

MACHINE LEARNING WITH PYTHON

Themistoklis Diamantopoulos

Contents

- Part 1:
 - Introduction to Machine Learning
 - Classification with Decision Trees
 - Classification with Naïve Bayes
 - Classification with SVMs
- Part 2
 - Classification and Regression with kNN
 - Linear and Polynomial Regression
 - Feature Selection and Feature Extraction
 - Centroid-based and Connectivity-based Clustering

INTRODUCTION TO MACHINE LEARNING

What is Machine Learning?

- Subfield of Artificial Intelligence
- Term coined in 1959 by Arthur Samuel

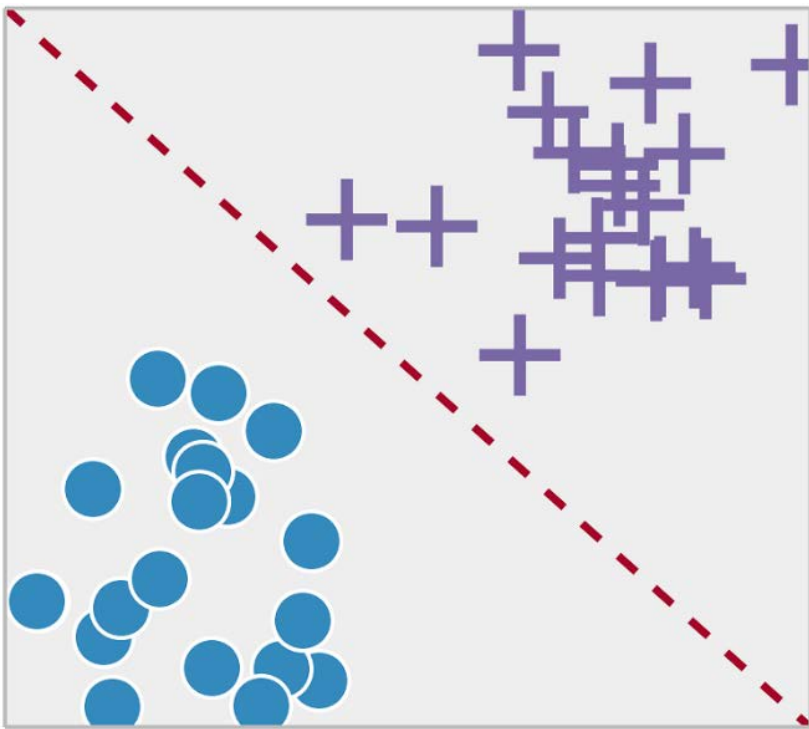
Progressively improve performance on a specific task with data, without being explicitly programmed

Types of Machine Learning tasks

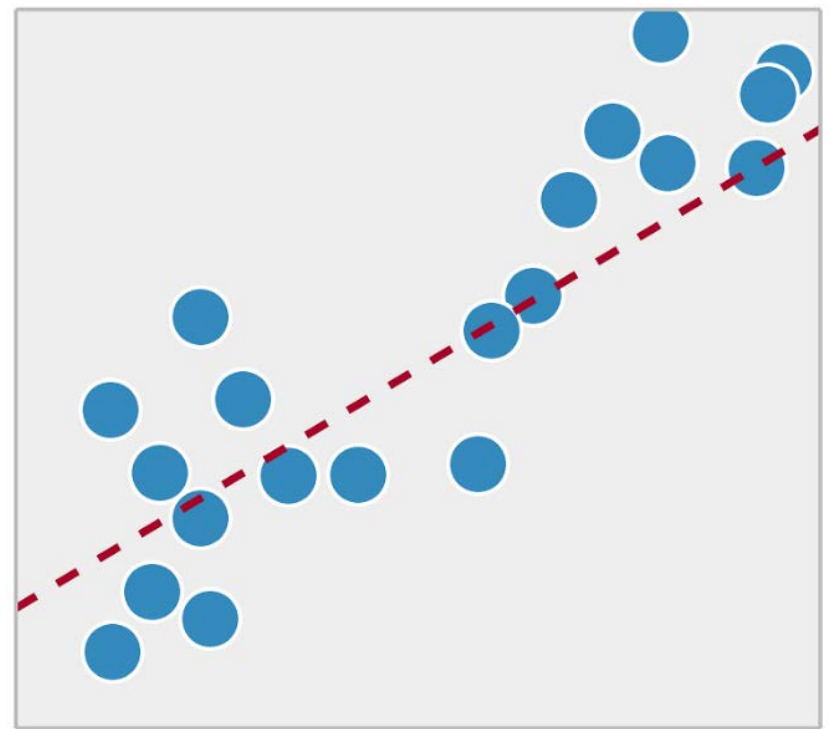
- Supervised Learning
 - Learn output based on input data
- Unsupervised Learning
 - Find structure in given data
- Reinforcement Learning
 - Learn from the environment

