

ARISTOTLE UNIVERSITY OF THESSALONIKI
FACULTY OF SCIENCES
SCHOOL OF INFORMATICS



POSTGRADUATE STUDIES PROGRAM ON INFORMATICS AND COMMUNICATIONS
SPECIALIZATION ON DIGITAL MEDIA AND COMPUTATIONAL INTELLIGENCE

Content-based Image Retrieval

Presenter: Efstathios Chatzikyriakidis

M.Sc. in Informatics and Communications
Specialization in Computational Intelligence and Digital Media
School of Informatics, Faculty of Sciences
Aristotle University of Thessaloniki, Hellas

B.Sc. in Informatics and Communications
Specialization in Software Engineering
Department of Informatics and Communications, Faculty of Applied Technology
Technological Educational Institute of Central Macedonia, Hellas

Email : contact@efxa.org

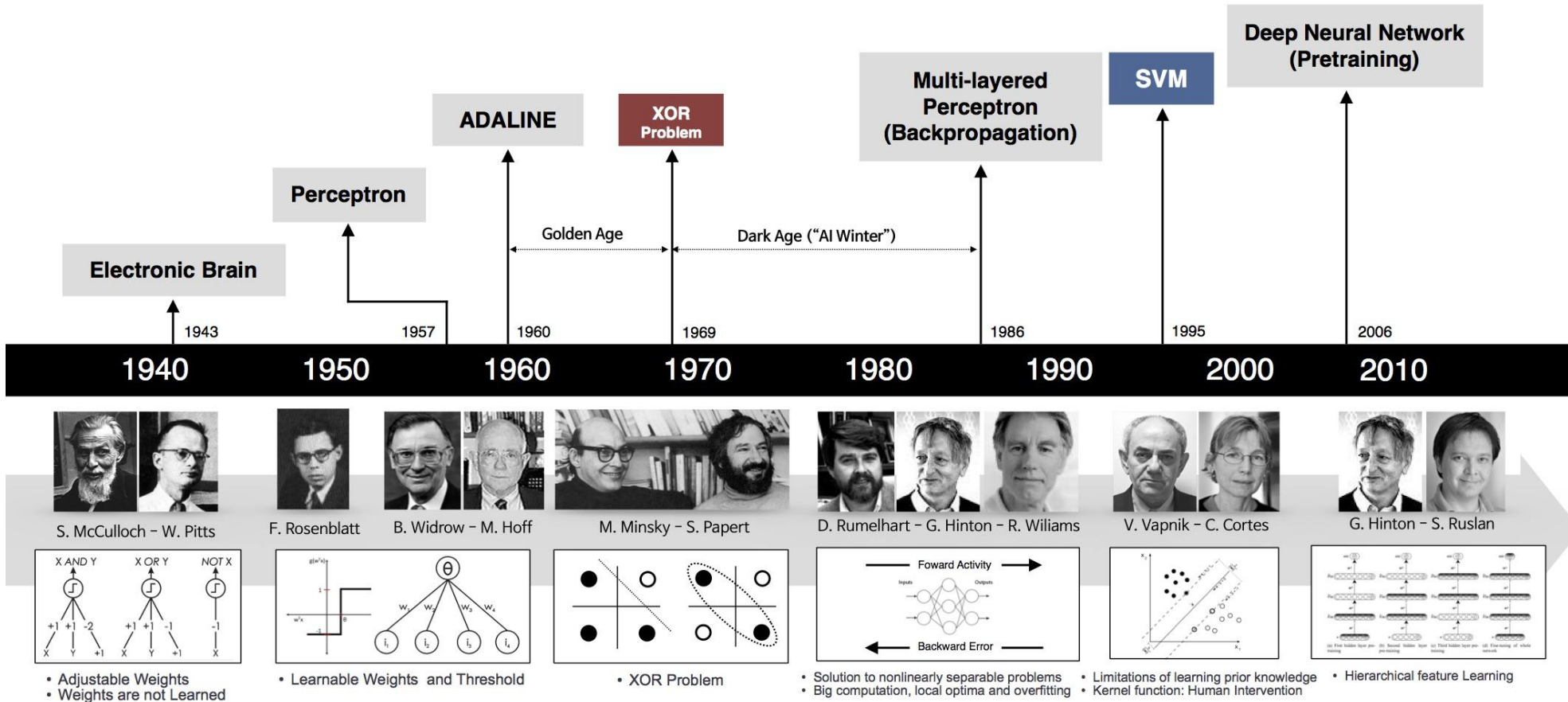
Website : <http://www.efxa.org/>

Introduction

Important history of “Artificial Intelligence”



We had a long journey... and we are still at the birth of it...

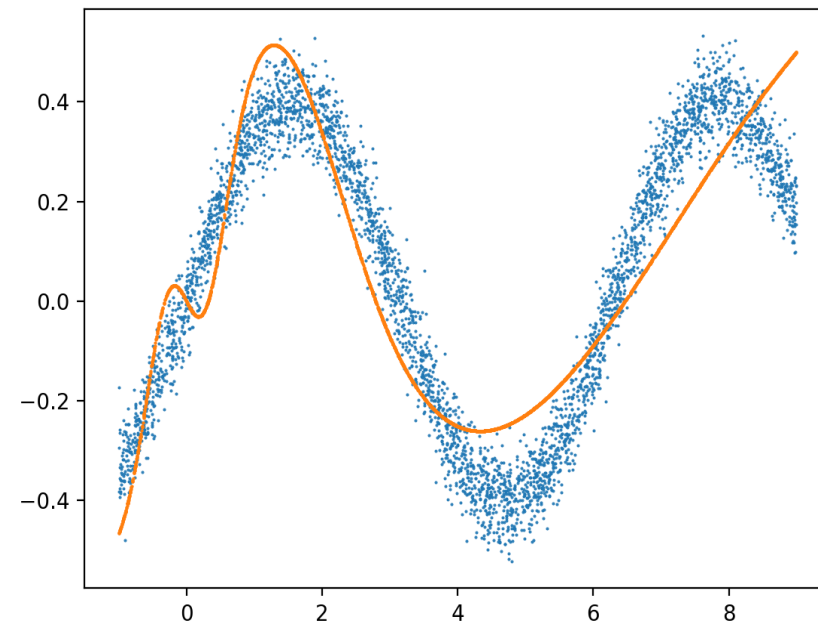
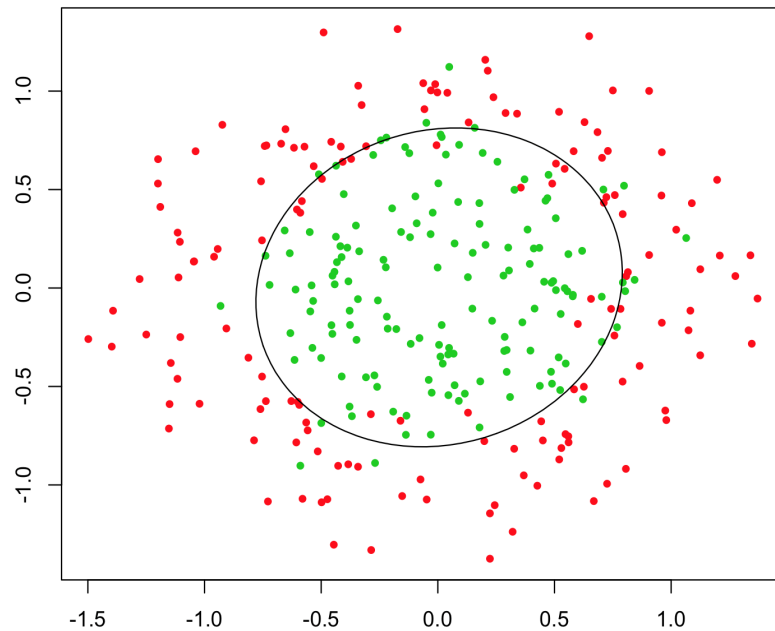


Introduction

Advantages of Artificial Neural Networks



- Satisfactory separation of non-linear separable input data
- Satisfactory function approximation using only input data

