

# shed

A FRAMEWORK FOR REPRODUCIBLE TEXT MINING RESEARCH

15 October, 2015

Chris Emmery, Walter Daelemans



- 3 days of author profiling with CLiPS a.o.\*.
- Standard tasks: prepare data, discuss features, try classifiers.
- Missing: general framework to uniformly report and compare pipeline performance across tasks and dataset combinations.

\*Walter Daelemans, Guy De Pauw, Mike Kestemont, Tom De Smedt, Ben Verhoeven, Janneke van de Loo, Chris Emmery, Enrique Manjavacas, Florian Kunneman.

# GENERAL IDEA

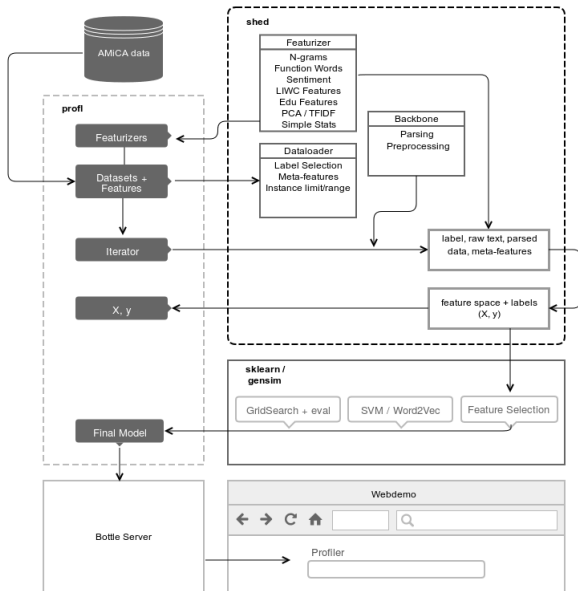
- Usable as a working environment as well as an (almost) end-to-end system.
- Flexible enough for extension and adaptation to different tasks.
- Open backbone, implementation for Frog (LaMachine) and spaCy.
- Report settings across the entire pipeline.
- Save backbone settings, fitted featurizers, and trained classifiers for easy deployment and reproduction.
- Well documented and open-source.

- **experiment.py** - wrapping and reporting standard procedures.
- **environment.py** - saving and loading progress.
- **datareader.py** - iterative loading of data in standard format.
- **backbone.py** - preprocessing according to own preference.
- **featurizer.py** - skeleton and implementations for feature extraction.

# FEATURIZERS - STANDARD IMPLEMENTATIONS

- **SimpleStats**
  - flooding, character types, number emoticons, mentions, hashtags
  - average word length, caps, start caps, urls, photos, videos
- **Ngrams**
  - *n* amount of: characters, tokens, POS (up to max features)
  - relative & cutoff on frequency
- **Function Words**
- **SentimentFeatures**
- **Readability**

# SCHEMA



# ENVIRONMENT

---

```
1 import shed
2 from shed.featurizer import Ngrams, SentimentFeatures, FuncWords
3 from sklearn.svm import SVC
4 from sklearn import metrics
5
6 data = ['/home/chris/data/netlog.csv', '/home/chris/data/otherstuff.csv']
7 features = [
8     Ngrams(level='token', n_list=[1]),
9     Ngrams(level='char', n_list=[2, 3]),
10    SentimentFeatures(),
11    FuncWords()
12 ]
13
14 env = shed.Environment(name, backbone='frog')
15
16 ftf_loader = env.load(data=data[0], target_label='age', proc='both', max_n=5000)
17 X, y = env.fit_transform(ftf_loader, features)
18
19 trf_loader = env.load(data=data[1], target_label='age', proc='both', max_n=1000)
20 Xi, yi = env.transform(trf_loader)
21
22 env.train(SVC(), X, y)
23 res = env.classify(clf, Xi, yi)
24
25 print(metrics.f1_score(yi, res, pos_label=1, average='binary'))
26 env.save()
```

---

# EXPERIMENT.PY

---

```
1 ...
2
3 if __name__ == '__main__':
4
5     features = [
6         Ngrams(level='token', n_list=[1]),
7         Ngrams(level='char', n_list=[2, 3])
8     ]
9
10    def dr(a): return '/home/chris/data/{0}.csv'.format(a)
11    exp = {
12        # data          instances    config  balance  params
13        'CV(X)':      ([dr('netlog')], (1000, 5000), False, False, ('rbf', 1e-3, 500)),
14        'X->Xi':      ([dr('netlog'),
15                       dr('stuff')], (-1, -1), 'trte', True, None),
16        'CVGrid(X)': ([dr('netlog')], (1000, 1000), 'grid', False, False)
17    }
18
19    for k, v in exp.items():
20        clf = svm.SVC
21        name, features, instances, config, balance, params = (k,) + v
22        experiment(name, features, instances, config, balance, clf, params)
```

---



---

```
1  ---- Shed ----
2
3  Config:
4
5      feature:  token_ngram
6      n_list:   [1, 2]
7      max_feat: None
8
9
10     feature:  char_ngram
11     n_list:   [3]
12     max_feat: None
13
14     name: somename
15     seed: 123
16     clasf: gender
17     pos, neg: 1000, 1000
18
19     Reading from /home/chris/data/netlog.csv...
20     Acquired (1000, 1000) from data.
21     Reading from /home/chris/data/stuff.csv...
22     Acquired (1000, 1000) from data.
```

---

---

```
23 Classes: [0 1]
24
25 Sparse shape: (2000, 8739)
26
27 Total Data instances: 2000
28 Total Label instances: 2000
29
30 Distribution: [(0, 1000), (1, 1000)]
31 Accuracy @ majority baseline: 0.5
32
33 Gridsearch:
34 {'svc__gamma': 0.01, 'svc__C': 100, 'svc__kernel': 'rbf'}
35
36 Tf-CV Result:
37 0.78565912214803
```

---