Supervised Learning tasks

Classification



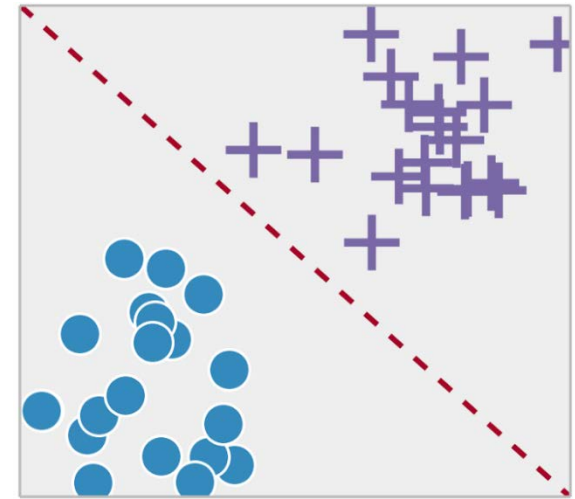
Regression



Classification

- Classify data to 1, 2 or more classes
- Confusion Matrix

		Actual Class	
		Positive	Negative
Predicted Class	Positive	TP	FP
	Negative	FN	TN
		$P = TP + FN$	$N = FP + TN$



- Evaluation Metrics
 - Accuracy = $(TP + TN) / (P + N)$
 - Precision = $TP / (TP + FP)$
 - Recall = $TP / (TP + FN)$

Regression

- Build a model that fits the data
- Actual (y_i) and predicted values (\hat{y}_i)

- Mean Absolute Error

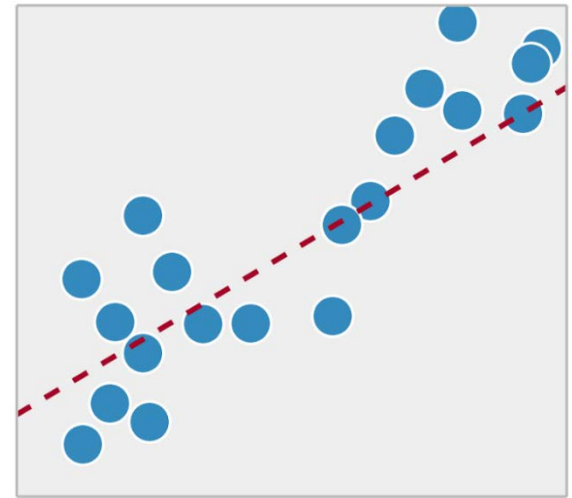
$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