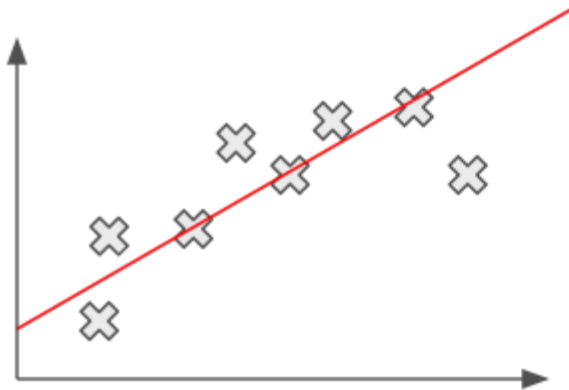


Introduction

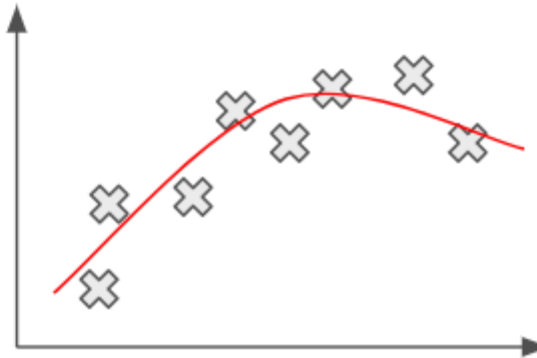
Advantages of Artificial Neural Networks



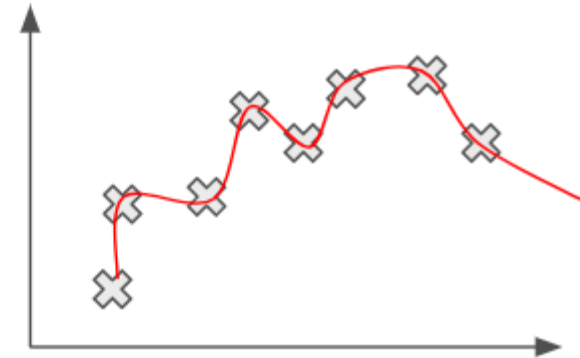
- Good generalization that captures the general manifold of data



Underfitting



Optimal



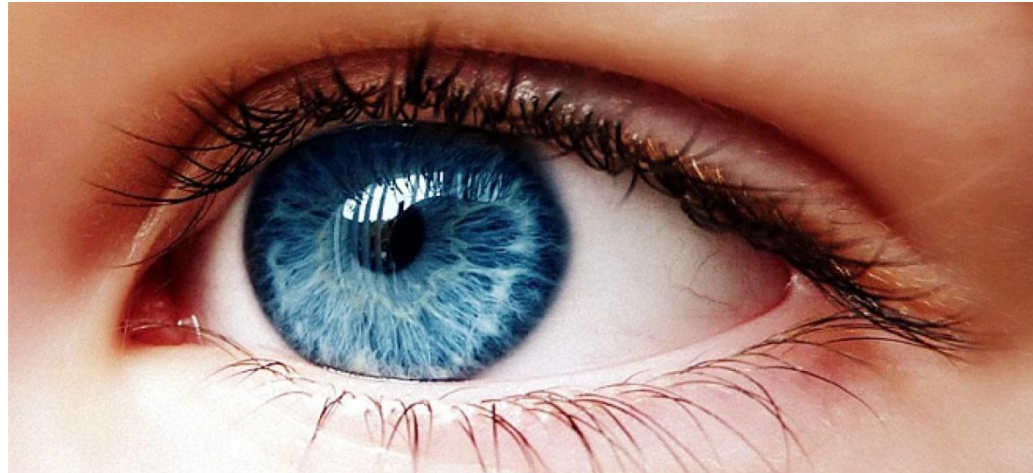
Overfitting

Introduction

Computer Vision and its future goal

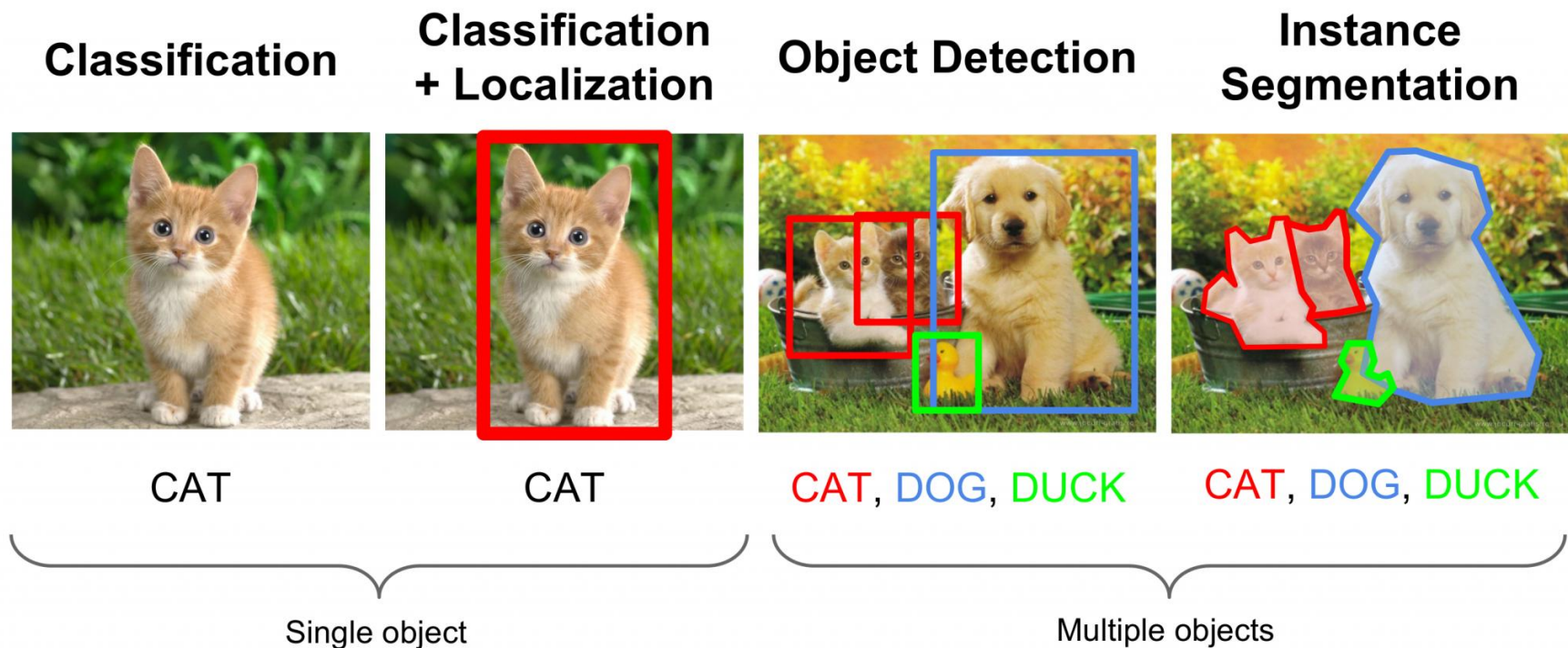


- The strongest sense of most animal species
- Simulate how brains see and understand the world through vision sense



