

**Enough Machine  
Learning To Make  
Hacker News  
Readable Again**

**Ned Jackson Lovely**

**@nedjl / njl@njl.us**

**slides @ [www.njl.us](http://www.njl.us)**

# **A Simple, Achievable Project**

**A personalized filter for Hacker News**

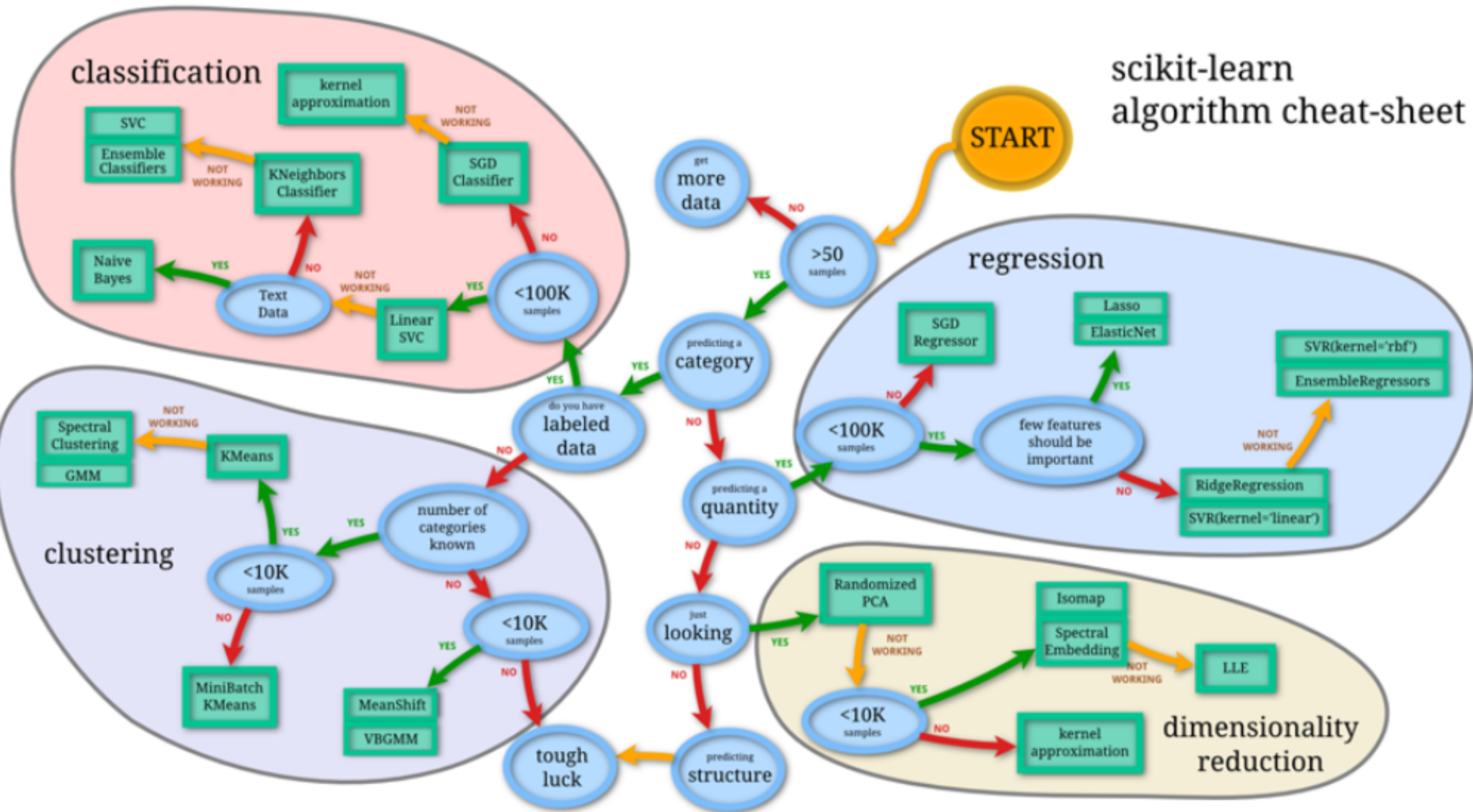
**I Can Machine  
Learn and You  
Can Too!**

