

Intro to Machine Learning

December BUG

Agenda

- Personal Background
- Defining machine learning
- Real World Applications
- Machine Learning / Data Science Tools
- Machine Learning Workflow
- Demo
- Resources

Personal Background

YouTube



CS:GO

BLIZZARD
ENTERTAINMENT



What is machine learning?

”Machine learning is the science of getting computers to act without being explicitly programmed.”

-Andrew Ng

- Machine learning today involves applying statistical models to datasets to predict data, discover patterns, etc.

ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



MACHINE LEARNING

Machine learning begins to flourish.



DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

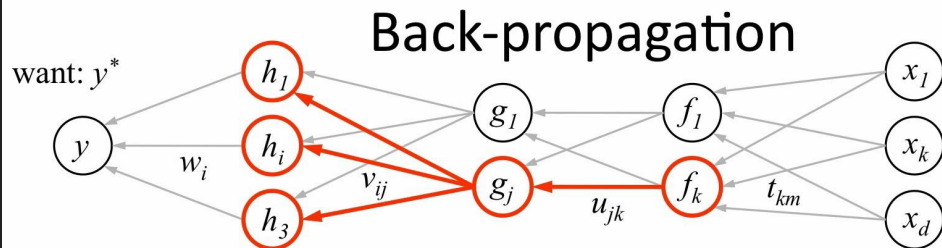
2000's

2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

Deep Learning - Neural Nets and Backpropagation

Backpropagation - Backward propagation of errors



1. receive new observation $\mathbf{x} = [x_1 \dots x_d]$ and target y^*
2. **feed forward:** for each unit g_j in each layer $1 \dots L$
compute g_j based on units f_k from previous layer: $g_j = \sigma \left(u_{j0} + \sum_k u_{jk} f_k \right)$
3. get prediction y and error $(y - y^*)$
4. **back-propagate error:** for each unit g_j in each layer $L \dots 1$

(a) compute error on g_j

$$\frac{\partial E}{\partial g_j} = \sum_i \underbrace{\sigma'(h_i)}_{\text{should } g_j \text{ be higher or lower?}} \underbrace{v_{ij}}_{\text{how } h_i \text{ will change as } g_j \text{ changes}} \underbrace{\frac{\partial E}{\partial h_i}}_{\text{was } h_i \text{ too high or too low?}}$$

