# pomegranate

fast and flexible probabilistic modelling in python

Jacob Schreiber
Paul G. Allen School of Computer Science & Engineering
University of Washington

jmschreiber91

@jmschrei

@jmschreiber91

pomegranate is **more flexible** than other packages, **fast**, is **intuitive to use**, and can do it all **in parallel**

**Probability Distributions**
**General Mixture Models**
**Hidden Markov Models**
**Naive Bayes / Bayes' Classifiers**
**Markov Chains**
**Bayesian Networks**
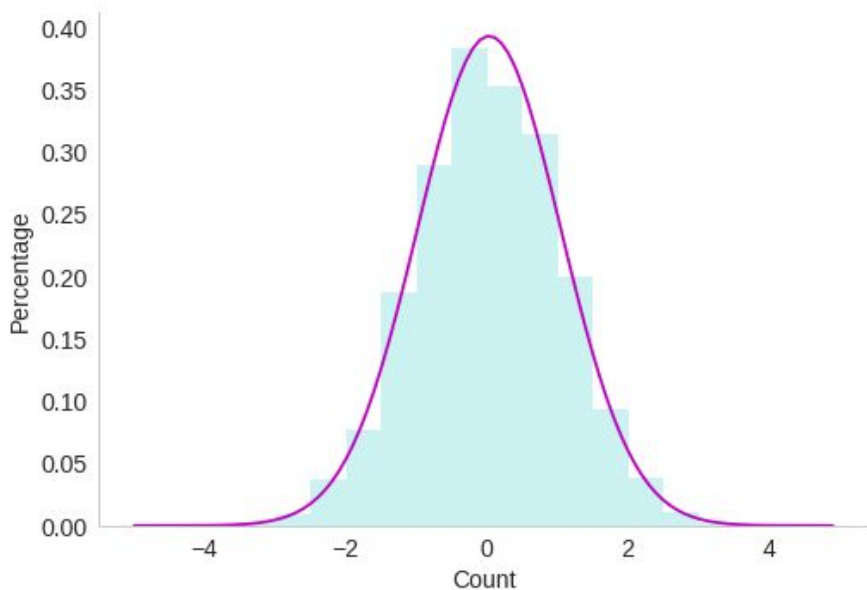

**k-means / kmeans++ / kmeans||**
**Factor graphs**
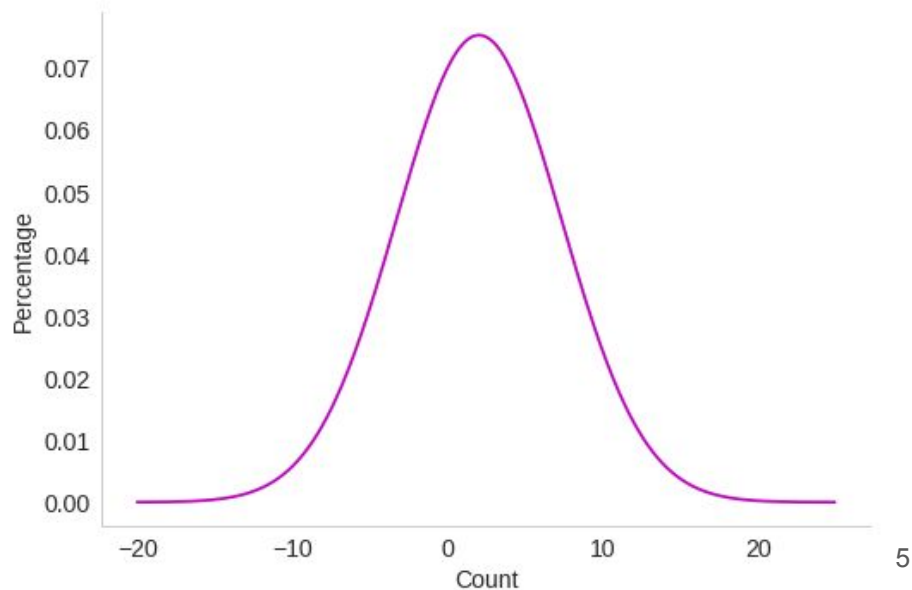
# Models can be made in two ways...

