# Hunting for Pulsars

Joe Davies

# Useful Links

- **Particle physics application:**
**https://ilmonteux.github.io/2018/10/15/jet-tagging-cnn.html**

- **Kaggle page for the pulsar dataset:**
**https://www.kaggle.com/shivam1901/pulsar-star**

- **Github repo with all the code from these tutorials:**
**https://github.com/adrianbevan/IntroToML/tree/master/Pulsars**

- **https://learn.datacamp.com/courses/introduction-to-deep-learning-with-keras**

- **Misc useful links for more information on machine learning:**

  - **https://towardsdatascience.com/**

  - **https://www.coursera.org/learn/machine-learning**

  - **https://elitedatascience.com/learn-machine-learning**
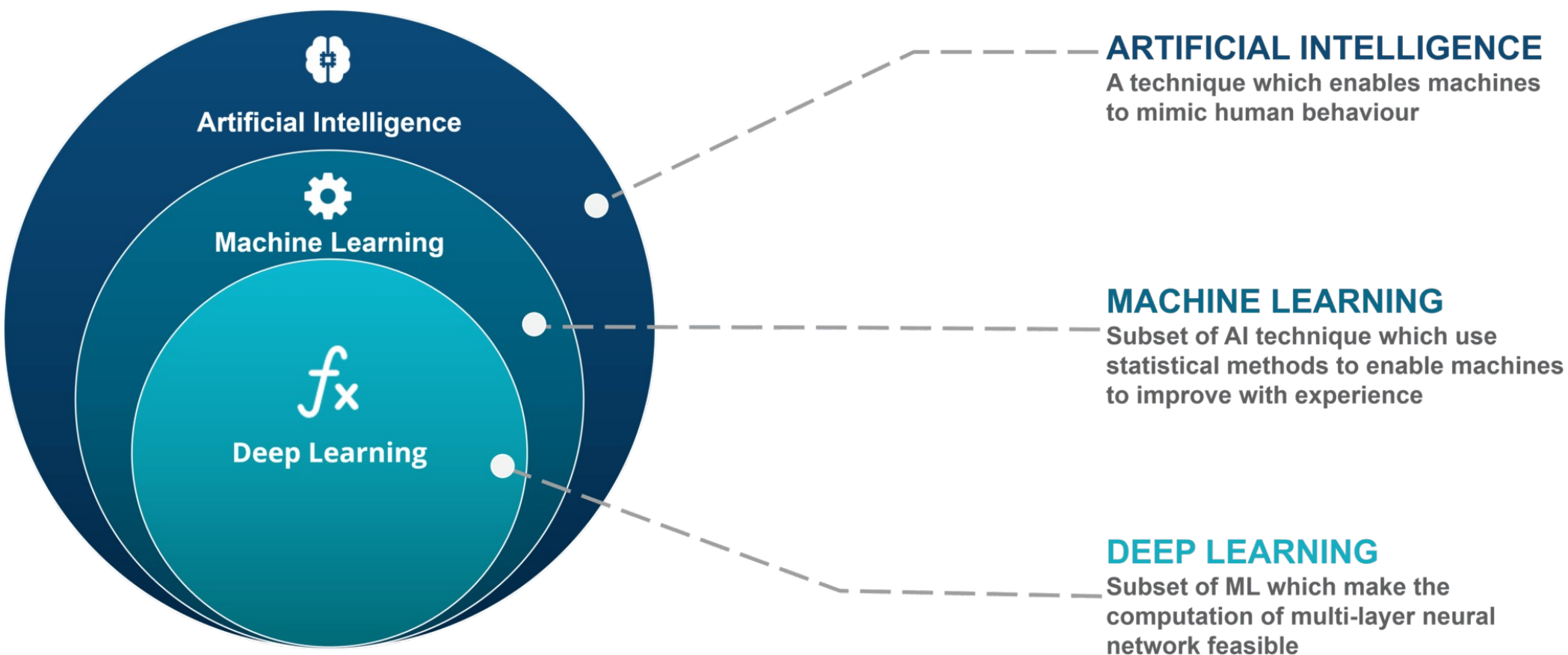
- **j.m.m.davies@qmul.ac.uk**

