

Artificial Intelligence

V08: Supervised Learning using Neural Networks



Introduction to learning agents
Artificial Neural Networks
Doing machine learning

Based on material by

- Stuart Russell, UC Berkeley
- Andreas Krause, ETH Zurich



Educational objectives



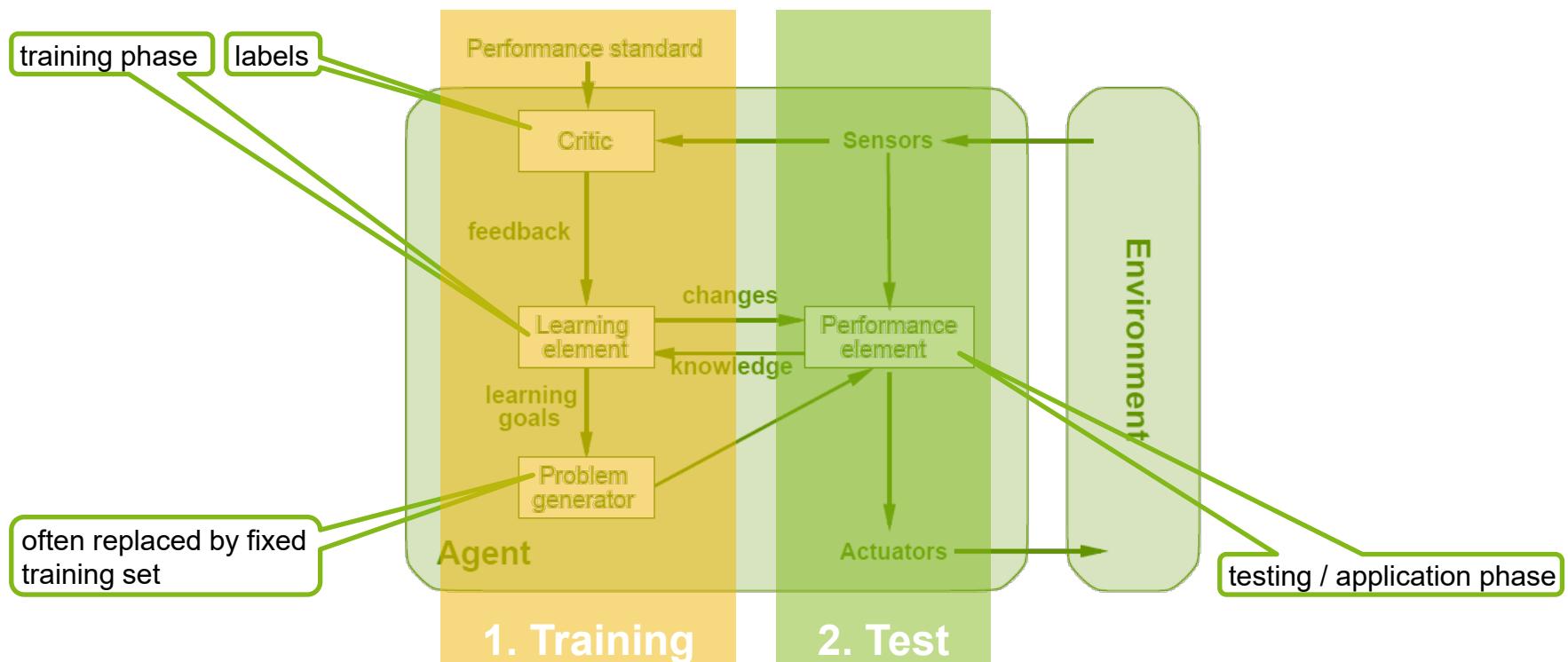
- **Remember** the correct **technical terms** to explain **machine learning**
- Have a **basic understanding** of the architecture and working **of neural networks**
- **Defend** your **own view** on the existence of good **general learners**



“In which we describe agents that can improve their behavior through diligent study of past experiences and predictions about the future.”