...from data

```
d = NormalDistribution.from_samples(X)
```

...from known values
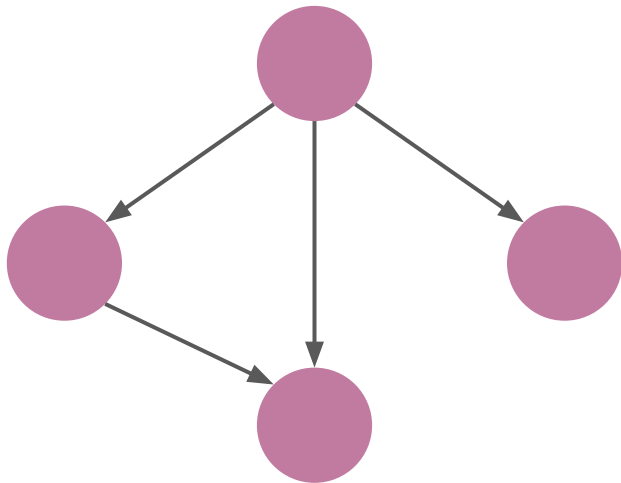
```
d = NormalDistribution(5, 2.3)
```

## ...from data

```
d = BayesianNetwork.from_samples(X)
```



## ...from known values

```
n1 = Node(...)
n2 = Node(...)
model = BayesianNetwork()
model.add_nodes(n1, n2...)
model.add_edges(...)
```

# The API is common to all models

model.log_probability(X) / model.probability(X)

model.sample()

model.fit(X, weights, inertia)

All models have these methods!

model.summarize(X, weights)

model.from_summaries(inertia)

Model.from_samples(X, weights)

model.predict(X)

model.predict_proba(X)

All models composed of distributions (like GMM, HMM...) have these methods too!

model.predict_log_proba(X)

```
GeneralMixtureModel.from_samples(NormalDistribution, n_components=3, X=X)

GeneralMixtureModel.from_samples(ExponentialDistribution, n_components=3,
X=X)

BayesClassifier.from_samples(MultivariateGaussianDistribution, X, y)


d1 = GeneralMixtureModel.from_samples...
d2 = GeneralMixtureModel.from_samples...
model = BayesClassifier([d1, d2])
```

# pomegranate can be faster than numpy

Fitting Multivariate Gaussian to 10,000,000 samples of 10 dimensions

```python
data = numpy.random.randn(10000000, 10)

print "numpy time:"
%timeit -n 10 data.mean(axis=0), numpy.cov(data, rowvar=False, bias=True)
print "\n" "pomegranate time:"
%timeit -n 10 MultivariateGaussianDistribution.from_samples(data)
```

```
numpy time:
10 loops, best of 3: 3.52 s per loop

pomegranate time:
10 loops, best of 3: 2.87 s per loop
```

pomegranate reduces data to sufficient statistics for updates and so only has to go datasets once (for all models).

Here is an example of the Normal Distribution sufficient statistics

$$\sum_{i=1}^{n} w_i \qquad \sum_{i=1}^{n} w_i x_i \qquad \sum_{i=1}^{n} w_i x_i^2 \longrightarrow$$

$$\mu = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i}$$

$$\sigma^2 = \frac{\sum_{i=1}^{n} w_i x_i^2}{\sum_{i=1}^{n} w_i} - \frac{\left(\sum_{i=1}^{n} w_i x_i\right)^2}{\left(\sum_{i=1}^{n} w_i\right)^2}$$

Batches from a dataset can be reduced to additive summary statistics, enabling exact updates from data that can't fit in memory.



On Disk     In Memory     Extract summaries     Sufficient Statistics     New Parameters

# pomegranate supports mini-batching

Instead of going through the full dataset before updating parameters, one could update parameters at each step.