# Introductions!



## The Dark Machines Anomaly Score Challenge: Benchmark Data and Model Independent Event Classification for the Large Hadron Collider

T. Aarrestad[CERN]   M. van Beekveld[Ox]   M. Bona[QMUL]   A. Boveia[OSU]
S. Caron[HEF, Nikhef]   J. Davies[QMUL]   A. De Simone[SISSA, INFN]   C. Doglioni[Lund]
J.M. Duarte[UCSD]   A. Farbin[UnivArlington]   H. Gupta[GSoC]   L. Hendriks[HEF, Nikhef]
L. Heinrich[CERN]   J. Howarth[Glasgow]   P. Jawahar[WPI, CERN]   A. Jueid[UnivKonkuk]
J. Lastow[Lund]   A. Leinweber[UnivAdelaide]   J. Mamuzic[IFIC]   E. Merényi[UnivRice]
A. Morandini[RWTH]   P. Moskvitina[HEF, Nikhef]   C. Nellist[HEF, Nikhef]
J. Ngadiuba[FNAL, Caltech]   B. Ostdiek[Harvard, AIFI]   M. Pierini[CERN]   B. Ravina[Glasgow]
R. Ruiz de Austri[IFIC]   S. Sekmen[KNU]   M. Touranakou[NKUA, CERN]
M. Vaškevičiūte[Glasgow]   R. Vilalta[UnivHouston]   J.-R. Vlimant[Caltech]   R. Verheyen[UCL]
M. White[UnivAdelaide]   E. Wulff[Lund]   E. Wallin[Lund]   K.A. Wozniak[UniVie, CERN]
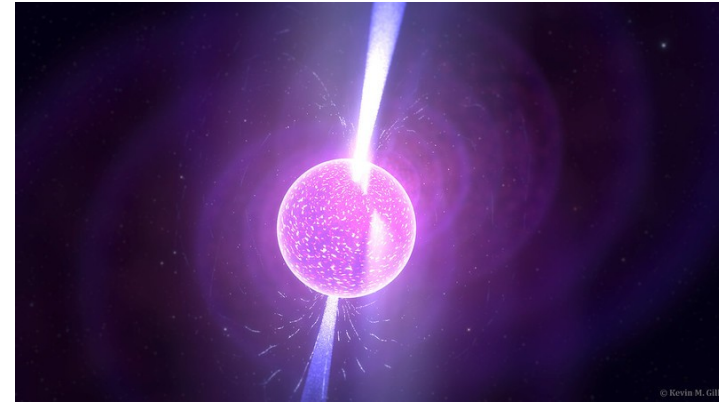Z. Zhang[HEF, Nikhef]

# Machine Learning vs Artificial Intelligence

**Artificial Intelligence**

**Machine Learning**

**Deep Learning**

**ARTIFICIAL INTELLIGENCE**
A technique which enables machines to mimic human behaviour

**MACHINE LEARNING**
Subset of AI technique which use statistical methods to enable machines to improve with experience

**DEEP LEARNING**
Subset of ML which make the computation of multi-layer neural network feasible

# The Data

- **Data is based on a Kaggle set investigating pulsars**

- **Pulsars are highly magnetized neutron stars that emit radiation from their magnetic poles**

- **Data contains 8 columns including metrics like kurtosis and dispersion of radiation measures**

- **The data also contains a target class: 0 or 1 depending on not-a-pulsar or pulsar**

# What are we doing?

- Create 3 machine learning algorithms that identify pulsars from the data

- Understand how to use specific python modules to do this

- Get an idea of when to use and not use each algorithm

- sklearn, scipy, pandas, matplotlib, numpy, keras, tensorflow

# Logistic Regression

- **Statistical model that uses a logistic function to (usually) model a binary variable (0 or 1)**
- **An example of this function is the sigmoid:**

$$f(x) = \frac{1}{1 + e^{-x}}$$

- **Uses gradient descent for the meat of the regression**
- **An issue with this approach is that extreme cases are presumed to become progressively rarer at a specific rate**
  - Not very good if one doesn't have much data



