# Distant Reading in R

## Machine Learning

Simone Rebora & Giovanni Pietro Vitali
simone.rebora@univr.it      giovannipietrovitali@gmail.com

# THE "TWO PARADIGMS" OF ARTIFICIAL INTELLIGENCE

## Top-down

Define a set of rules which model the human cognition
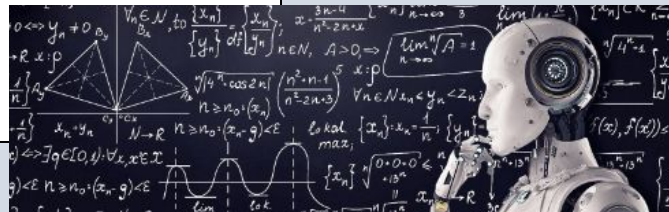Apply those rules to new subjects and situations

Dominating paradigm in the "first wave" of AI (2nd half XX century)

Main con
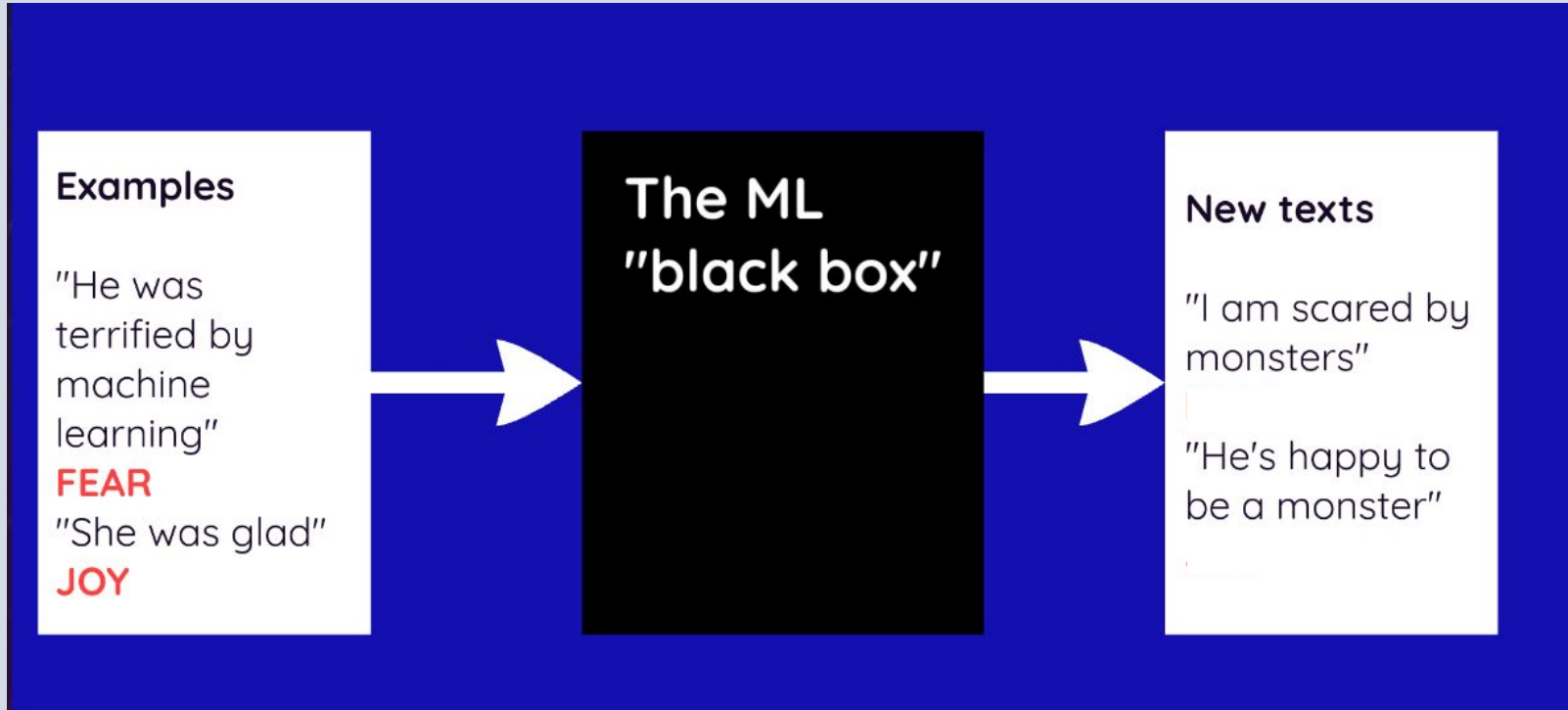- absence of flexibility

## Bottom-up

Define a system that is able to "learn" a task from a set of examples (i.e. imitate the human)
Apply the system to new samples
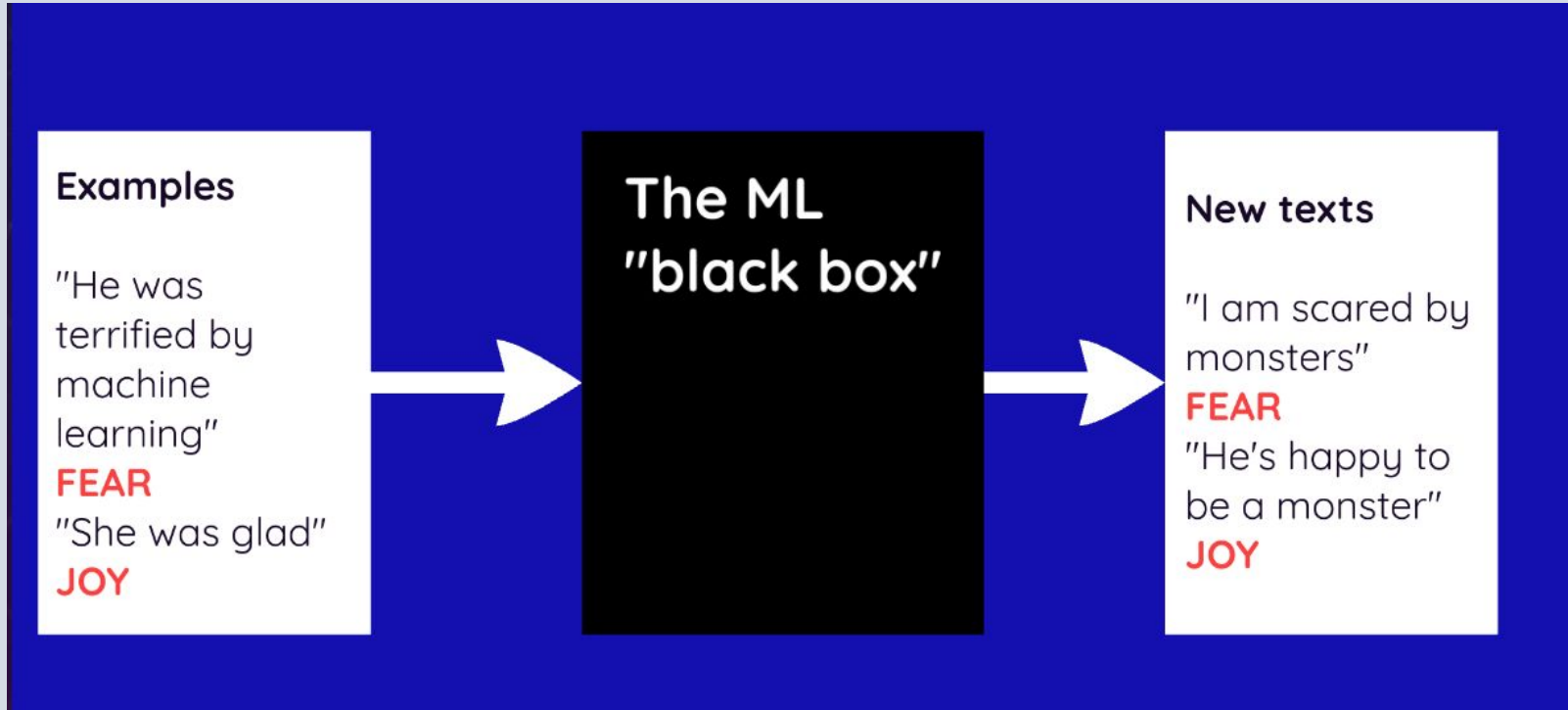
This is the machine learning approach!





