

# Affective User Research & Human-AI Interaction

Seminar Summer 2024, Karlsruhe  
Institute of Technology  
Dr. Ivo Benke, BioNTech  
Dr. Lennard Schmidt, Google



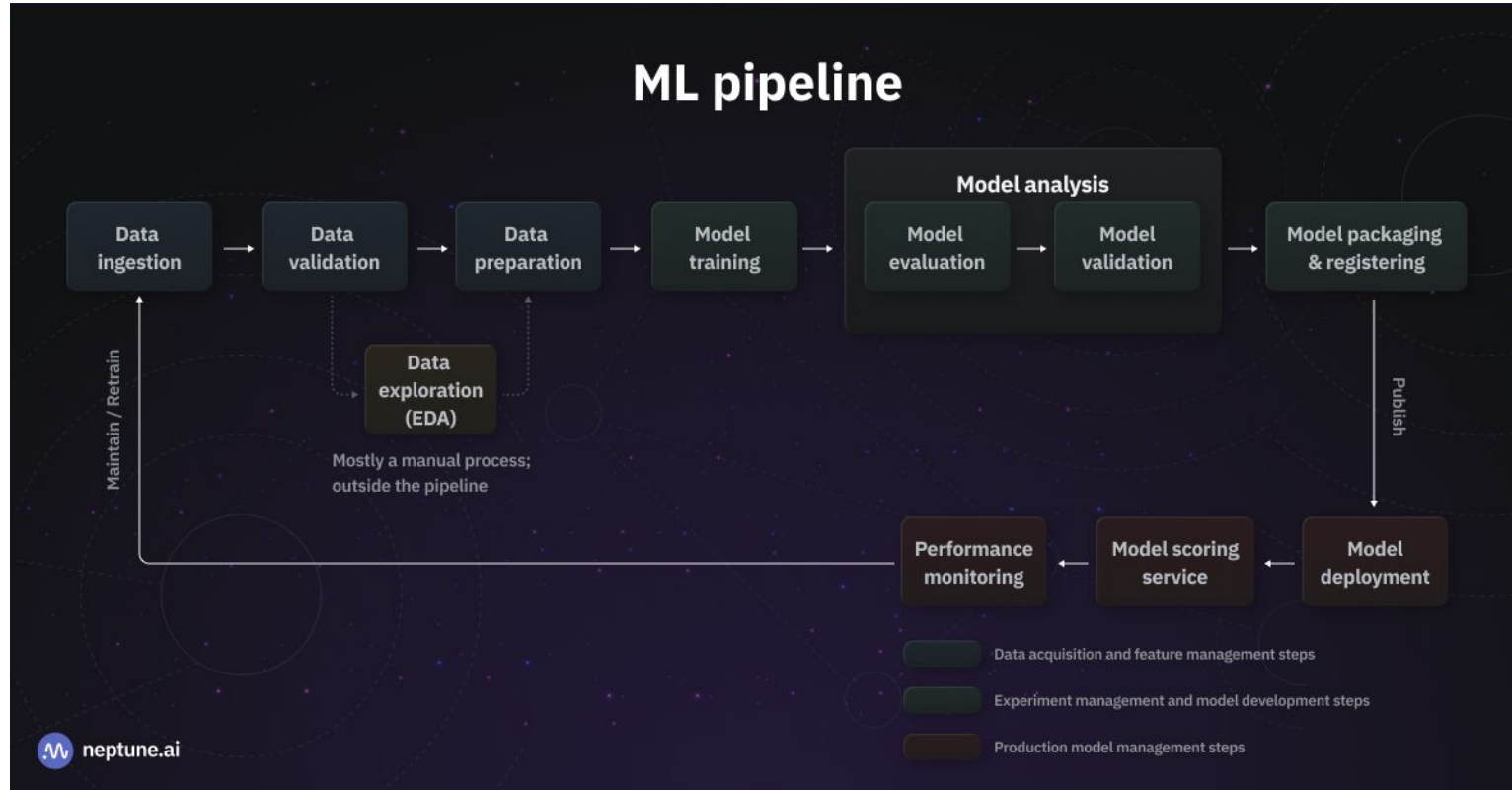
# Affective User Research & Human-AI Interaction