Computer Vision

Well-known tasks of high-level image understanding

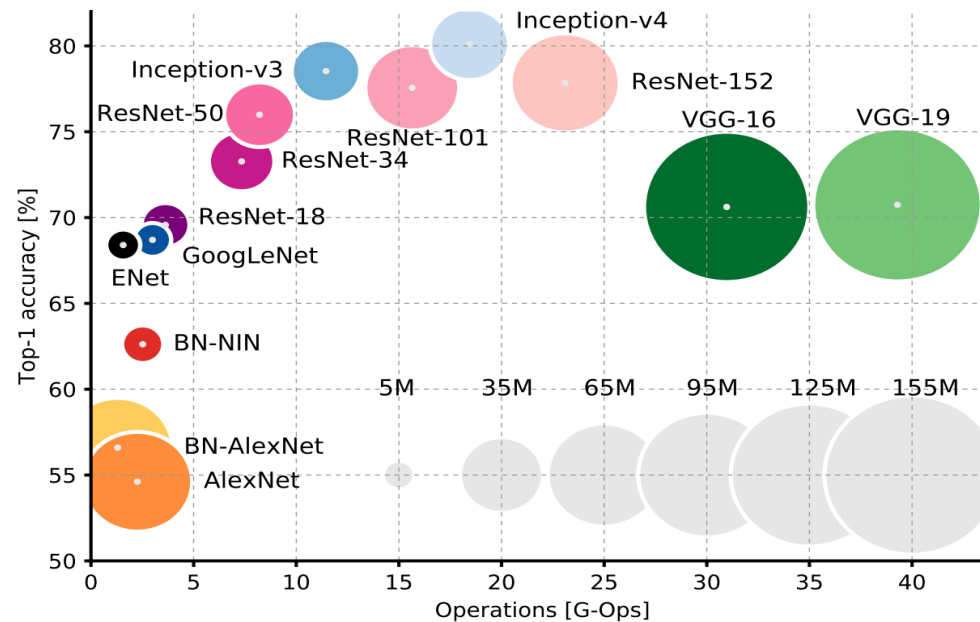
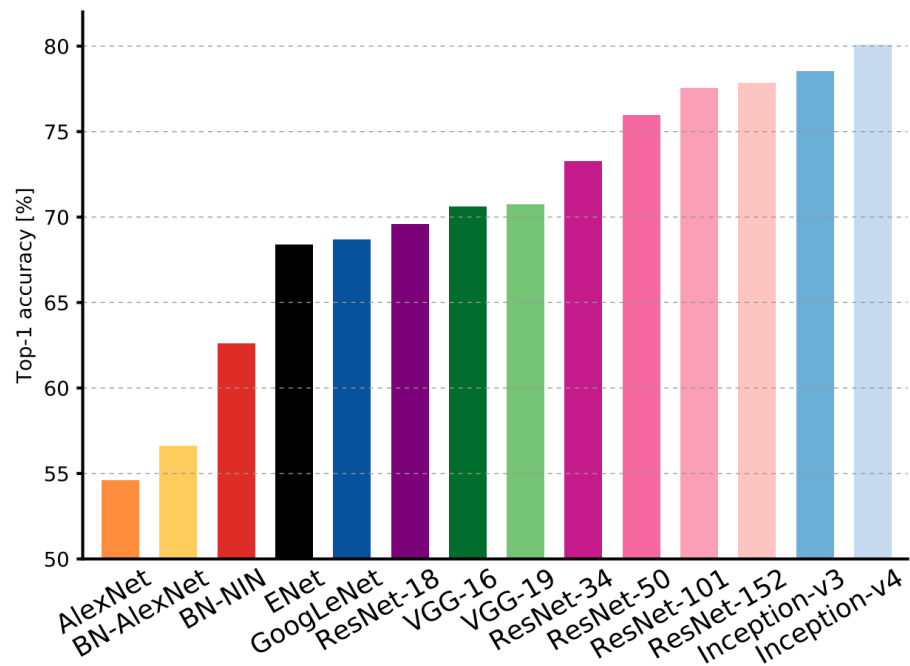


Computer Vision

State-of-the-art Image Classification



ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Dataset

Fashion MNIST



Dataset

Fashion MNIST



Characteristics:

- 70000 total grayscale images
- 28x28 pixels image resolution
- 10 total class labels
- 60000 training images (6000 images per class label)
- 10000 testing images (1000 images per class label)

Dataset

Fashion MNIST

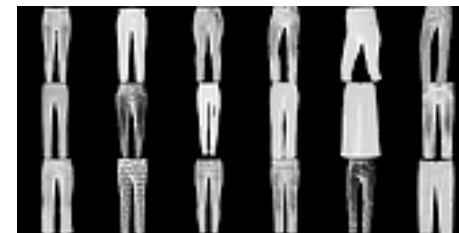
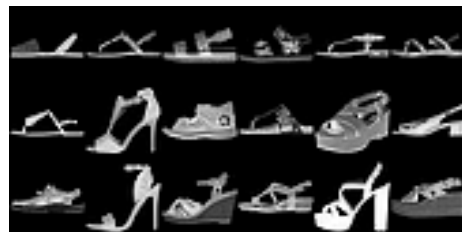


Class labels:

- T-shirt/top
- Trouser
- Pullover
- Dress
- Coat
- Sandal
- Shirt
- Sneaker
- Bag
- Ankle boot

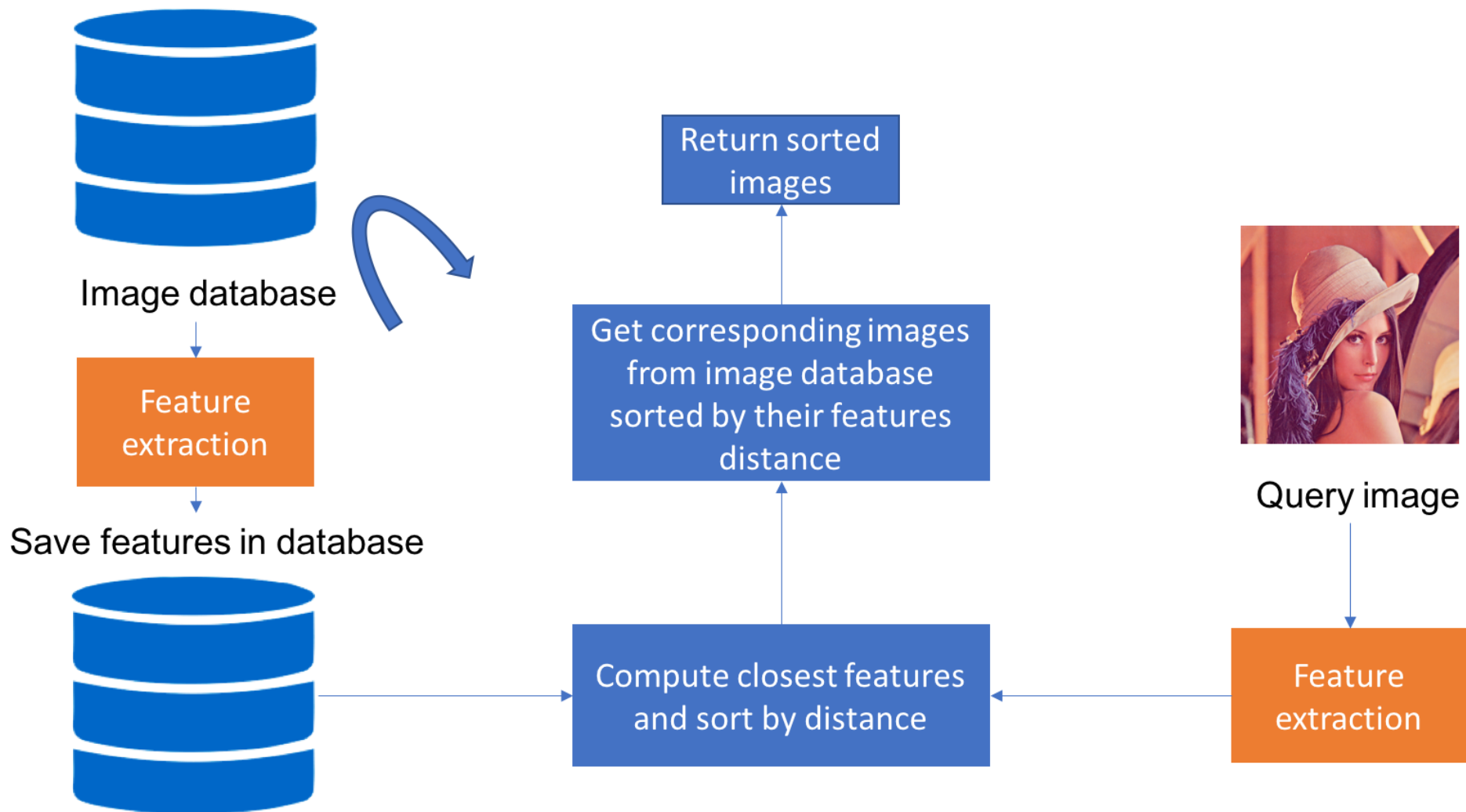
Dataset

Fashion MNIST



Content-based Image Retrieval (CBIR)

The general idea



A simple CBIR prototype

Using a convolutional deep neural network and the KNN algorithm



Steps:

1. Train a convolutional deep neural network with the Fashion MNIST for multi-label classification
2. Extract the feature vectors of all dataset images as represented internally in the neural network
3. Store the extracted image feature vectors in a feature database
4. Fit a KNN model using the extracted image feature vectors stored in the feature database
5. For an input query image:
 - Extract its image feature vector using the same technique used in step 2
 - Use the KNN model to find the closest K neighbor images using a distance metric
 - Present the closest K images sorted by distance