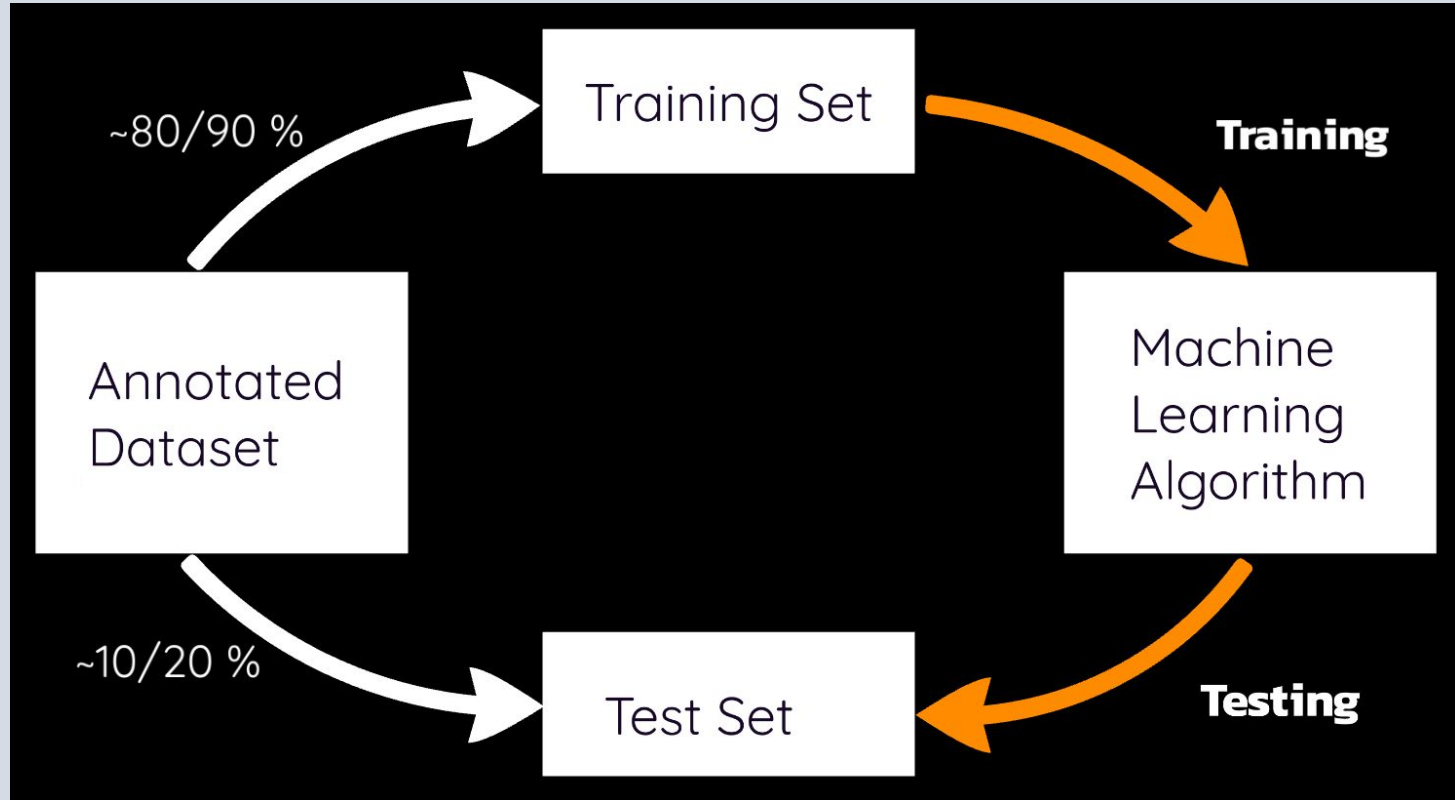Distant Reading in ®

# ML - THE WORKING LOGIC

**Examples**

"He was terrified by machine learning"
**FEAR**
"She was glad"
**JOY**

**The ML "black box"**

**New texts**

"I am scared by monsters"

"He's happy to be a monster"