Joe Davies, QMUL                                        23/11/21

# Logistic Regression

Example file: pulsar_logistic_sklearn.py

**Data Preperation**

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from scikitplot.estimators import plot_learning_curve
df=pd.read_csv('~/Documents/ML/pulsars/pulsar_stars.csv')
```

```python
x_data=df.drop(columns='target_class')
x=(x_data-np.min(x_data))/(np.max(x_data)-np.min(x_data)) #scaling
y=df.target_class.values

compare_score=[]

#training and testing split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=123)
```

**Model Creation**

```python
lr=LogisticRegression()
lr.fit(x_train, y_train)
```

**Visualization**

```python
lr_score=lr.score(x_test, y_test) * 100
compare_score.append(lr_score)

print('Test accuracy: {}%'.format(lr_score))

plot_learning_curve(lr, x_test, y_test)
plt.show()
```

# Decision Trees





Decision Tree Classification (Test set)

| Advantages | Disadvantages |
|---|---|
| Simple to understand and interpret | Not very robust |
| Useful for classification and regression | Make locally optimal decisions (info gain max, entropy min) |
| Requires little data preparation and can handle large data sets | Prone to over-fitting |

# Boosted Decision Tree



- Boosting involves running multiple trees one after the other and minimizing error
- Each tree output, h(x), is given a weight, w, based on accuracy: $\hat{y}(x) = \sum_t w_t h_t(x)$
- Each data sample is given a weight based on misclassification
- We minimize the objective function: $O(x) = \sum_i l(\hat{y}_i, y_i) + \sum_t \Omega(f_t)$
- First term is loss function, second is regularization (penalizes complexity of the t-th tree)

https://towardsdatascience.com/understanding-adaboost-2f94f22d5bfe
https://indico.fnal.gov/event/15356/contributions/31377/attachments/19671/24560/DecisionTrees.pdf
https://scikit-learn.org/stable/auto_examples/ensemble/plot_adaboost_regression.html

# Boosted Decision Tree

Example file: pulsar_bdt.py

**Data Preperation**

```python
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.model_selection import train_test_split
from sklearn import metrics
#Vizualization
import matplotlib.pyplot as plt
from scikitplot.estimators import plot_learning_curve

data=pd.read_csv('~/Documents/ML/pulsars/pulsar_stars.csv')
```

```python
features=data.columns[:-1]
X=data[features]
#output:
y=data.target_class

#Now we need to split the data using train_test_split
X_train, X_test, y_train, y_test=train_test_split(X, y,
                              test_size=0.2, random_state=42)
```

**Model Creation**

```python
classifier=AdaBoostClassifier(DecisionTreeClassifier())

classifier=classifier.fit(X_train, y_train)
#Preducting the response for the test dataset
y_pred=classifier.predict(X_test)
```

**Visualization**

```python
score=round(metrics.accuracy_score(y_test, y_pred)*100, 2)
print('Accuracy = {}%'.format(score))