- Mean Squared Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

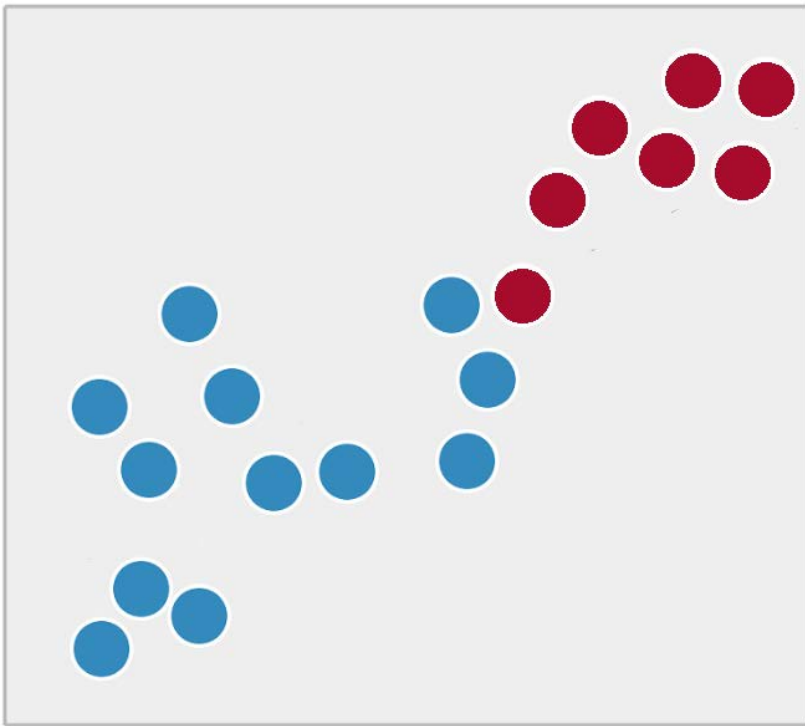
- Coefficient of Determination

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} \quad \text{where} \quad SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{and} \quad SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$$

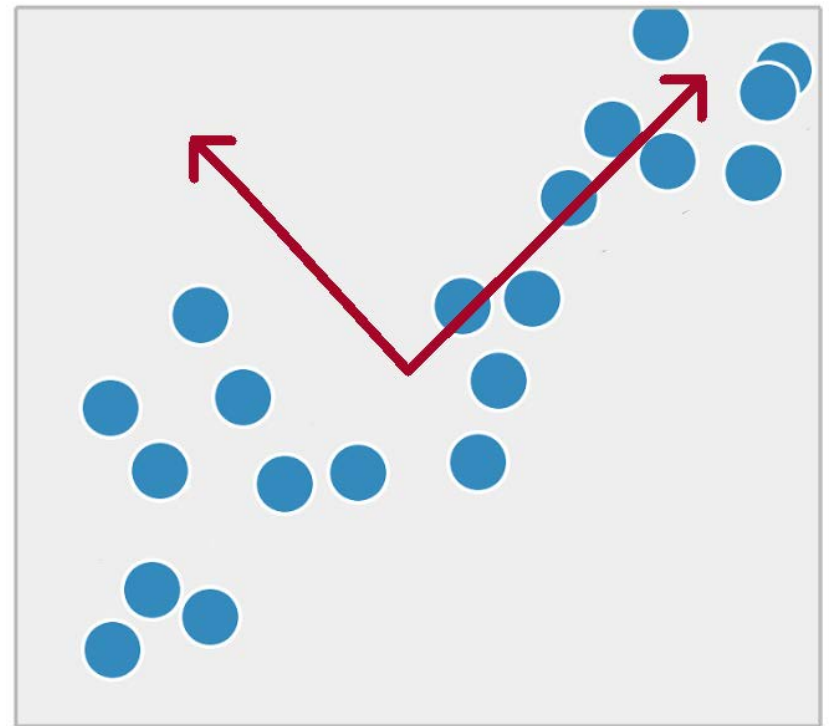


Unsupervised Learning tasks

Clustering



Dimensionality Reduction



Clustering

- Cluster data according to their features
- Evaluation Metrics

- Cohesion (Within Cluster Sum of Squares)

$$SSE = WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

- Separation (Between Cluster Sum of Squares)

