# Introduction to Deep Learning
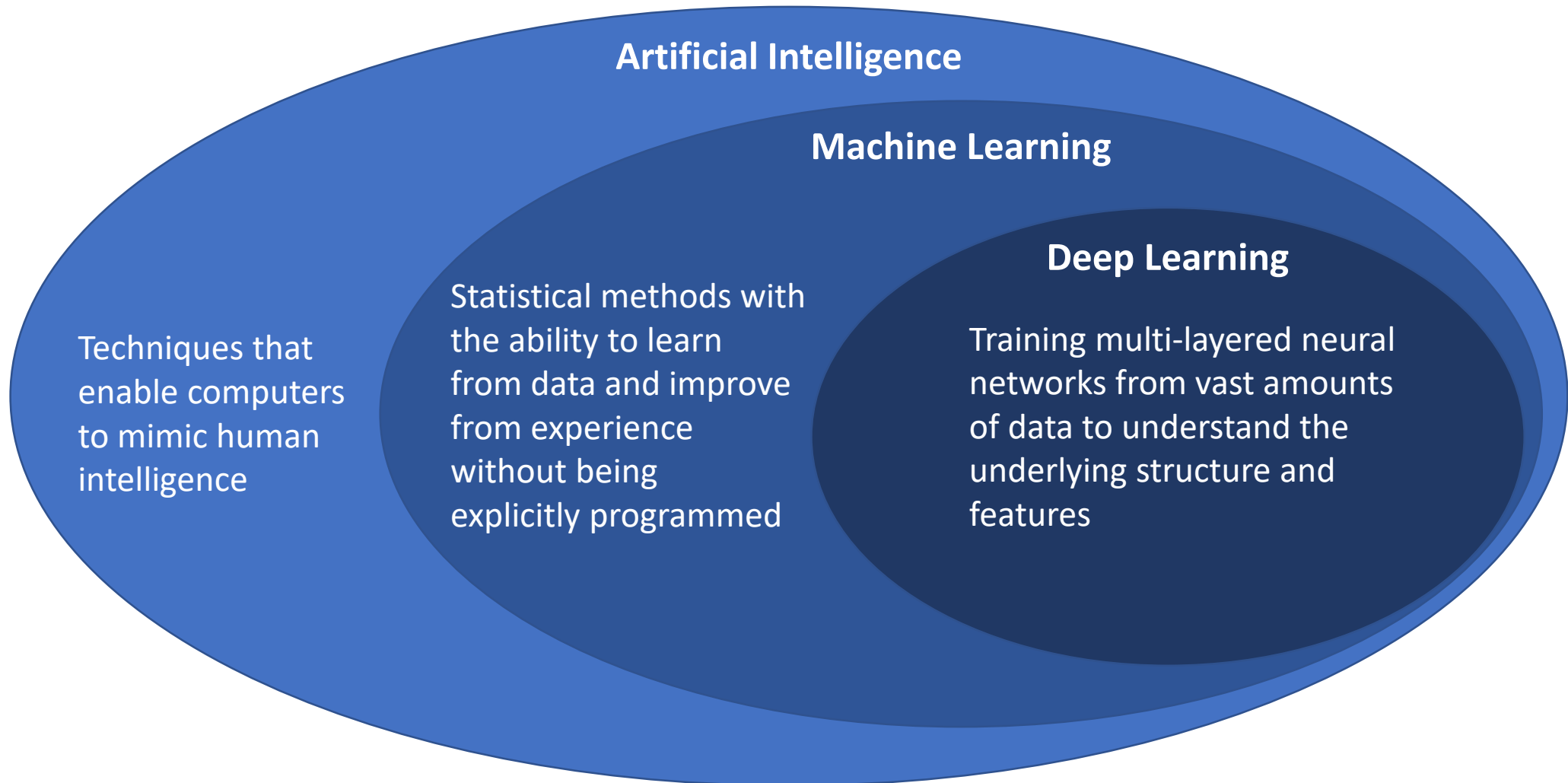
Prof. Dr. Martin Kada
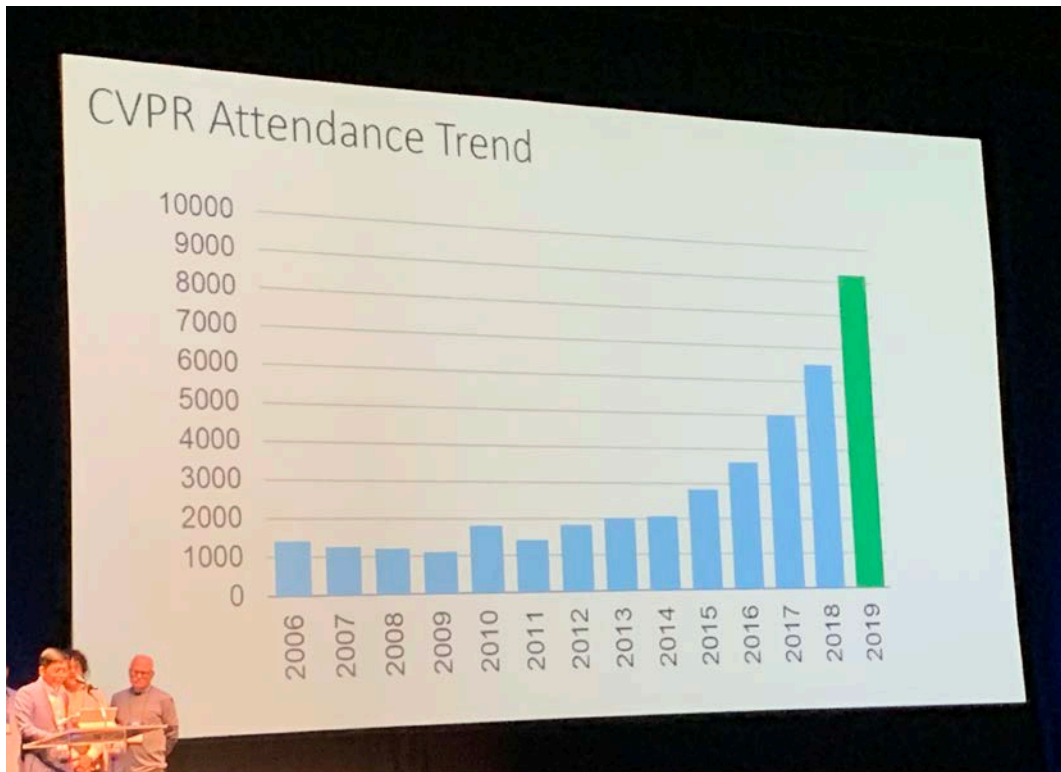
# Artificial Intelligence (AI)
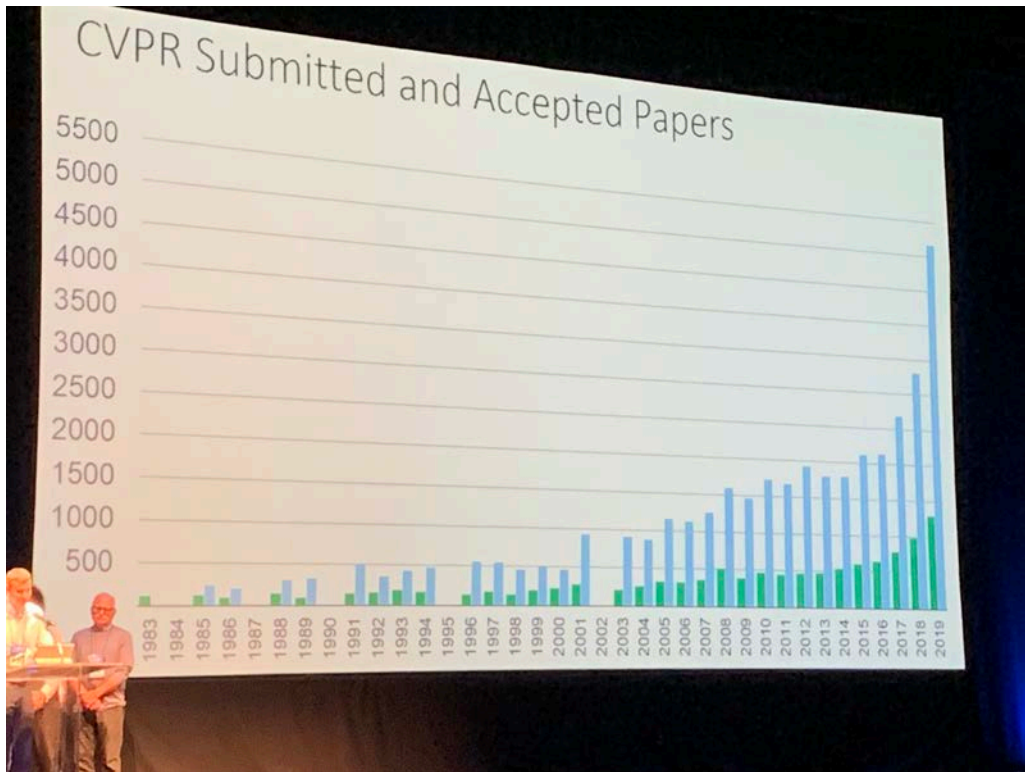


**Artificial Intelligence**

**Machine Learning**

**Deep Learning**

Techniques that enable computers to mimic human intelligence

Statistical methods with the ability to learn from data and improve from experience without being explicitly programmed

Training multi-layered neural networks from vast amounts of data to understand the underlying structure and features

# AI in Computer Vision

# AI in Computer Vision

# AI in Computer Vision

# AI in Computer Vision



Autonomes Fahren: Softbank investiert zwei Milliarden Euro in GM-Tochter Cruise

US-Behörde genehmigt eine Zwei-Milliarden-Euro-Investition des japanischen Softbank-Konzerns in Cruise, der Roboterwagen-Tochter von General Motors.