Convolutional Neural Network Classifier

Summary and architecture of the neural network



| Layer (type) | Output Shape | Param # |
|--------------------------------|--------------------|---------|
| conv2d_1 (Conv2D) | (None, 28, 28, 64) | 320 |
| max_pooling2d_1 (MaxPooling2D) | (None, 14, 14, 64) | 0 |
| dropout_1 (Dropout) | (None, 14, 14, 64) | 0 |
| conv2d_2 (Conv2D) | (None, 14, 14, 32) | 8224 |
| max_pooling2d_2 (MaxPooling2D) | (None, 7, 7, 32) | 0 |
| dropout_2 (Dropout) | (None, 7, 7, 32) | 0 |
| flatten_1 (Flatten) | (None, 1568) | 0 |
| dense_1 (Dense) | (None, 256) | 401664 |
| dropout_3 (Dropout) | (None, 256) | 0 |
| dense_2 (Dense) | (None, 10) | 2570 |

Total params: 412,778
Trainable params: 412,778
Non-trainable params: 0

Conv(64, Kernel(2, 2), Padding(Same))+Relu
MaxPooling(PoolSize(2, 2), Strides(2, 2))
Dropout(0.3)
Conv(32, Kernel(2, 2), Padding(Same))+Relu
MaxPooling(PoolSize(2, 2), Strides(2, 2))
Dropout(0.3)
FC(256)+Relu
Dropout(0.5)
FC(10)+Softmax

Convolutional Neural Network Classifier

Training information



| Dataset | Fashion MNIST |
|----------------------|-----------------|
| Total Classes | 10 |
| Total Images | 70000 |
| Training Images | 55000 |
| Validation Images | 5000 |
| Testing Images | 10000 |
| Images Resolution | 28x28 |
| Images Normalization | MinMax |
| Learning Rate | 1e-3 |
| Training Algorithm | Backpropagation |
| Optimization Method | Adam |
| Loss Function | Cross-entropy |
| Batch Size | 64 |
| Training Epochs | 10 |

Example Queries

Example 1 (closest K=30 neighbors)



T-shirt/Top

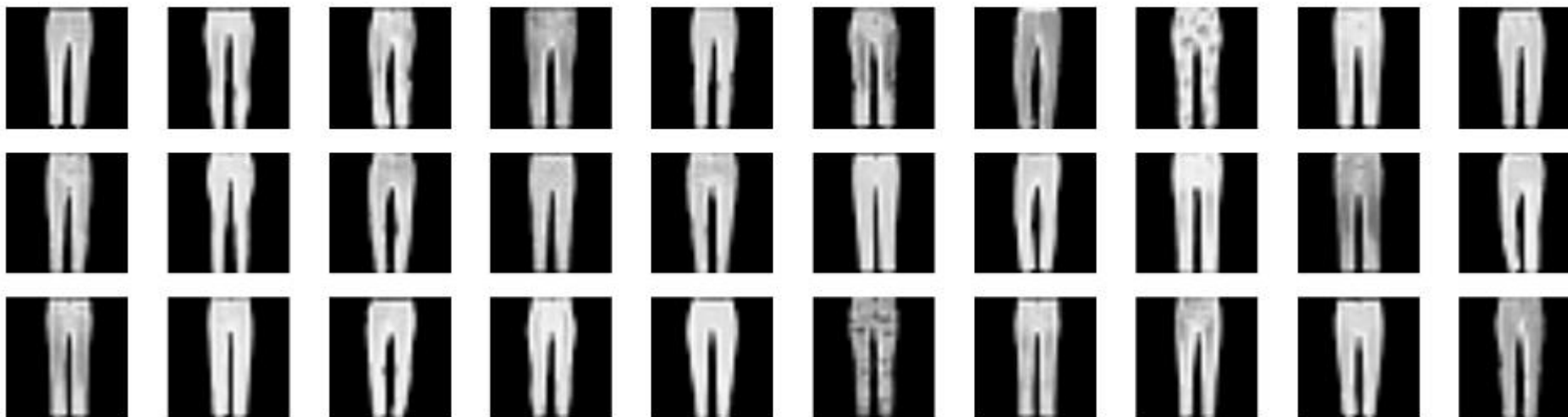


Example Queries

Example 2 (closest $K=30$ neighbors)



Trouser



Example Queries

Example 3 (closest K=30 neighbors)



Pullover

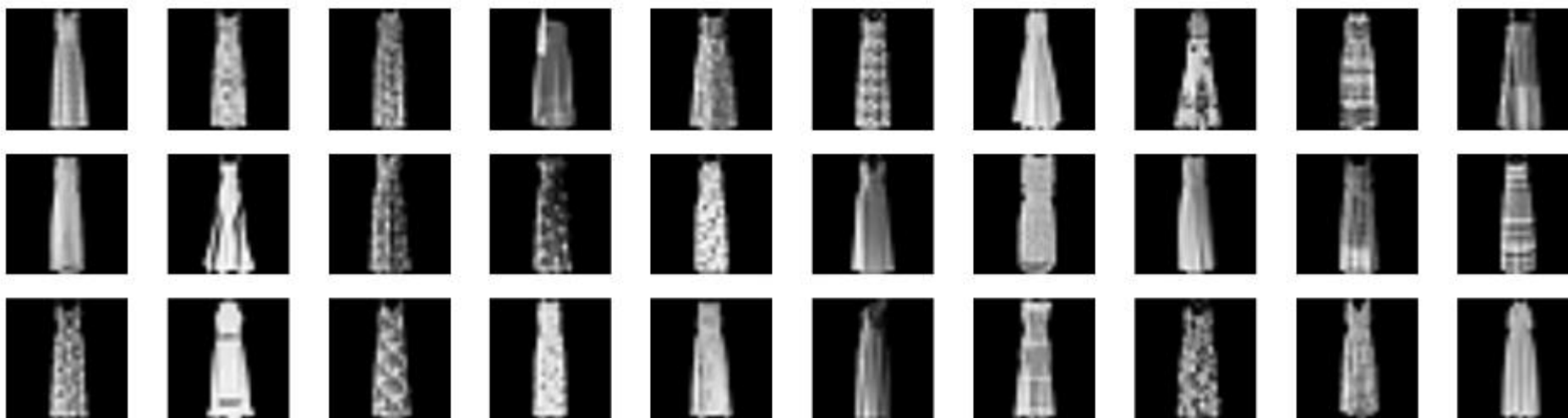
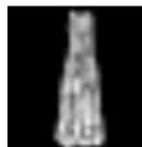


Example Queries

Example 4 (closest K=30 neighbors)



Dress



Example Queries

Example 5 (closest K=30 neighbors)



Coat



Example Queries

Example 6 (closest $K=30$ neighbors)



Sandal



Example Queries

Example 7 (closest K=30 neighbors)



Shirt

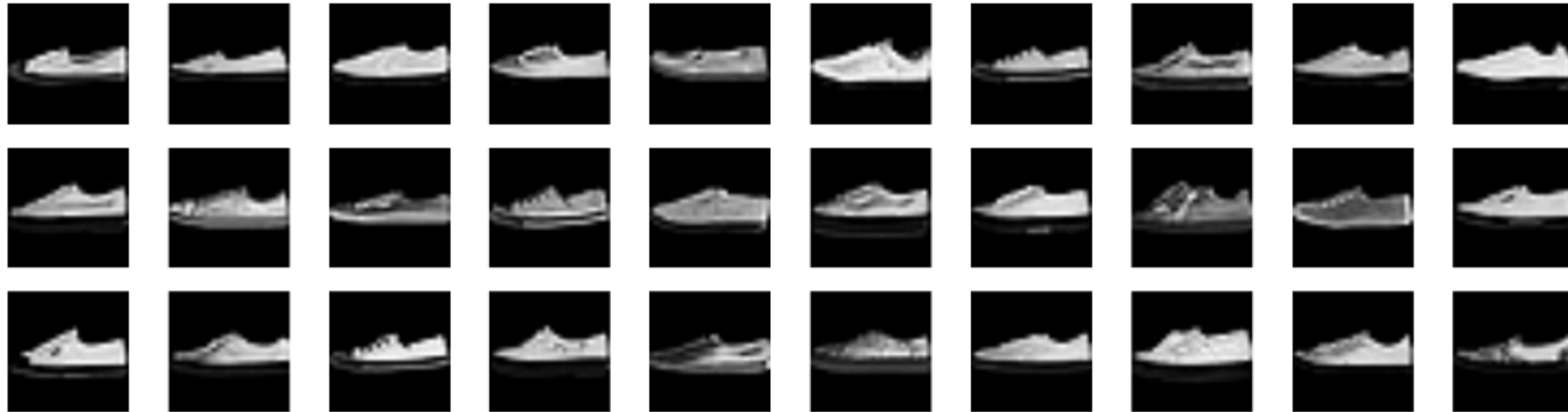


Example Queries

Example 8 (closest K=30 neighbors)