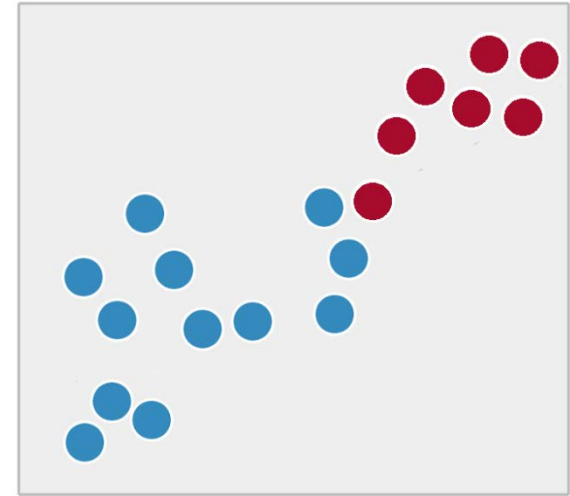
$$BSS = \sum_i |C_i| (m - m_i)^2$$

- Silhouette

$$s = (b - a) / \max(a, b)$$

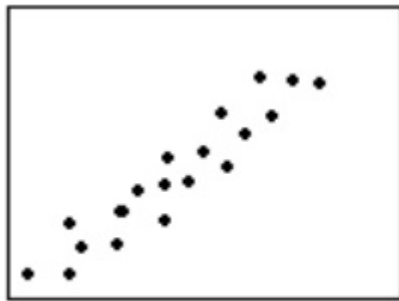
where a = average distance of i to the points in its cluster

b = min(average distance of i to points in another cluster)

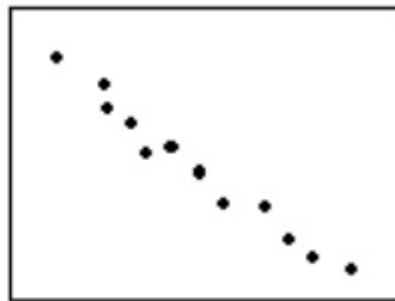


Dimensionality Reduction

- Transform the data to extract useful information
 - Measure correlation
 - Maximize variance

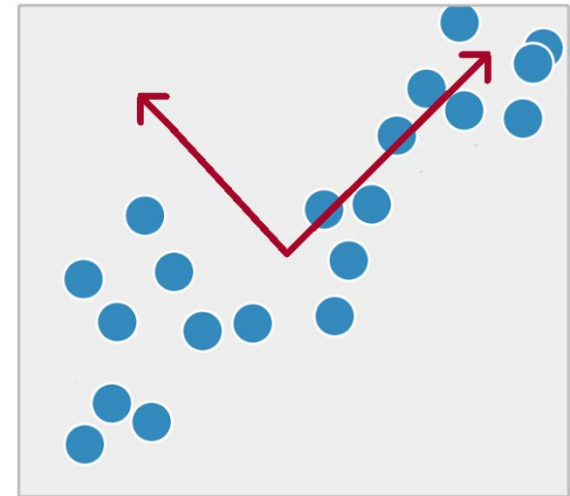


Positive



Negative

- Evaluation Metrics
 - Percentage of Variance
 - Cumulative Percentage of Variance