Lesezeit: 1 Min.　In Pocket speichern　　　　　9



(Bild: General Motors)

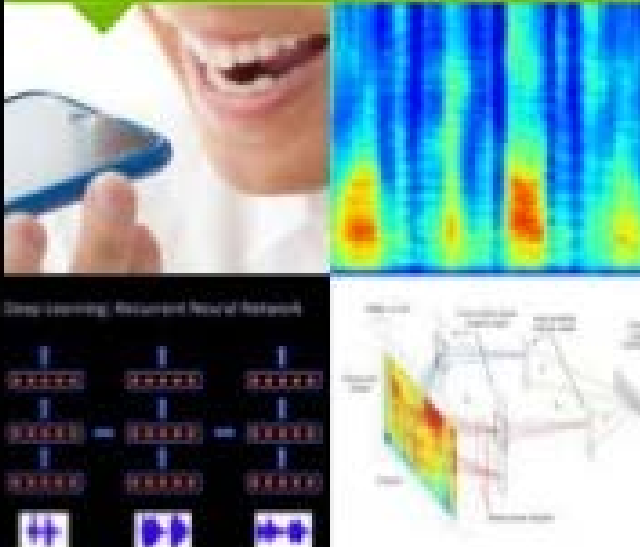07.07.2019　17:28 Uhr

Von Bernd Mewes

# AI Applications

# Goals for Today

- What we want to achieve today:
  - Theory of (Convolutional) Neural Networks
  - Arouse interest for further studies

- What we cannot achieve today:
  - Cover all details of Deep Learning
  - Go deep into Deep Learning for precision farming

# Supervised Learning

**Labelled Data**

Tom    Jerry    Tom

Tom    Jerry    Jerry

**Algorithm**

training

**Model**

**Unlabelled Data**

prediction

Jerry

**Predicted Labels**

Tom

# Unsupervised Learning

# Reinforcement Learning



reward

Environment

State

Agent

action

observation

# Categories of Machine Learning Algorithms

- Supervised learning
  - Given (training) data, which contains the correct answer for each dataset, the learning algorithm tries to find a hypothesis (model) that allows to predict the outcome for unseen datasets

- Unsupervised learning
  - The learning algorithm finds structure in the given data based on similarity and groups the data elements into clusters
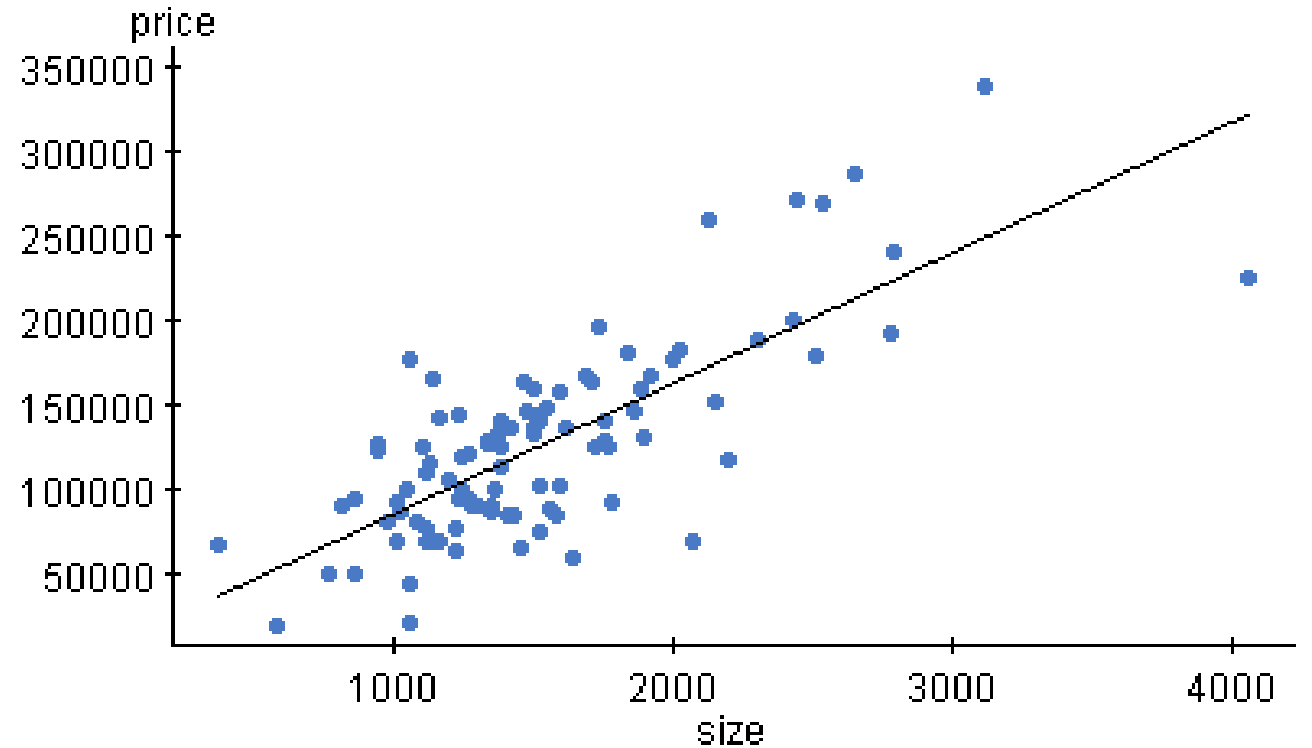
- Reinforcement learning
  - The learning algorithm learns from rewards of previous decisions

# Regression



- Learn a model by fitting a (straight) line through all (training) examples
- Predict the outcome for an unseen dataset by substituting the input values into the model
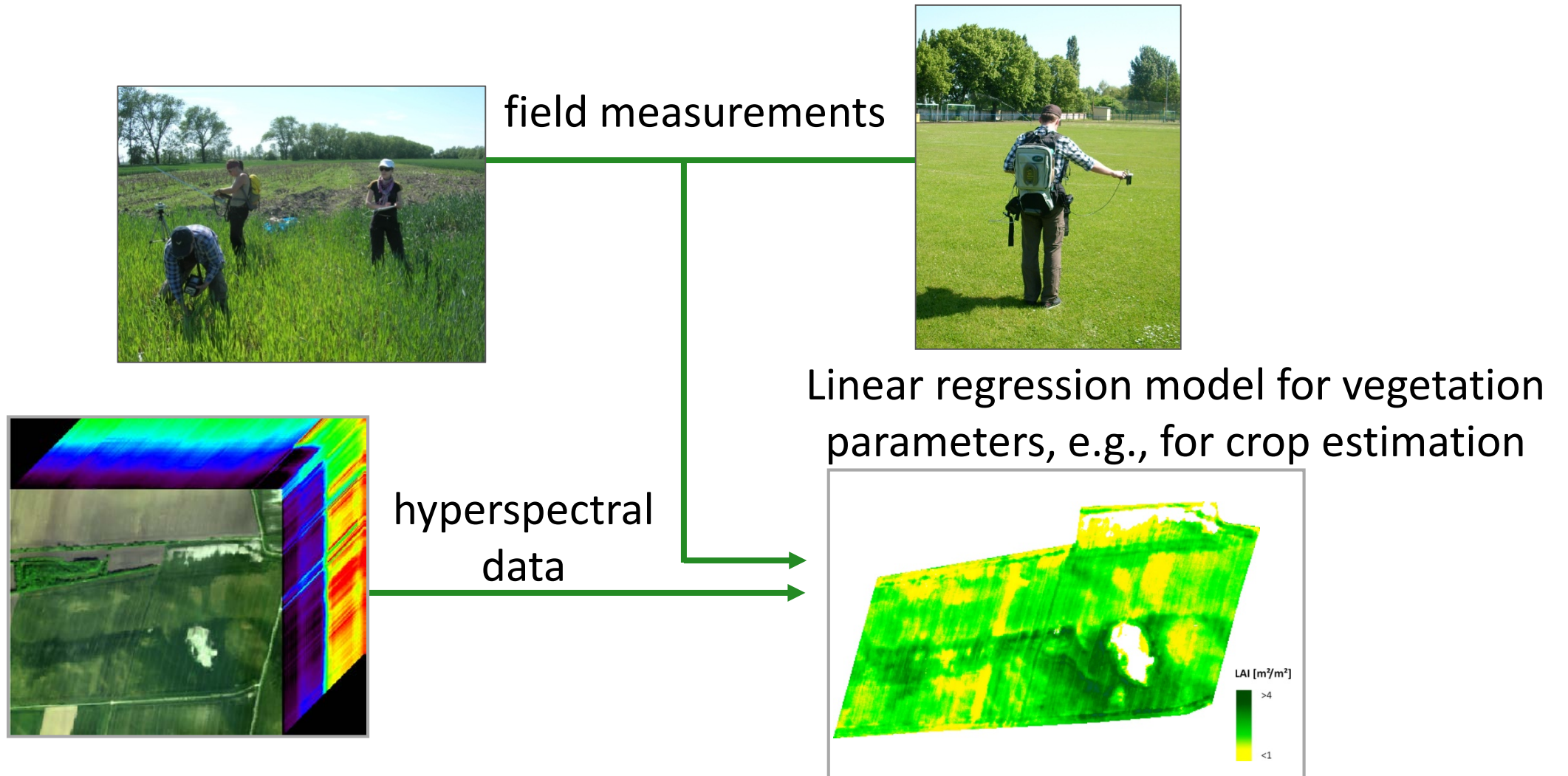
# Linear Regression