(b) for each u_{jk} that affects g_j

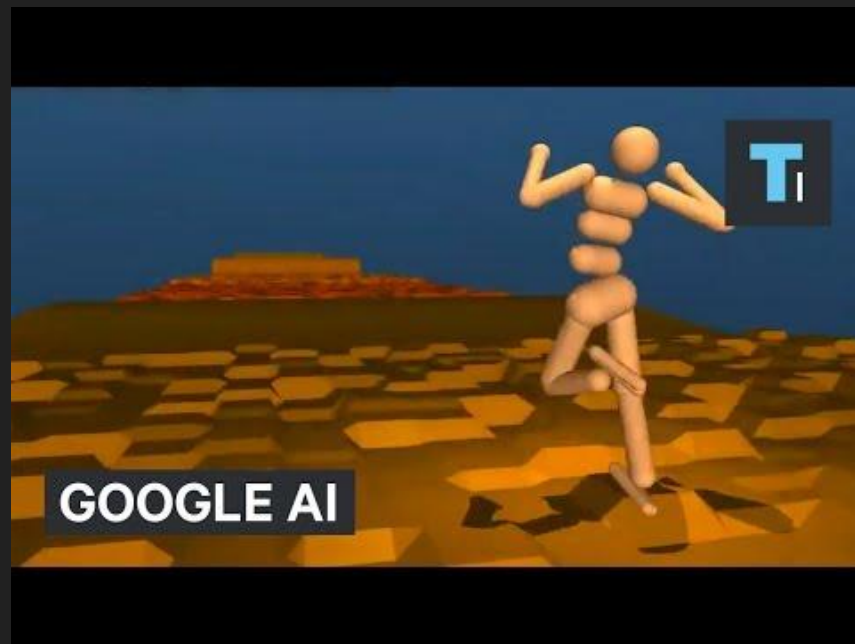
(i) compute error on u_{jk}

$$\frac{\partial E}{\partial u_{jk}} = \frac{\partial E}{\partial g_j} \underbrace{\sigma'(g_j)}_{\text{do we want } g_j \text{ to be higher/lower}} \underbrace{f_k}_{\text{how } g_j \text{ will change if } u_{jk} \text{ is higher/lower}}$$

(ii) update the weight

$$u_{jk} \leftarrow u_{jk} - \eta \frac{\partial E}{\partial u_{jk}}$$

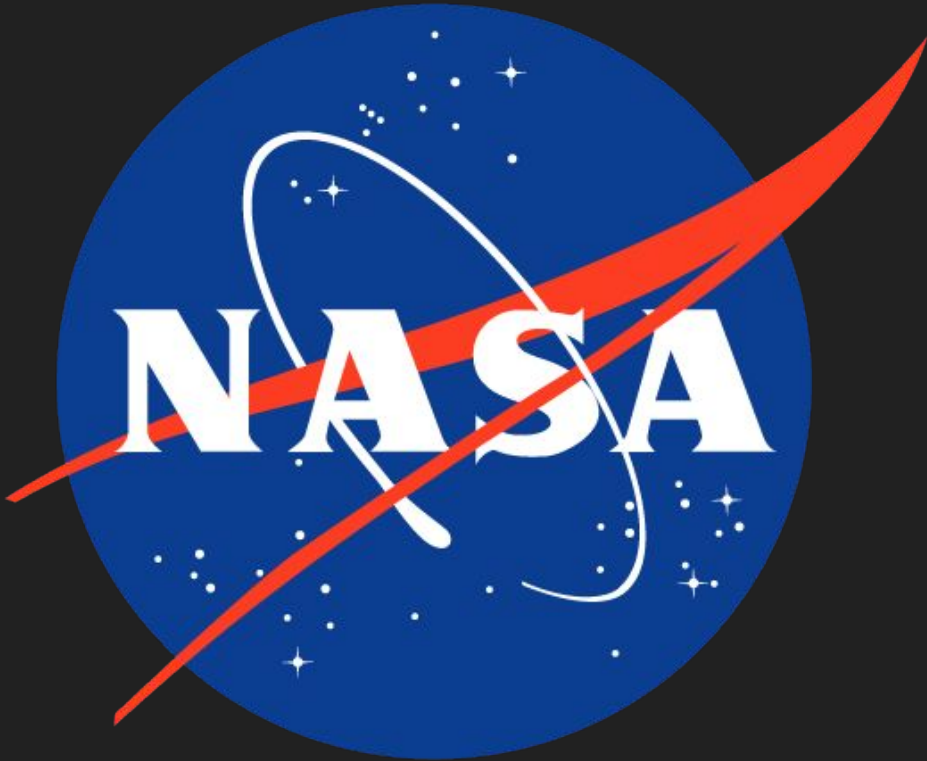
Real World Applications



Real World Applications

The Netflix logo is displayed in a white rectangular box. It consists of the word "NETFLIX" in a bold, red, sans-serif typeface. The letters are slightly slanted to the right, giving it a dynamic feel.