Sneaker



Example Queries

Example 9 (closest K=30 neighbors)



Bag

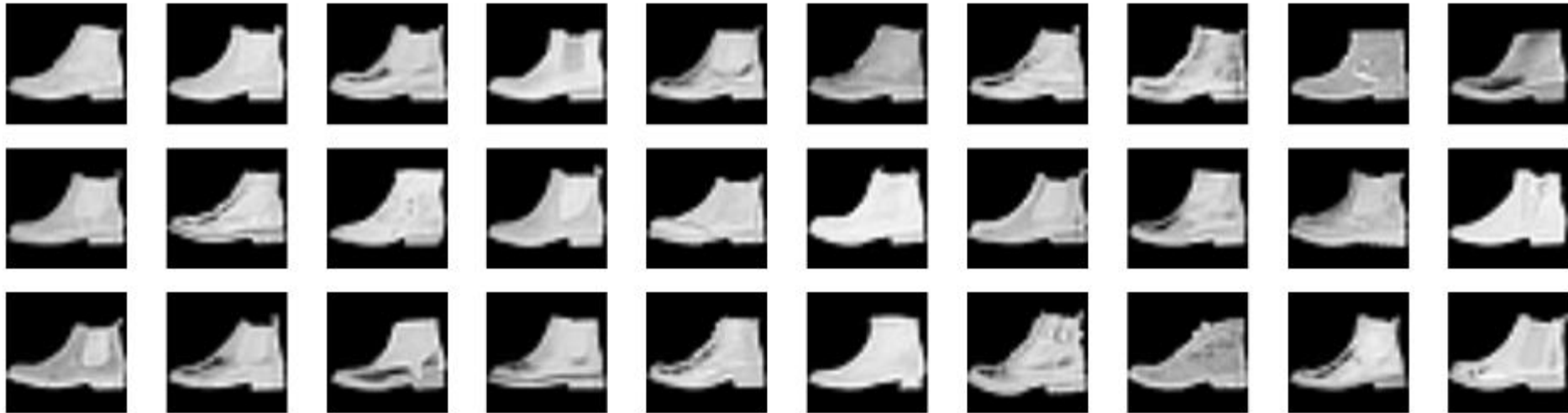


Example Queries

Example 10 (closest $K=30$ neighbors)

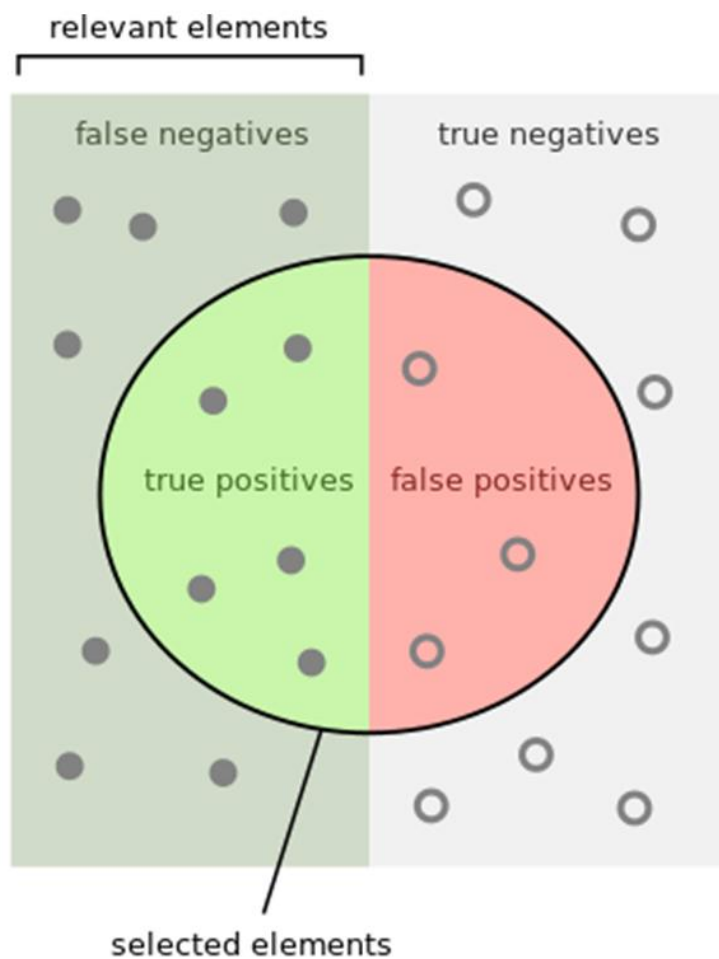


Ankle Boot



Metrics

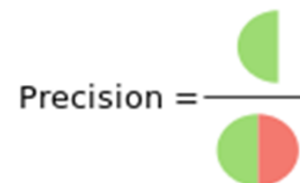
Precision & Recall



$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

How many selected items are relevant?



How many relevant items are selected?



Experiments

Using various similarity metrics, neural network layers, K sizes



We have calculated average precision and recall metrics for experiments using:

- 100 different random test images as search input queries
- various K sizes ranging from 1 to 500
- features vectors extracted from various neural network layers
- various similarity metrics (Cosine, Euclidean, etc)

Results

The experimental results are stored in a CSV file



CSV file format:

Neural Network Layer, Similarity Metric, K size, Avg Precision (%), Avg Recall (%)

flatten_1, euclidean, 1, 87.0, 0.015

...

...

flatten_1, cosine, 500, 71.6, 5.970

...

...

dense_1, euclidean, 1, 94.0, 0.016

...

...

dense_1, cosine, 500, 84.8, 7.074

Results

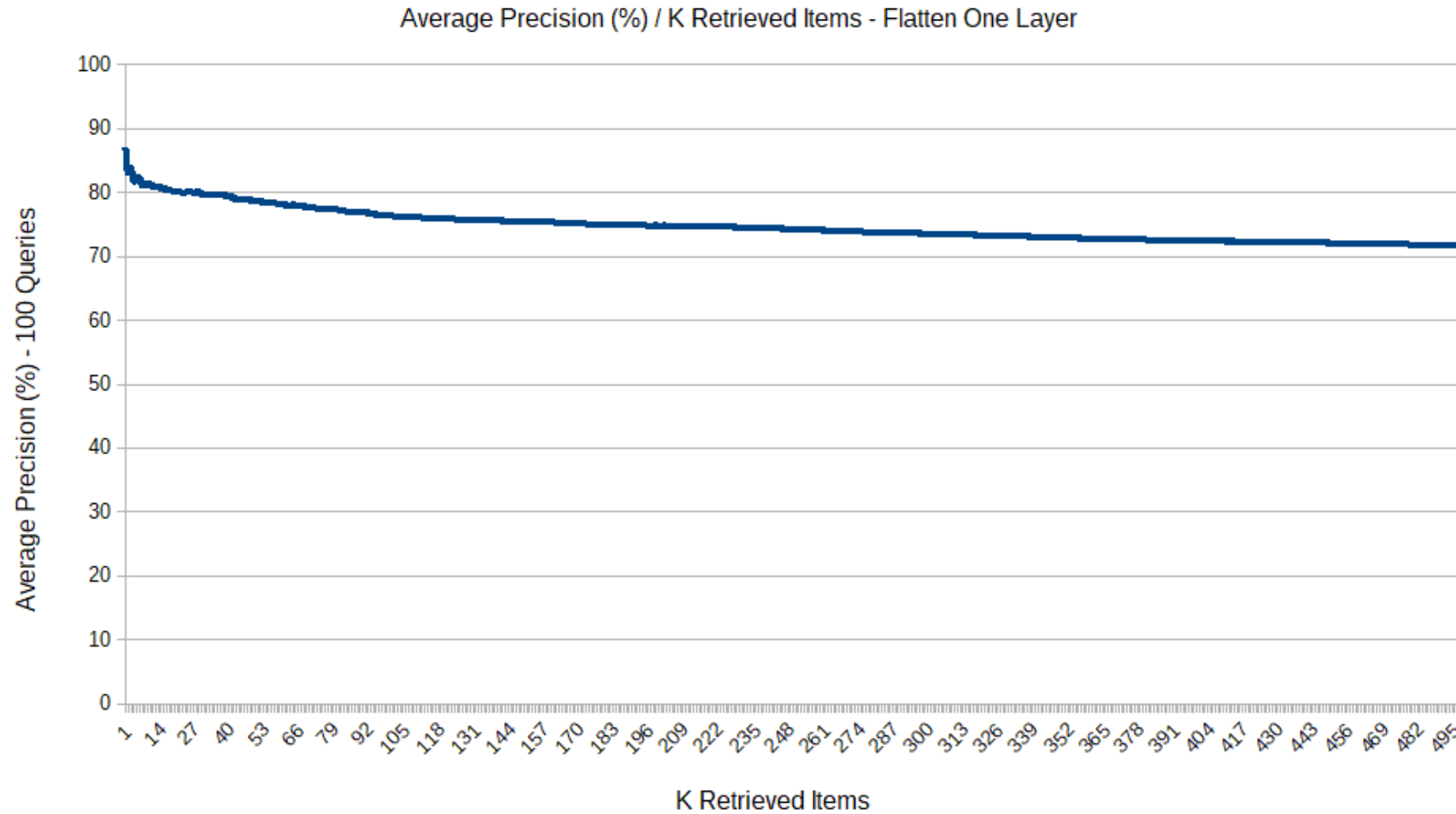
Conclusions



- When K size increases the precision decreases
- When K size increases the recall increases
- The feature vectors extracted from 'dense_1' layer give better precision, recall results
- The cosine similarity gives better precision, recall results
- In general there is a trade-off between precision and recall metrics