➔ Reading: AIMA, ch. 19-19.6

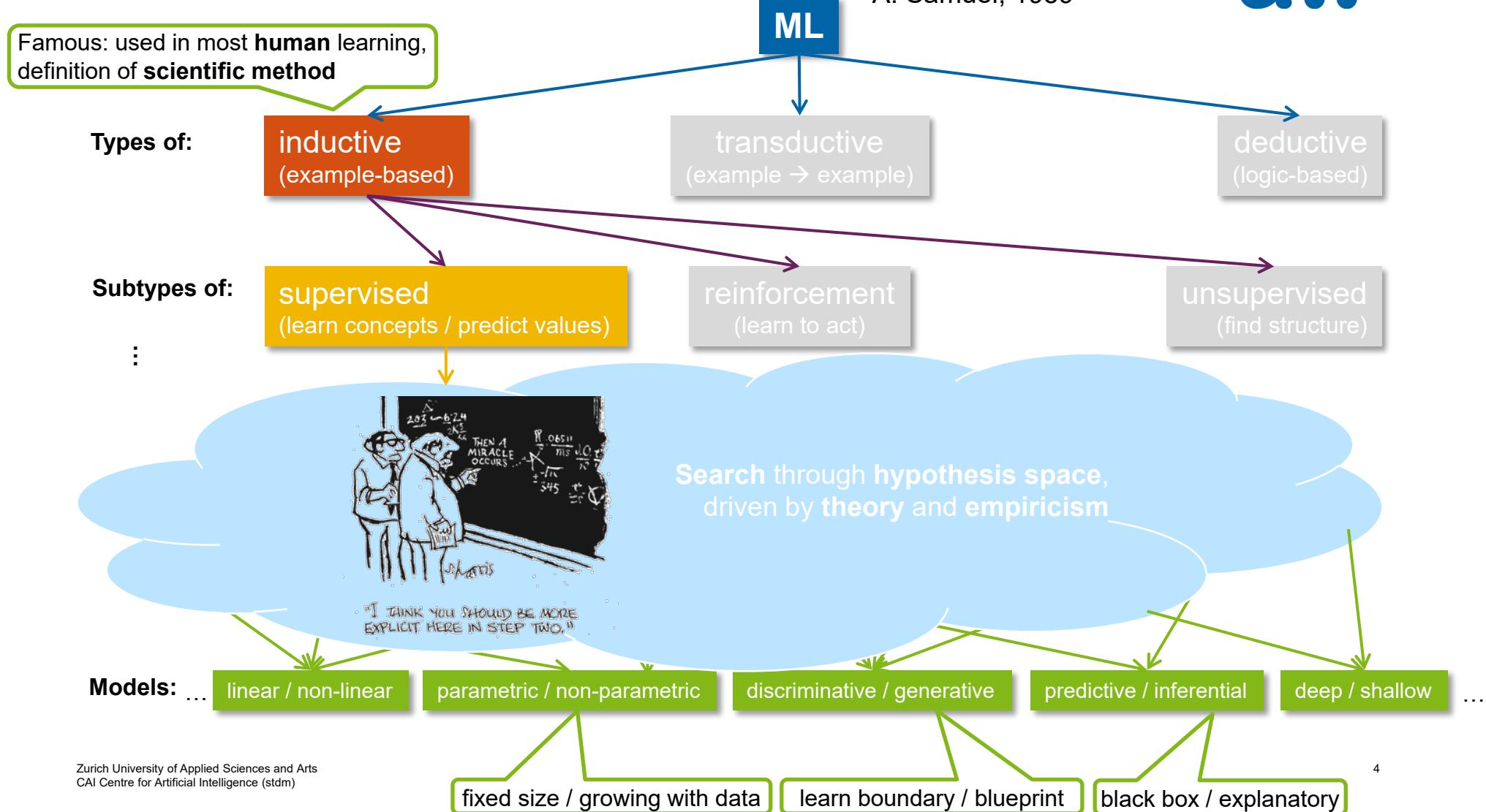
1. INTRODUCTION TO LEARNING AGENTS



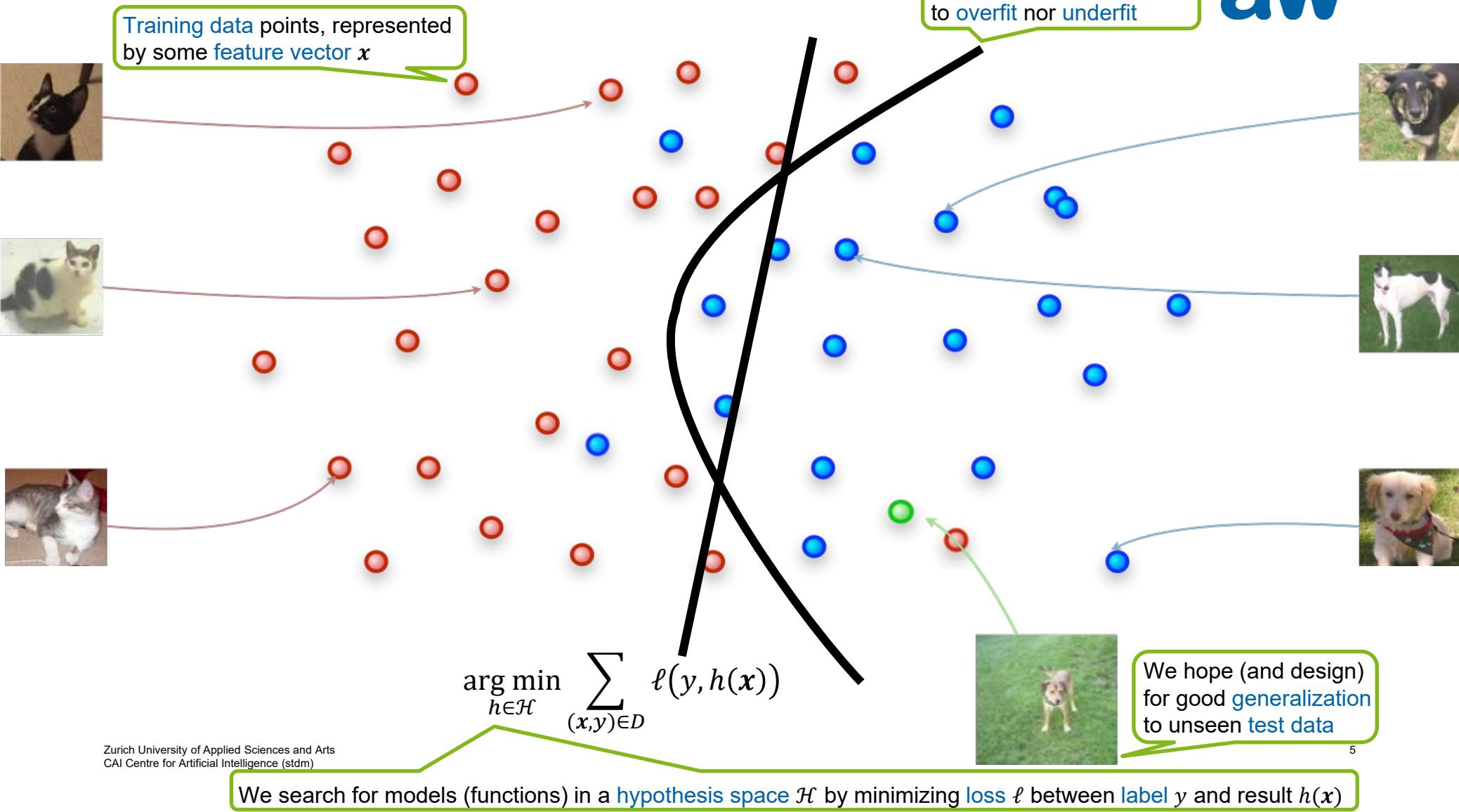
The discipline of machine learning – mapped

«...gives computers the ability to [act] without being explicitly programmed.»

A. Samuel, 1959

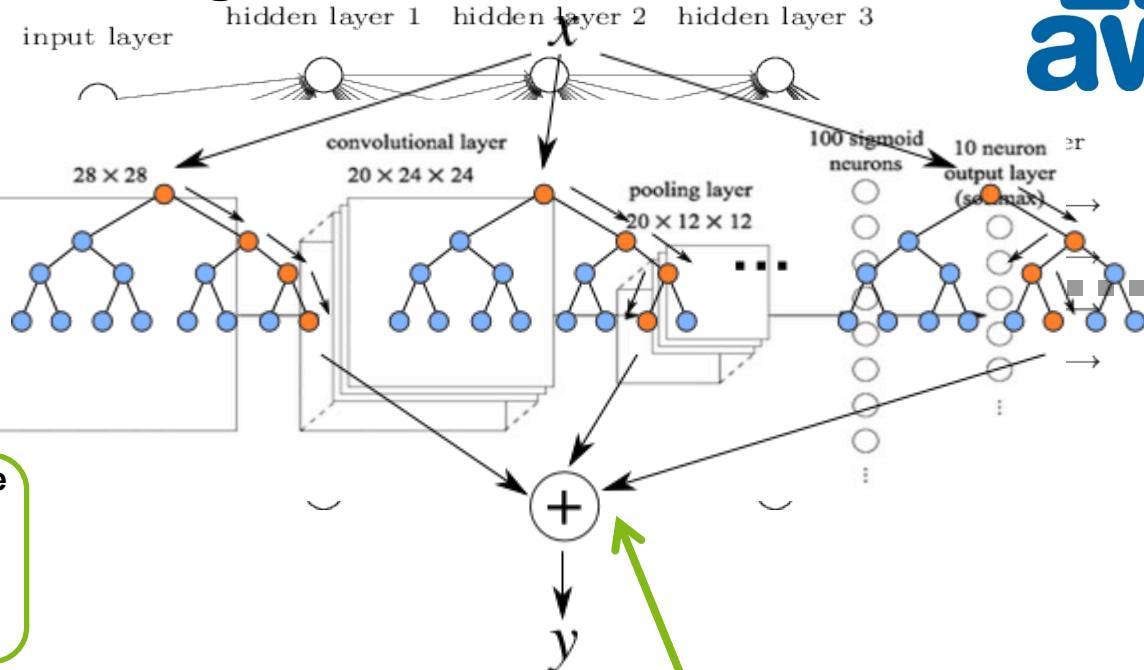


Supervised machine learning in a nutshell



Learning as search through \mathcal{H}

$$\mathcal{H} = \{$$



Success is largely determined by **choosing the correct hypothesis space** for the problem:

- Linear? Polynomial?
- Deep neural network? CNN?
- Ensemble of decision trees? ...

$$h(x) = h(x, w)$$

A **good model** complies with **Ockham's razor**: Maximize a combination of **consistency** and **simplicity**

Learning then means finding good **parameters** (sometimes called θ)

Shallow vs. deep learning

Add depth (layers → capability) to learn features automatically

Classic computer vision

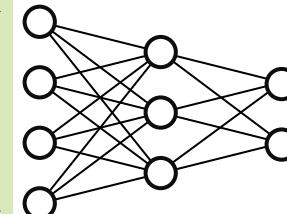


Feature extraction
(SIFT, SURF, LBP, HOG, etc.)

(0.2, 0.4, ...)

(0.4, 0.3, ...)

Classification
(SVM, neuronal net, etc.)



container ship

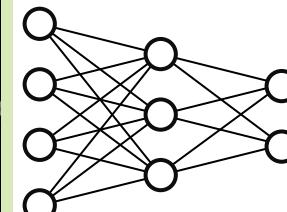
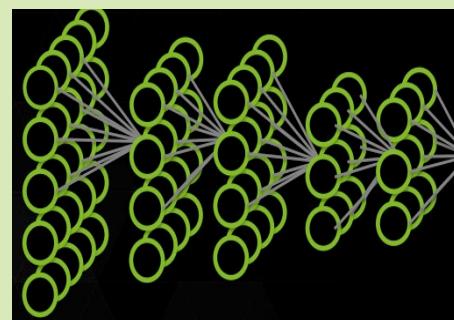
tiger

...

Convolutional neural networks (CNNs)



Takes raw pixels as input, learns good features automatically!



container ship

tiger

...

Why study machine learning in general?