plot_learning_curve(classifier, X_test, y_test)
plt.show()
```

# Artificial Neural Networks

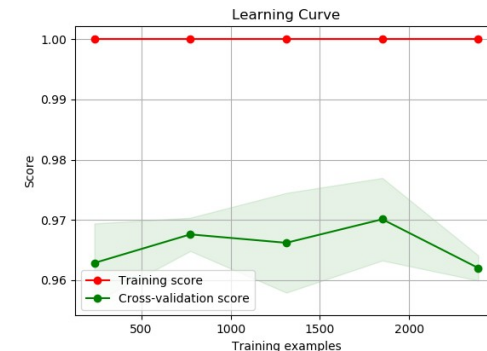$$y = f\Big(\sum_{i=1}^{N} w_i x_i + \theta\Big)$$
$$= f(w^T x + \theta)$$



input layer     hidden layer     output layer

output

$$w_{r+1} = w_r - \gamma \sum_{i=1}^{N} \frac{\partial E_i}{\partial w}$$

$$\theta_{r+1} = \theta_r - \gamma \sum_{i=1}^{N} \frac{\partial E_i}{\partial \theta}$$

| Advantages | Disadvantages |
|---|---|
| Among the most accurate of modelling approaches | Computationally intensive |
| Useful for classification and regression | Easy to over- or under-training data: |
| Makes few assumptions about relationships in the data | Results in a complex black box model |

# Artificial Neural Network

Example file: pulsar_ann.py

**Data Preperation**

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from keras import Sequential
from keras.layers import Dense
from sklearn.metrics import confusion_matrix
from scikitplot.estimators import plot_learning_curve
df=pd.read_csv('~/Documents/ML/pulsars/pulsar_stars.csv')
```

```python
x_data=df.drop(columns='target_class')
X=StandardScaler().fit_transform(x_data)
y=df.target_class.values

x_train, x_test, y_train, y_test=train_test_split(X, y,
                          test_size=0.2, random_state=42)
```

**Model Creation**

```python
classifier=Sequential()

#first hidden layer
#we have 8 input features, 1 output and the kernel_initializer uses a normal distribution to
#function
classifier.add(Dense(8, activation='relu', kernel_initializer='random_normal', input_dim=8))

#second
classifier.add(Dense(8, activation='relu', kernel_initializer='random_normal'))

#output
classifier.add(Dense(1, activation='sigmoid', kernel_initializer='random_normal'))

#compiling the network
classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

#fitting the data to the training set
history=classifier.fit(x_train, y_train, validation_split=0.33, batch_size=10, epochs=15)

#evaluate the loss value and metrics values for the model in test mode using evaluate funcn.
eval_model=classifier.evaluate(x_train, y_train)
print('eval model: ', eval_model)
```

- 1 input layer, 1 hidden layer, 1 output layer
- 8 features, 1 output
- Minimum 8 neurons in the hidden layer
- Validation set to get an idea of accuracy per epoch
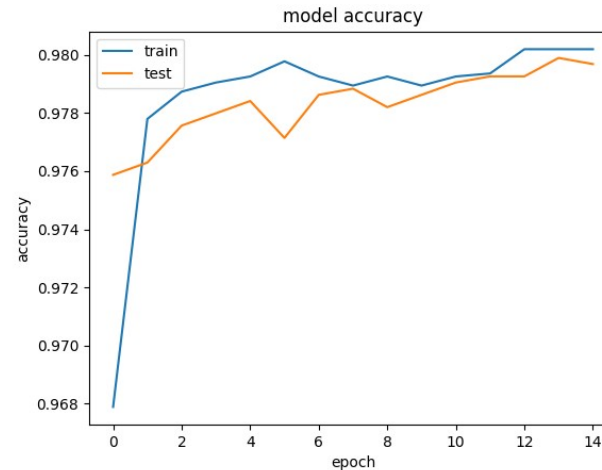
# Artificial Neural Network
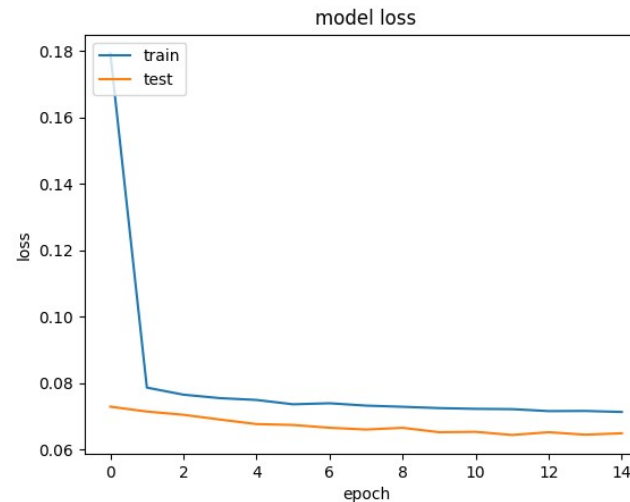
Example file: pulsar_ann.py