## Seminar #6 Part: Data Science Methods

Dr. Ivo Benke, Dr. Lennard Schmidt

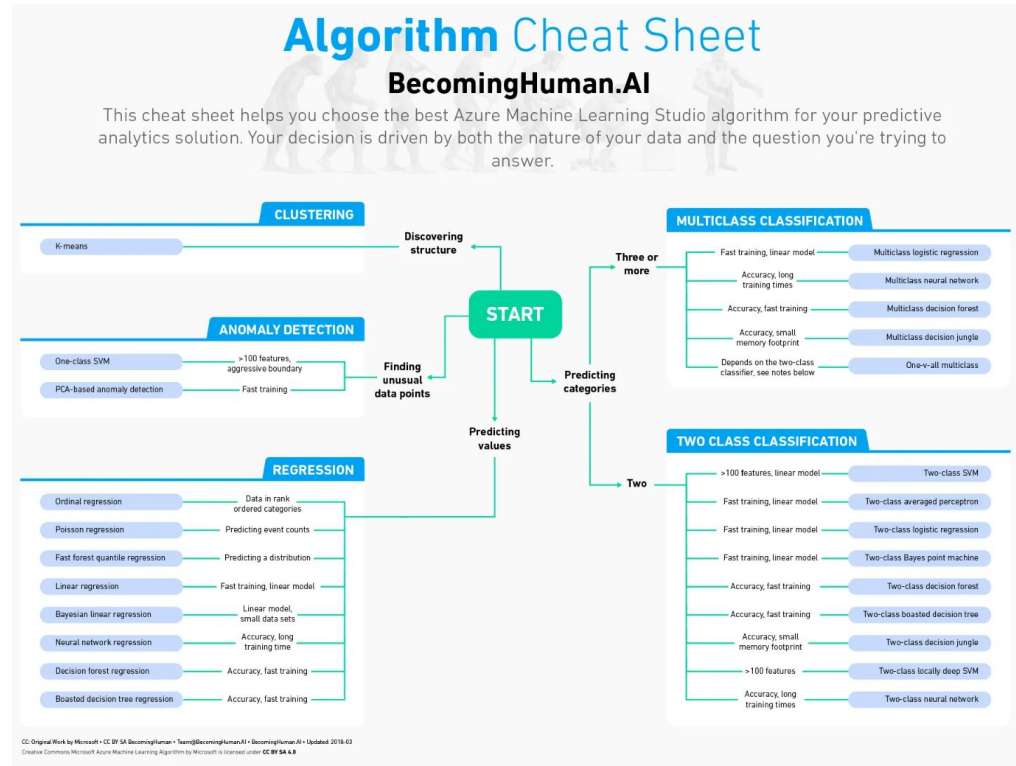


# ML Pipeline



# Machine Learning Helpful Packages

# Model Refresher





- A popular Python library for machine learning that provides simple and efficient tools for data mining and data analysis.
- Built on NumPy, SciPy, and matplotlib.
- Supports supervised and unsupervised learning with a rich documentation and active community.
- Main sections in [User Guide](#)
  - 1. Supervised learning
  - 2. Unsupervised learning
  - 3. Model selection and evaluation
  - 5. Visualizations
  - 6. Dataset transformations
  - 10. Common pitfalls and recommended practices

# Scikit Learn Cheat Sheet

## Cheat-Sheet Scikit learn Python For Data Science Becoming Human.AI



### Scikit Learn

Scikit Learn is an open source Python library that implements a range of machine learning, processing, cross validation and visualization algorithm using a unified

#### A basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.cross_validation import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
>>> scaler = preprocessing.StandardScaler(100, train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