Dataset

Extract summaries

On Disk     In Memory     Sufficient Statistics     New Parameters

# pomegranate supports parallelization

Multiple batches can be loaded at the same time and processed by different threads using `n_jobs` in either fitting or prediction methods



On Disk  In Memory    Extract summaries   Sufficient Statistics  New Parameters

Dataset

Time to Train HMM Mixture

```
model.fit(X, n_jobs=n)
```

For many tasks, there is limited labeled data but a deluge of unlabeled data, and one wants to utilize both.



15

Summaries from MLE on the labeled data can be added to summaries from EM on the unlabeled data

# pomegranate allows semisupervised learning



Supervised Accuracy: 0.93

Semisupervised Accuracy: 0.96

Many real world tasks involve missing data, but common approaches aren't sufficient for tackling the problem.



Bimodal Distribution

# pomegranate supports missing data

Many real world tasks involve missing data, but common approaches aren't sufficient for tackling the problem.

# pomegranate supports missing data

Many real world tasks involve missing data, but common approaches aren't sufficient for tackling the problem.

## Error in Multivariate Gaussian Covariance Matrix

# pomegranate supports missing data

Pomegranate supports **model fitting**, **structure learning**, and **model predictions** on data sets that include missing values, no matter how complicated the model or sparse the data set.

You can **fit a Gaussian mixture model** to incomplete data sets.

You can run the **Viterbi or forward-backward algorithm** using a HMM on incomplete data sets.

You can **learn the structure of a Bayesian network** on incomplete data sets.

All without having to change your command, simply by including np.nan in the place of the missing value

Multivariate Gaussian with GPU Acceleration

# pomegranate can be faster than scipy

```python
mu, cov = numpy.random.randn(2000), numpy.eye(2000)
d = MultivariateGaussianDistribution(mu, cov)
X = numpy.random.randn(2000, 2000)
print "scipy time: ",
%timeit multivariate_normal.logpdf(X, mu, cov)
print "pomegranate time: ",
%timeit MultivariateGaussianDistribution(mu, cov).log_probability(X)
print "pomegranate time (w/ precreated object): ",
%timeit d.log_probability(X)
```

```
scipy time: 1 loop, best of 3: 1.67 s per loop
 pomegranate time: 1 loop, best of 3: 801 ms per loop
 pomegranate time (w/ precreated object): 1 loop, best of 3: 216 ms per loop
```

$$P(X|\mu,\sigma) = \frac{1}{\sqrt{2\pi}\sigma}exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$logP(X|\mu,\sigma) = -\log\left(\sqrt{2\pi}\sigma\right) - \frac{(x-\mu)^2}{2\sigma^2}$$

$$logP(X|\mu,\sigma) = \alpha - \frac{(x-\mu)^2}{\beta}$$

GOSSIPGIRL

Spotted: Lonely Boy. Can't believe the love of his life has returned. If only she knew who he was. But everyone knows Serena. And everyone is talking. Wonder what Blair Waldorf thinks. Sure, they're BFF's, but we always thought Blair's boyfriend Nate had a thing for Serena.

Why'd she leave? Why'd she return? Send me all the deets. And who am I? That's the secret I'll never tell. The only one. —XOXO. Gossip Girl.

Better lock it down with Nate, B. Clock's ticking.

+1 Nate
-1 Blair

This just in: S and B committing a crime of fashion. Who doesn't love a five-finger discount. Especially if it's the middle one.

-1 Blair
-1 Serena

# Simple summations don't work well

# Beta distributions can model uncertainty

# Bayesian networks

Bayesian networks are powerful inference tools which define a dependency structure between variables.

# Bayesian networks

Two main difficult tasks:
(1) Inference given incomplete information
(2) Learning the dependency structure from data

**???**

Three primary ways:
- "Search and score" / Exact
- "Constraint Learning" / PC
- Heuristics

???

pomegranate supports:
- "Search and score" / Exact
- "Constraint Learning" / PC
- Heuristics

Time to Learn Structure

$P(D|M)$ with Resulting Model

Easy! Tractable!