**Visualization**

```
cm=confusion_matrix(y_test, y_pred)
print(cm)

# list all data in history
print(history.history.keys())
# summarize history for accuracy
plt.plot(history.history['accuracy'])
plt.plot(history.history['val_accuracy'])
plt.title('model accuracy')
plt.ylabel('accuracy')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()
# summarize history for loss
plt.plot(history.history['loss'])
plt.plot(history.history['val_loss'])
plt.title('model loss')
plt.ylabel('loss')
plt.xlabel('epoch')
plt.legend(['train', 'test'], loc='upper left')
plt.show()
```
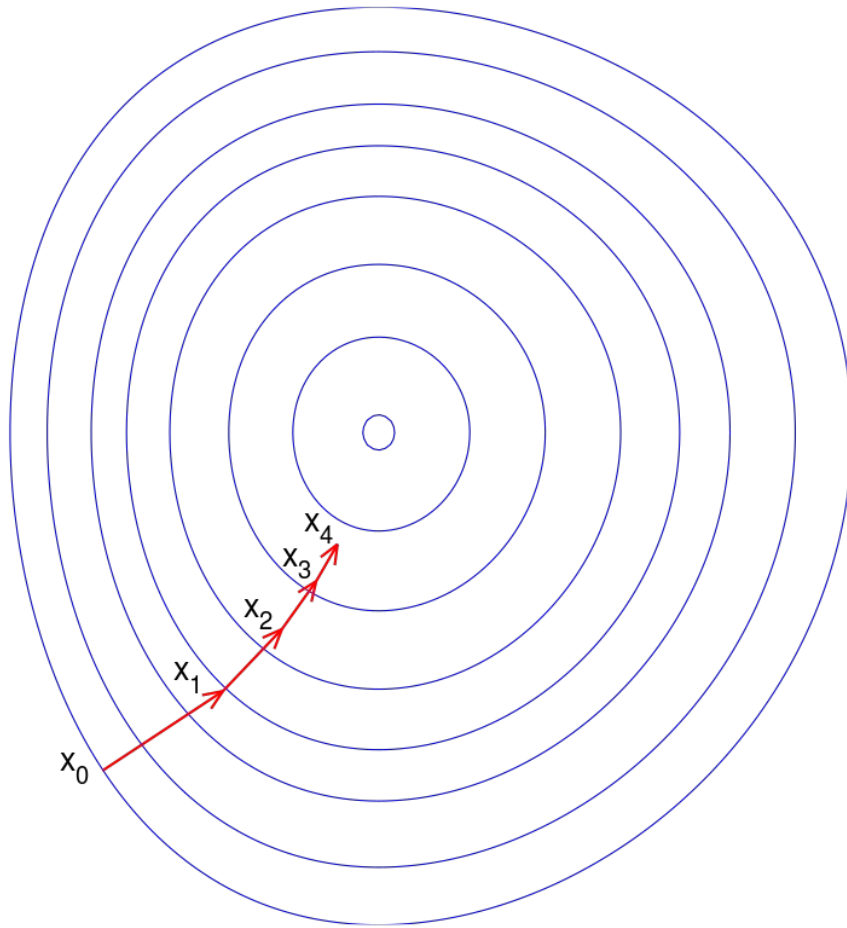
```
[[3233    26]
 [  42   279]]
```

# Something cool to end: GANs

# Backup: Gradient Descent

- F(x) is differentiable around point theta
- F will decrease fastest if going in the direction of negative gradient
- Gamma is some real, positive number
- Converges to a local minimum

$$\theta_{i+1} = \theta_i - \gamma \nabla F(\theta_i)$$