### Prediction

#### Supervised Estimators

```
>>> y_pred = model.predict(X_test)
>>> y_pred = model.predict(X_test)
```

Predict Labels  
Predict Labels  
Estimate probability of labels

#### Unsupervised Estimators

```
>>> y_pred = k_means.predict(X_test)
```

Predict labels in clustering algorithms

### Loading the Data

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrix, other types that they are convertible to numeric arrays, such as Pandas Dataframe, are also acceptable

```
>>> import numpy as np
>>> X = np.random.random(100)
>>> y = np.array(['M', 'F', 'M', 'F', 'M', 'F', 'M', 'F', 'M', 'F'])
>>> X(X < 0.7) = 0
```

### Preprocessing The Data

#### Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler(100, train)
>>> standardized_X = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

#### Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> normalizer = Normalizer(100, train)
>>> normalized_X = normalizer.transform(X_train)
>>> normalized_X_test = normalizer.transform(X_test)
```

#### Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0)
>>> binary_X = binarizer.transform(X)
```

#### Encoding Categorical Features

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit(X_train, y_train)
```

#### Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit(X_train, y_train)
```

#### Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit(X_train, y_train)
```

### Evaluate Your Model's Performance

#### Classification Metrics

```
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)
>>> accuracy_score(y_test, y_pred)
```

Estimator score method  
Metric scoring functions

#### Classification Report

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

Precision, recall, F1 score and support

#### Confusion Matrix

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

#### Regression Metrics

##### Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
>>> y_train = X_train * 0.5 + 1
>>> mean_absolute_error(y_train, y_pred)
```

##### Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
>>> mean_squared_error(y_train, y_pred)
```

##### R<sup>2</sup> Score

```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_train, y_pred)
```

#### Clustering Metrics

##### Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_train, y_pred)
```

##### Homogeneity

```
>>> from sklearn.metrics import homogeneity_score
>>> homogeneity_score(y_train, y_pred)
```

##### V-measure

```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_train, y_pred)
```

#### Cross-Validation

```
>>> from sklearn.cross_validation import cross_val_score
>>> print(cross_val_score(X_train, y_train, cv=5))
>>> print(cross_val_score(X_test, y_test, cv=5))
```

### Model Fitting

#### Supervised learning

```
>>> from sklearn import svm
>>> svm.fit(X_train, y_train)
```

Fit the model to the data

#### Unsupervised Learning

```
>>> from sklearn import svm
>>> svm.fit(X_train, y_train)
```

Fit the model to the data  
Fit to data, then transform it

### Create Your Model

#### Supervised Learning Estimators

##### Linear Regression

```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

##### Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

##### Naive Bayes

```
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
```

##### KNN

```
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

#### Unsupervised Learning Estimators

##### Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

##### K Means

```
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=3, random_state=0)
```

### Training And Test Data

```
>>> from sklearn.cross_validation import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
```

### Tune Your Model

#### Grid Search

```
>>> from sklearn.grid_search import GridSearchCV
>>> params = {'n_neighbors': np.arange(1, 31)}
>>> grid = GridSearchCV(estimator=knn, param_grid=params)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

#### Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {'n_neighbors': range(1, 31),
>>>           'weights': ['uniform', 'distance']}
>>> search = RandomizedSearchCV(estimator=knn, param_distributions=params, n_iter=10, random_state=5)
>>> print(search.best_score_)
```

Content Copyright by DataCamp.com. Design Copyright by BecomingHuman.AI. See Original here.

# NLTK

## What is NLTK?

- A powerful Python library for working with human language data (text).
- Provides easy-to-use interfaces to over 50 corpora and lexical resources.
- Contains a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning.

## Features:

- **Corpora and Lexical Resources:** Access to large text corpora, such as Gutenberg, Brown, and many others.
- **Text Processing:** Tools for tokenizing text, finding n-grams, part-of-speech tagging, and named entity recognition.
- **Classification:** Implements various machine learning algorithms for text classification.
- **Text Analysis:** Tools for language modeling, text generation, and sentiment analysis.
- **Educational:** Comes with comprehensive documentation and tutorials for learning.