Real World Applications



TYPE



PLAYER



VISION



HEALTH



Machine Learning & Data Science Tools

Data Science / Machine Learning Tools

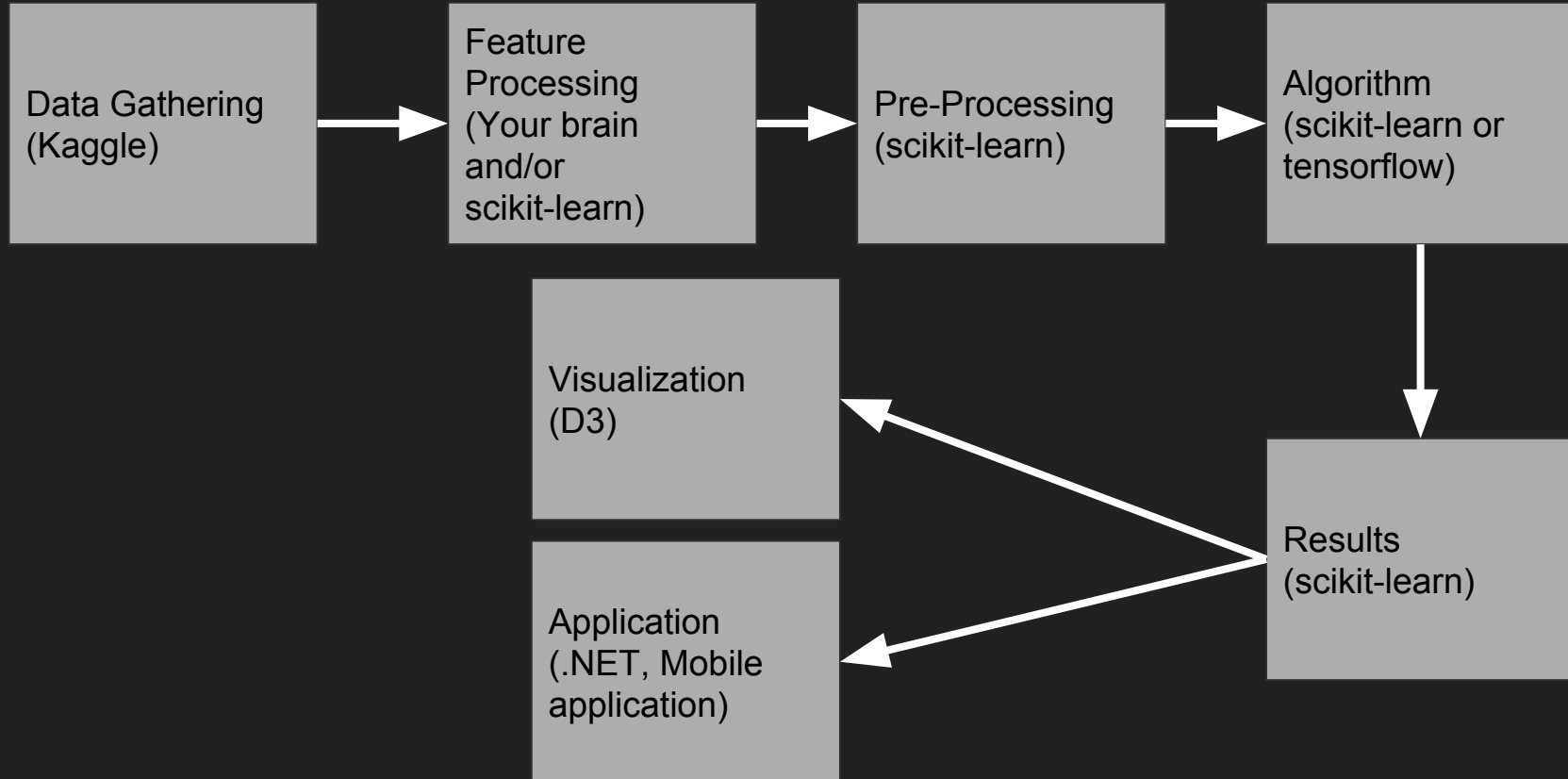


Data Science / Machine Learning Tools

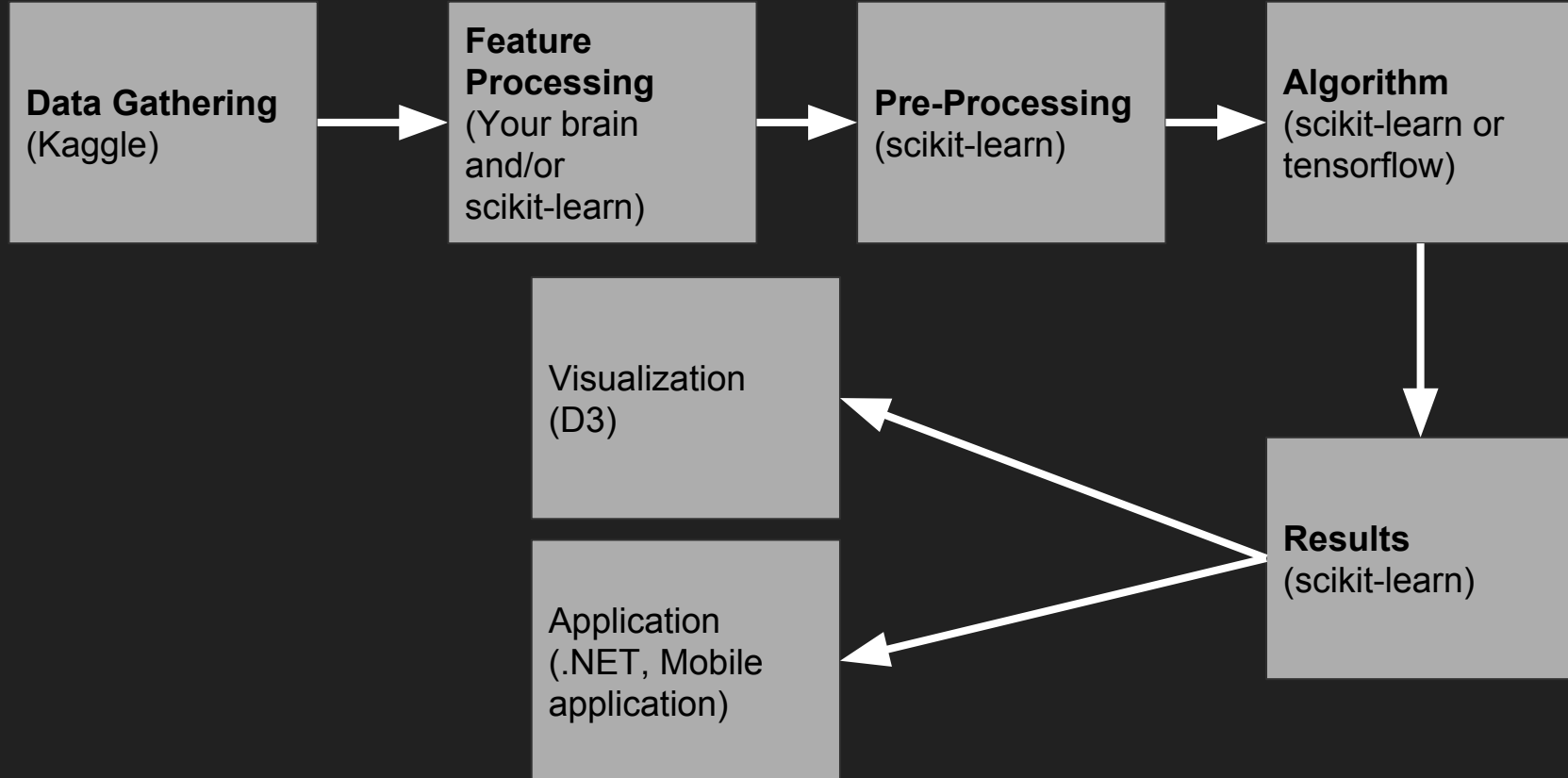


Machine Learning Workflow

Machine Learning Workflow



Machine Learning Workflow



Data Gathering



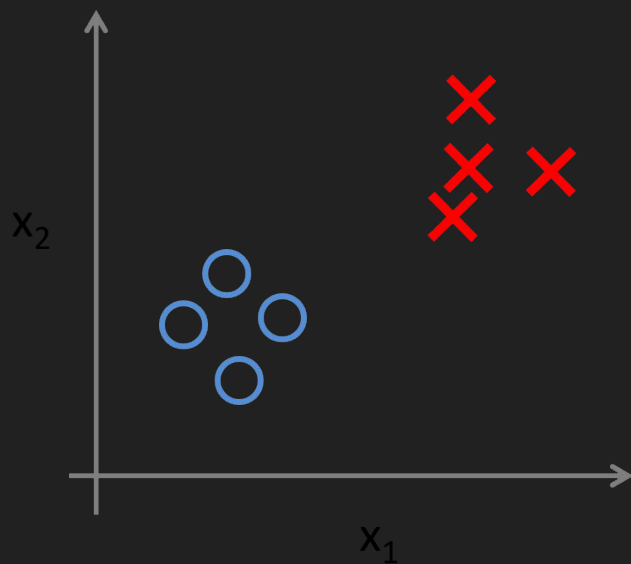