- Linear model that predicts a target value $\hat{y}$ by computing a weighted sum of the input features $(x_1, x_2, \ldots, x_n)$ plus a bias term $b$

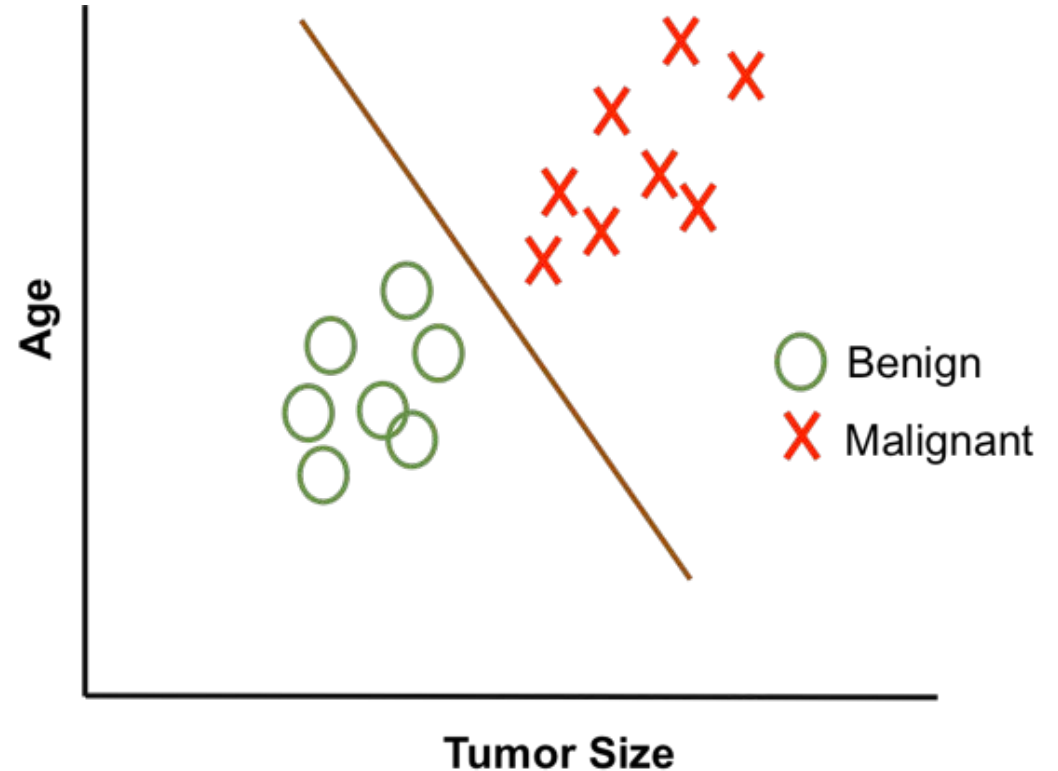$$\hat{y} = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

$$\hat{y} = \sum_{i=1}^{n} w_i x_i + b$$

# Linear Regression for Crop Estimation



field measurements

hyperspectral data

Linear regression model for vegetation parameters, e.g., for crop estimation
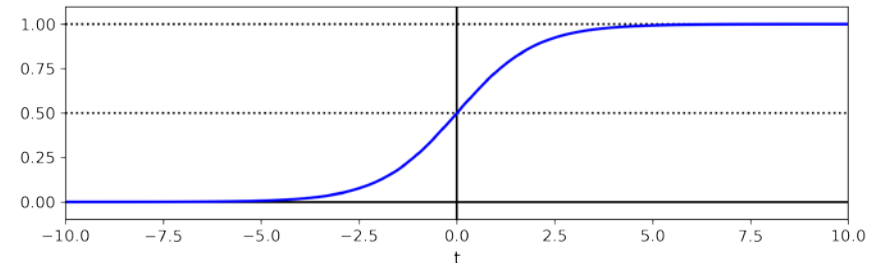
LAI [m²/m²]

>4

<1

# Classification

- Learn a model by finding a (straight) line that separates the (two) classes
- Predict the class by determining in which region your unseen input dataset lies
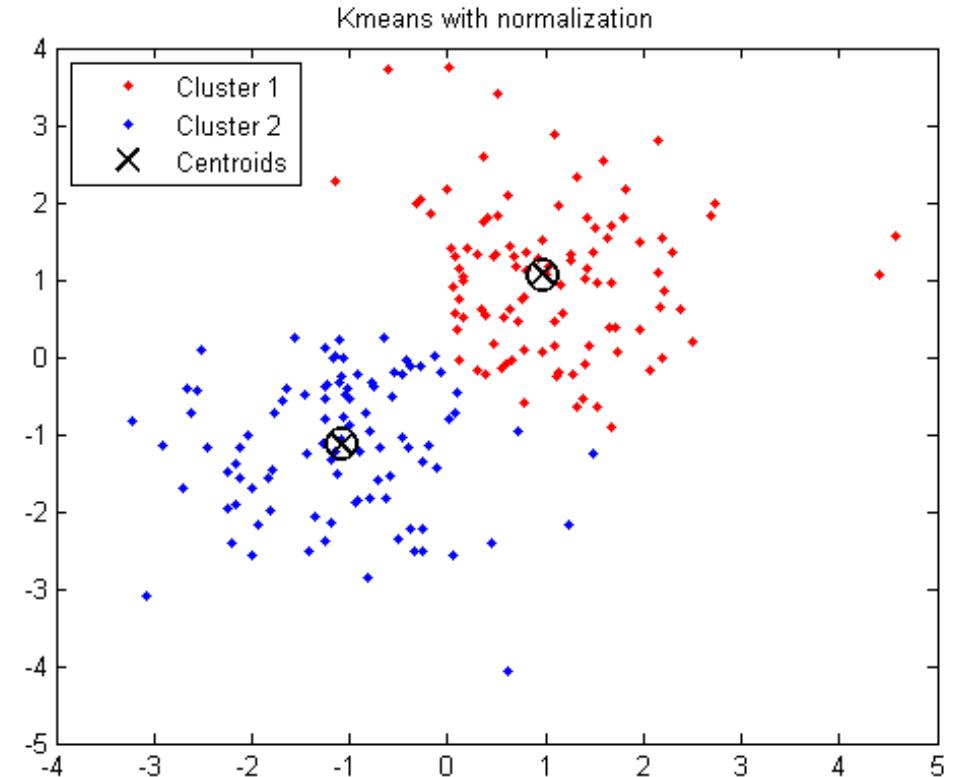
# Logistic Regression

- Computes a probability that the object with the given input features belongs to the class (or does not belong to the class) by using the sigmoid logistic function

$$\hat{p} = \sigma\left(\sum_{i=1}^{n} w_i x_i + b\right)$$
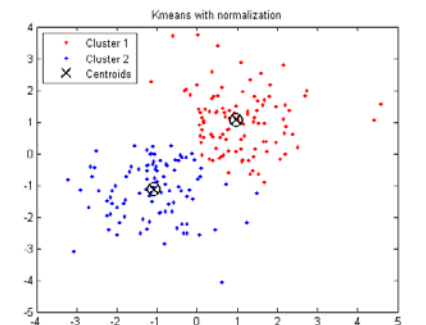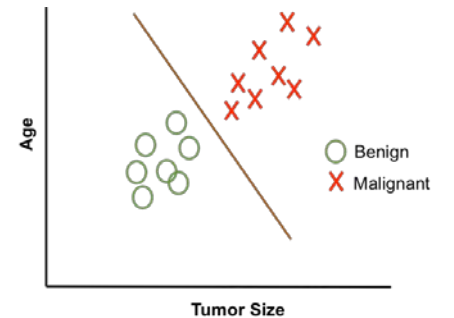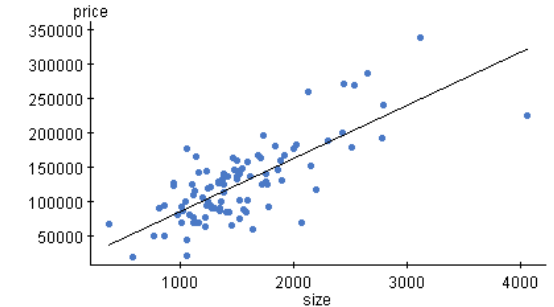
$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