«A learner that makes ***no a priori assumptions*** regarding the identity of the target concept has ***no rational basis for classifying*** any unseen instances»

[Mitchell, 1997, ch. 2.7.3]

There's no single best algorithm

- **No free lunch theorem** (NFL) regarding the general equivalence of learners [Wolpert, 1996]:
When all hypotheses h are equally likely, the probability of observing an arbitrary sequence of cost values during training does not depend upon the learning algorithm \mathcal{L}
→ there's **no universally best learner** (across problems)
- Empirical study [Caruana et al., 2006]:
«Even the best models sometimes perform poorly, and models with poor average performance occasionally perform exceptionally well»
→ All learning algorithms **have** advantages & **disadvantages**, depending on the current data



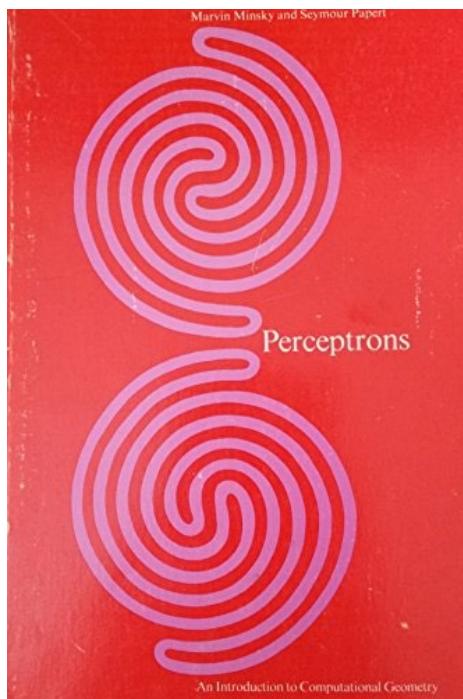
Examples of sensor data for pattern recognition tasks («**Labeled faces in the wild**» dataset) and tabular data («**Iris**» dataset)

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

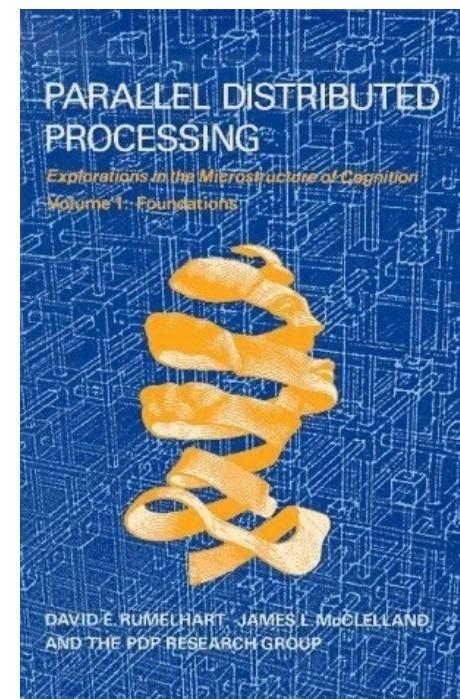
Ascertainment from kaggle.com

- **Tabular** data: do **handcrafted** feature engineering, followed by an **ensemble of decision trees**
 - **Sensor** data (images, speech, ...): use a **suitable** deep neural network
- See <https://www.import.io/post/how-to-win-a-kaggle-competition/>

2. ARTIFICIAL NEURAL NETWORKS



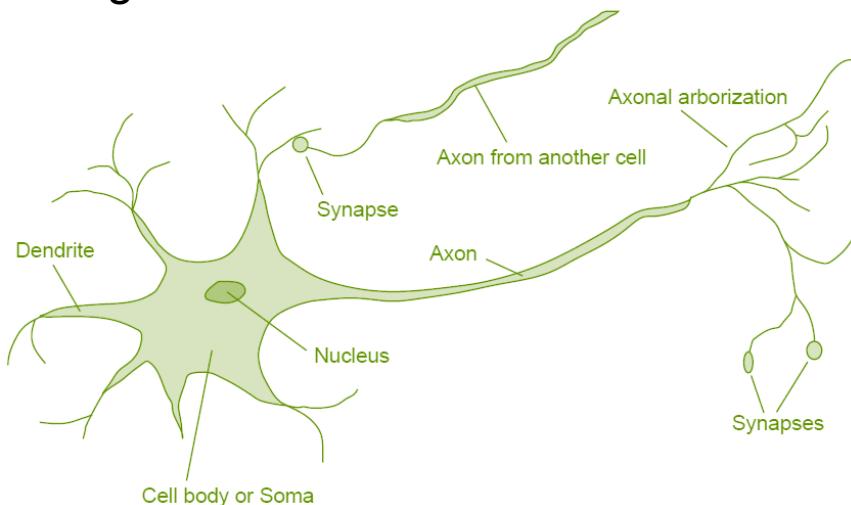
The critique of the first neural network hype:
Minsky & Papert, "Perceptrons", MIT Press Ltd, 1969



...and its resurrection:
Rumelhart et al., "Parallel Distributed Processing, Volume 1", MIT Press Ltd, 1987

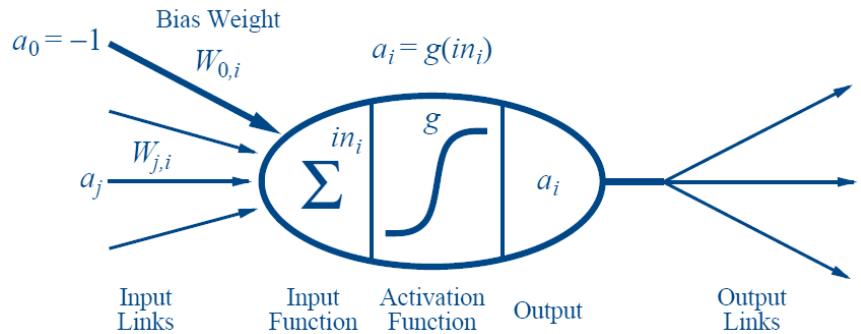
Neurons

Biological model



- 10^{11} neurons of > 20 types
- 10^{14} synapses
- 1ms – 10ms cycle time
- Signals are noisy “spike trains” of electrical potential
- Organized in layers to form a brain

McCulloch-Pitts “unit” (1943)

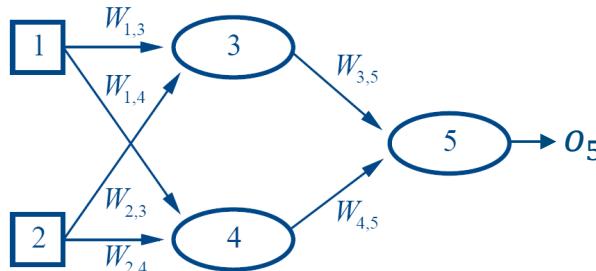


- Output is a **thresholded linear function** of the inputs: $a_i = g(in_i) = g(\sum_j W_{j,i} \cdot a_j)$
- Changing the **bias weight** $W_{0,i}$ moves the threshold location
- A **gross oversimplification** of real neurons!
- Purpose: develop understanding of what **networks of simple units** can do

Feed-forward network example

FNN: a parameterized family of nonlinear functions $g()$

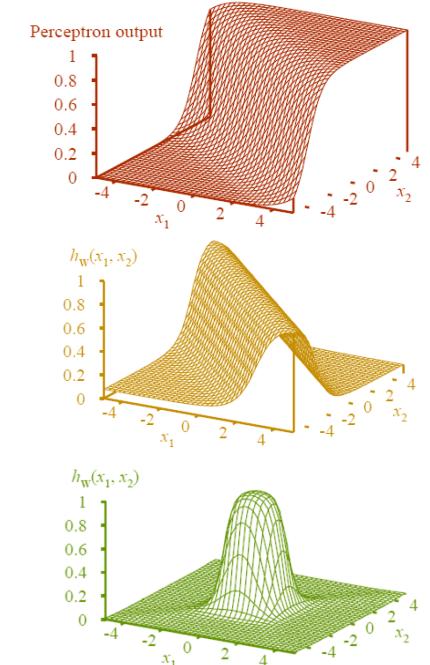
- $$\begin{aligned} o_5 &= g(W_{3,5} \cdot o_3 + W_{4,5} \cdot o_4) \\ &= g\left(W_{3,5} \cdot g(W_{1,3} \cdot o_1 + W_{2,3} \cdot o_2) + W_{4,5} \cdot g(W_{1,4} \cdot o_1 + W_{2,4} \cdot o_2)\right) \end{aligned}$$