MACHINE LEARNING METHODOLOGY

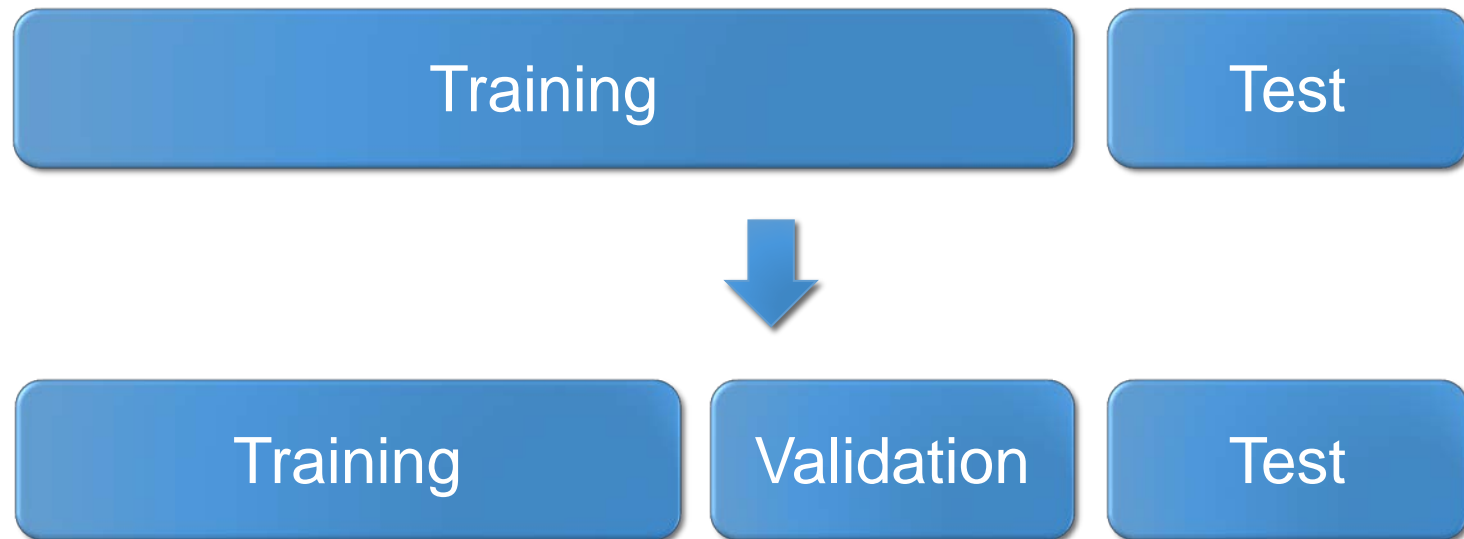
Machine Learning Steps

- Data cleaning and preprocessing equally important with model selection and training



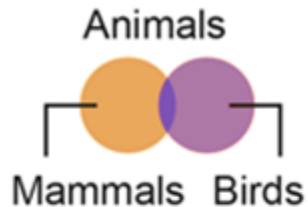
Data Splitting

- Use training data to train the model
 - Some data can be used to validate the model → validation set
 - Use folds of training data for validation → Cross-validation
- Evaluate the model on test data
 - Test set must not overlap with training data



Tribes of Machine Learning

Symbolists



Use symbols, rules, and logic to represent knowledge and draw logical inference

Favored algorithm
Rules and decision trees

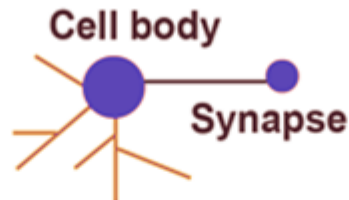
Bayesians



Assess the likelihood of occurrence for probabilistic inference

Favored algorithm
Naive Bayes or Markov

Connectionists



Recognize and generalize patterns dynamically with matrices of probabilistic, weighted neurons