# Segmentation



- Learn the structure of data by grouping similar examples into a set of clusters
- Predict the properties of an unseen dataset by its closeness to a cluster

# Categories of Machine Learning Problems

- Regression – supervised learning problem where the answer to be learned is a continuous value



- Classification – supervised learning problem where the answer is discreet (one of finitely many) values



- Segmentation – unsupervised learning problem where the structure to be learned is a set of clusters of similar examples
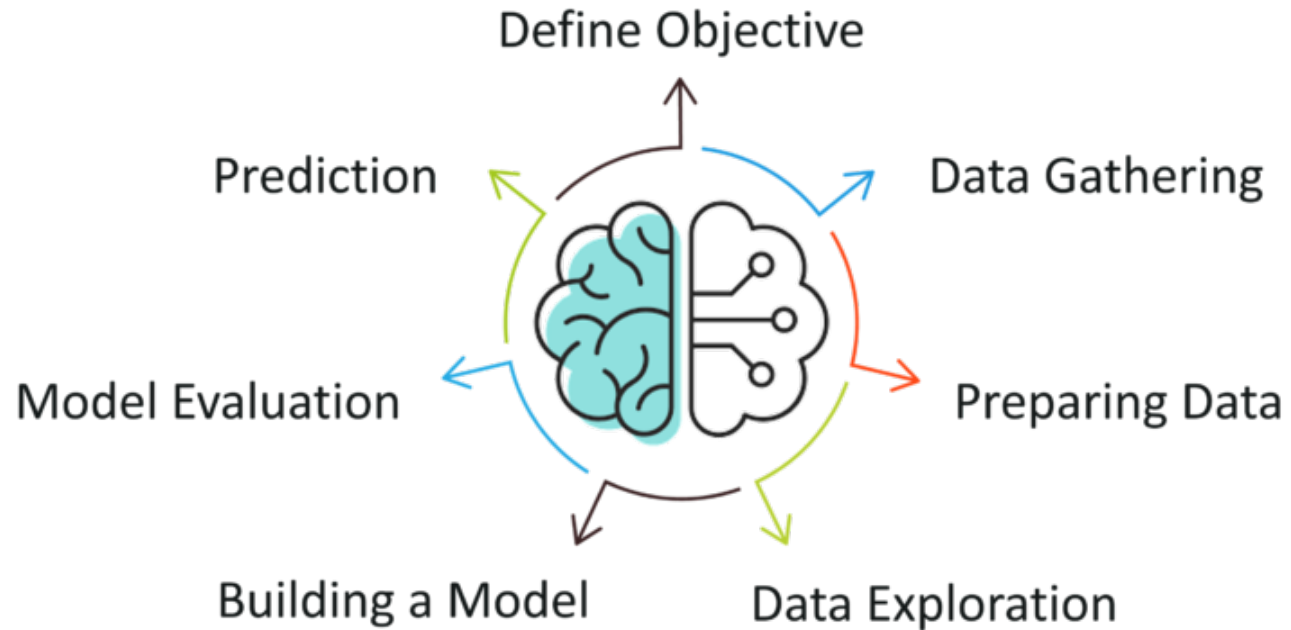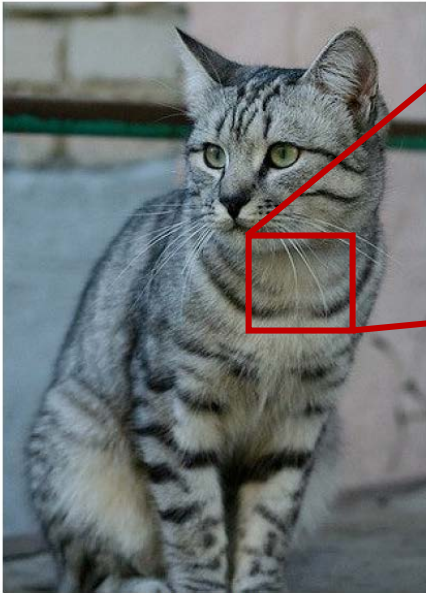
# Machine Learning Process



Define Objective

Data Gathering

Preparing Data

Data Exploration

Building a Model

Model Evaluation

Prediction

# Image Classification



```
[[105 112 108 111 104  99 106  99  96 103 112 119 104  97  93  87]
 [ 91  98 102 106 104  79  98 103  99 105 123 136 110 105  94  85]
 [ 76  85  90 105 128 105  87  96  95  99 115 112 106 103  99  85]
 [ 99  81  81  93 120 131 127 100  95  98 102  99  96  93 101  94]
 [106  91  61  64  69  91  88  85 101 107 109  98  75  84  96  95]
 [114 108  85  55  55  69  64  54  64  87 112 129  98  74  84  91]
 [133 137 147 103  65  81  80  65  52  54  74  84 102  93  85  82]
 [128 137 144 140 109  95  86  70  62  65  63  63  60  73  86 101]
 [125 133 148 137 119 121 117  94  65  79  80  65  54  64  72  98]
 [127 125 131 147 133 127 126 131 111  96  89  75  61  64  72  84]
 [115 114 109 123 150 148 131 118 113 109 100  92  74  65  72  78]
 [ 89  93  90  97 108 147 131 118 113 114 113 109 106  95  77  80]
 [ 63  77  86  81  77  79 102 123 117 115 117 125 125 130 115  87]
 [ 62  65  82  89  78  71  80 101 124 126 119 101 107 114 131 119]
 [ 63  65  75  88  89  71  62  81 120 138 135 105  81  98 110 118]
 [ 87  65  71  87 106  95  69  45  76 130 126 107  92  94 105 112]
 [118  97  82  86 117 123 116  66  41  51  95  93  89  95 102 107]
 [164 146 112  80  82 120 124 104  76  48  45  66  88 101 102 109]
 [157 170 157 120  93  86 114 132 112  97  69  55  70  82  99  94]
 [130 128 134 161 139 100 109 118 121 134 114  87  65  53  69  86]
 [128 112  96 117 150 144 120 115 104 107 102  93  87  81  72  79]
 [123 107  96  86  83 112 153 149 122 109 104  75  80 107 112  99]
 [122 121 102  80  82  86  94 117 145 148 153 102  58  78  92 107]
 [122 164 148 103  71  56  78  83  93 103 119 139 102  61  69  84]]
```
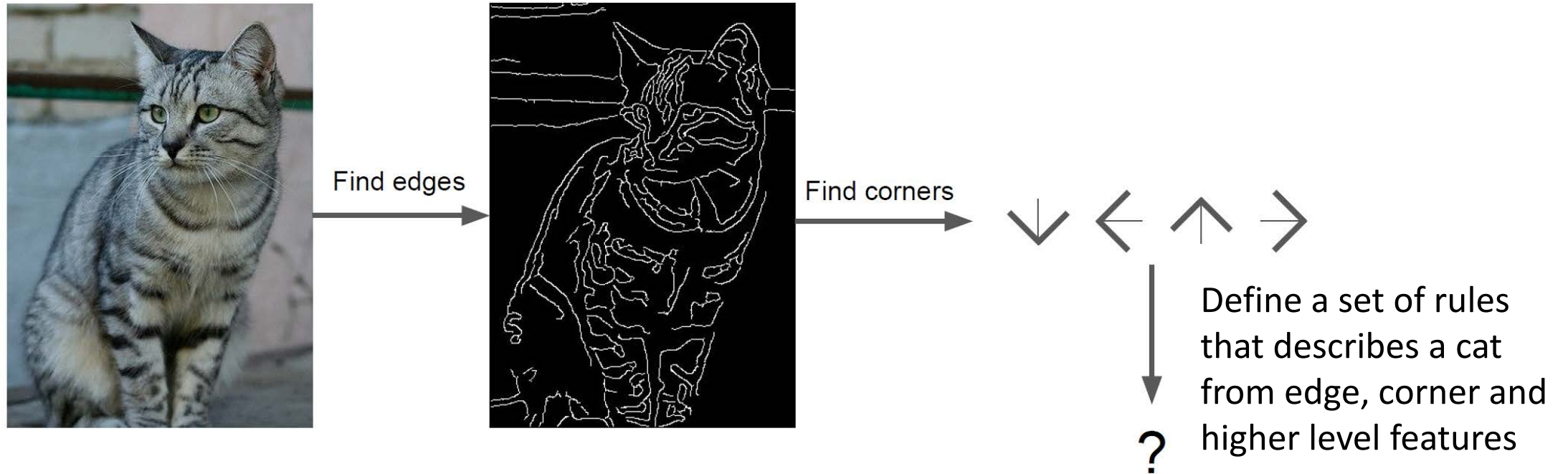
(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

→ cat

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

# Classical Approaches to Image Classification



Find edges

Find corners

Define a set of rules that describes a cat from edge, corner and higher level features

?