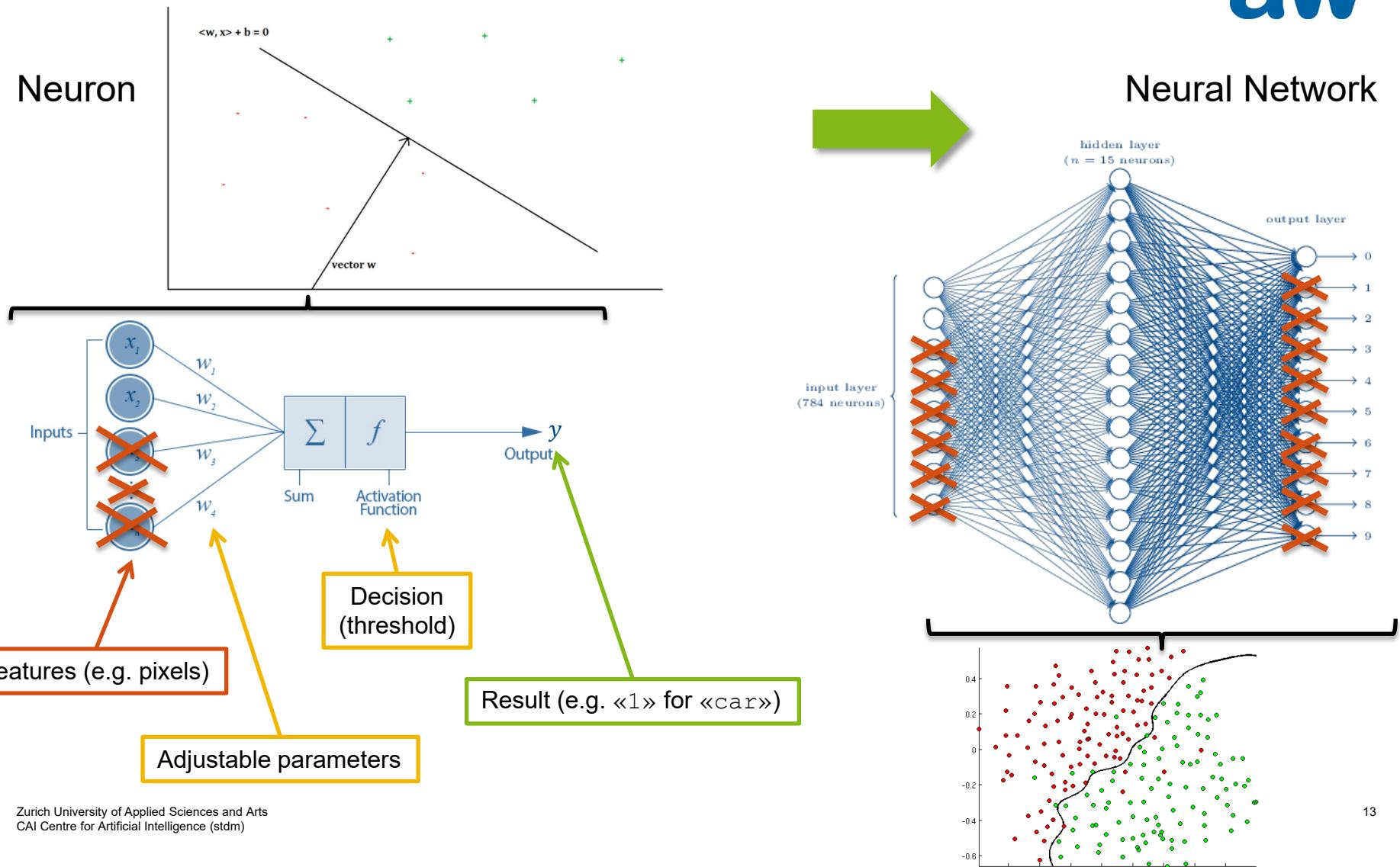
- Adjusting weights** changes the function: **learning** works this way!
(→ see appendix for first ideas)

Expressiveness of multilayer networks (**multilayer perceptrons**)

- All continuous functions w/ 2 layers, all functions w/ 3 layers
 - Combine two opposite-facing threshold functions to make a ridge
 - Combine two perpendicular ridges to make a bump
 - Add bumps of various sizes and locations to fit any surface



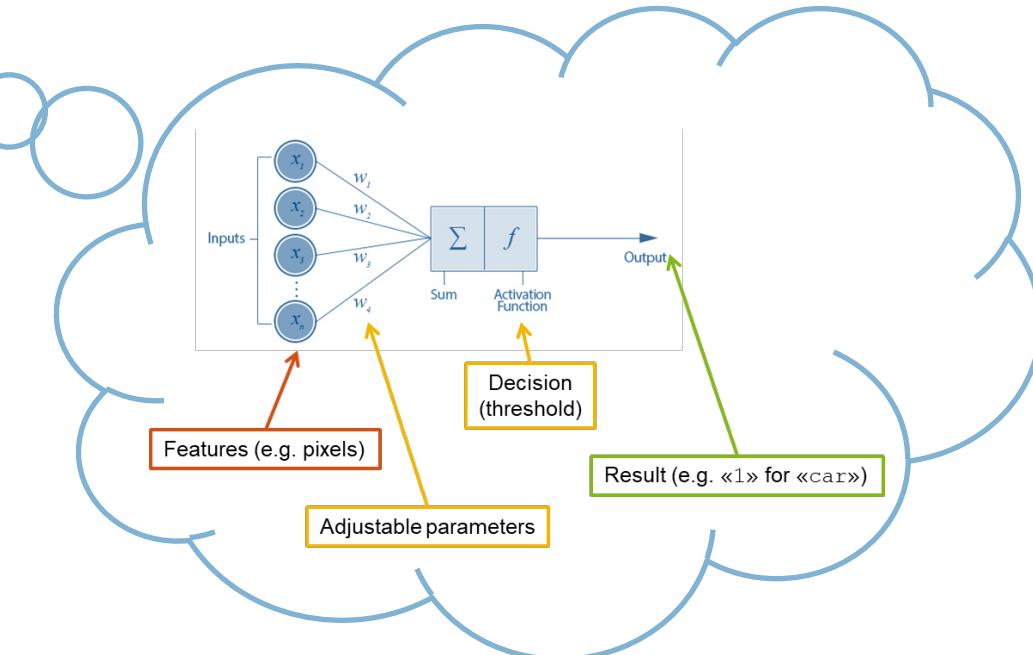
What is the effect of weight adjustment?



How are the weights adjusted?

First intuition

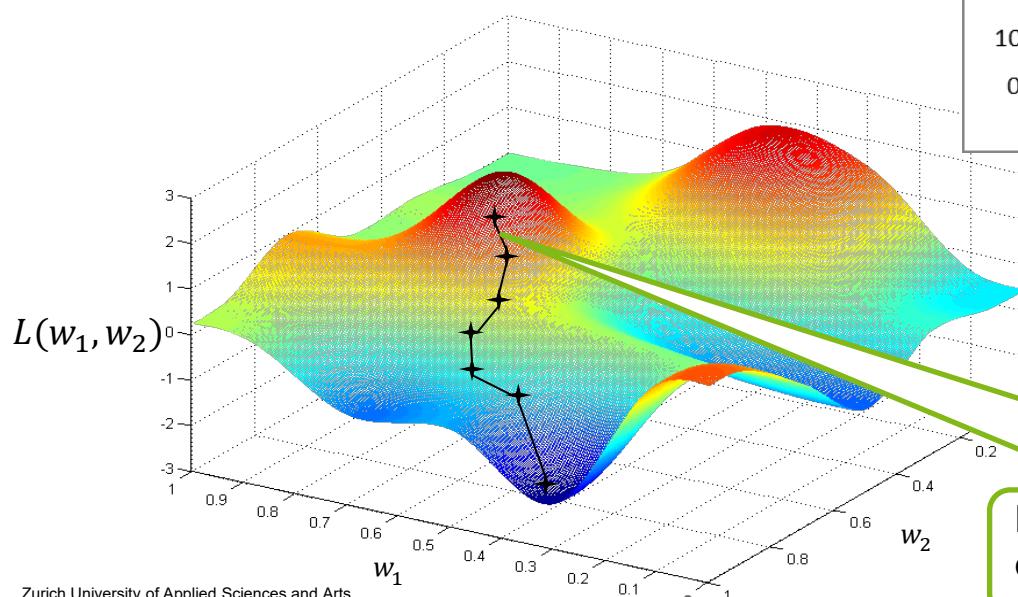
- Our example neural network: $f_W(x) = y$ with **image x** , **ground truth y** und **parameters W** ($W = \{w_1, w_2, \dots\}$ initialized randomly)