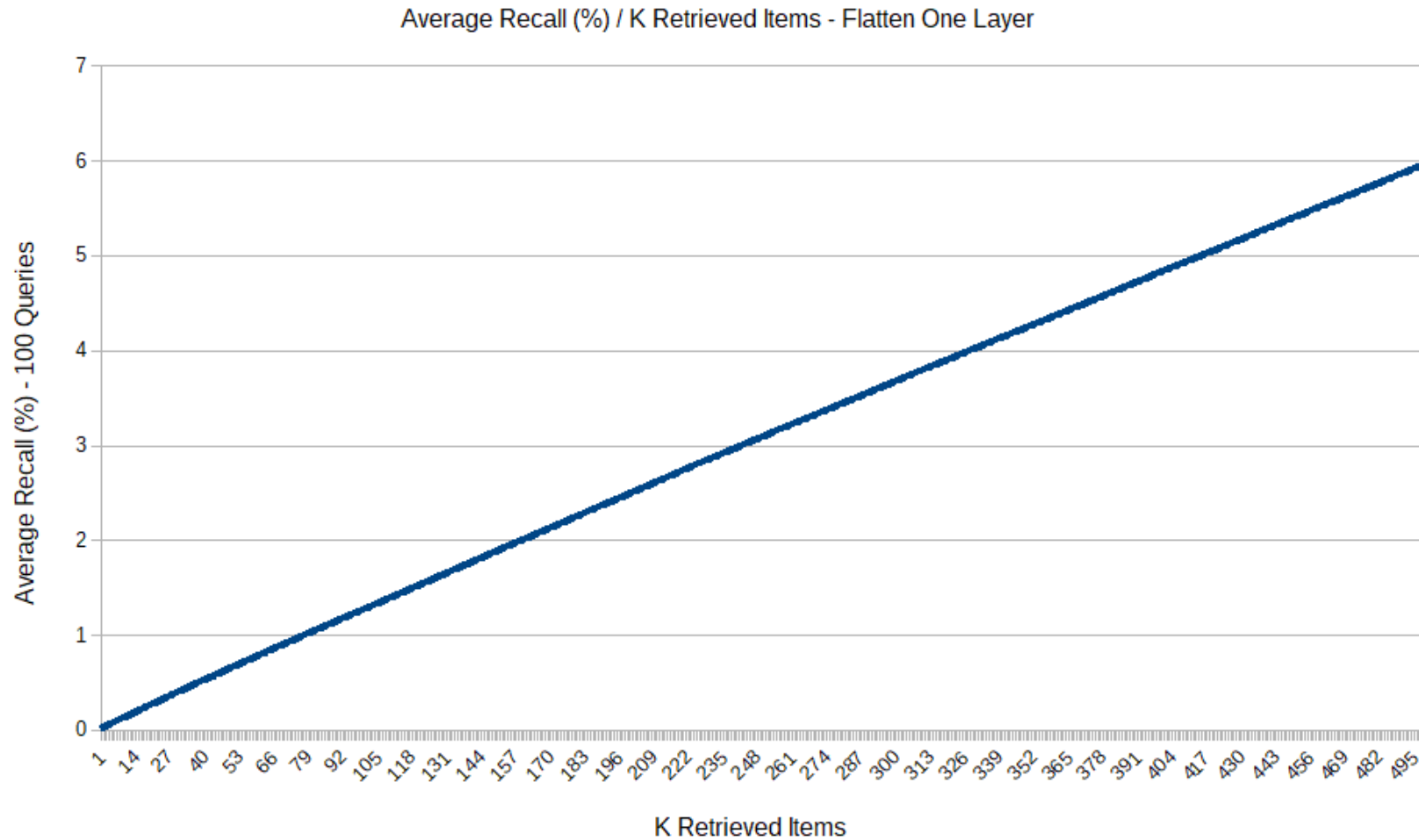
Curves

Average Precision / K Retrieved Items (for flatten_1 layer)



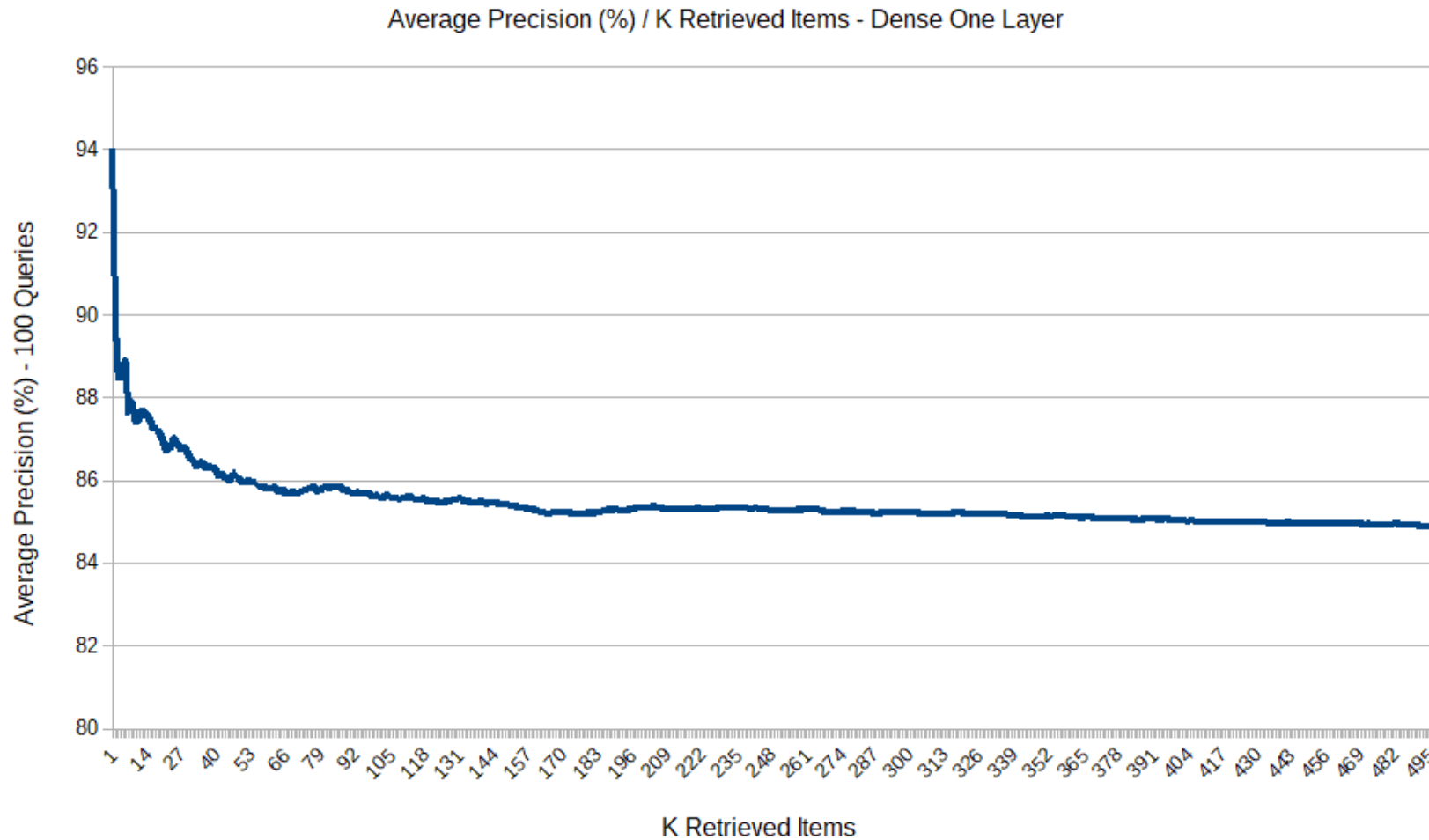
Curves

Average Recall / K Retrieved Items (for flatten_1 layer)