John Canny, "A Computational Approach to Edge Detection", IEEE TPAMI 1986

# Data-Driven Approach to Image Classification

1. Collect a (huge) dataset of labelled images

2. Train a classifier using machine learning techniques

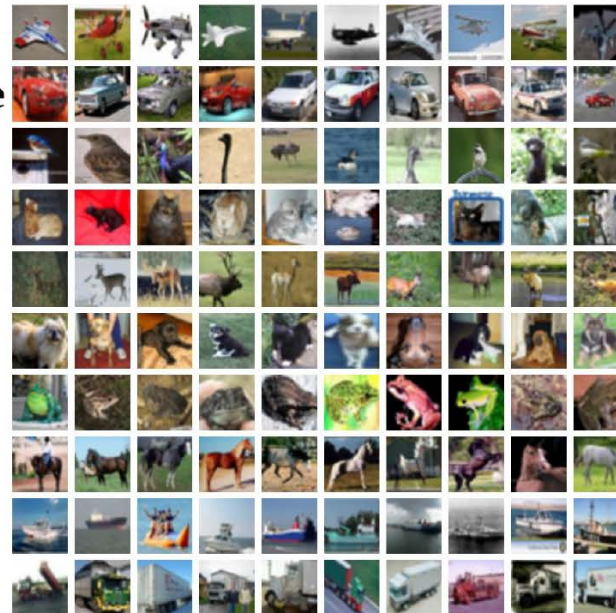3. Predict the class with the trained classifier

**CIFAR10**



**50,000** training images
each image is **32x32x3**

**10,000** test images.

# Neuron

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho
is licensed under CC-BY 3.0

$x_0$

$w_0$

synapse

axon from a neuron

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

$$\sum_i w_i x_i + b$$

$f$

$$f\left(\sum_i w_i x_i + b\right)$$

output axon

activation function

$w_2 x_2$

sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

# Perceptron

- Supervised learning algorithm for a binary classifier
- Takes vector data $\mathbf{x}$ as input and computes a single output value $y$

input



Input image

output

Probability that the input image shows a cat

# Mark I Perceptron

- First implementation of the perceptron algorithm by Frank Rosenblatt ~ 1957

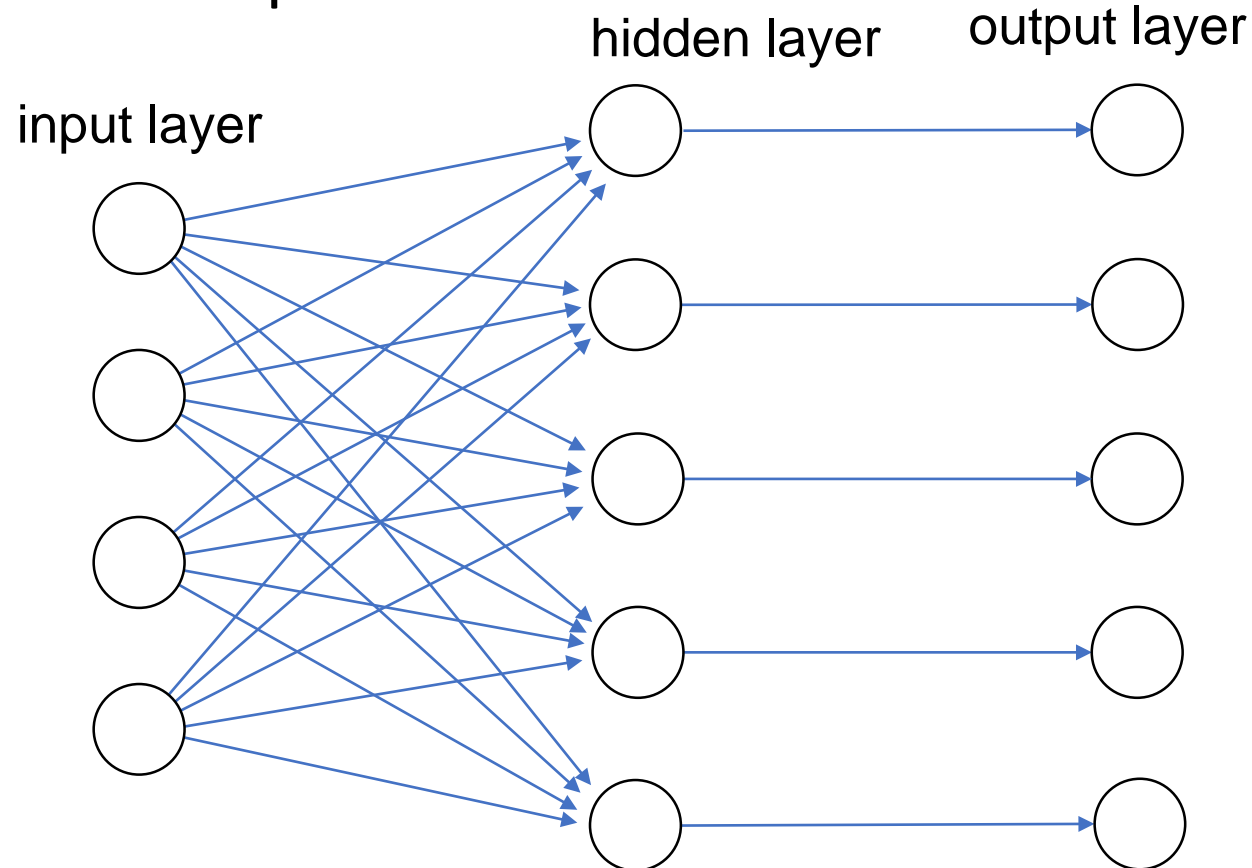- 20 × 20 cadmium sulfide photocells (400 pixel image)

$$f(x) = \begin{cases} 1 & if \sum_{i=1}^{n} w_i x_i + b > 0 \\ 0 & otherwise \end{cases}$$

- Recognizes letters of the alphabet

# Multilayer Perceptron

- Many perceptron are grouped so that the output is a vector output instead of a scalar output value
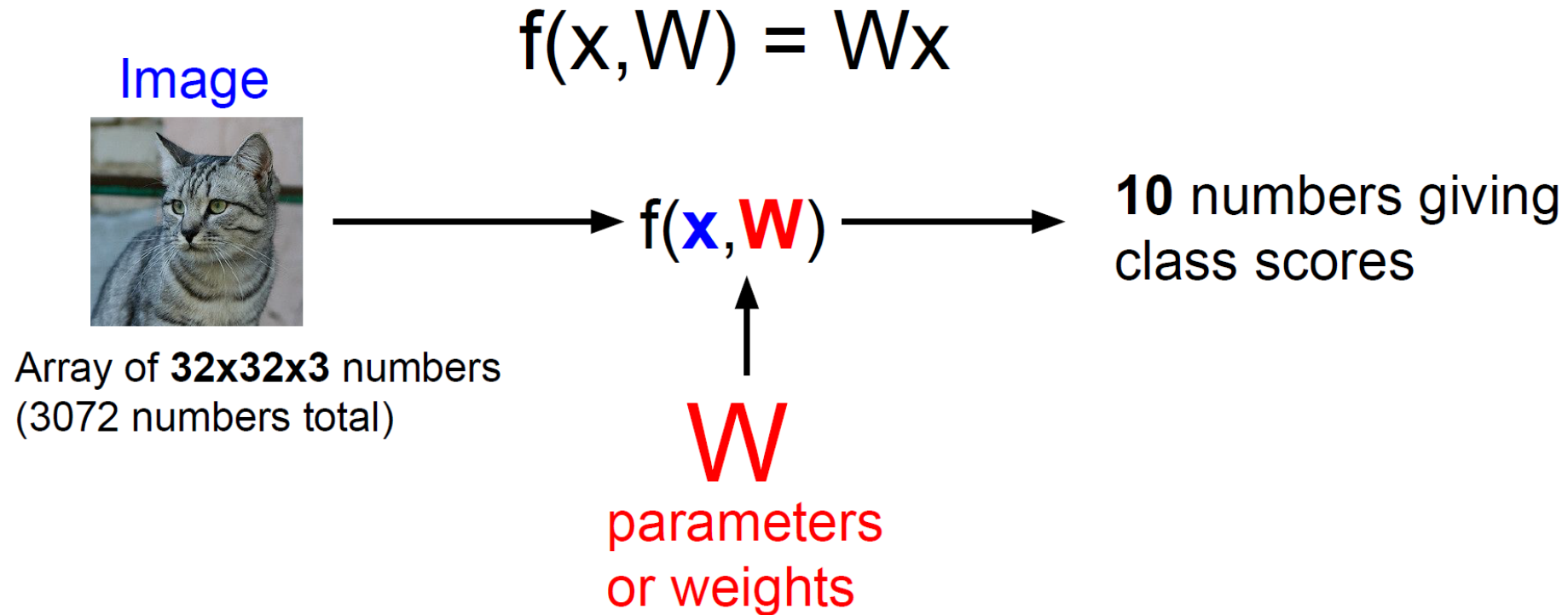
hidden layer    output layer

