --> 0.46

Class 2: 0.33 --> 0.36

Class 3: 0.50 --> 0.51

Class 4: 0.41 --> 0.40

Class 5: 0.58 --> 0.56

Predictions

89170 predictions

Class 0: 13140 (15%)

Class 1: 4230 (5%)

Class 2: 3050 (3%)

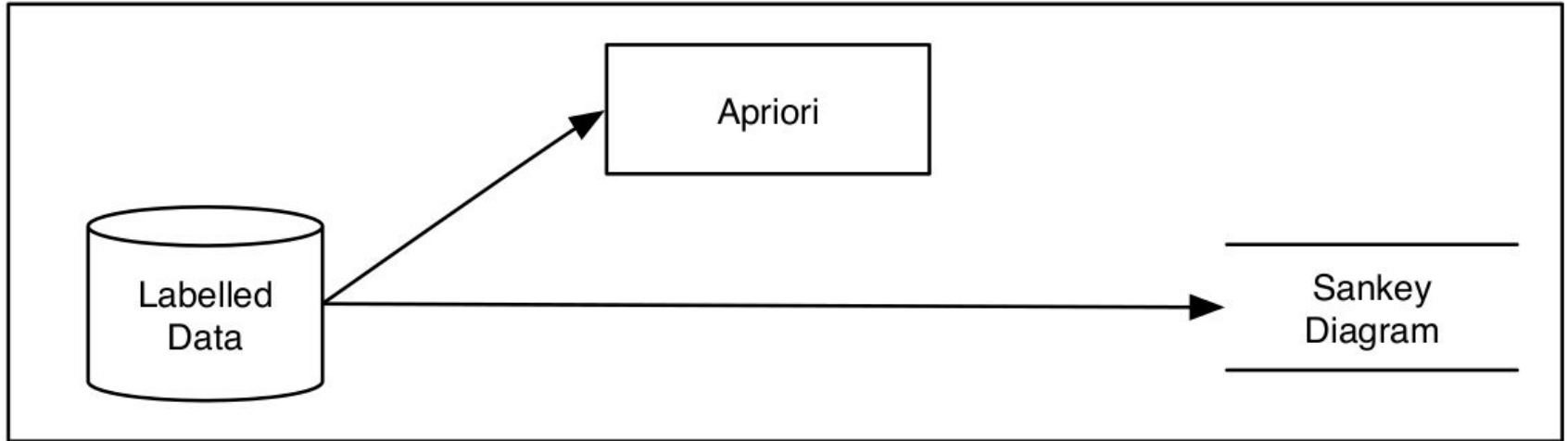
Class 3: 15200 (17%)

Class 4: 15300 (17%)

Class 5: 18560 (21%)

