# Easy Review Sentiment Analysis with `pandas` and `scikit-learn`

Arnout Devos

indy.epfl.ch

## Problem

Predict whether a new review has a rather positive (+) or negative (-) sentiment, given the experience from known labeled (+/-) reviews.

## Dataset

Around 400,000 Amazon product reviews
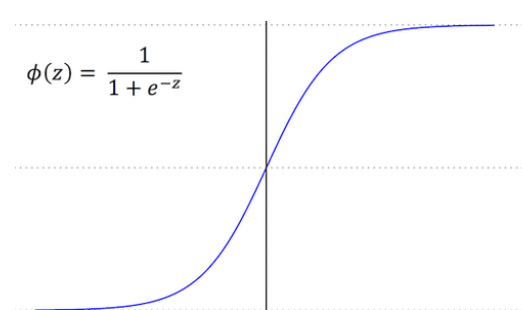Sentiment data: 50% positive, 50% negative

## Preprocessing

Python: `pandas` and `scikit-learn`
80% Train, 20% Test (every $5^{th}$ sample $\in$ test)
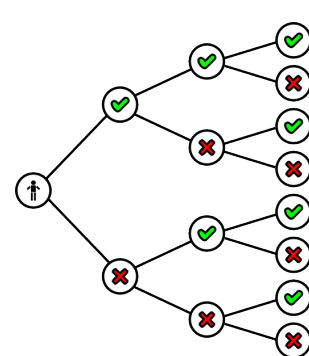Convert *textual reviews* into *numerical vectors* with `CountVectorizer()`

## Algorithms

### 1. Logistic Regression



$\phi(z) = \frac{1}{1+e^{-z}}$

- Natural choice for binary classification problems with vectors
- Simple interpretable model (coefficients)
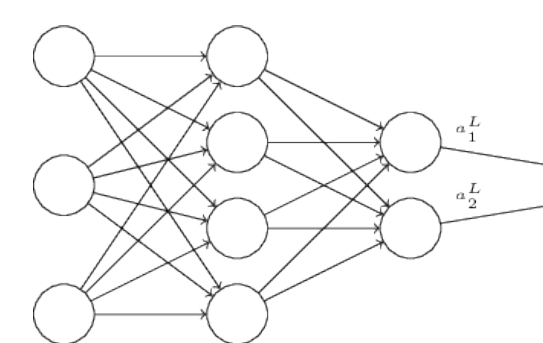
`LogisticRegression()`

### 2. Decision Tree



- Natural choice for decision problems in general (low-dimensional, human-size)
- Simple interpretable model (tree)
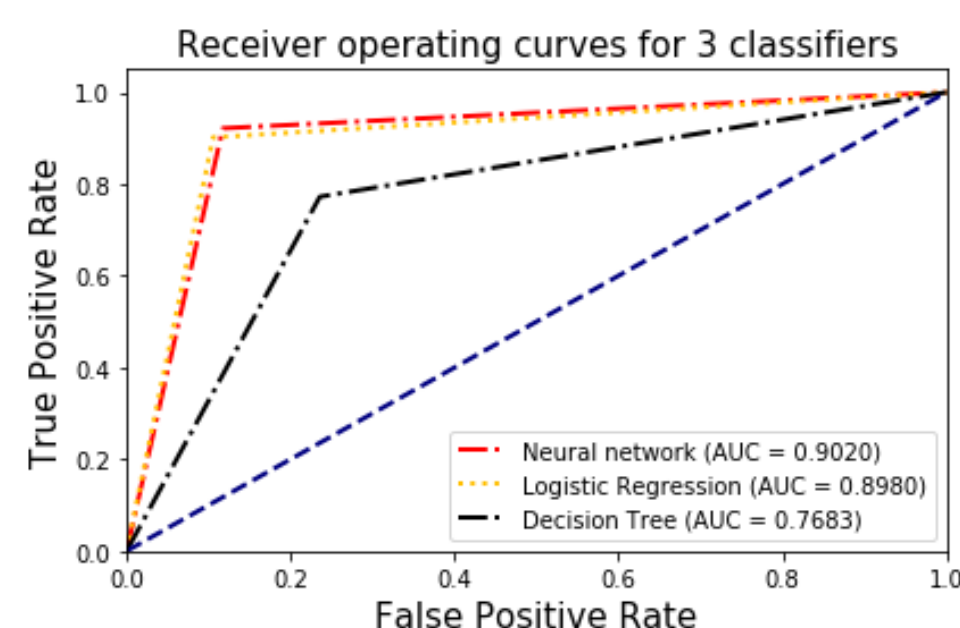
`DecisionTreeClassifier()`