input layer

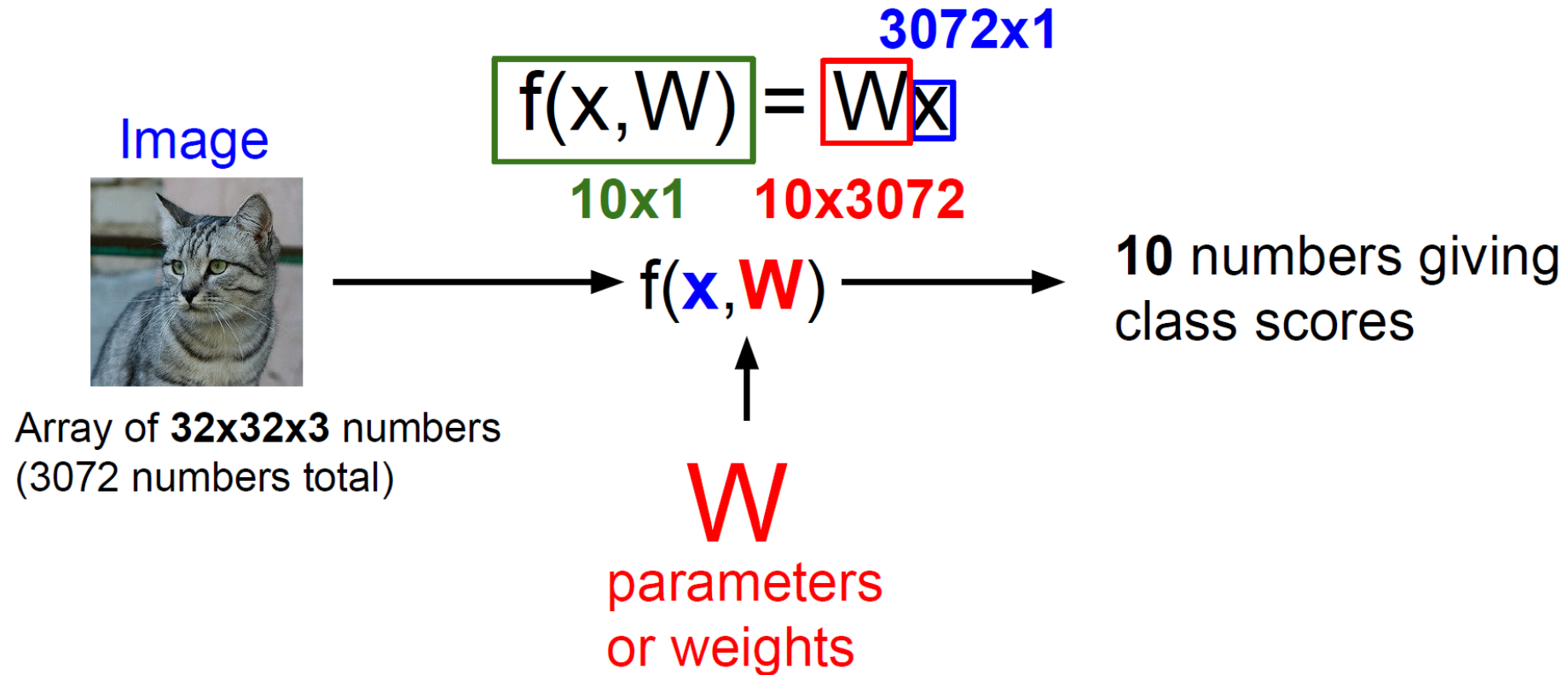each neuron of the output layer stands for a certain class

the output value of a neuron is the score for its class

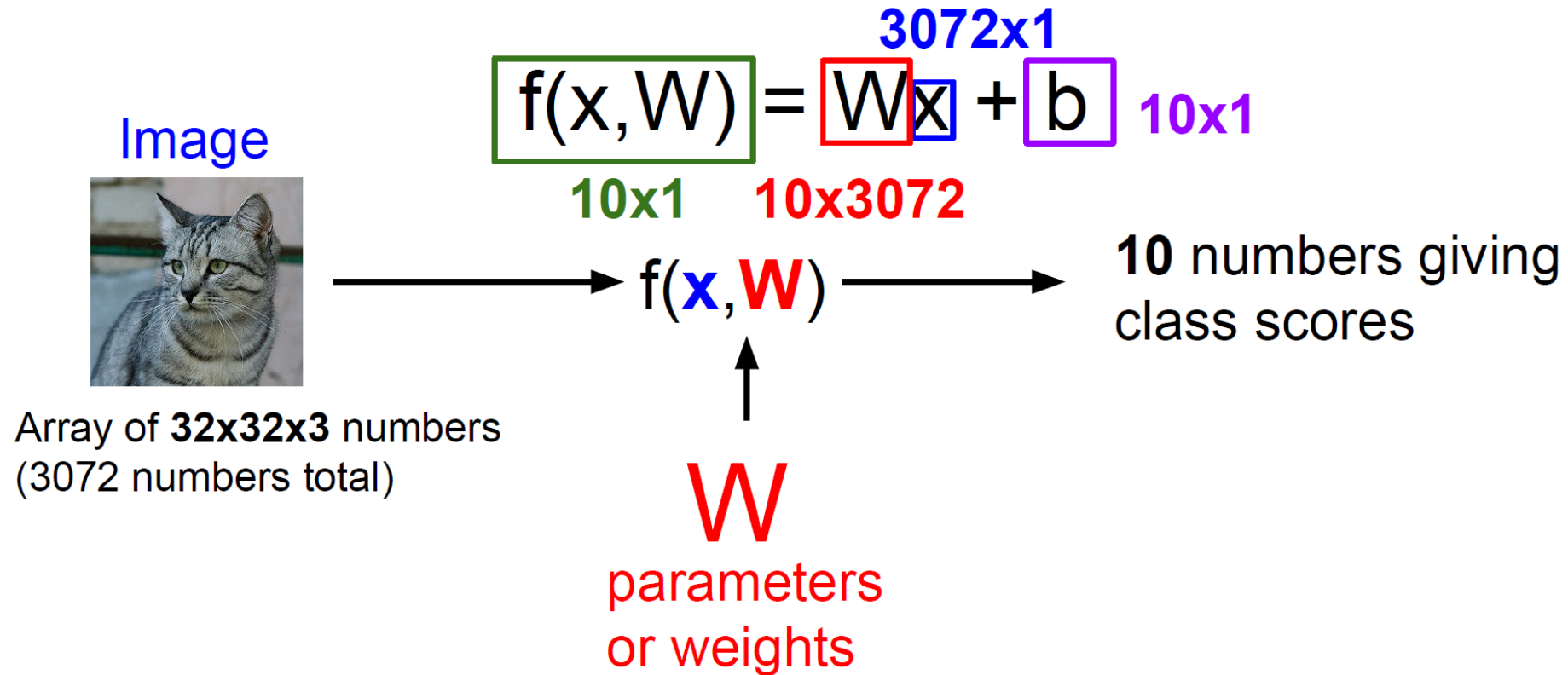the neuron with the highest score defines the predicted class

# Parametric Approach to Linear Classifier

$$f(x,W) = Wx$$

**Image**



Array of **32x32x3** numbers
(3072 numbers total)

$f(\textcolor{blue}{x},\textcolor{red}{W})$

$\textcolor{red}{W}$
parameters
or weights

**10** numbers giving
class scores

# Parametric Approach to Linear Classifier

$$f(x,W) = Wx$$

3072x1

10x1   10x3072

Image



Array of **32x32x3** numbers
(3072 numbers total)

f(**x**,**W**)

**10** numbers giving
class scores

W
parameters
or weights

# Parametric Approach to Linear Classifier

$$f(x,W) = Wx + b$$

3072x1

10x1    10x3072    10x1

Image



Array of **32x32x3** numbers
(3072 numbers total)

$f(\mathbf{x},\mathbf{W})$

**W**
parameters
or weights

**10** numbers giving
class scores

# Parametric Approach to Linear Classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Stretch pixels into column

| | | | |
|---|---|---|---|
| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

W

Input image

56, 231, 24, 2

| 56 |
|---|
| 231 |
| 24 |
| 2 |

x

+

| 1.1 |
|---|
| 3.2 |
| -1.2 |

b

=

| -96.8 | Cat score |
|---|---|
| 437.9 | Dog score |
| 61.95 | Ship score |

# Interpreting a Linear Classifier



$$f(x,W) = Wx + b$$

Example trained weights of a linear classifier trained on CIFAR-10:

# Interpreting a Linear Classifier

# Interpreting a Linear Classifier



$$f(x,W) = Wx + b$$

Array of **32x32x3** numbers
(3072 numbers total)

car classifier

airplane classifier

deer classifier

0

# German Traffic Sign Recognition Benchmark

- Single-image, multiple classes
- More than 40 classes
- More than 50,000 images
- Best recognition rates:
  1. 99.46% (committee of CNNs)
  2. 98.84% (human performance)
  3. 98.31% (multi-scale CNNs)

# Training
# Neural Networks

Prof. Dr. Martin Kada

# Training Neural Networks

- Initialize the weights of the network

- Evaluate how good the network is

- (Stepwise) improve the network
  - Gradient descent
  - Backpropagation
  - Learning rate

- Activation functions

# Weight Initialization

- Initialize all values of weight matrix $W$ with random gaussian noise with zero mean and a user-defined (e.g. 0.01) variance
  - Works only good for shallow networks
  - Weights initialized too small,
    then the signal shrinks as it passes through each layer until it vanishes
  - Weights initialized too large,
    then the signal grows as is passes through each layer until it explodes

- Xavier initialization
  - Makes sure the weights are just right, keeping the signal in a reasonable range of values through many layers