*example: sentiment analysis

Distant Reading in

# ML - THE WORKING LOGIC

**Examples**

"He was terrified by machine learning"
**FEAR**
"She was glad"
**JOY**

**The ML "black box"**

**New texts**

"I am scared by monsters"
**FEAR**
"He's happy to be a monster"
**JOY**

*example: sentiment analysis
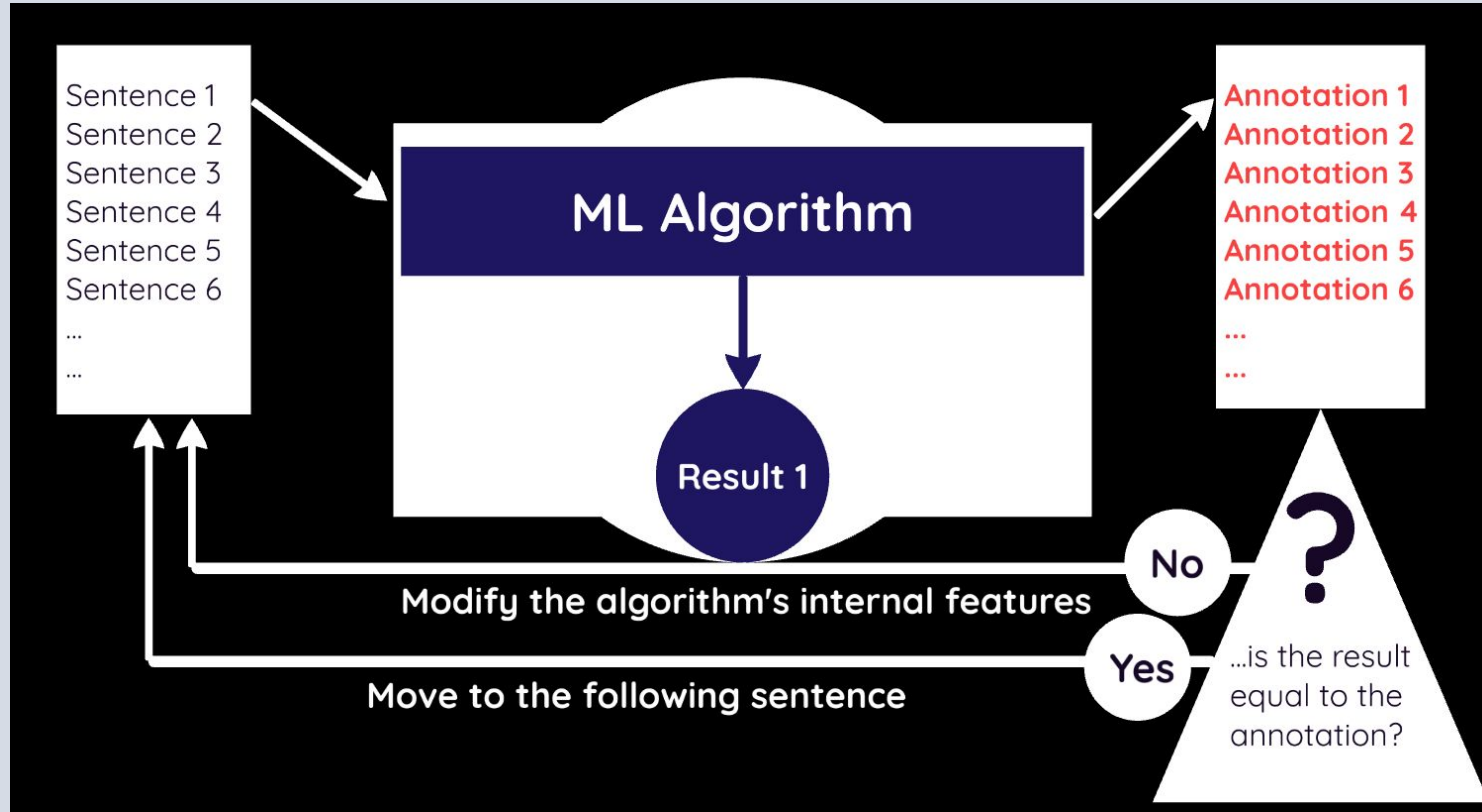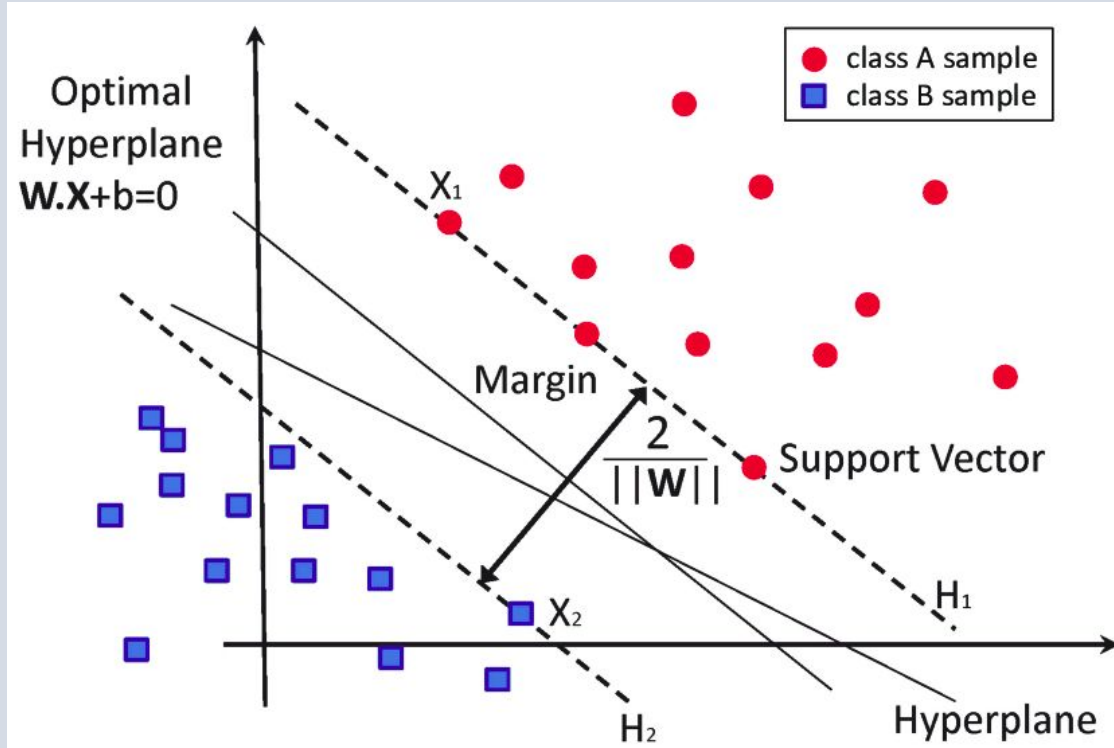
Distant Reading in