**Global Parameter Independence:** The parents of some variable A are independent of the parents of some variable B given that they don't form a cycle in the resulting graph

Hard! Exponential Time!

Variables to be Modeled

Constraint Graph

Constraint Graph

**Global Parameter Independence:** The parents of some variable A are independent of the parents of some variable B given that they don't form a cycle in the resulting graph

Constraint Graph



Time To Learn Bayesian Network

# Finding the optimal Bayesian network given a constraint graph

Jacob M. Schreiber[1] and William S. Noble[2]

[1] Department of Computer Science, University of Washington, Seattle, WA, United States of America
[2] Department of Genome Science, University of Washington, Seattle, WA, United States of America

## ABSTRACT

Despite recent algorithmic improvements, learning the optimal structure of a Bayesian network from data is typically infeasible past a few dozen variables. Fortunately, domain knowledge can frequently be exploited to achieve dramatic computational savings, and in many cases domain knowledge can even make structure learning tractable. Several methods have previously been described for representing this type of structural prior

# Naive Bayes produces ellipsoid boundaries



pomegranate naive Bayes

sklearn naive Bayes

```
model = NaiveBayes.from_samples(NormalDistribution, X, y)
```

$$Posterior = \frac{Likelihood * Prior}{Normalization}$$

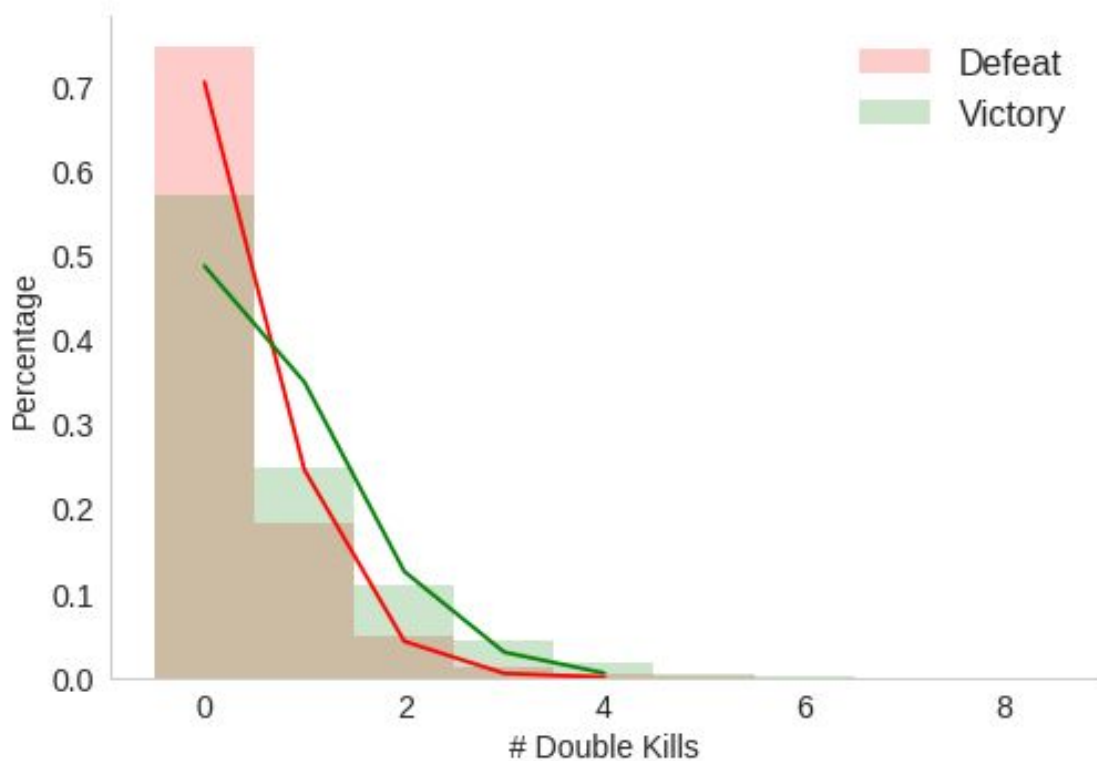$$P(M|D) = \frac{\prod_{i=1}^{d} P(D_i|M)P(M)}{\sum_{M} \prod_{i=1}^{d} P(D_i|M)P(M)}$$

```
dists = [LogNormalDistribution, PoissonDistribution,
ExponentialDistribution, PoissonDistribution]

model1 = NaiveBayes.from_samples(NormalDistribution, X, y)
model2 = NaiveBayes.from_samples(dists, X, y)
model3 = GaussianNB().fit(X, y)
```
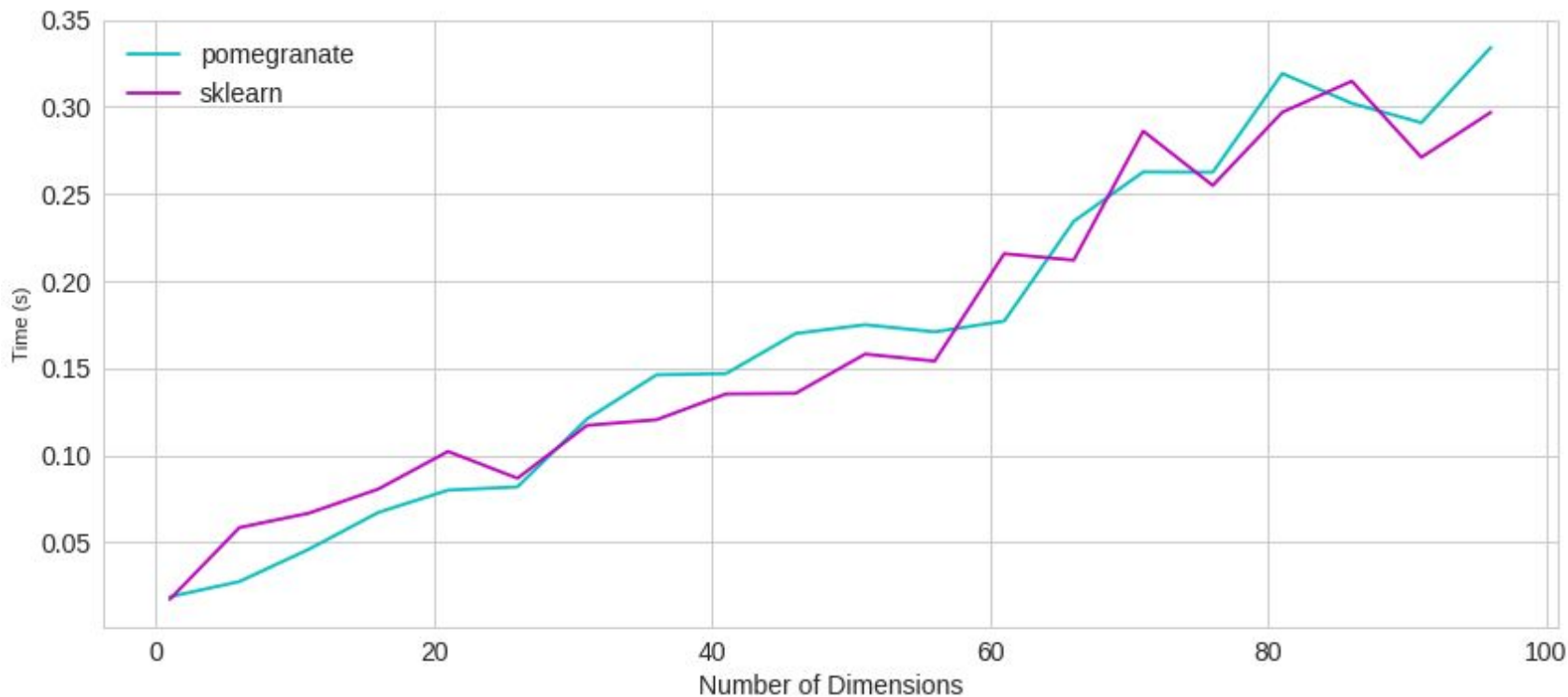
pomegranate Gaussian Naive Bayes:  0.711
sklearn Gaussian Naive Bayes:  0.711
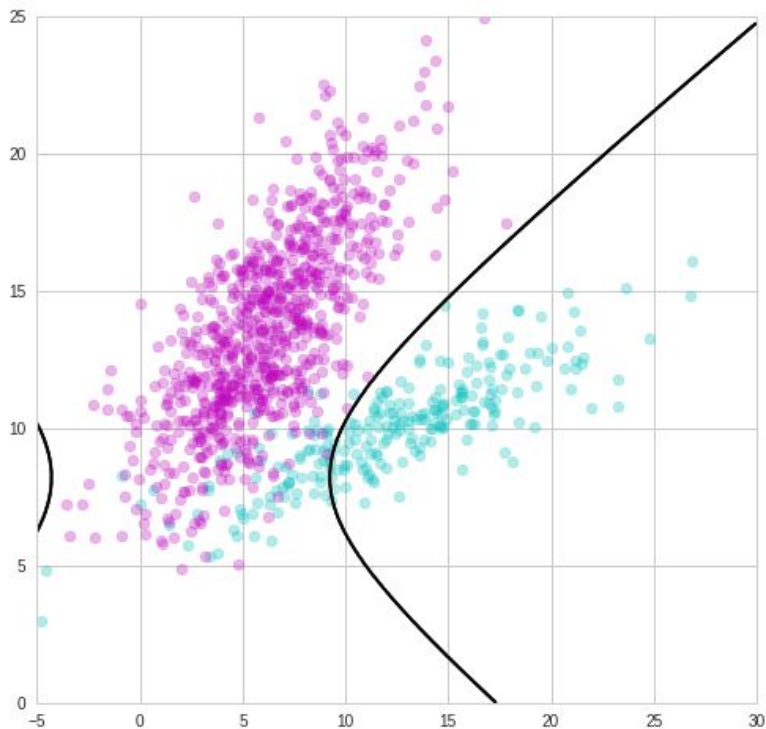Heterogeneous Naive Bayes:  0.726