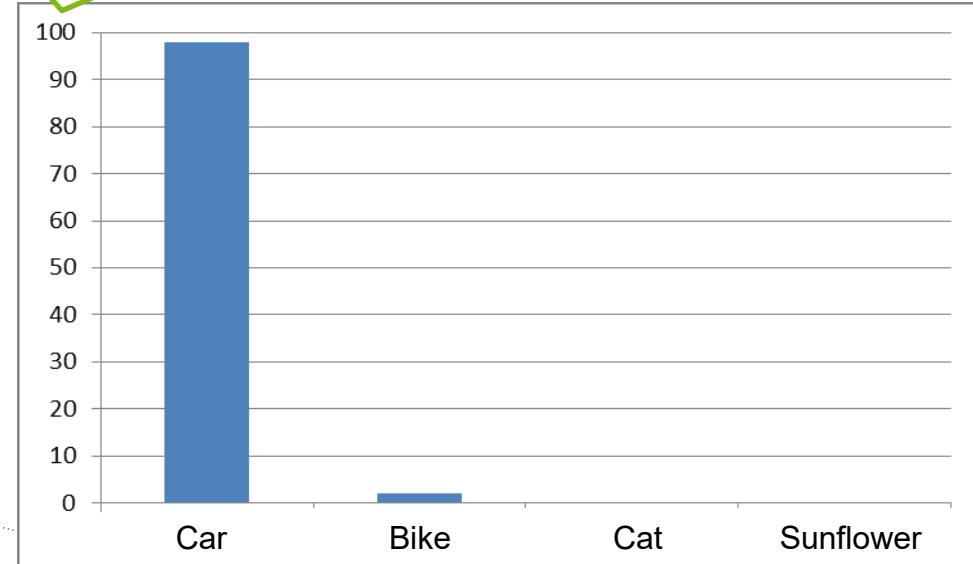
How are the weights adjusted? (contd.)

First intuition

- Our example neural network: $f_{\mathbf{W}}(\mathbf{x}) = \mathbf{y}$ with **image \mathbf{x}** , **ground truth \mathbf{y}** und **parameters \mathbf{W}** ($\mathbf{W} = \{w_1, w_2, \dots\}$ initialized randomly)
- Error measure: $L(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N (f_{\mathbf{W}}(\mathbf{x}_i) - \mathbf{y}_i)^2$
Average of quadratic difference on all images (loss function)



Likelihood [%] of certain event



← error «landscape»

Method: adapt weights of f in the direction of the steepest descent (downwards) of L

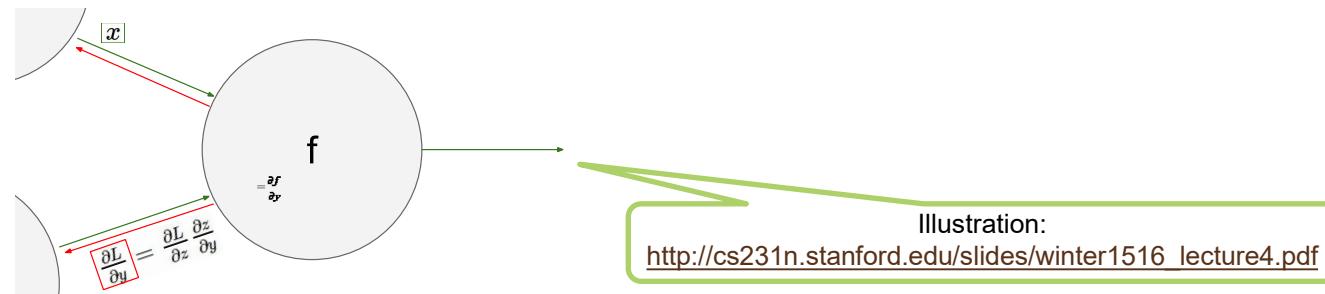
How are the weights adjusted?

Neural network training ideas

→ see also https://stdm.github.io/downloads/papers/ADS_2019_DeepLearning.pdf

Trained by gradient descent (complete network is differentiable)

- **Forward pass:** calculation of loss function L for a **mini batch** of training examples
- **Backward pass:** calculation of $\frac{\partial L}{\partial W_{l,i}}$ for each weight $W_{l,i}$ on overall loss
 - Efficiently computable by layer-wise application of chain rule (**backpropagation** algorithm)



Many details to be considered for training to work in practice

- Weight initialization: choose random initial weights according to the magnitude of the inputs
- Gradient flow: secure sufficient gradient magnitude for fast training convergence via **batchnorm**
- Learning rate: choose adaptive learning rates, e.g. using the **ADADELTA** optimizer
- Batch composition: care for sufficient randomness in the presentation order
- Regularization: use **dropout** to overcome the problem of more parameters than input data

Convolutional Neural Networks (CNNs)

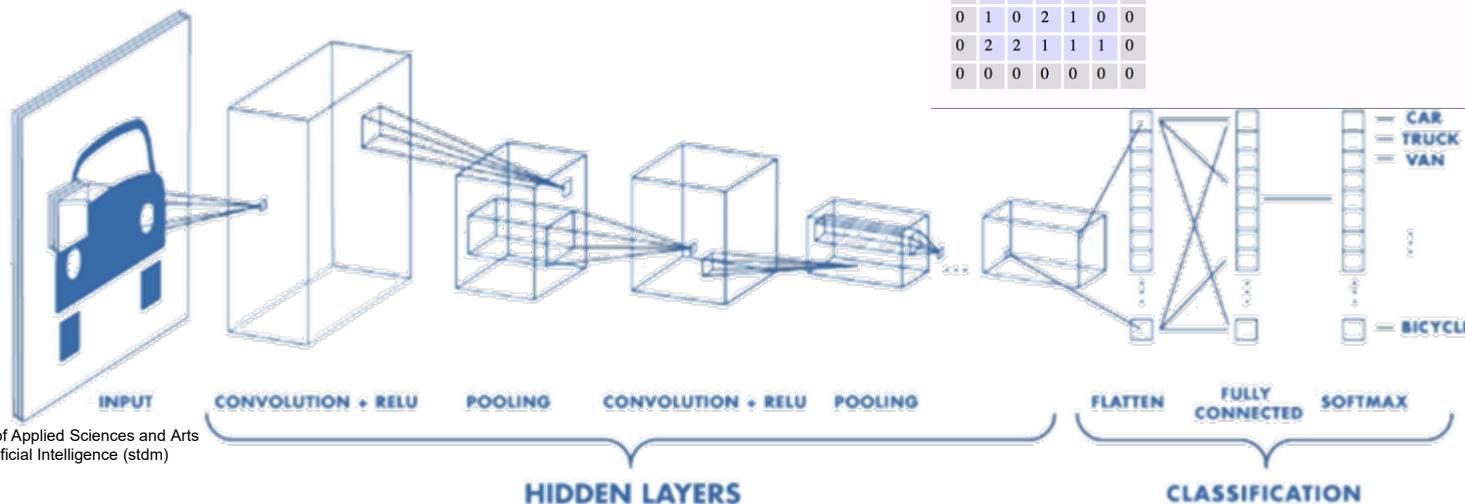
Intuition: cp. <https://medium.freecodecamp.org/an-intuitive-guide-to-convolutional-neural-networks-260c2de0a050>

Goal: fewer free parameters → eases learning

Idea: exploit 2D-correlated local structure in (image) input data
→ inspired by mammal visual cortex