**Machine learning is just applying statistics to big piles of data**

**Get Data**  
**Engineer the Data**  
**Train and Tune Models**  
**Apply Model**

**The scikit-learn Docs are Fantastic**

# scikit-learn algorithm cheat-sheet

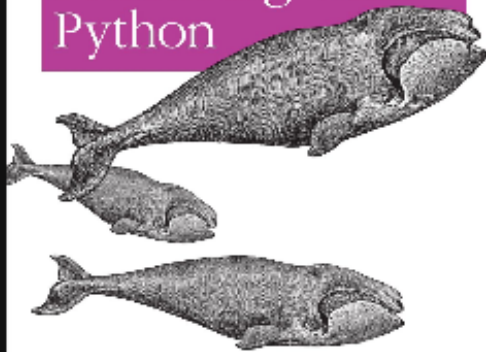




# Supervised vs Unsupervised

Analyzing Text with the Natural Language Toolkit

# Natural Language Processing with Python



O'REILLY\*

Steven Bird, Ewan Klein  
with Richard Loper



Community Experience Distilled

# Building Machine Learning Systems with Python

Master the art of machine learning with Python and build effective machine learning systems with this intensive hands-on guide

Willi Richert  
Luis Pedro Coelho

PACKT open source\*



Community Experience Distilled

# Learning scikit-learn: Machine Learning in Python

Experience the benefits of machine learning techniques by applying them to real-world problems using Python and the open source scikit-learn library

Raúl Garriga  
Guillermo Moncecchi

PACKT open source\*

Building Smart Web 2.0 Applications

Programming

# Collective Intelligence



O'REILLY\*

2,000 chapters  
Reviewed by 116,000+ users

# scikit-learn Patterns

# Parallel Arrays

$X, y$

# Set aside a validation set

```
1 | from sklearn.cross_validation import train_test_split  
2 | X, X_val, y, y_val = train_test_split(X_full, y_full)
```

**How Do You Work?**

**Build Pipelines**

**Optimize Parameters**

# Pipeline

**A sequence of operations on  
data**

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.feature_extraction import text
3 from sklearn.svm import LinearSVC
4
5 p = Pipeline([ ('hv', text.HashingVectorizer()),
6               ('svm', LinearSVC()) ])
```



# Hyper-Parameters

**Tune the Magic**

```
1 from sklearn.grid_search import GridSearchCV
2 params = { 'svm__C':[0.5, 1.0, 2.0, 4.0],
3           'svm__loss':['l1', 'l2'],
4           'hv__ngram_range':[ (1,1), (1,2)],}
5 gs = GridSearchCV(p, params, verbose=2, n_jobs=-1)
6 gs = gs.fit(X,y)
```

# Functions You'll See a Lot Of

`transform()`

`fit()`

`predict()`

# transform(X [,y])

```
1 hv = HashingVectorizer()  
2 hv.transform(['Simple is better than complex.',  
3             'Sparse is better than dense.'])  
4 <2x1048576 sparse matrix of type '<type 'numpy.float64'>'  
5 with 10 stored elements ...
```

# fit(X,y)

```
1 svm = LinearSVC()  
2 svm = svm.fit(X, y)
```

# predict(X)

```
1 | y_predictions = svm.predict(X_new)
```

# Dealing with the Super-Messy Real World



# Get the Data

requests & lxml





	1852	73.0860299921%
	677	26.7166535122%
	5	0.197316495659%
	2534	

Bring The Pain

rd data from Nieman-Marcus.	Framed	Y	Good
In 5 out of 6 stars for protecting your privacy	Framed	Y	Good
Course	Framed	Y	Good
neurship: How to Avoid Becoming a Stressed Out Loner	Framed	Y	Good
roducer Multi-consumer Queue on Ring Buffer	Framed	Y	Good
olution claimed for Navier-Stokes equations, one of the millennium problems	Framed	Y	Good
ne Day We Fight Back: MIT website hacked by Anonymous	Framed	Y	Good
onymous hacks MIT website on anniversary of Aaron Swartz suicide	Framed	Y	Good

# Classify Dreck and Non-Dreck

**Title, URL, Submitter, Content of Link,  
Rank, Votes, Comments, Time of Day,  
Dreck or Not**



# **Messy Data to Normalized Numpy Arrays**



**Bag Of Words**

**n-grams**

**Normalization**

**Stop Words**

**TF-IDF**



# Bag of Words

"Time Flies Like An Arrow,  
Fruit Flies Like Bananas"

Flies	2
Like	2
Time	1
An	1
Arrow	1
Fruit	1
Bananas	1
Bitcoin	0
...	

# n-grams

"Time Flies Like An Arrow,  
Fruit Flies Like Bananas"

Time Flies
Flies Like
Like An
An Arrow
Arrow Fruit
...



# Normalization