naive accuracy: 0.929

bayes classifier accuracy: 0.966

# pomegranate: fast and flexible probabilistic modeling in python

**Jacob Schreiber**
Paul G. Allen School of Computer Science
University of Washington
Seattle, WA 98195
jmschr@cs.washington.edu

## Abstract

We present pomegranate, an open source machine learning package for probabilistic modeling in Python. Probabilistic modeling encompasses a wide range of methods that explicitly describe uncertainty using probability distributions. Three widely used probabilistic models implemented in pomegranate are general mixture models, hidden Markov models, and Bayesian networks. A primary focus of pomegranate is to abstract away the complexities of training models from their definition. This allows users to focus on specifying the correct model for their

57

# pomegranate is NumFOCUS affiliated

## NUMF⊙CUS
### OPEN CODE = BETTER SCIENCE

Home    About ⌄    **Open Source Projects** ⌄    Community ⌄    Programs ⌄    Blog

## pomegranate

**pomegranate**

pomegranate is a Python module for fast and flexible probabilistic modeling inspired by the design of scikit-learn. A primary focus of pomegranate is to abstract away the intricacies of a model from its definition, allowing users to easily prototype with complex models and training strategies. Its modular implementation allows for probability distributions to be swapped in or out for each other with ease and for models to be stacked within each other, yielding such delights as a mixture of Bayesian networks or a Gaussian mixture model Bayes classifier.

https://www.numfocus.org/open-source-projects/affiliated-projects/

# Documentation available at Readthedocs



https://pomegranate.readthedocs.io/en/latest/

# Tutorials available on GitHub

Branch: master ▾     **pomegranate** / **tutorials** /

Create new file | Upload files | Find file | History

jmschrei ENH NB/BC notebook                                     Latest commit 5cd8d68 5 days ago

.. 