Principle

- A “**filter**” moves over every input pixel and calculates a feature that **describes** the **pixel’s local context**
→ map result to same spatial location
→ filter weights (i.e., feature meaning) is trainable
- Have **several such** “filters” to encode different features
- After each filtering layer, **sub-sample** result to reduce spatial resolution and increase “field of vision”

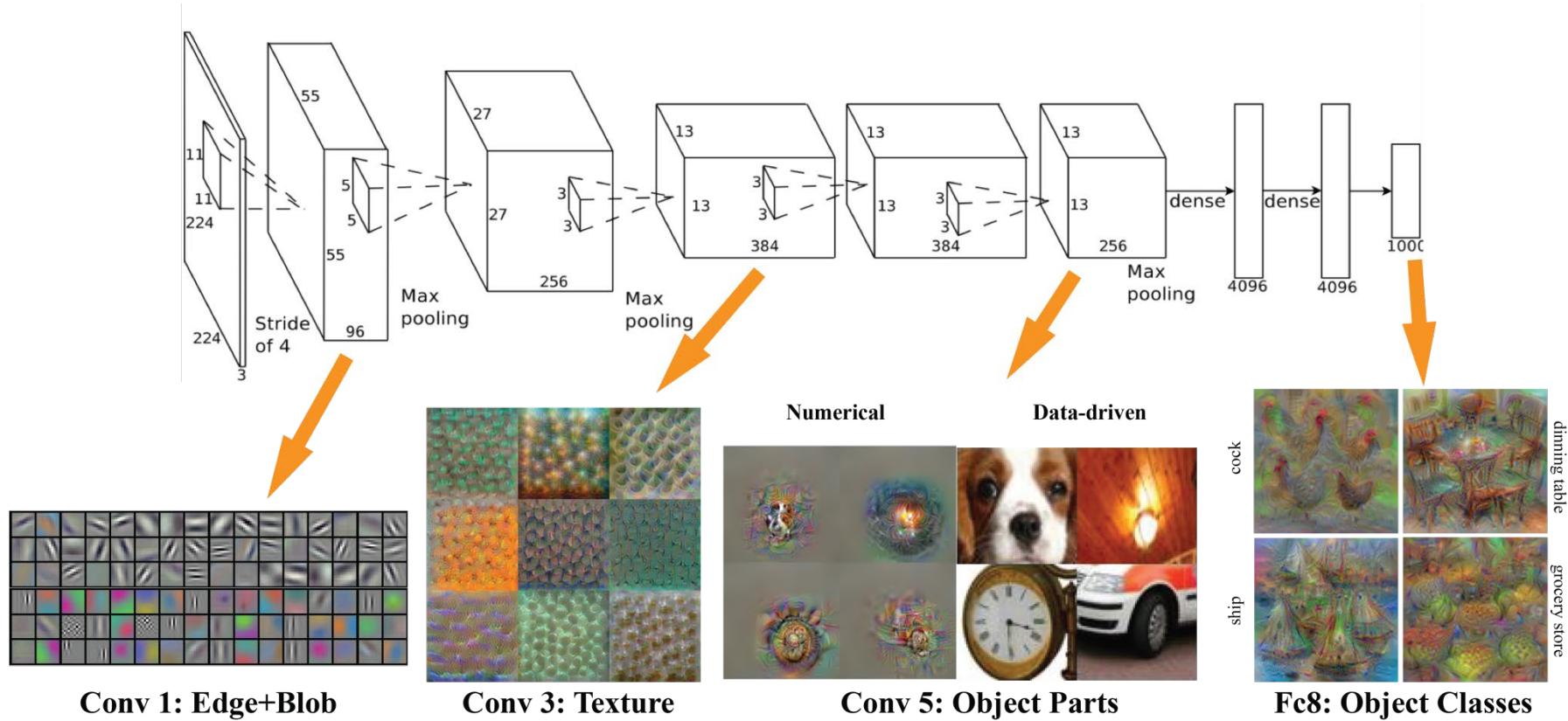


Input Volume (+pad 1) (7x7x3)	Filter W0 (3x3x3)	Output Volume (3x3x2)
$x[:, :, 0]$	$w0[:, :, 0]$	$o[:, :, 0]$
0 0 0 0 0 0 0 0 0 0 1 0 2 0 0 1 0 2 0 1 0 0 1 0 2 2 0 0 0 2 0 0 2 0 0 0 2 1 2 2 0 0 0 0 0 0 0 0 0	-1 0 1 0 0 1 1 -1 1 -1 0 1 1 -1 1 0 1 0 -1 1 1 1 1 0 0 -1 0 1 0 0 1 1 -1 0 -1 -1 1 0 0	2 3 3 3 7 3 8 10 -3 -1 0 0 1 -1 0 1 -1 0 -1 1 0 0 -1 -1 1 0 0 -3 1 0 -3 -8 -5
$x[:, :, 1]$	$w0[:, :, 1]$	$o[:, :, 1]$
0 0 0 0 0 0 0 0 2 1 2 1 1 0 0 2 1 2 0 1 0 0 0 2 1 0 1 0 0 1 2 2 2 2 0 0 0 1 2 0 1 0 0 0 0 0 0 0 0	-1 1 1 1 1 0 0 -1 0 -1 1 1 1 1 0 0 -1 0 1 1 -1 0 -1 -1 1 0 0	-8 -8 -3
$x[:, :, 2]$	$w0[:, :, 2]$	$o[:, :, 2]$
0 0 0 0 0 0 0 0 2 1 1 2 0 0 0 1 0 0 1 0 0 0 0 1 0 0 0 0 0 1 0 2 1 0 0 0 2 2 1 1 1 0 0 0 0 0 0 0 0	1 1 1 0 1 0 0 -1 0 -1 1 1 1 1 0 0 -1 0 -1 1 -1 0 -1 -1 1 0 0	1 0 0
	Bias $b0 (1 \times 1 \times 1)$ $b0[:, :, 0]$	$b1 (1 \times 1 \times 1)$ $b1[:, :, 0]$
	1	0

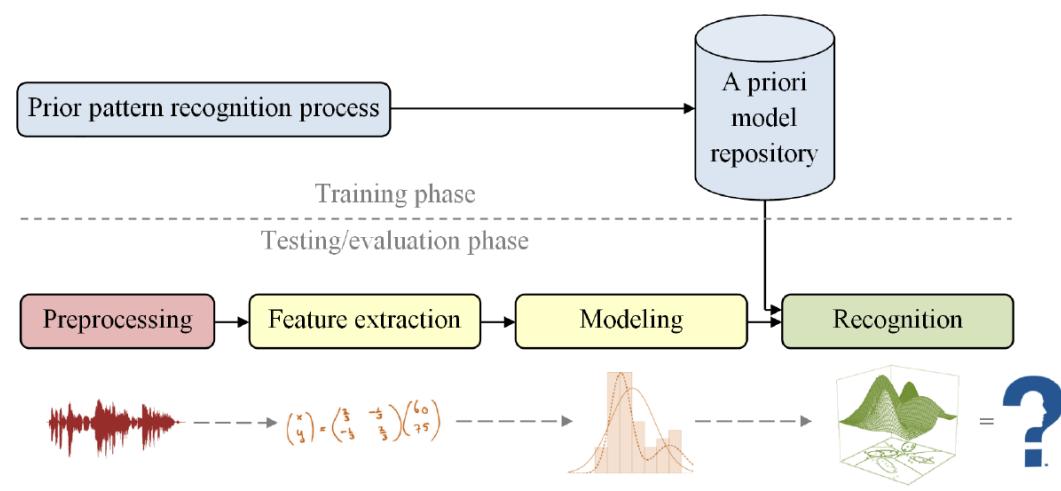
toggle movement

What does a convnet «see»?

A hierarchy of progressively complex features, visualized



3. DOING MACHINE LEARNING



Performance measurement

The ML development process being an empirical science

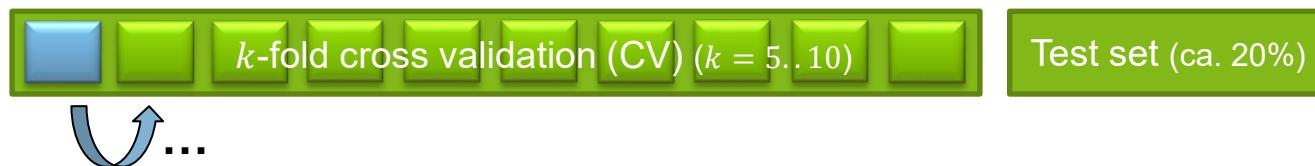
Hume's "Problem of Induction" (1740): when is generalization admissible?