Probably the most important step in machine learning.

Labeled Data - Includes a meaningful “tag” to data. It’s typically what you’re trying to predict.

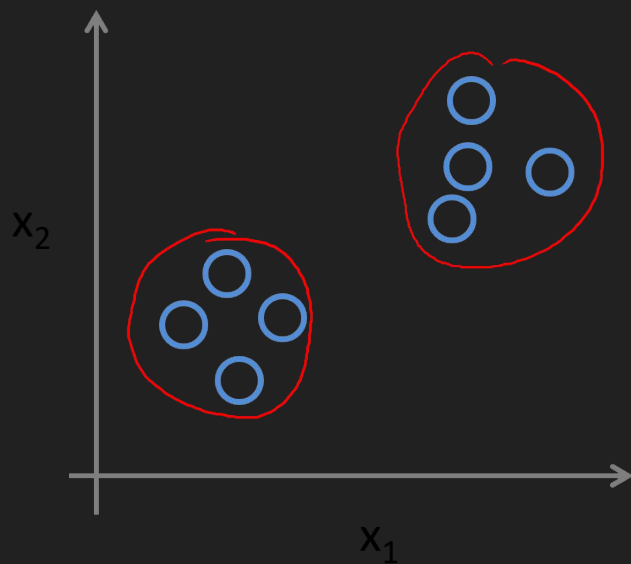
Unlabeled data - Just the data, no classification. (photos, videos)

Unsupervised vs. Supervised Learning

Supervised Learning



Unsupervised Learning

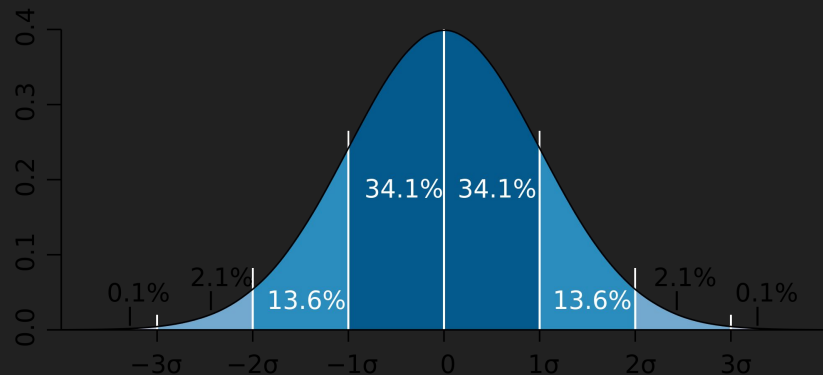


Feature Processing

- Data-Domain Research
- Filling in missing data
- Cartesian product
 - Population Density: urban, suburban, rural
 - State: Washington, Oregon, California
 - Product: urban_Washington, suburban_Washington, rural_Washington, urban Oregon, etc.
- Non-linear transformations
 - Binning
- Domain-specific features
 - Ex. $\text{length} * \text{width} * \text{height} = \text{volume}$
- Variable-specific features
 - Text features -