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

# Loss Function

- Quantifies how good the model is at the intended task

# Softmax Loss Function

- Generalization of the logistic function to multiple classes

- Assumption: scores are unnormalized log probabilities of the classes

| 2. normalize to get probabilities | 1. take the exponential of the scores |
|---|---|

- Minimize the loss (for a given image $i$) w.r.t. the correct class $k$
  - Categorical Cross-Entropy loss (also called Softmax Loss) is a Softmax activation plus a Cross-Entropy loss

$$L_i = -\log\left(\frac{e^{(s_k(\mathbf{x}))}}{\sum_{j=1}^{K} e^{(s_j(\mathbf{x}))}}\right)$$

# Softmax Loss Function

- Quantifies how good the model is at the intended task



$$L_i = -\log\left(\frac{e^{(s_k(\mathbf{x}))}}{\sum_{j=1}^{K} e^{(s_j(\mathbf{x}))}}\right)$$

# Softmax Loss Function

- Quantifies how good the model is at the intended task



input layer

hidden layer

output layer

N images

loss value $L$

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i)$$

# Softmax Loss Function



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

unnormalized probabilities

|       |       |          |        |          |       |
|-------|-------|----------|--------|----------|-------|
| cat   | **3.2**  |          | **24.5**  |          | **0.13** |
| car   | 5.1   | exp →    | 164.0  | normalize → | 0.87  |
| frog  | -1.7  |          | 0.18   |          | 0.00  |

L_i = -log(0.13)
    = **0.89**

unnormalized log probabilities          probabilities

# Gradient Descent

- Finds optimal weights by iteratively tweaking the model parameters $W$ in order to minimize the cost function $L$

- Idea:
  - Go downhill in the direction of the steepest slope until you reach a valley
  - Measure the local gradient of the cost function with regard to the parameter vector $W$ and go in the direction of descending gradient until a minimum is reached

# Gradient Descent

- Algorithm:
  - Initialize $W$ with random values (random initialization)
  - Gradually improve $W$ by <span style="color:red">backpropagation</span> to decrease the cost function
  - Stop when $W$ converges to a minimum

# Gradient Descent

- The loss $L$ is a function of $W$

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i)$$

- Calculate analytical $\boxed{\text{gradient vector } \nabla_w L}$ ⟵ backpropagation

- Update $W$ by computing $W' = W - \eta \cdot \nabla_w L$ where $\eta$ is the learning rate

# Gradient Descent

# Backpropagation

- Forward pass – compute the aggregated output of all neurons and the loss

# Backpropagation

- Backwards pass – update the weights $W$ w.r.t. to the loss

# Backpropagation

## Forward pass

$x$

$w$

$f(x, w)$

$z$

$z = x \cdot w$

$$\frac{dL}{dx} = \frac{dL}{dz}\frac{dz}{dx}$$

$df$

$$\frac{dL}{dz}$$

$$\frac{dL}{dw} = \frac{dL}{dz}\frac{dz}{dy}$$

$$\frac{dz}{dx} = w$$

$$\frac{dz}{dw} = x$$

# Backpropagation

Forward pass

$t \longrightarrow f(t) \longrightarrow z$

$\dfrac{dL}{dt} = \dfrac{dL}{dz}\dfrac{dz}{dt} \longleftarrow df \longleftarrow \dfrac{dL}{dz}$

$z = \sigma(t) = \dfrac{1}{1 + e^{-t}}$

$\dfrac{dz}{dt} = (1 - \sigma(t))\sigma(t)$

# Backpropagation



$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \quad \Big| \quad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \quad \Big| \quad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

# Learning Rate

- The hyperparameter learning rate $\eta$ determines how well the algorithm converges



If the learning rate is too small, then the algorithm will take many iterations to converge

If the learning rate is too high, then the algorithm might overjump the minimum, possibly ending up even further away from the minimum than before → algorithm diverges and fails to find a good solution

# Learning Rate

- Cost function might have valleys, ridges, plateaus and other irregular shapes → makes it difficult to converge to the minimum

- Challenges:
    - Getting stuck in a local minimum, which is not as good as the global minimum
    - Taking very long to cross a plateau and unwillingly stopping too early before the global minimum is reached

# Learning Rate

- How to get a good learning rate

# Optimization Algorithms

- Adapt the learning rate to find the global minimum
- Build up "velocity" to overcome local minima and plateaus

(Stochastic) Gradient Descent

$$W_{t+1} = W_t - \eta \cdot \nabla_{w_t} L$$

(Stochastic) Gradient Descent
+ Momentum

$$v_{t+1} = \rho v_t + \nabla_{w_t} L$$
$$W_{t+1} = W_t - \eta \cdot v_{t+1}$$

$\rho$ gives "friction" (e.g. 0.9 or 0.99)

# Optimization Algorithms

# Activation Functions

**Sigmoid**

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

**tanh**

$\tanh(x)$

**ReLU**

$\max(0, x)$

**Leaky ReLU**

$\max(0.1x, x)$

**ELU**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

# Convolutional Neural Networks

Prof. Dr. Martin Kada

# ImageNet



IM GENET                                    www.image-net.org

**22K** categories and **14M** images

- Animals
    - Bird
    - Fish
    - Mammal
    - Invertebrate
- Plants
    - Tree
    - Flower
- Food
- Materials
- Structures
- Artifact
    - Tools
    - Appliances
    - Structures
- Person
- Scenes
    - Indoor
    - Geological Formations
- Sport Activities

Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009

# ImageNet



IMAGENET **Large Scale Visual Recognition Challenge**

steel drum

The Image Classification Challenge:

1,000 object classes

1,431,167 images

**Output:**
Scale
T-shirt
Steel drum
Drumstick
Mud turtle

✔

**Output:**
Scale
T-shirt
Giant panda
Drumstick
Mud turtle

✘

Russakovsky et al. arXiv, 2014

# ImageNet



Russakovsky et al. arXiv, 2014

# AlexNet

# VGG

[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

8 layers (AlexNet)
-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13
(ZFNet)
-> 7.3% top 5 error in ILSVRC'14



AlexNet

VGG16

VGG19

# GoogLeNet

*[Szegedy et al., 2014]*

**Deeper networks, with computational efficiency**

- 22 layers
- Efficient "Inception" module
- No FC layers
- Only 5 million parameters!
  12x less than AlexNet
- ILSVRC'14 classification winner
  (6.7% top 5 error)



Inception module

Hint: only the convolutional layers and the fully connected layer (at the head of the network) that predicts the class scores are counted into the 22 layers

# ResNet

[He et al., 2015]

**Very deep networks using residual connections**

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!



Residual block

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$Wx$

10 x 3072
weights

**activation**

1

10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

# Convolution Layer

32x32x3 image -> preserve spatial structure

5x5x3 filter

32 height

32 width

3 depth

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer

7×7 input with 3 × 3 filter  ⇒  5×5 output

# Convolution Layer

32x32x3 image

Filters always extend the full depth of the input volume

32

32

3

5x5x3 filter

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer



32x32x3 image

5x5x3 filter $w$

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation map**

28

28

1

# Convolution Layer

consider a second, green filter



32x32x3 image
5x5x3 filter

activation maps

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Activation Maps

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



activation maps

32
32
3

Convolution Layer

28
28
6

We stack these up to get a "new image" of size 28x28x6!

# Convolutional Neural Network

- A Convolutional Neural Network (CNN) is a sequence of Convolutional Layers, interspersed with activation functions

# Convolutional Neural Network



VGG-16 Conv1_1     VGG-16 Conv3_2     VGG-16 Conv5_3

# Stride



32x32x3 image

5x5x3 filter

convolve (slide) over all spatial locations

**activation map**

# Stride



7x7 input (spatially)
assume 3x3 filter

applied **with stride 1**

=> **5x5 output**

# Stride



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

**=> 3x3 output!**

# Stride



7

7

7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

**doesn't fit!**
cannot apply 3x3 filter on
7x7 input with stride 3.

# Stride



N

F

F

N

Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

# Padding



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

(recall:)                              **7x7 output!**
(N - F) / stride + 1

in general, common to see CONV layers with
stride 1, filters of size FxF, and zero-padding with
(F-1)/2. (will preserve size spatially)
e.g. F = 3 => zero pad with 1
     F = 5 => zero pad with 2
     F = 7 => zero pad with 3

# Examples



Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Output volume size:
(32+2*2-5)/1+1 = 32 spatially, so
**32x32x10**

# Examples



Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?

each filter has 5*5*3 + 1 = 76 params          (+1 for bias)
=> 76*10 = **760**

# 1×1 Convolution



1x1 CONV
with 32 filters

(each filter has size 1x1x64, and performs a 64-dimensional dot product)

# Pooling Layer

- makes the representations smaller and more manageable
- operates over each activation map independently:

# Max Pooling

Single depth slice



max pool with 2x2 filters
and stride 2

# CNN for Image Classification

# What's going on inside ConvNets?

Input Image:
3 x 224 x 224

Class Scores:
1000 numbers

What are the intermediate features looking for?

Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.
Figure reproduced with permission.