How do we know that $h \approx f$ (the true function)?

1. Use theorems of computational/statistical learning theory

2. Try h on a new test set of examples

- Prerequisite for inductive learning: generalizes (only) to same distribution as seen in training set!
- Best practice: use cross-validation to train & validate on different sets before final test



3. Report performance using recognized figures of merit

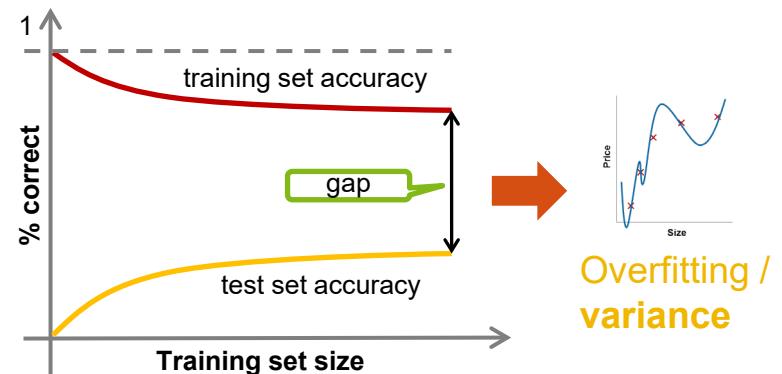
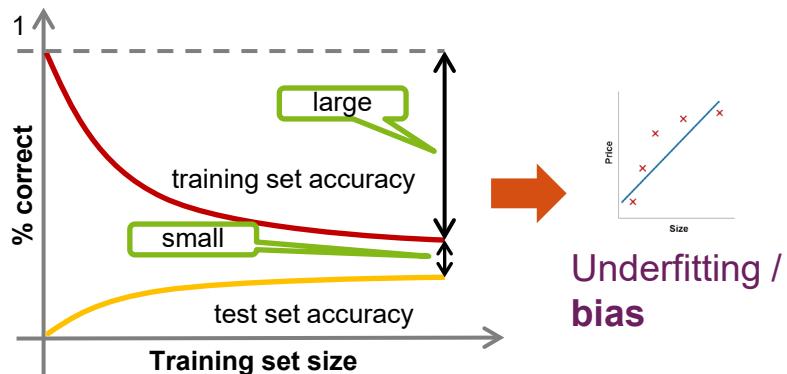
- E.g. accuracy (or test set error) if all errors are equally costly: $\text{accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$
- E.g. recall/precision if false alarms and misses differ in cost: $\text{recall} = \frac{TP}{TP+FN}$, $\text{precision} = \frac{TP}{TP+FP}$
- Conduct repeatable experiments (i.e., fully scriptable, full documentation of inputs and results)

classification → ↓ label	1	0
1	true positive (TP, "hit")	false negative (FN, "miss")
0	false positive (FP, "false alarm")	true negative (TN)

Debugging machine learning models

Learning curve: %correct on train & test set as a function of training set size

- Diagnostic: reveals over- and underfitting as well as realizability (→ see appendix)



What to try next when a given model generalizes poorly?

- Get **more training** examples → fixes **overfitting**
- Try **smaller sets of features** → fixes **overfitting**
- Try getting **additional features** → fixes **underfitting**
- Try **adding polynomial** features $x_1, x_2, x_1^2, x_2^2, \dots$ → fixes **underfitting**
- Try **less regularization** → fixes **underfitting**
- Try **more regularization** → fixes **overfitting**
- Build **ensembles** → fixes **overfitting**, uses limited data best (→ see V09)

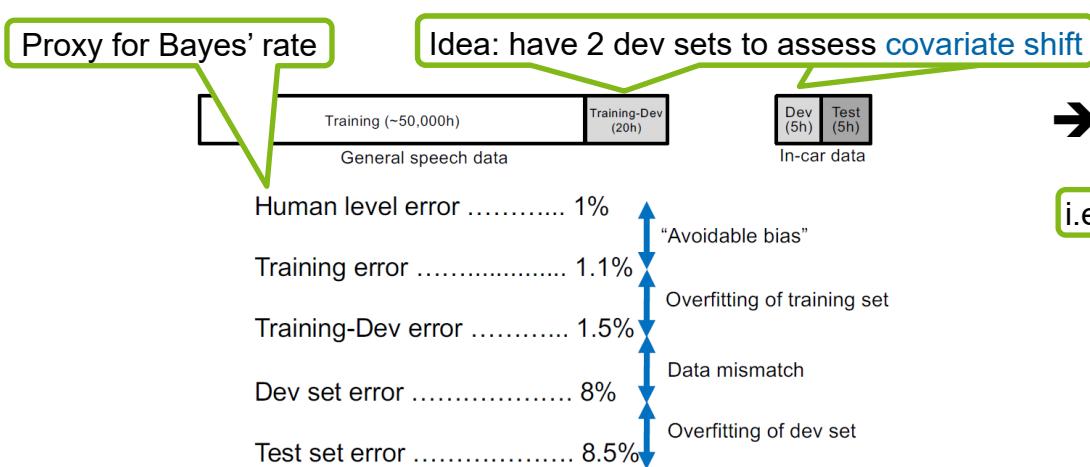
Regularization: Any method that limits the expressiveness of the hypothesis space by adding constraints to learning; e.g., pruning decision trees.

Bias & variance in the deep learning era

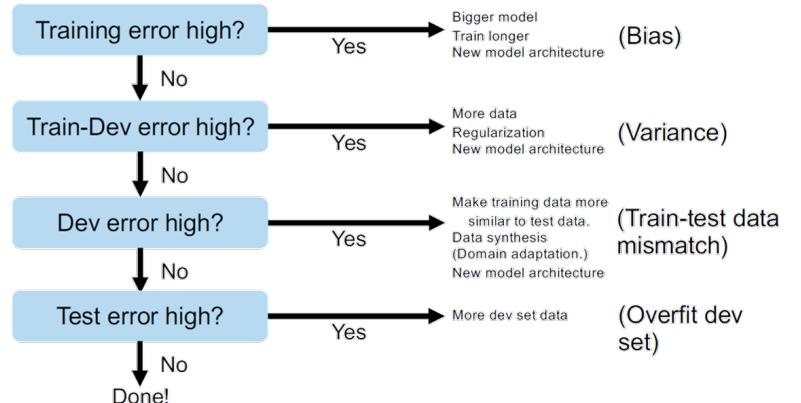
A. Ng's tutorial „Nuts and bolts of building AI...“ @ NIPS 2016

Today's situation

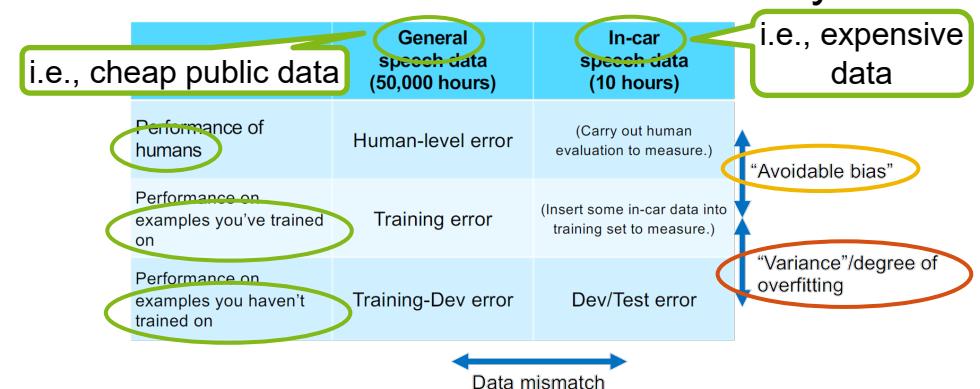
- We usually combine data sets from different distributions → **make validation and test data similar to obtain insight**
- Example in-car speech recognition:
Spread special in-car data to dev (validation) and test set!