Favored algorithm
Neural networks

Evolutionaries



Generate variations and then assess the fitness of each for a given purpose

Favored algorithm
Genetic programs

Analogizers

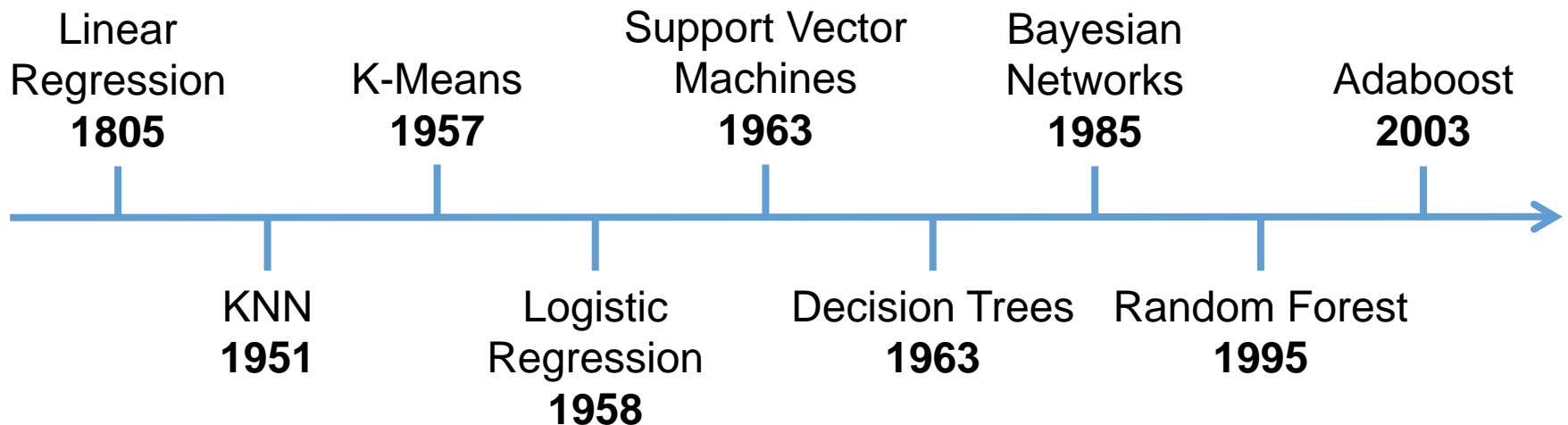


Optimize a function in light of constraints ("going as high as you can while staying on the road")

Favored algorithm
Support vectors

A brief history course

- 1960s: golden age of AI
- Predictive Statistical Algorithms



Machine Learning Applications

- Predictive maintenance or condition monitoring
- Warranty reserve estimation
- Propensity to buy
- Demand forecasting
- Process optimization
- Telematics

Manufacturing



- Predictive inventory planning
- Recommendation engines
- Upsell and cross-channel marketing
- Market segmentation and targeting
- Customer ROI and lifetime value

Retail



- Alerts and diagnostics from real-time patient data
- Disease identification and risk stratification
- Patient triage optimization
- Proactive health management
- Healthcare provider sentiment analysis

Healthcare and Life Sciences



- Aircraft scheduling
- Dynamic pricing
- Social media – consumer feedback and interaction analysis
- Customer complaint resolution
- Traffic patterns and congestion management

Travel and Hospitality



- Risk analytics and regulation
- Customer Segmentation
- Cross-selling and up-selling
- Sales and marketing campaign management
- Credit worthiness evaluation