# INSIDE THE ML "BLACK BOX"

# TRAINING

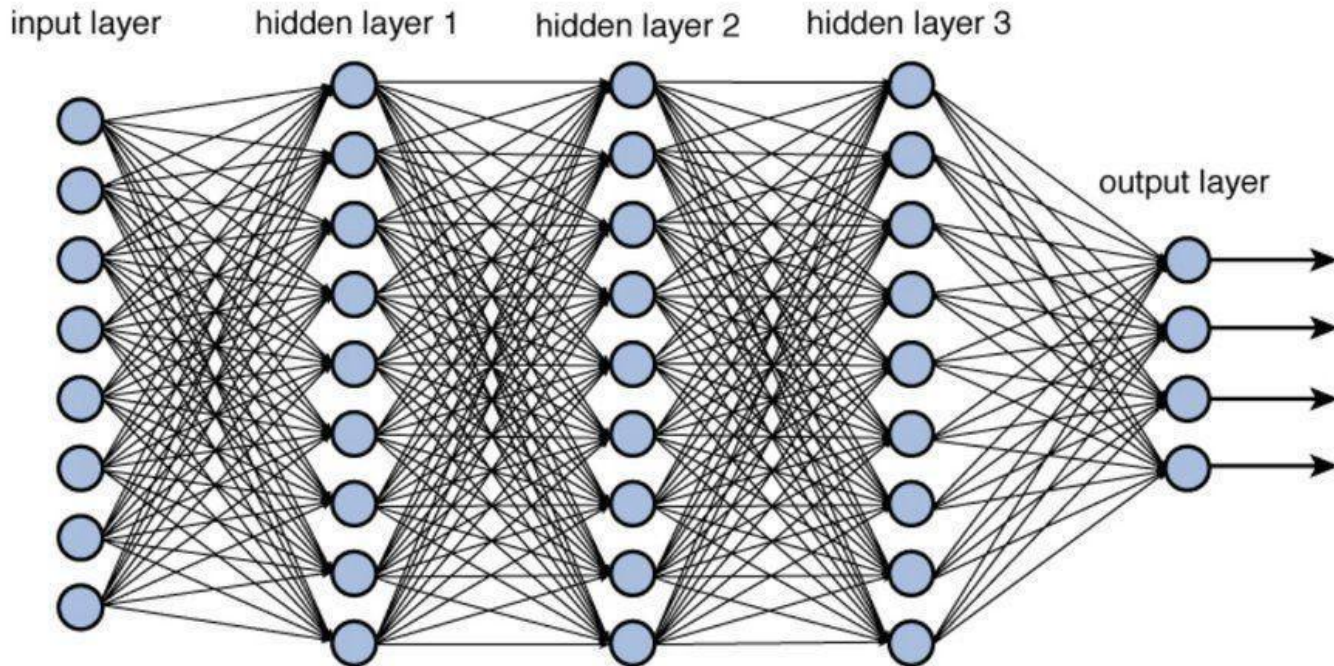# ML ALGORITHMS



Support Vector Machines (SVM)
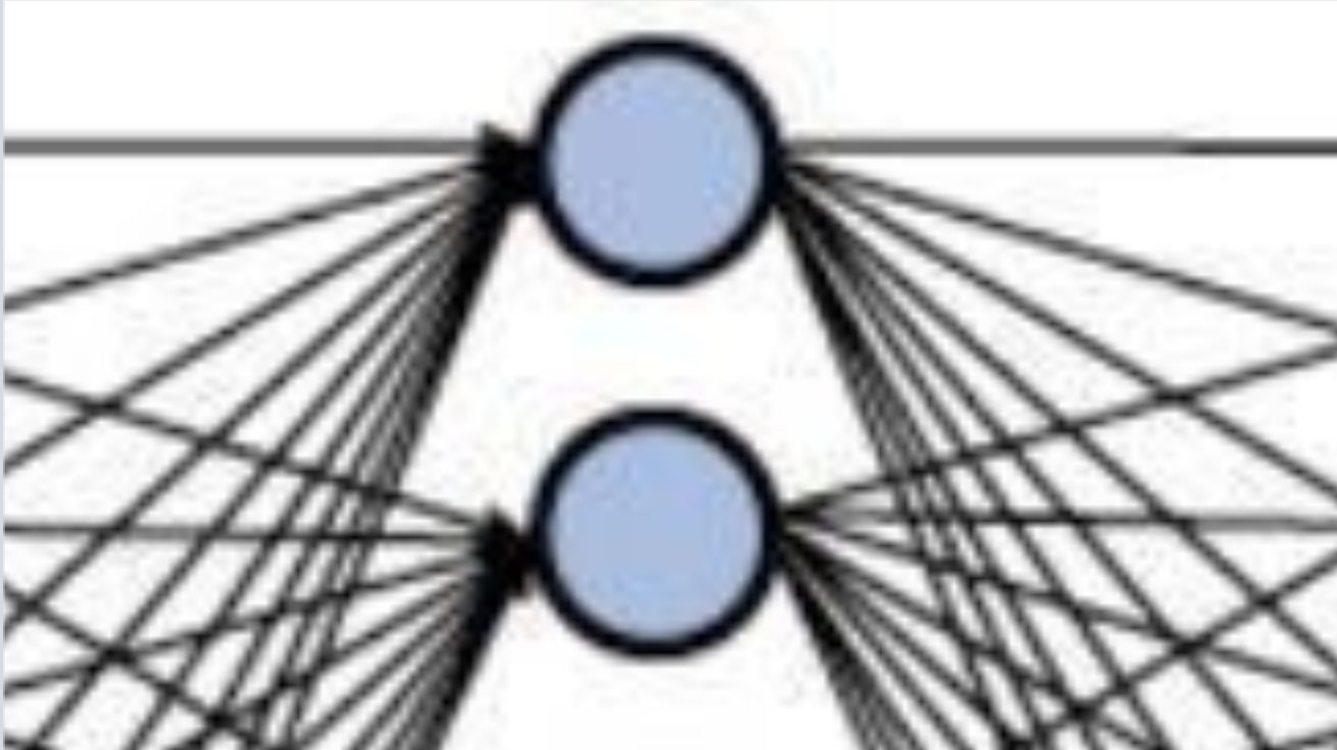
# ML ALGORITHMS