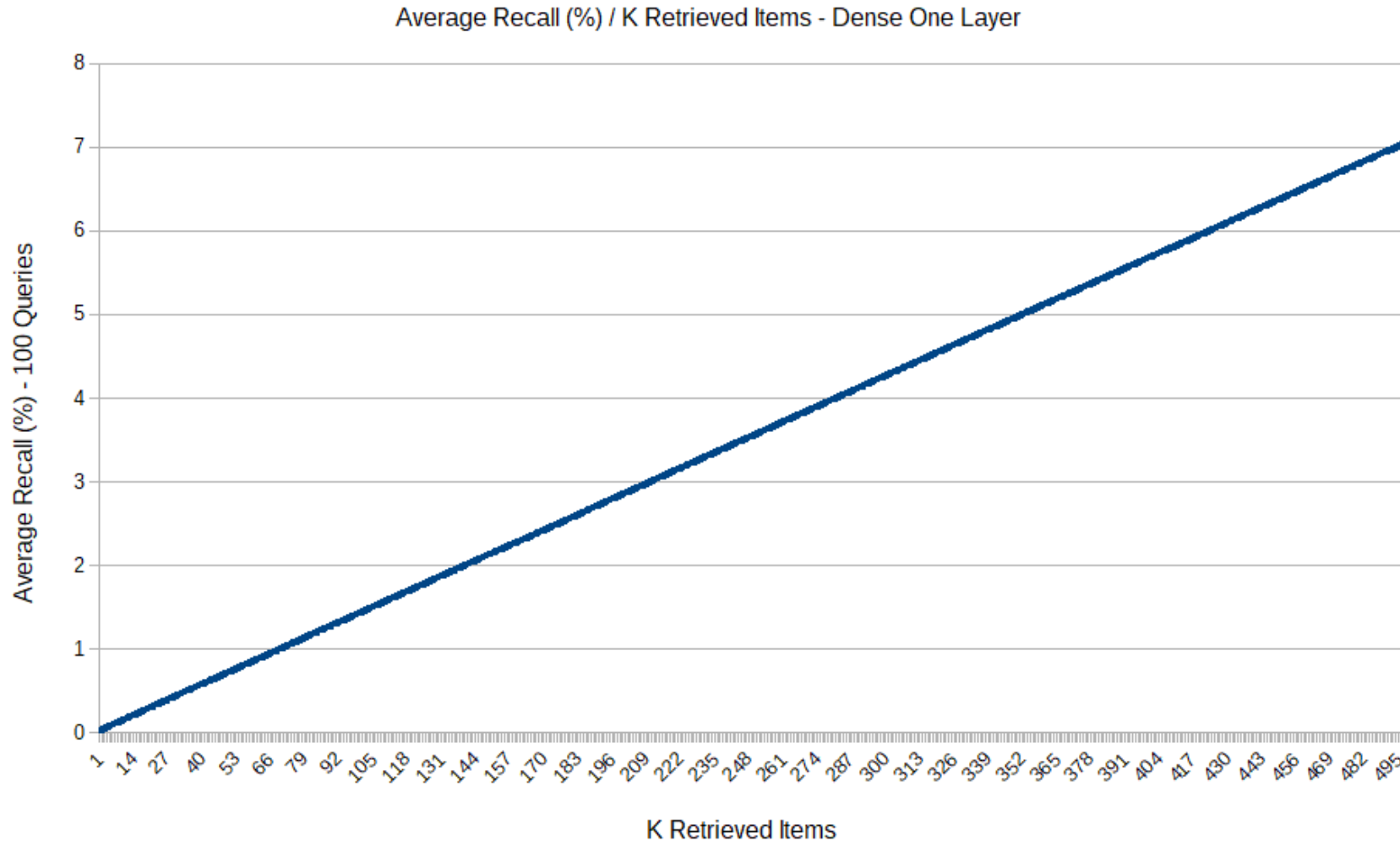
Curves

Average Precision / K Retrieved Items (for dense_1 layer)



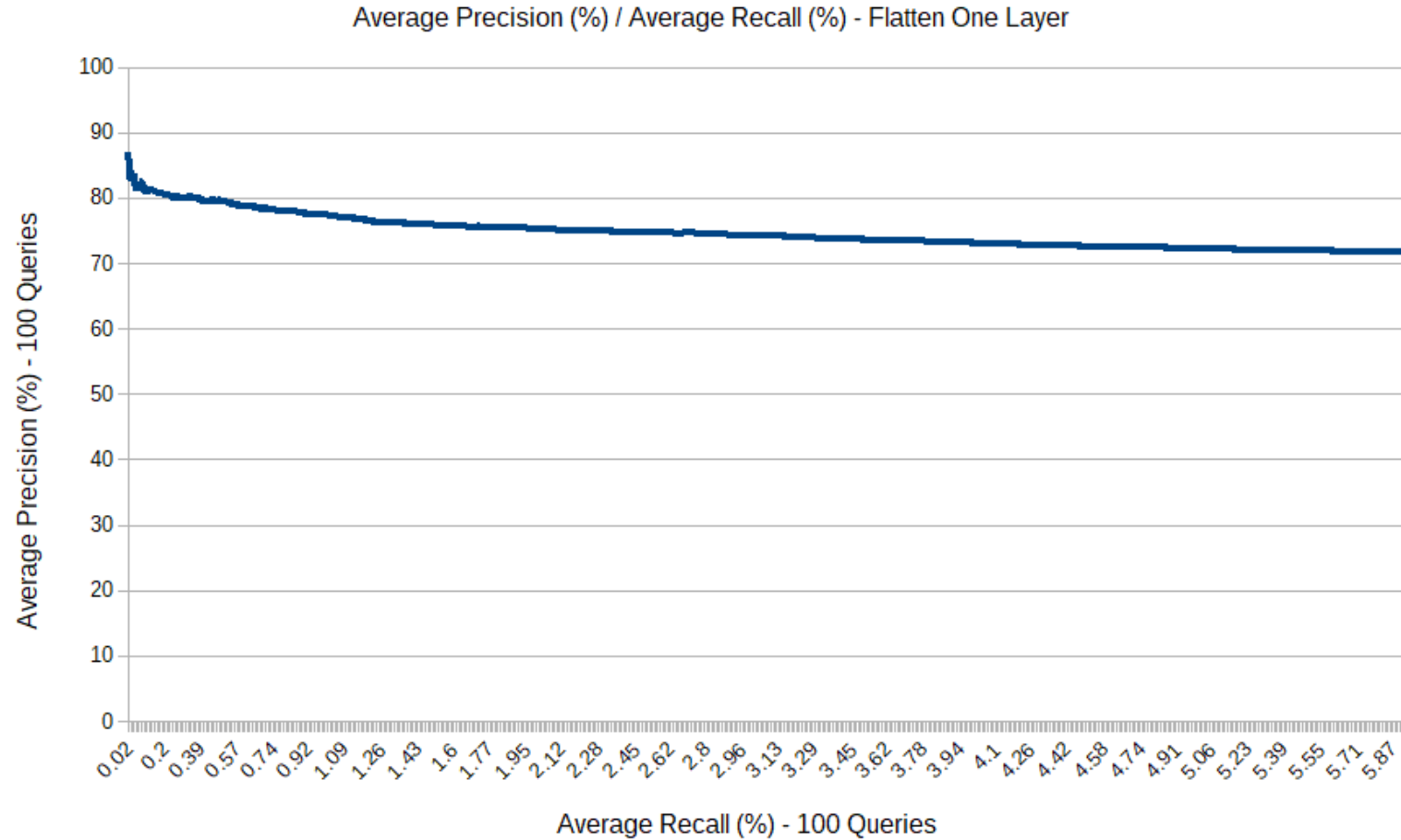
Curves

Average Recall / K Retrieved Items (for dense_1 layer)