Financial Services



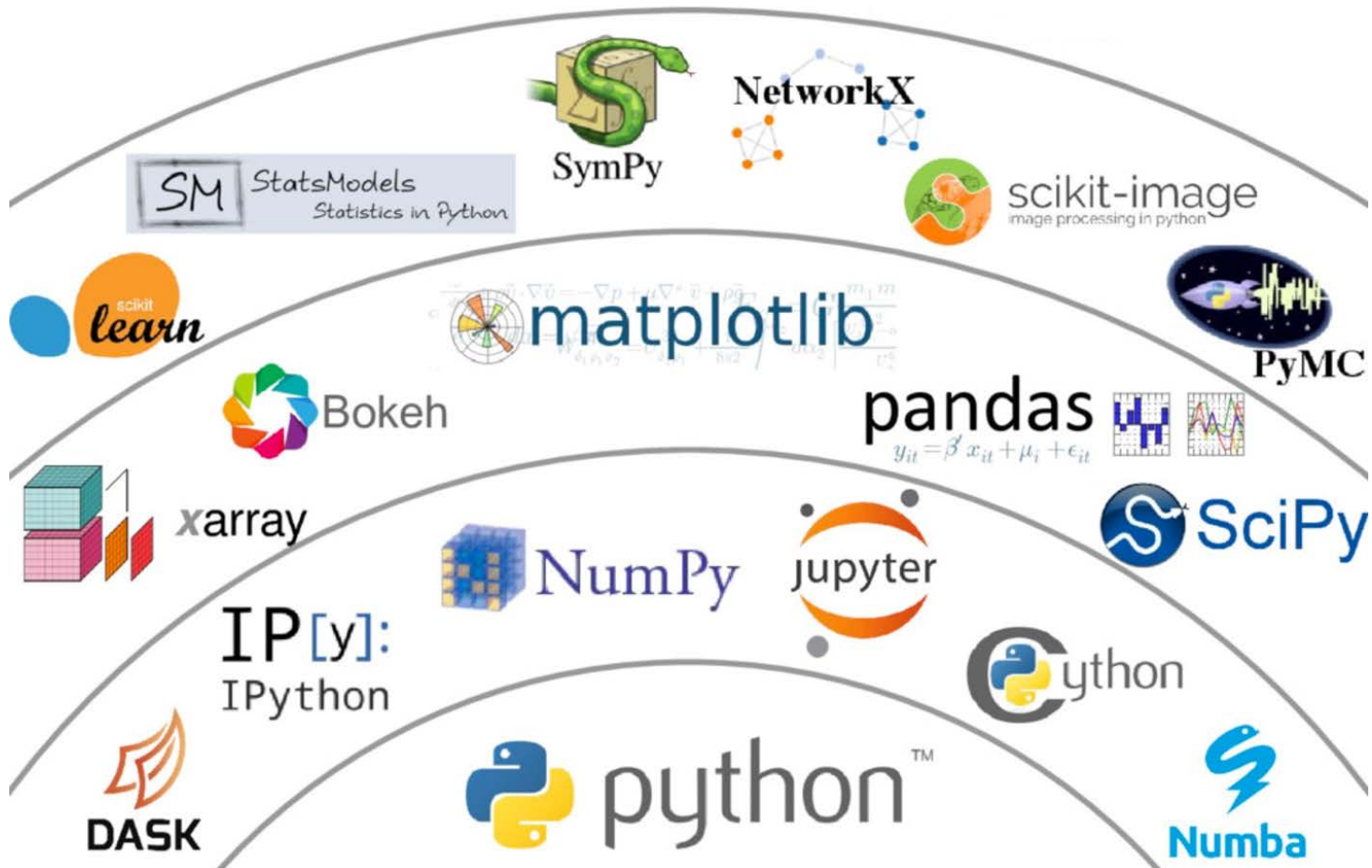
- Power usage analytics
- Seismic data processing
- Carbon emissions and trading
- Customer-specific pricing
- Smart grid management
- Energy demand and supply optimization

Energy, Feedstock, and Utilities



MACHINE LEARNING WITH PYTHON

Set of Powerful Libraries

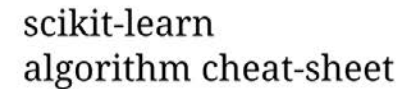


Machine Learning Libraries

- numpy
 - Arrays: universal point of reference in the python ML world
- pandas
 - Data manipulation made easy
- scipy
 - Basis of scientific computing
- scikit-learn
 - (Almost) all machine learning algorithms you will ever need
- matplotlib
 - Plot all of the above

... and all of these are seamlessly connected!

- Almost any model you will ever need



Source: http://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

Scikit-learn Examples

- Very easy to use

```
from sklearn.svm import SVC  
model = SVC(gamma=0.1)  
model.fit(x_train, y_train)  
y_pred = model.predict(x_test)
```

} Create, Train,
and Run Model

- Supports multiple data transformations
- And multiple evaluation metrics

```
from sklearn import metrics  
print(metrics.classification_report(y_test, y_pred))  
print(metrics.confusion_matrix(y_test, y_pred))  
print(metrics.accuracy_score(y_test, y_pred))
```

} Classification
Evaluation Metrics

Time for hands-on!