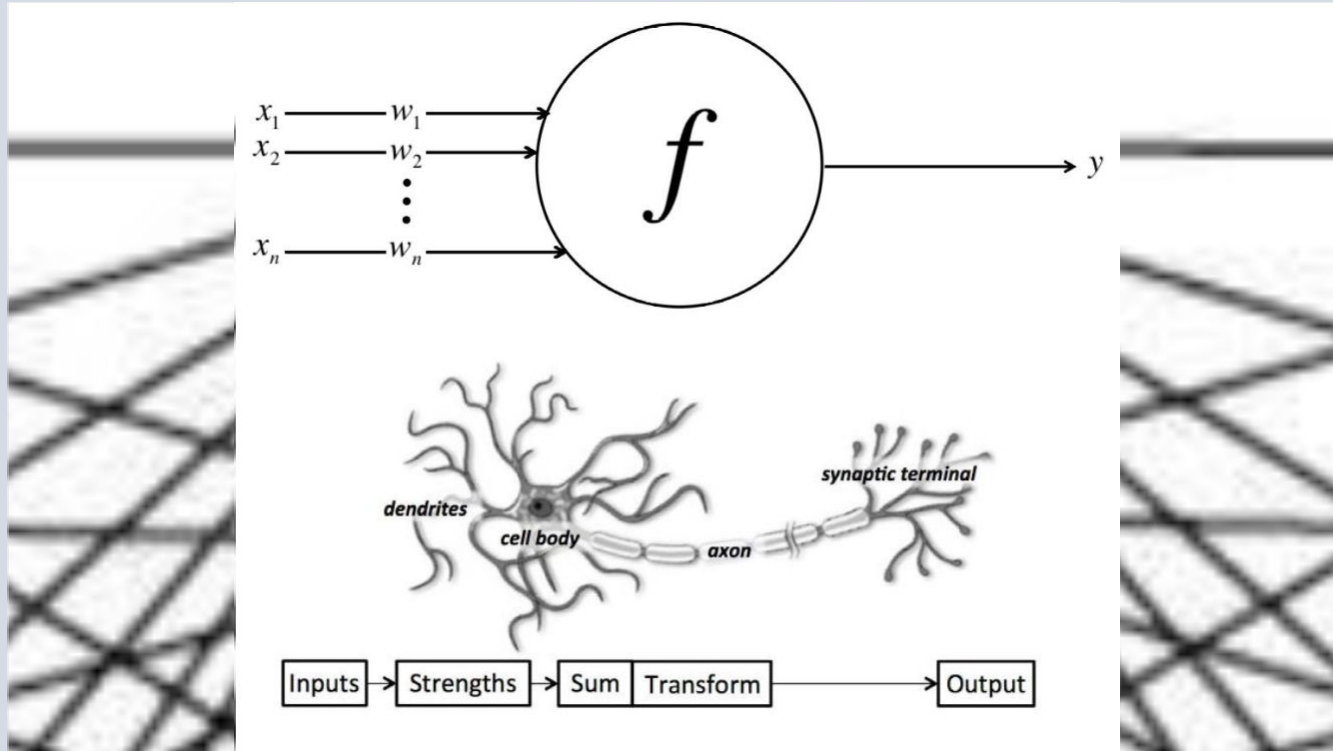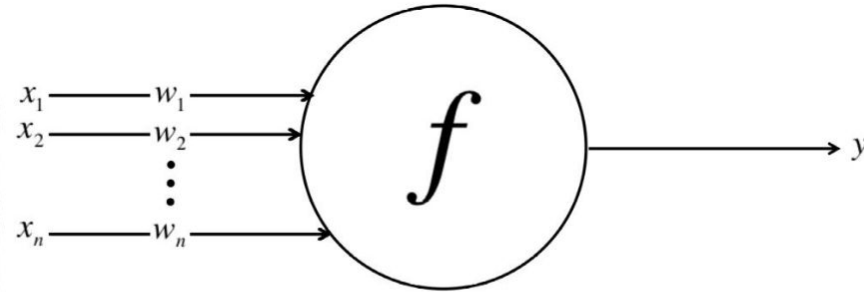Figure 12.2 Deep network architecture with multiple layers.