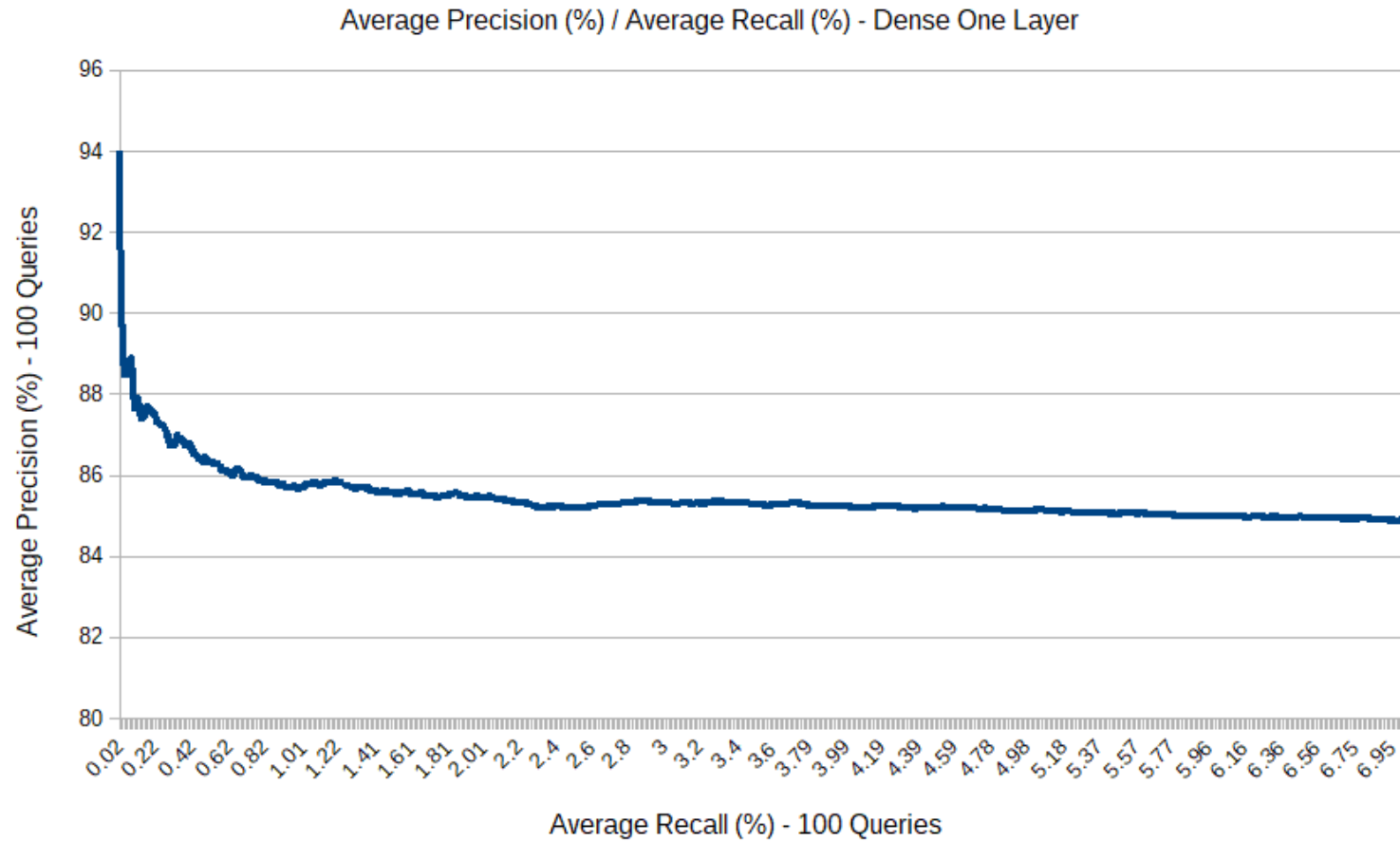
Curves

Average Precision / Average Recall (for flatten_1 layer)



Curves

Average Precision / Average Recall (for dense_1 layer)



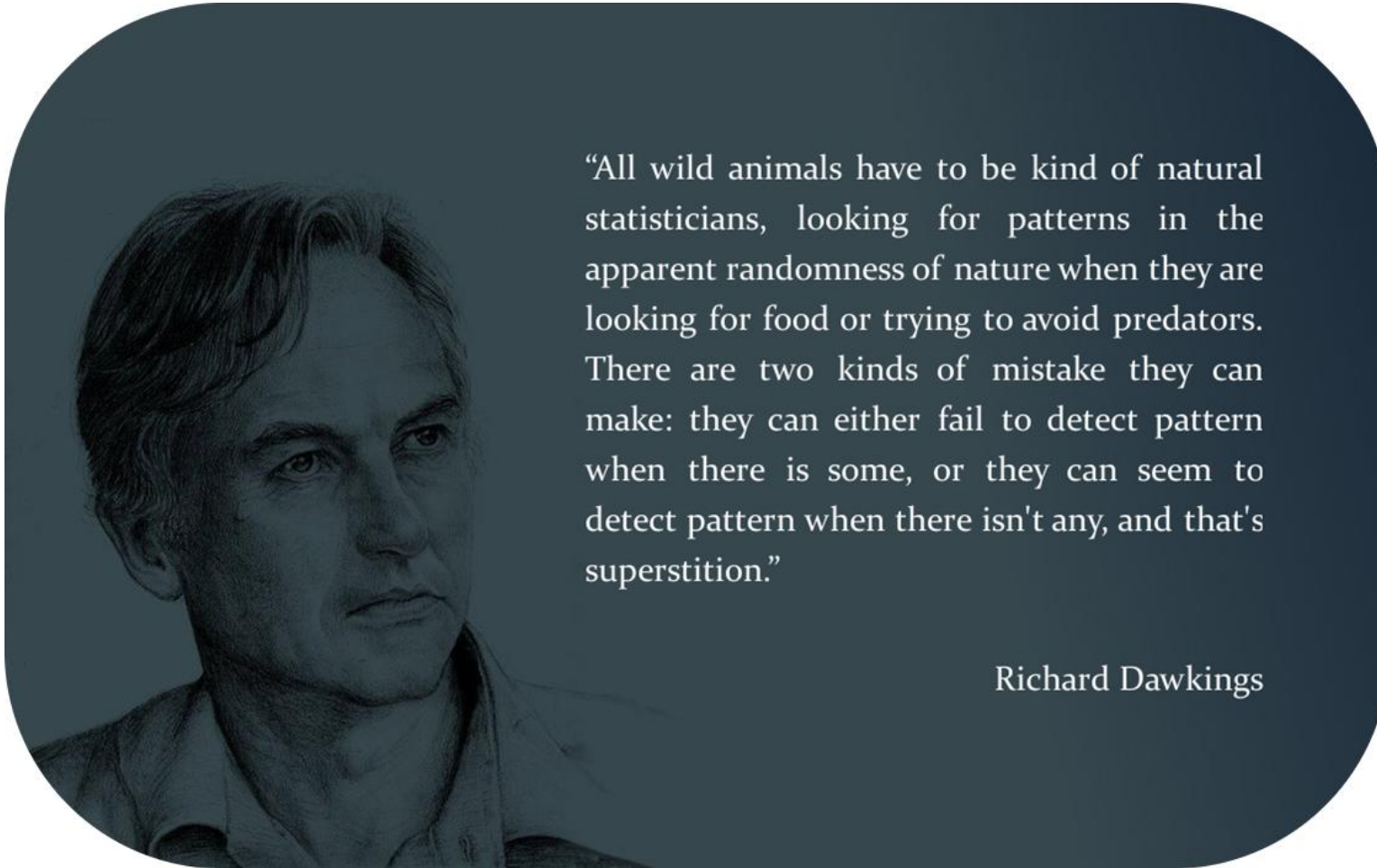
Software Requirements



- Programming Languages : Python
- IDE : Spyder, JetBrains PyCharm
- Application-level Dependency Manager: Anaconda
- Python Modules :

| Module | Version |
|------------|---------|
| sklearn | 0.19.1 |
| numpy | 1.14.3 |
| matplotlib | 2.2.2 |
| keras | 2.1.6 |
| tensorflow | 1.8.0 |

Thank you a lot and have a nice day!



“All wild animals have to be kind of natural statisticians, looking for patterns in the apparent randomness of nature when they are looking for food or trying to avoid predators. There are two kinds of mistake they can make: they can either fail to detect pattern when there is some, or they can seem to detect pattern when there isn't any, and that's superstition.”

Richard Dawkins