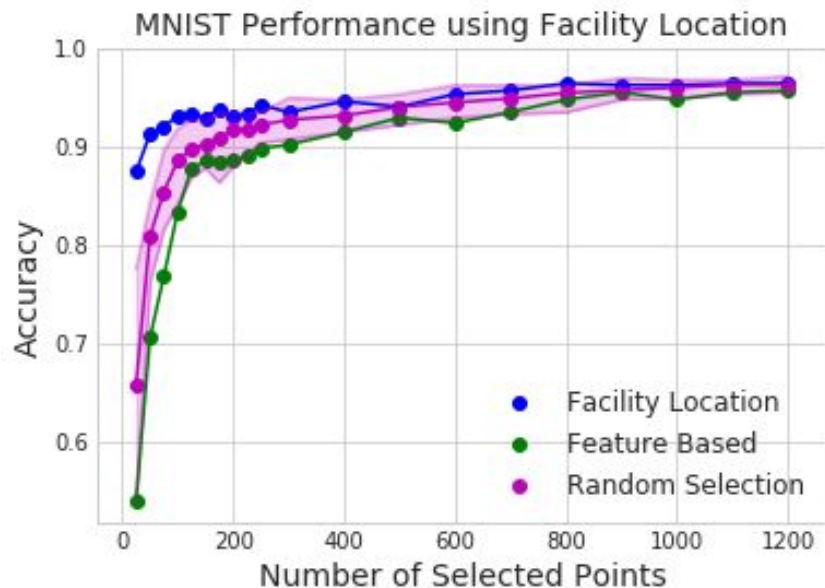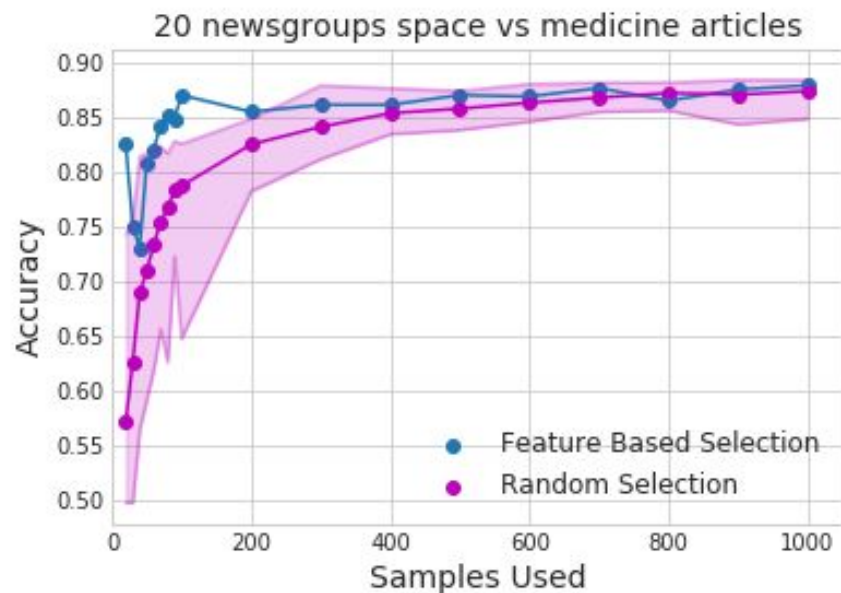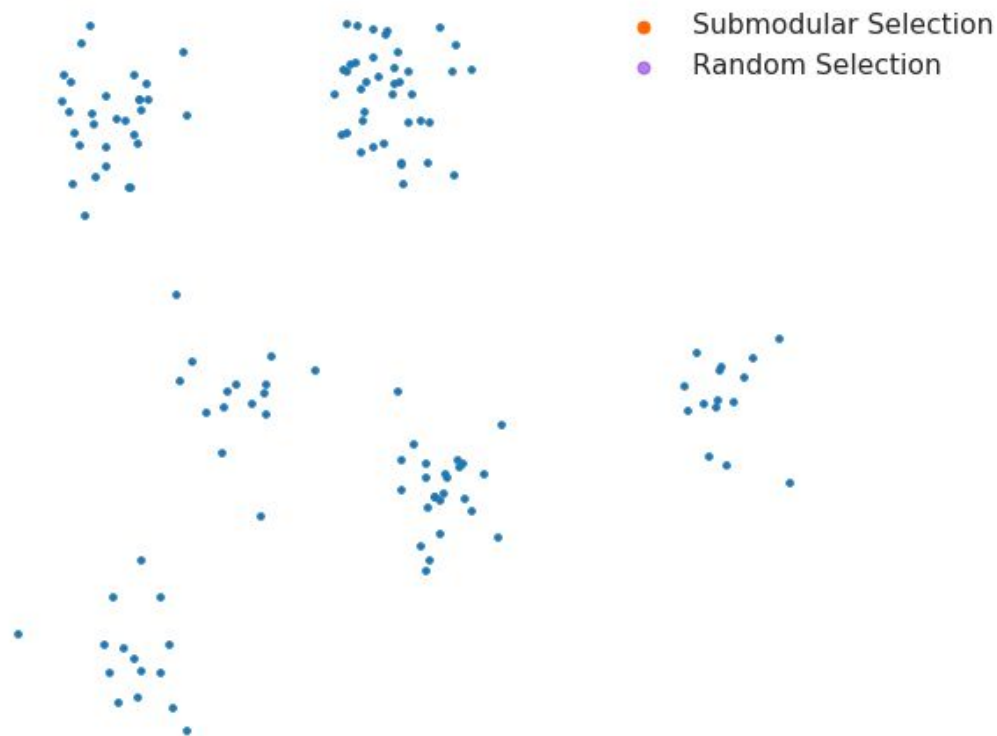| 📁 old | ADD new overview tutorial | a month ago |
| 📄 A_Overview.ipynb | ADD new notebook features | 12 days ago |
| 📄 B_Model_Tutorial_1_Distributions.ipynb | ENH NB/BC notebook | 5 days ago |
| 📄 B_Model_Tutorial_2_General_Mixture_Models.ipynb | ADD new notebook features | 12 days ago |
| 📄 B_Model_Tutorial_3_Hidden_Markov_Models.ipynb | ADD new notebook features | 12 days ago |
| 📄 B_Model_Tutorial_4_Bayesian_Networks.ipynb | ENH NB/BC notebook | 5 days ago |
| 📄 B_Model_Tutorial_4b_Bayesian_Network_Structure_Learning.ip... | ADD new notebook features | 12 days ago |
| 📄 B_Model_Tutorial_5_Bayes_Classifiers.ipynb | ENH NB/BC notebook | 5 days ago |
| 📄 B_Model_Tutorial_6_Markov_Chain.ipynb | ADD new notebook features | 12 days ago |
| 📄 C_Feature_Tutorial_1_Parallelization_and_GPUs.ipynb | ADD new notebook features | 12 days ago |
| 📄 C_Feature_Tutorial_8_Semisupervised_Learning.ipynb | ADD new notebook features | 12 days ago |
| 📄 C_Feature_Tutorial_9_Missing_Values.ipynb | ADD new notebook features | 12 days ago |
| 📄 GGBlasts.xlsx | PyData Chicago 2016 | 2 years ago |
| 📄 README.md | Update README.md | 3 years ago |

https://github.com/jmschrei/pomegranate/tree/master/tutorials

60

# apricot implements submodular selection for training machine learning models faster



20 newsgroups space vs medicine articles

MNIST Performance using Facility Location

https://github.com/jmschrei/apricot

# apricot implements submodular selection for training machine learning models faster



Submodular Selection
Random Selection

https://github.com/jmschrei/apricot

# pomegranate

fast and flexible probabilistic modelling in python

Jacob Schreiber
Paul G. Allen School of Computer Science & Engineering
University of Washington

jmschreiber91

@jmschrei

@jmschreiber91