Pre-Processing

- Standardization
 - Some algorithms make assumptions about the data
 - SVM - assume that all features are centered around zero and have the same variance.
 - Standard normally distributed data (the bell curve)
- Scaling Features
- Normalization
- Binarization



Training vs Test Data

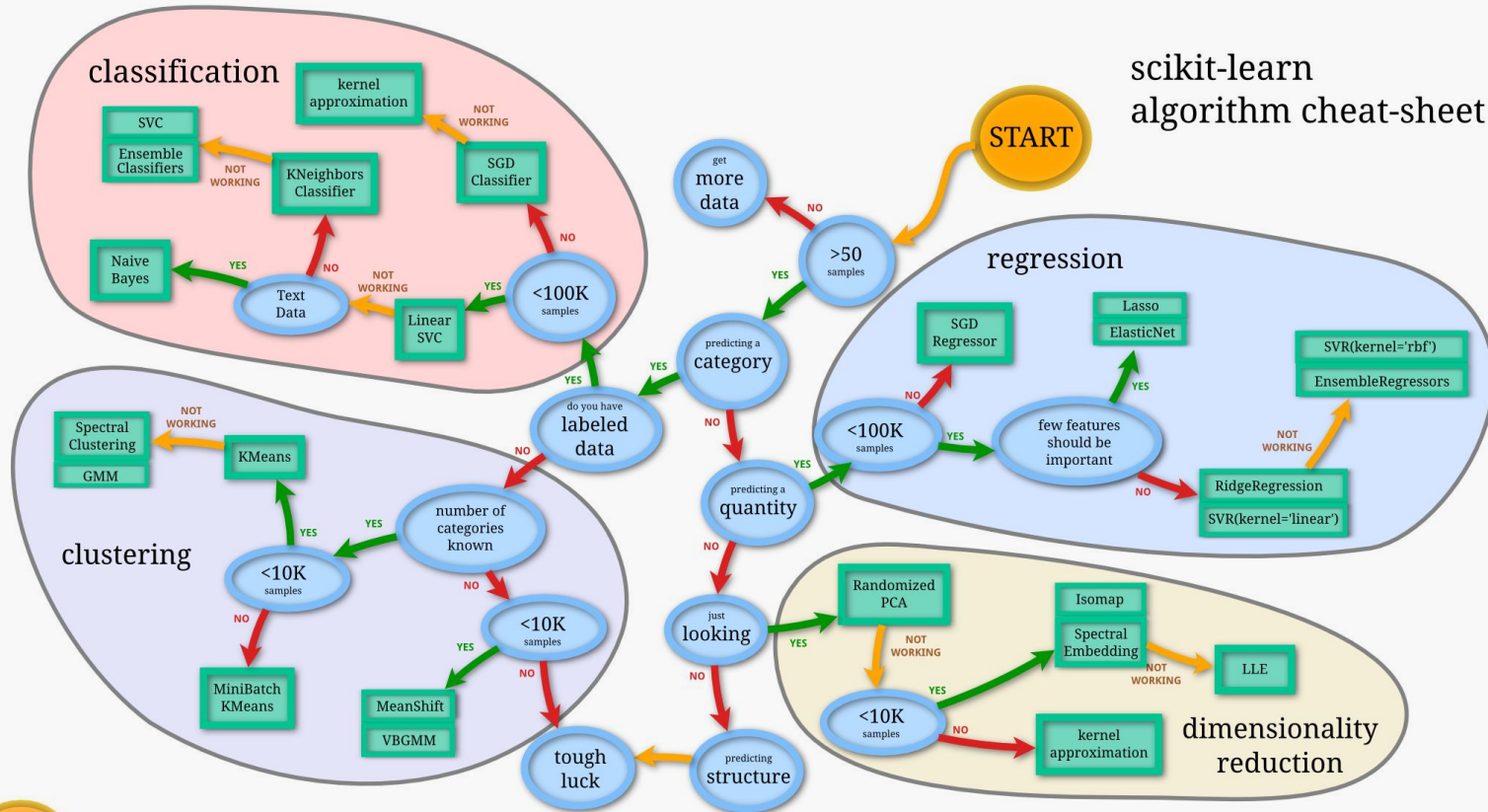
“The fundamental goal of ML is to generalize beyond the data instances used to train models.”