# ML ALGORITHMS

# ML ALGORITHMS

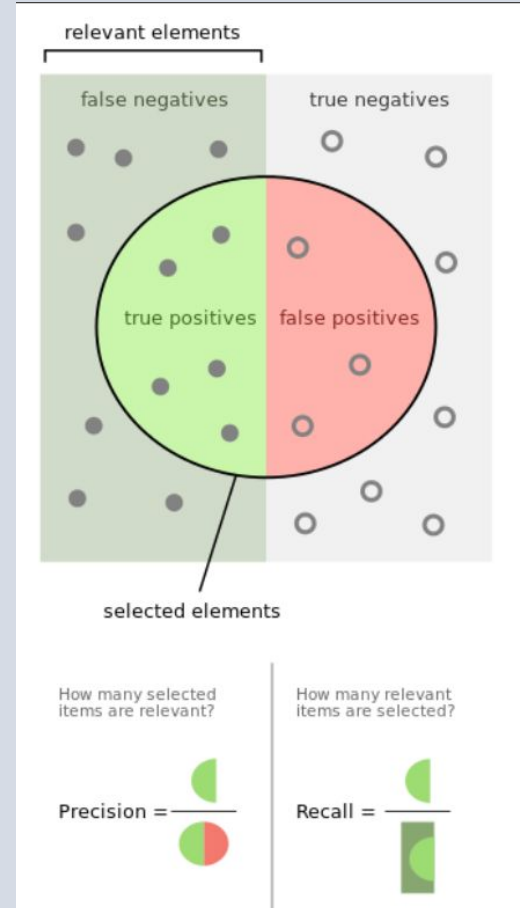# ML ALGORITHMS

# INSIDE THE ML "BLACK BOX"

# TESTING

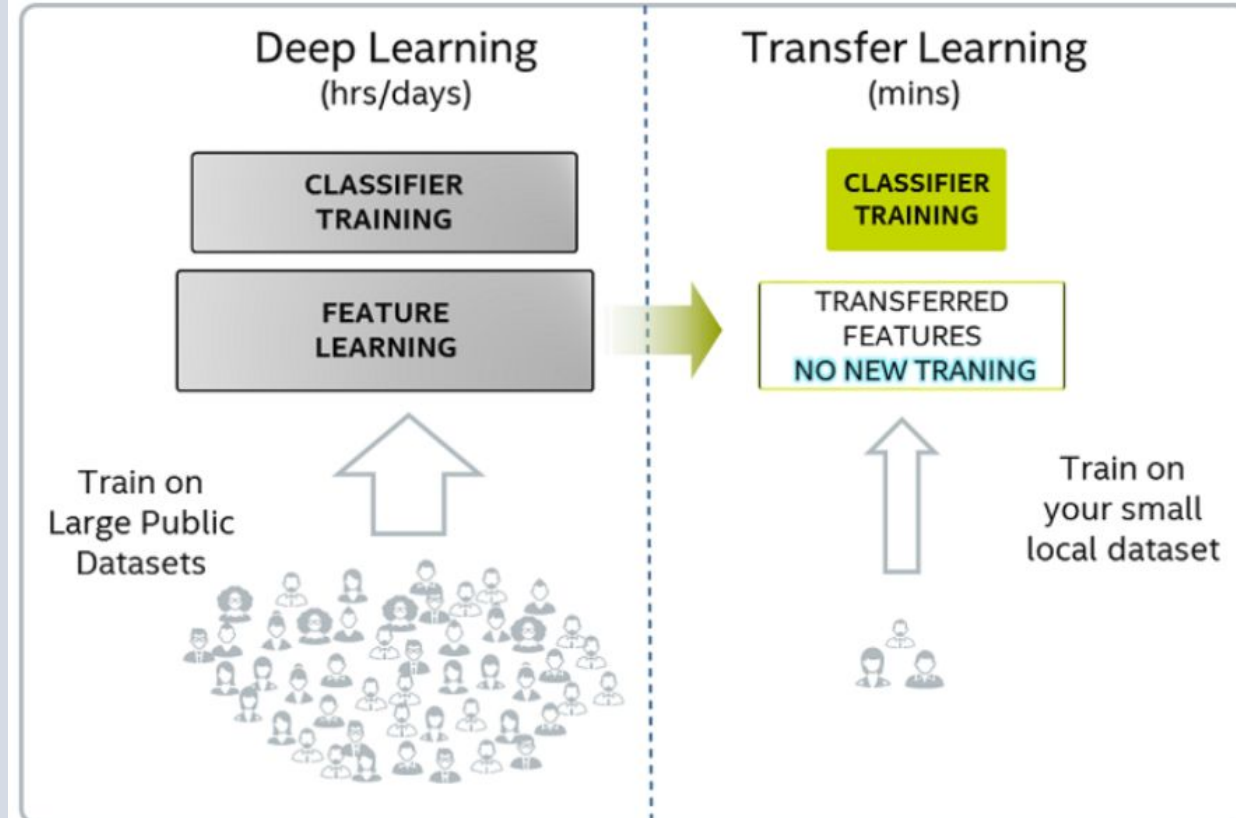Once the training is complete, the trained algorithm is "tested" on never-seen documents

Main reason: the algorithm might have learned how to work just on the training data, but not on the task in general (overfitting)

When the task is that of assigning a label (e.g. an emotion), precision/recall are generally used