Association

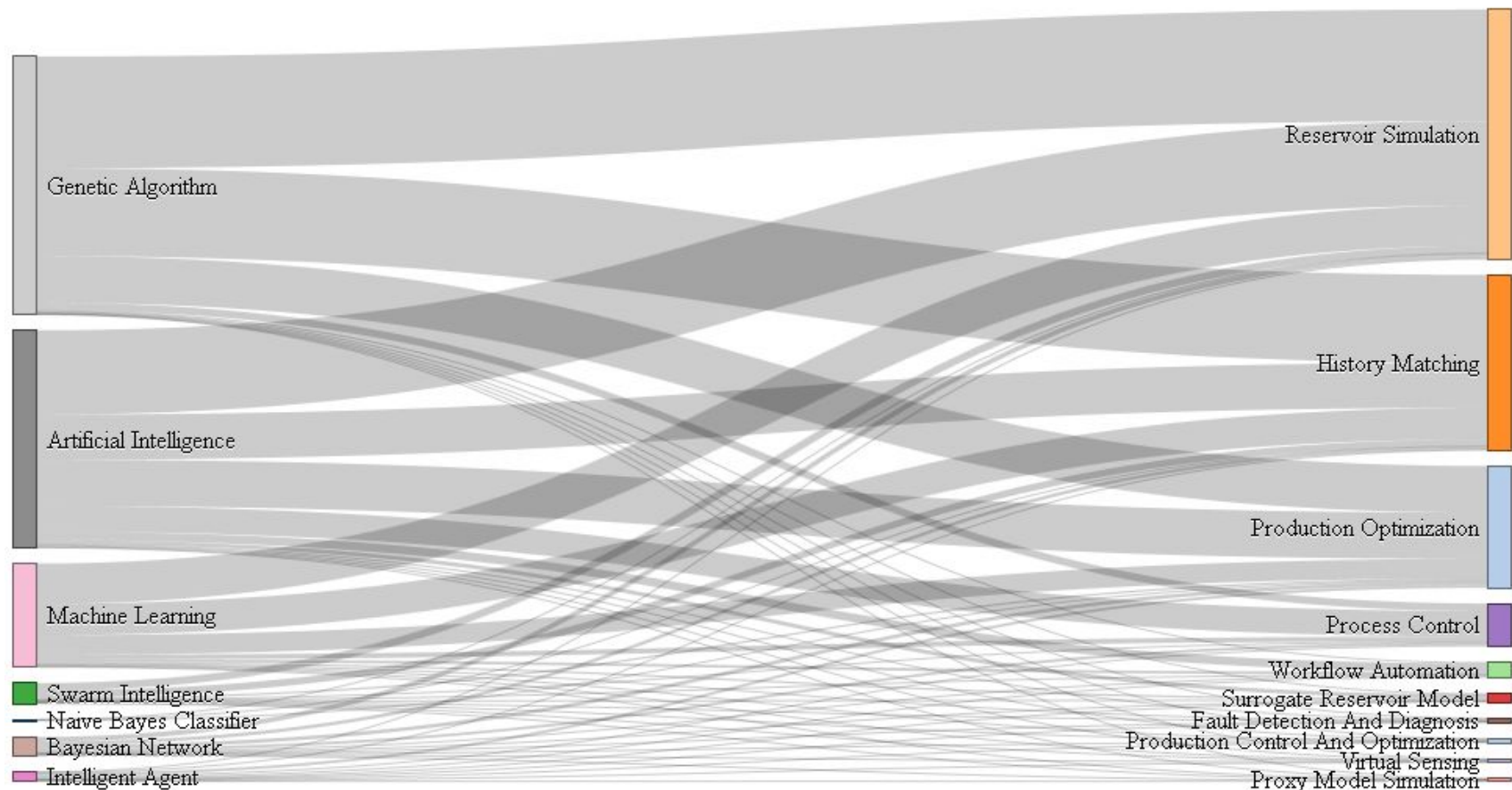
Association



Apriori

	rules	support	confidence	lift
{}	=> {class_0_prediction=1}	0.1190476190	0.119047619	1.000000000
{regress=1}	=> {class_0_prediction=1}	0.0430597771	0.144680851	1.21531915
{neural_network=1}	=> {class_0_prediction=1}	0.0215298886	0.122478386	1.02881844
{mont_carlo=1}	=> {class_0_prediction=1}	0.0197568389	0.100905563	0.84760673
{regress_analysi=1}	=> {class_0_prediction=1}	0.0141843972	0.200000000	1.680000000
{regress_analysi=1,regress=1}	=> {class_0_prediction=1}	0.0141843972	0.200000000	1.680000000
{data_manag=1}	=> {class_0_prediction=1}	0.0134245187	0.187279152	1.57314488
{expert_system=1}	=> {class_0_prediction=1}	0.0116514691	0.211981567	1.78064516
{genet=1}	=> {class_0_prediction=1}	0.0078520770	0.064049587	0.53801653
{bayesian=1}	=> {class_0_prediction=1}	0.0058257345	0.076411960	0.64186047

[illegible]



Thanks!