NLTK

Documentation

Search

Natural Language

NLTK Documentation  
API Reference  
Example Usage  
Module Index

NLTK is a leading platform  
provides easy-to-use interf  
with a suite of text process  
and semantic reasoning, w  
forum.



# NLTK Cheat Sheet

## Cheatography

richardmurray\_nlp\_python\_nltk\_cheatsheet Cheat Sheet  
by murenei via [cheatography.com/58736/cs/15485/](https://cheatography.com/58736/cs/15485/)

<b>Handling Text</b>		<b>Lemmatization &amp; Stemming</b>		<b>Sentence Parsing</b>	
<code>text="Some words"</code>	assign string	<code>input = "List listed lists listing listings"</code>	Different suffixes	<code>g=nltk.data.load('grammar.cfg')</code>	Load a grammar from a file
<code>list(text)</code>	Split text into character tokens	<code>words = input.lower().split('')</code>	Normalize (lowercase) words	<code>g=nltk.CFG.fromstring('...', '')</code>	Manually define grammar
<code>set(text)</code>	Unique tokens	<code>porter=nltk.PorterStemmer</code>	Initialise Stemmer	<code>parser=nltk.ChartParser(g)</code>	Create a parser out of the grammar
<code>len(text)</code>	Number of characters	<code>[porter.stem(t) for t in words]</code>	Create list of stems	<code>trees=parser.parse_all(text)</code>	
<b>Accessing corpora and lexical resources</b>		<code>WNL.WordNetLemmatizer()</code>	Initialise WordNet lemmatizer	for tree in trees: ... print tree	
<code>from nltk.corpus import brown</code>	import CorpusReader object	<code>WNL.lemmatize(t) for t in words]</code>	Use the lemmatizer	<b>Sentence Parsing</b>	
<code>brown.words(text_id)</code>	Returns pretokenised document as list of words	<b>POS Tagging</b>		<code>g=nltk.data.load('grammar.cfg')</code>	Load a grammar from a file
<code>brown.fileids()</code>	Lists docs in Brown corpus	<code>nltk.help.upenn_tagset('MD')</code>	Lookup definition for a POS tag	<code>g=nltk.CFG.fromstring('...', '')</code>	Manually define grammar
<code>brown.categories()</code>	Lists categories in Brown corpus	<code>nltk.pos_tag(words)</code>	nltk in-built POS tagger	<code>parser=nltk.ChartParser(g)</code>	Create a parser out of the grammar
<b>Tokenization</b>		<use an alternative tagger to illustrate ambiguity>		<code>trees=parser.parse_all(text)</code>	
<code>text.split(" ")</code>	Split by space			for tree in trees: ... print tree	
<code>nltk.word_tokenizer(text)</code>	nltk in-built word tokenizer			from nltk.corpus import treebank	
<code>nltk.sent_tokenizer(doc)</code>	nltk in-built sentence tokenizer			<code>treebank.parsed_sents('wsj_0001.mrg')</code>	Treebank parsed sentences

C By murenei  
[cheatography.com/murenei/](https://cheatography.com/murenei/)

Published 28th May, 2018.  
Last updated 28th May, 2018.  
Page 1 of 2.

Sponsored by [Readability-Score.com](https://readability-score.com)  
Measure your website readability!  
<https://readability-score.com>

# Tensorflow

- **TensorFlow** ([www.tensorflow.org/api](https://www.tensorflow.org/api)) **a set of tools and functions that to build and deploy machine learning models.** Allows for different levels of abstraction, from low-level operations on tensors to high-level APIs like Keras, which simplifies the process of defining and training neural networks.
- **TensorFlow Hub** ([www.tensorflow.org/hub](https://www.tensorflow.org/hub)) **is a repository of trained machine learning models** ready for fine-tuning and deployable anywhere. Reuse trained models like BERT and Faster R-CNN with just a few lines of code.



# Questions, Comments, Observations



# Thank You



**Dr. Ivo Benke**

ivolucasbenke@gmail.com



**Dr. Lennard Schmidt**

lennard.schmidt@hhl.de