# TRANSFER LEARNING
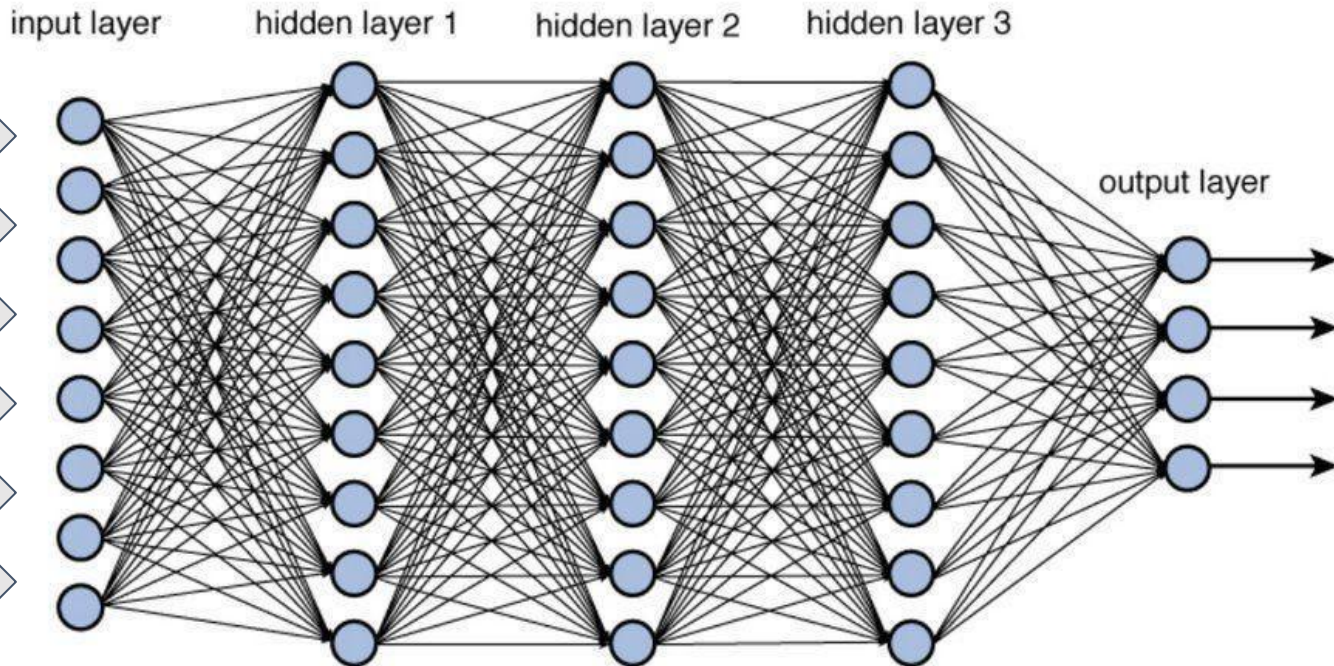


...see BERT

# EMBEDDINGS

Deep Neural Network

input layer    hidden layer 1    hidden layer 2    hidden layer 3

output layer

To
be
or
not
to
be

Figure 12.2 Deep network architecture with multiple layers.

Distant Reading in

# EMBEDDINGS

0.123

1.743

0.325

1.143

0.463

1.153

Distant Reading in
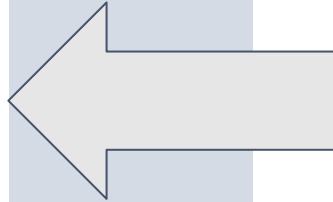
# EMBEDDINGS

To ⟹ 0.123

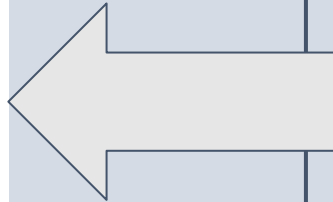be ⟹ 1.743

or ⟹ 0.325

not ⟹ 1.143

to ⟹ 0.463

be ⟹ 1.153

word embeddings

# COMPUTATIONAL MODELING

"1. A model is a model *of something.* A model is always a kind of mapping. It represents something, an object, a concept, and so on, by representing it using something else like clay, words, images, and so forth.

2. A model is *not the original* and it is not a *copy of the original*. Unlike a copy, a model doesn't capture all features of the entity it represents, only some of them. The choice of features selected to be present in the model is usually based on assumptions by the creator of the model concerning which features are relevant for the intended use of the model."
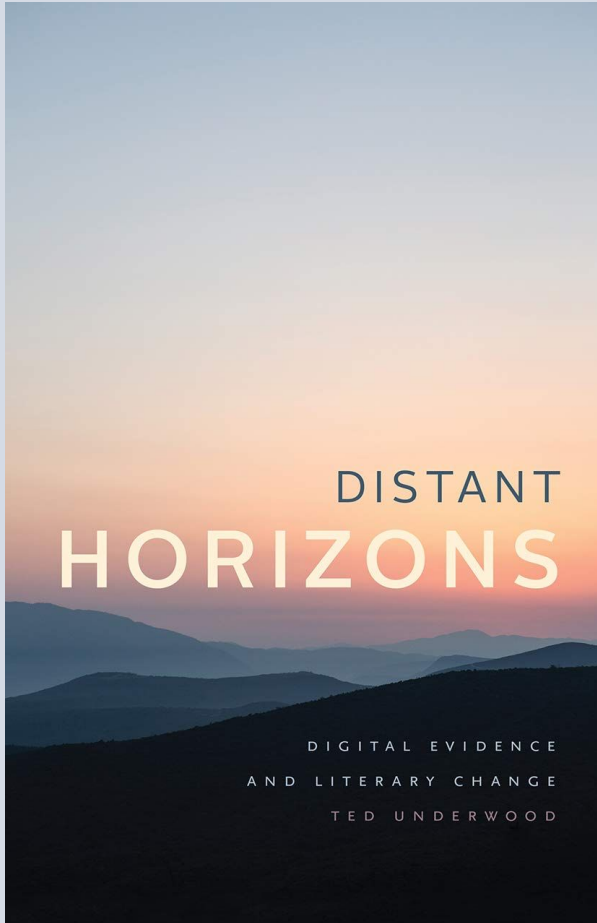(Jannidis and Flanders, 2019)

# APPLICATIONS

- using a combination of sentiment analysis, topic modeling, et al. (i.e. creating embeddings)
- to train a ML classifier that predicts the commercial success of novels

# APPLICATIONS



DISTANT
HORIZONS

DIGITAL EVIDENCE
AND LITERARY CHANGE
TED UNDERWOOD

- using a simple ML classifier (logistic regression) to predict various phenomena (like genre, literariness, etc.)
- and then looking at the features that made the classifier successful (e.g. the presence of certain words, etc.)

# CRITICAL ASPECTS

**Opaqueness of the most advanced ML algorithms**

...so you just get the results, but cannot interpret them

**In ML, intelligence is intended as a form of imitation**

...so can it really make predictions?

...so how can it be "creative"?

**ML (and in particular transfer learning) depends heavily on the quality of the training materials**

...ML models built on wide (and uncontrolled) datasets can embed strong biases