# What's going on inside ConvNets?



AlexNet:
64 x 3 x 11 x 11

ResNet-18:
64 x 3 x 7 x 7

ResNet-101:
64 x 3 x 7 x 7

DenseNet-121:
64 x 3 x 7 x 7

Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016
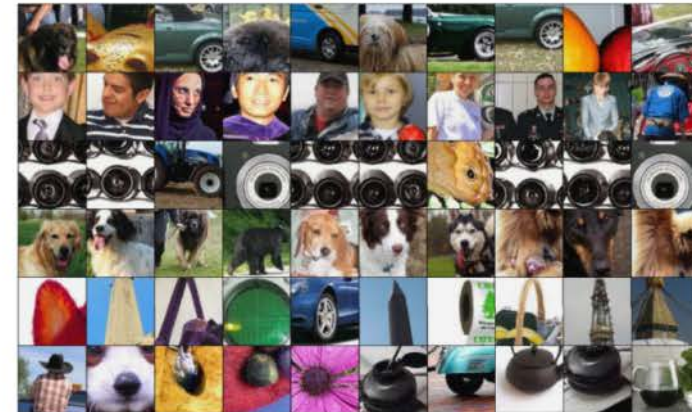Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# What's going on inside ConvNets?
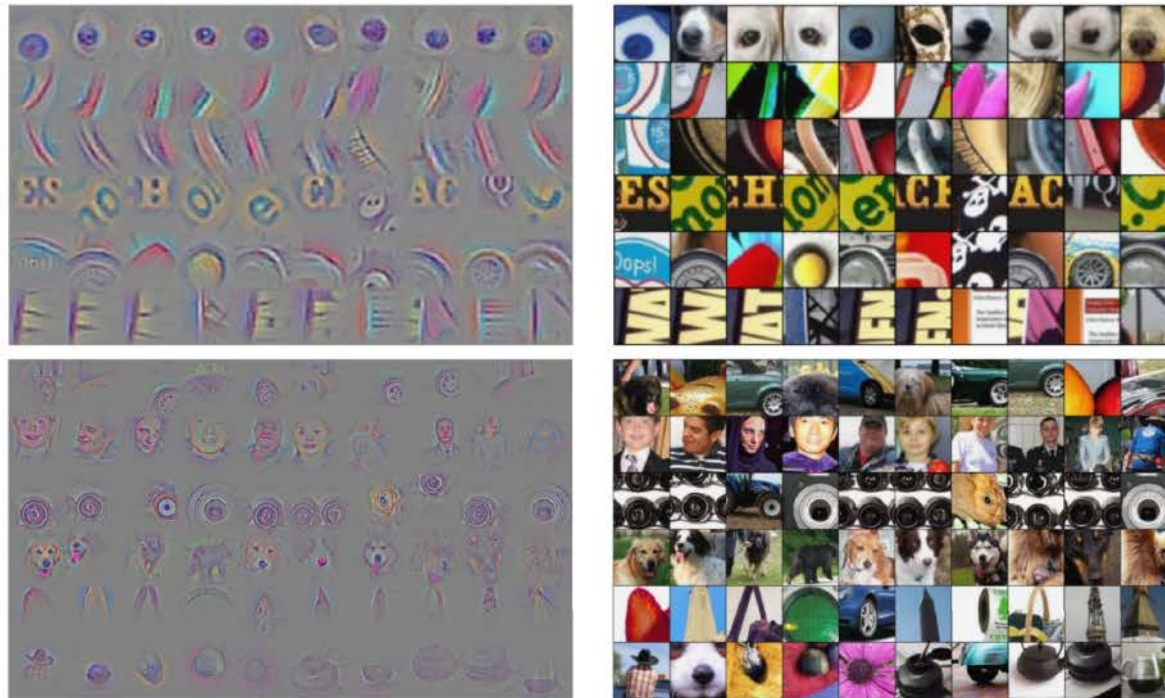


Pick a layer and a channel; e.g. conv5 is 128 x 13 x 13, pick channel 17/128

Run many images through the network, record values of chosen channel

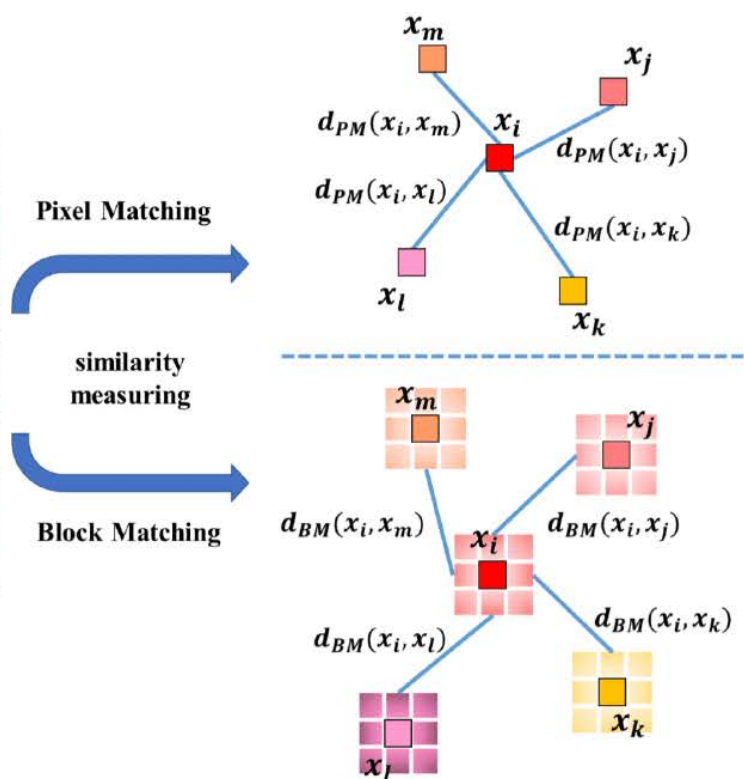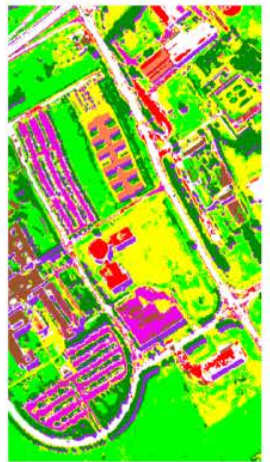Visualize image patches that correspond to maximal activations

# What's going on inside ConvNets?

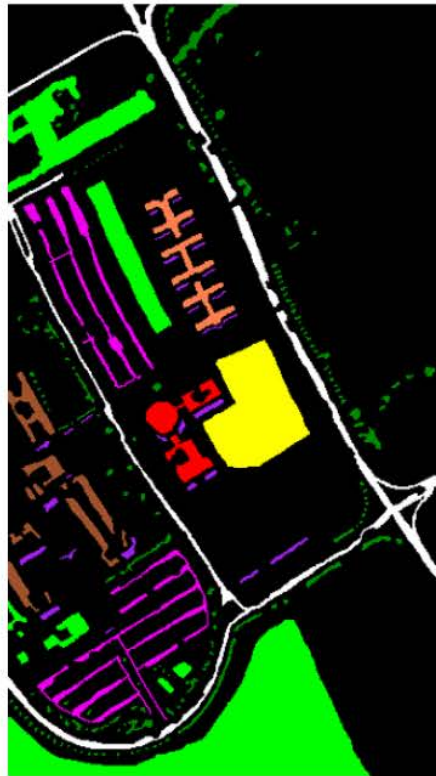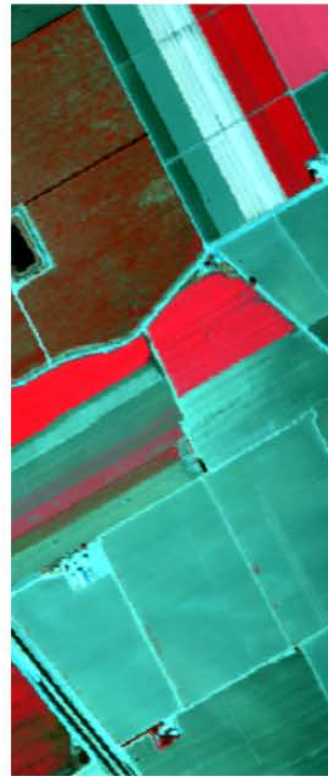# What to do now?

# What to do now?



**Class**
- No Ground Reference
- Asphalt
- Meadow
- Gravel
- Trees
- Painted Metal Sheets
- Bare Soil
- Bitumen
- Self-Blocking Bricks
- Shadows

**Class**
- No Ground Reference
- Brocoli_green_weeds_1
- Brocoli_green_weeds_2
- Fallow
- Fallow_rough_plow
- Fallow_smooth
- Stubble
- Celery
- Grapes_untrained
- Soil_vinyard_develop
- Corn_senesced_green_weeds
- Lettuce_romaine_4wk
- Lettuce_romaine_5wk
- Lettuce_romaine_6wk
- Lettuce_romaine_7wk
- Vinyard_untrained
- Vinyard_vertical_trellis

# What to do now?

- Stanford course CS231n on
  "Convolutional Neural Networks for Visual Recognition"
  - PDF lecture presentation & YouTube lecture videos
    http://cs231n.stanford.edu/

- Deep Learning Book by Goodfellow, Bengio, Courville

- Machine Learning book by Géron