```
1 from nltk.stem.snowball import EnglishStemmer
2 stemmer = EnglishStemmer()
3 input='when he flies he likes to fly upon early flights'
4 print(' '.join(stemmer.stem(x) for x in input.split()))
5
6 "when he fli he like to fli upon earli flight"
```

# Stop Words

```
1 from sklearn.feature_extraction import text
2
3 print(', '.join(x for x in
4                 list(text.ENGLISH_STOP_WORDS)[:8]))
5
6 all, six, less, being, indeed, over, move, anyway
```



# TF-IDF

Term-Frequency, Inverse Document Frequency

```
1 from sklearn.feature_extraction import text
2
3 tfidf = text.TfidfVectorizer()
4 tfidf = tfidf.fit(X, y)
```

# Engineering Features



## Pulling Out the Relevant Text

```
1 from readability.readability import Document
2 from lxml.html import fragment_fromstring
3
4 def get_relevant(raw):
5     try:
6         cleaned = Document(raw).summary(True)
7         et = fragment_fromstring(cleaned)
8         return ''.join(x.text for
9                         x in et.getiterator()
10                        if x.text)
11     except Exception as e:
12         print(e)
13     return ''
```





## My Data Starts as Dictionaries

```
1 from sklearn import base
2
3 class DictFeatureExtractor(base.BaseEstimator,
4                             base.TransformerMixin):
5     def __init__(self, prop_name):
6         self.prop_name = prop_name
7     def fit(self, X, y=None):
8         return self
9     def transform(self, X, y=None):
10        return [x[self.prop_name] for x in X]
```

## Maybe Hostnames Are Relevant

```
1 import urlparse
2
3 class HostnameExtractor(BaseEstimator, TransformerMixin):
4     def fit(self, X, y=None):
5         return self
6     def transform(self, X, y=None):
7         return [{'host': urlparse.urlparse(x).netloc}
8                 for x in X]
```

## A Hostname Pipeline

```
1 from sklearn.feature_extraction import DictVectorizer
2
3 hosts = Pipeline([
4     ('dfe', DictFeatureExtractor('url')),
5     ('he', HostnameExtractor()),
6     ('dv', DictVectorizer()), ])
```



# The Actual Appli- cation



# Predict

```
1 from pickle import load
2 classifier = load(open('ridiculous_pipeline.pickle'))
3 classifier.predict(X_input_data)
```

	1 [Full-disclosure] Administrivia: The End	171 points by Morgawr 2 hours ago comments
ti triangle page to end most sierpinski triangle pages	56 points by deede 2 hours ago comments	
	3 Needy robotic toaster sells itself if neglected	20 points by Baustin 1 hour ago comments
	4 Lessons from a Silicon Valley job search	11 points by padolsey 41 minutes ago comments
	5 Nodemailer: Easy as cake e-mail sending from your Node.js applications	49 points by tilt 4 hours ago comments
from an IMAP server failure	28 points by Aoyagi 2 hours ago comments	
	7 Results of operation against child exploitation enterprise on tor network	21 points by alterrize 2 hours ago comments
	8 Why air travel hasn't gotten any faster since the 1960s	9 points by rmorell 1 hour ago comments
	9 What Happens to Older Developers?	192 points by dsirijus 4 hours ago comments
	10 Moto 360	587 points by uptown 21 hours ago comments
	11 Sony announces Project Morpheus, a virtual reality headset coming to PS4	176 points by nzonbi 12 hours ago comments
	12 The new Expats.Stackexchange Public Beta has launched	47 points by nsaparanoid 6 hours ago comments
	13 Storage Pod 4.0: Direct Wire Drives – Faster, Simpler and Less Expensive	8 points by nuriaion 1 hour ago comments
	14 Model View Whatever	3 points by kharlou 7 minutes ago comments
onfiguration Patterns	9 points by frist45 2 hours ago comments	
JS: Divide and conquer complex web apps	44 points by porker 7 hours ago comments	
	17 Android coming to wearables	374 points by deepblueocean 22 hours ago comments
	18 I don't want to be a Real Programmer	54 points by verybadalloc 6 hours ago comments

**hn.njl.us**



**If You Have The  
Data...**



# Unsupervised Learning

Cluster Hacker News Articles



# Regression

Predict article scores



**Machine Learning is  
Becoming Engineering  
Go Engineer!**





# Thank You!

**Ned Jackson Lovely**

**@nedjl / njl@njl.us**

**slides @ [www.njl.us](http://www.njl.us)**