### 3. Neural Net



- Popular choice with Word2Vec features
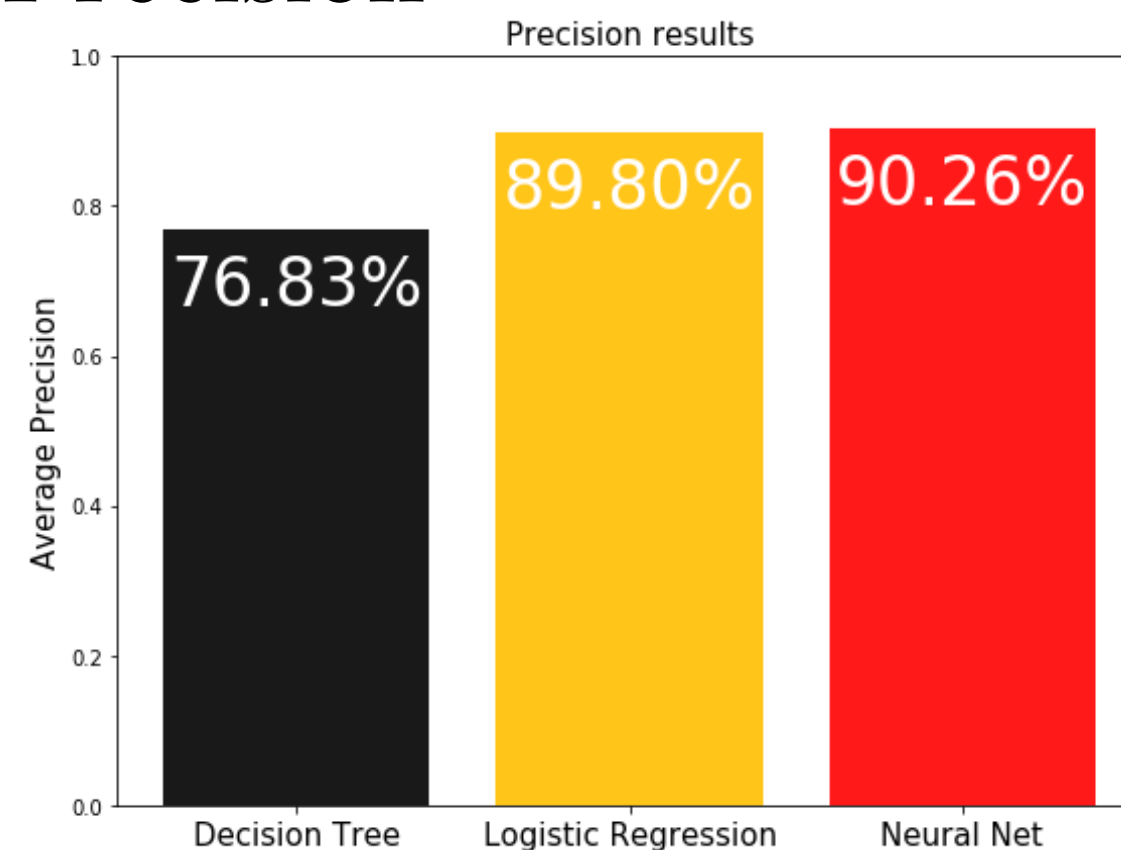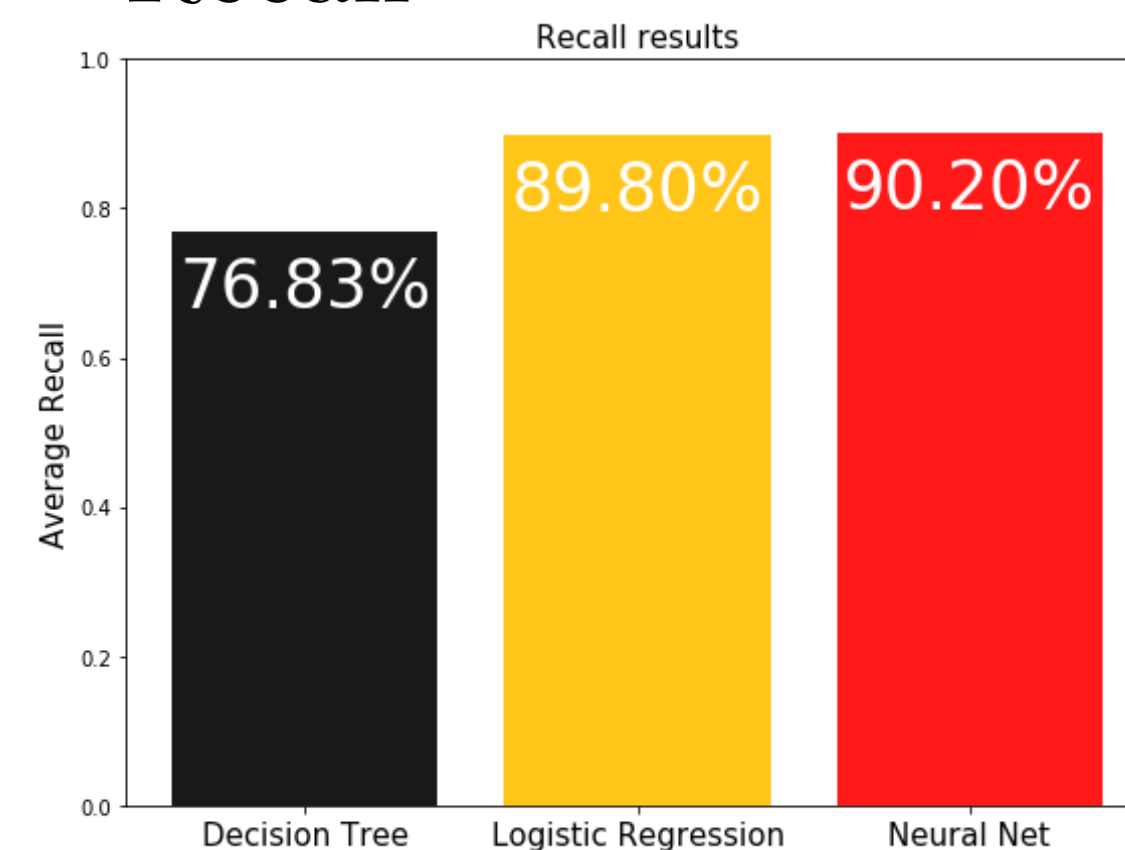- Can fit highly nonlinear functions/ distributions