Training Data (70% - 80%) - Use this data to train the model.

Test Data (20% - 30%) - Use this data to evaluate the model.

- When splitting data, keep in mind that you want to accurately express your dataset in both your training and test dataset.
- Cross-Validation

scikit-learn algorithm cheat-sheet



Algorithm - Train model

- Several different algorithms to choose from to model your data.
 - Few different options:
 - Implement the algorithm yourself.
 - Use an existing implementation.

```
>>> from sklearn import svm
>>> from sklearn import datasets
>>> clf = svm.SVC()
>>> iris = datasets.load_iris()
>>> X, y = iris.data, iris.target
>>> clf.fit(X, y)
SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)
```

Algorithm - Evaluate model

```
>>> import pickle
>>> s = pickle.dumps(clf)
>>> clf2 = pickle.loads(s)
>>> clf2.predict(X[0:1])
array([0])
>>> y[0]
0
```


Results

The `sklearn.metrics` module includes score functions, performance `metrics` and pairwise `metrics` and distance computations.

Model Selection Interface

See the [The scoring parameter: defining model evaluation rules](#) section of the user guide for further details.

<code>metrics.get_scorer</code> (scoring)	Get a scorer from string
<code>metrics.make_scorer</code> (score_func[, ...])	Make a scorer from a performance metric or loss function.

Classification `metrics`

See the [Classification `metrics`](#) section of the user guide for further details.

<code>metrics.accuracy_score</code> (y_true, y_pred[, ...])	Accuracy classification score.
<code>metrics.auc</code> (x, y[, reorder])	Compute Area Under the Curve (AUC) using the trapezoidal rule
<code>metrics.average_precision_score</code> (y_true, y_score)	Compute average precision (AP) from prediction scores
<code>metrics.brier_score_loss</code> (y_true, y_prob[, ...])	Compute the Brier score.
<code>metrics.classification_report</code> (y_true, y_pred)	Build a text report showing the main classification <code>metrics</code>
<code>metrics.cohen_kappa_score</code> (y1, y2[, labels, ...])	Cohen's kappa: a statistic that measures inter-annotator agreement.
<code>metrics.confusion_matrix</code> (y_true, y_pred[, ...])	Compute confusion matrix to evaluate the accuracy of a classification
<code>metrics.f1_score</code> (y_true, y_pred[, labels, ...])	Compute the F1 score, also known as balanced F-score or F-measure
<code>metrics.fbeta_score</code> (y_true, y_pred, beta[, ...])	Compute the F-beta score
<code>metrics.hamming_loss</code> (y_true, y_pred[, ...])	Compute the average Hamming loss.
<code>metrics.hinge_loss</code> (y_true, pred_decision[, ...])	Average hinge loss (non-regularized)
<code>metrics.jaccard_similarity_score</code> (y_true, y_pred)	Jaccard similarity coefficient score
<code>metrics.log_loss</code> (y_true, y_pred[, eps, ...])	Log loss, aka logistic loss or cross-entropy loss.
<code>metrics.matthews_corrcoef</code> (y_true, y_pred[, ...])	Compute the Matthews correlation coefficient (MCC)
<code>metrics.precision_recall_curve</code> (y_true, ...)	Compute precision-recall pairs for different probability thresholds
<code>metrics.precision_recall_fscore_support</code> (...)	Compute precision, recall, F-measure and support for each class
<code>metrics.precision_score</code> (y_true, y_pred[, ...])	Compute the precision
<code>metrics.recall_score</code> (y_true, y_pred[, ...])	Compute the recall
<code>metrics.roc_auc_score</code> (y_true, y_score[, ...])	Compute Area Under the Receiver Operating Characteristic Curve (ROC AUC) from prediction scores.
<code>metrics.roc_curve</code> (y_true, y_score[, ...])	Compute Receiver operating characteristic (ROC)
<code>metrics.zero_one_loss</code> (y_true, y_pred[, ...])	Zero-one classification loss.