`MLPClassifier()`

## Results

### ROC



Receiver operating curves for 3 classifiers
Neural network (AUC = 0.9020)
Logistic Regression (AUC = 0.8980)
Decision Tree (AUC = 0.7683)

### Accuracy



Accuracy results
76.82%  89.80%  90.20%
Decision Tree  Logistic Regression  Neural Net

### Precision



Precision results
76.83%  89.80%  90.26%
Decision Tree  Logistic Regression  Neural Net

### Recall



Recall results
76.83%  89.80%  90.20%
Decision Tree  Logistic Regression  Neural Net
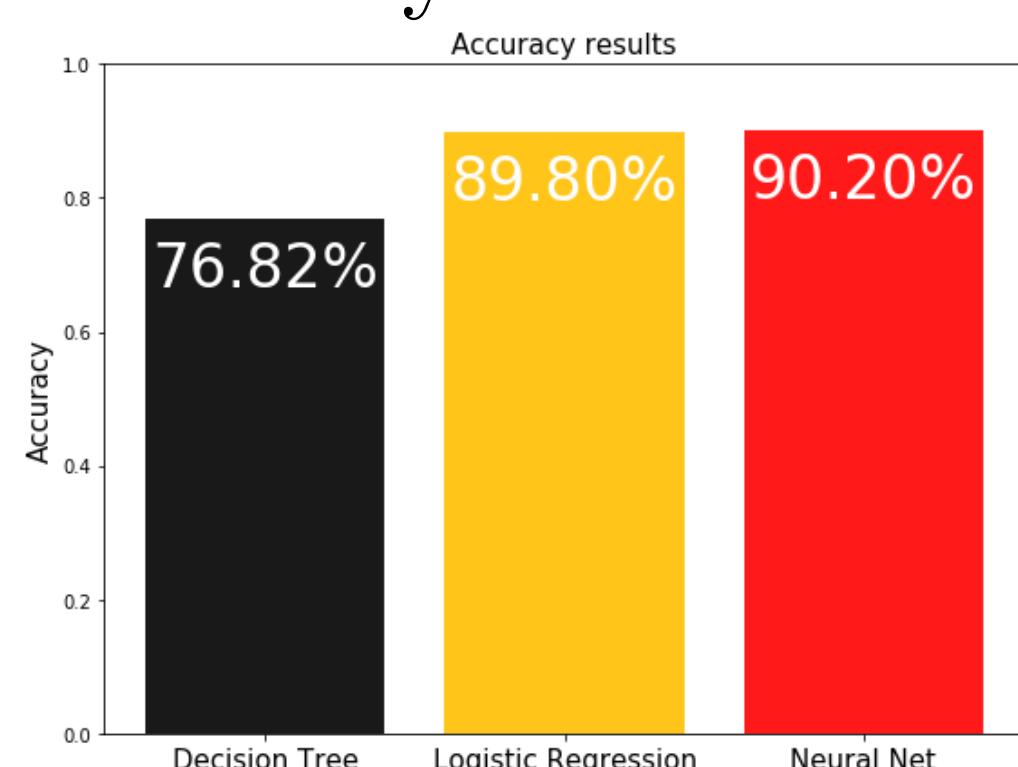
- *Neural Net* shows the best overall performance
- *Decision Tree* has trouble with high dimensionality and performs the worst
- *Logistic Regression* benefits from vectorization, but performs slightly worse than the Neural Net