Demo

Tutorials

<http://scikit-learn.org/stable/tutorial/basic/tutorial.html>

https://www.tensorflow.org/get_started/get_started

<http://docs.aws.amazon.com/machine-learning/latest/dg/building-machine-learning.html>

<https://www.kaggle.com/c/titanic>

Want to learn more?

<https://openai.com/>

<https://www.tensorflow.org/>

<https://ipython.org/notebook.html>

<https://www.kaggle.com/>

[Jason Mayes - Machine Learning 101](#)

Brush up on your math skills -> Statistics / Linear Algebra / Calculus

Questions???????

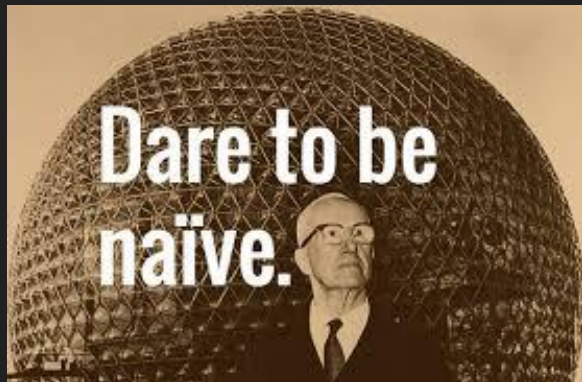
Extras

Machine Learning Tool Basics

IPython Notebook

Python / Pandas

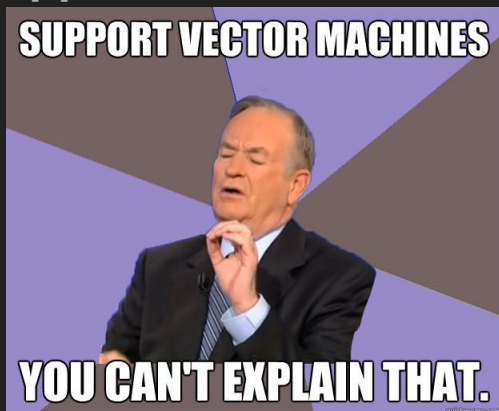
Naive Bayes



Random Forest



Support Vector Machine



K-Nearest Neighbors



The logo is set against a dark grey background. It features a blue square on the left containing a white circular wave-like shape. On the right is a green square containing a white silhouette of a fish. The text 'ALPHAZERO' is in white bold sans-serif font on the blue background, with 'vs' in a smaller font between two white horizontal bars. 'STOCKFISH' is in white bold sans-serif font on the green background.

ALPHAZERO vs STOCKFISH

The Chess.com logo, featuring a small green chess piece icon to the left of the text 'Chess.com' in a white sans-serif font.

 Chess.com

Real World Applications

The Netflix logo is displayed in a white rectangular box. It consists of the word "NETFLIX" in a bold, red, sans-serif typeface. The letters are slightly slanted to the right, giving it a dynamic feel.



Applications



AlphaGo

BILZZARD[®]
ENTERTAINMENT



DeepMind