→ New recipe for machine learning

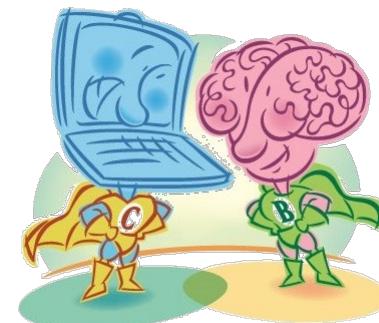


→ General human/bias/variance analysis



Where's the intelligence? Man vs. machine

- Machine learning offers **general function approximations purely learned** from examples
- But: **Success depends on** a good fit of the algorithm's inductive bias to problem at hand
→ i.e., **clever algorithm choice** based on experience
- Learning is a **powerful principle of self-optimization, applicable to all components** of previously seen agent designs
- But: **General** (domain crossing, knowledge-linking) **learning must be** based on way better inclusion of **unsupervised** learning principles (besides general inductive biases) and reasoning
→ current avant-garde deep learning research explores this route (→ see e.g., GANs in V11)



Review

- Learning needed for unknown environments, “lazy designers”
- Learning agent = performance element (**testing** / application **phase**)
+ learning element (**training phase**)
- **Learning method** (algorithm) **depends** on...
 - type of performance element (classify? regress? control?),
 - available feedback (labels),
 - type of component to be improved (representation? utility function? action?),
 - and data representation (numerical or categorical data, logical clauses, raw pixels, ...)
- For supervised learning, the **aim is** to find a **simple hypothesis** that is **approximately consistent** with training examples and **generalizes well**
- **Artificial neural networks** are inspired by biology, but **no “brain in a jar”**
 - **Popular** models because of feature-learning capabilities
 - Work especially well for pattern recognition
 - Difficult to interpret and maintain → give rise to new fields of **XAI** and **MLOps**
- Learning performance = prediction **accuracy** measured **on separate test set**
 - Development using 5-fold cross validation (without ever looking at test set!)
 - Systematic and repeatable experiments are paramount (e.g. using UNIX-style scripts)



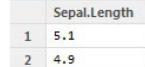
APPENDIX

Decision Trees Deep Learning Success Story



DECISION TREES

Attribute-based representations of data

Valid for all kinds of data ( , )

Examples described by **features**

- Possible attribute values: Boolean, discrete, continuous, etc.
- Example: “*Situations where I will/won't wait for a table*”

The diagram illustrates a dataset with 12 examples (X_1 to X_{12}) and 13 attributes. The attributes are categorized into two groups: **Attributes** (8 columns) and **Target** (1 column). The **Attributes** group includes *Alt*, *Bar*, *Fri*, *Hun*, *Pat*, *Price*, *Rain*, *Res*, *Type*, and *Est*. The **Target** group includes *WillWait*. The **Target** column is labeled *The label*.

Attributes:

	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>
X_1	T	F	F	T	Some	\$\$\$	F	T	French	0–10
X_2	T	F	F	T	Full	\$	F	F	Thai	30–60
X_3	F	T	F	F	Some	\$	F	F	Burger	0–10
X_4	T	F	T	T	Full	\$	F	F	Thai	10–30
X_5	T	F	T	F	Full	\$\$\$	F	T	French	>60
X_6	F	T	F	T	Some	\$\$	T	T	Italian	0–10
X_7	F	T	F	F	None	\$	T	F	Burger	0–10
X_8	F	F	F	T	Some	\$\$	T	T	Thai	0–10
X_9	F	T	T	F	Full	\$	T	F	Burger	>60
X_{10}	T	T	T	T	Full	\$\$\$	F	T	Italian	10–30
X_{11}	F	F	F	F	None	\$	F	F	Thai	0–10
X_{12}	T	T	T	T	Full	\$	F	F	Burger	30–60

Target:

<i>WillWait</i>
T
F
T
F
T
F
T
F
T
F
F
T

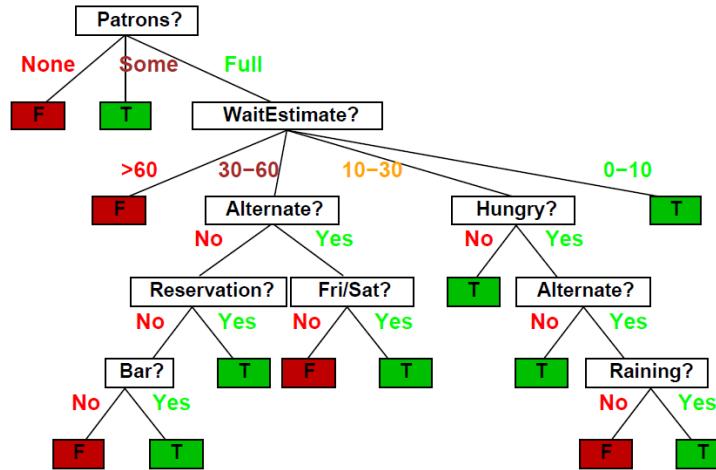
Labels:

<i>The label</i>

- Goal: **classification** of examples into positive (*T*) or negative (*F*) **class**

Decision tree representation of hypotheses

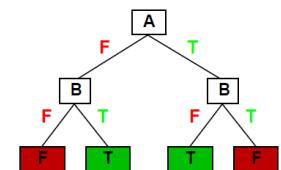
Example: Stuart Russell's “true” tree to decide whether to wait in a restaurant



Expressiveness

- Decision trees can express any function of the input attributes
E.g. for Boolean functions: truth table row → path to leaf
- Trivial tree ∨ training sets: one path to leaf for each example
But probably won't generalize to new examples
→ Prefer to find more compact decision trees

A	B	$A \oplus B$
F	F	F
F	T	T
T	F	T
T	T	F



Hypothesis spaces

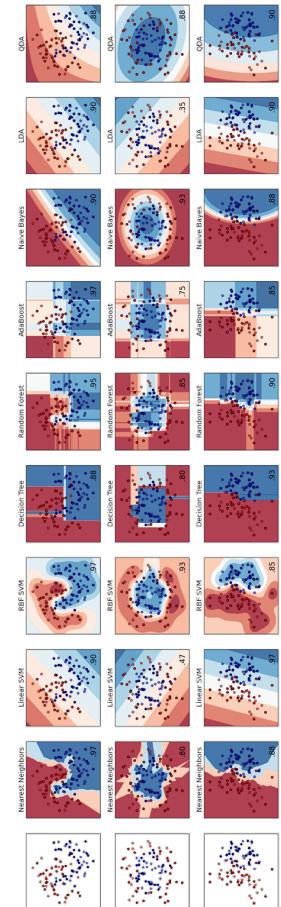
Even a constrained hypothesis space is large

- **How many distinct decision trees** with n Boolean attributes?
 - = number of Boolean functions
 - = number of distinct truth tables with 2^n rows = 2^{2^n}
 - Example: **6 Boolean attributes** → 18'446'744'073'709'551'616 possible trees
- How many purely conjunctive hypotheses (e.g., *Hungry* \wedge \neg *Rain*)
 - Each attribute can be either positive, negative, or out of the hypothesis
→ 3^n

More expressive hypothesis spaces

- ...increase chance that **target** function can be **expressed** 😊
- ...increases **number** of hypotheses **consistent** w/ training set
→ **may get worse** predictions ☹

Due to overfitting we have seen earlier



Decision tree learning

Goal: find a **small** tree **consistent** with the training examples

Idea: (**recursively**) choose “most significant” attribute as root of (sub)tree

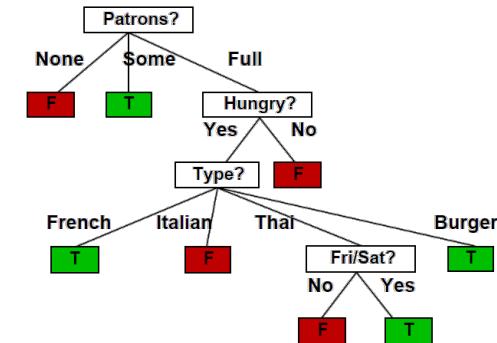
Algorithm

```

• function LearnDecisionTree(examples, attributes) returns a tree
    return DecisionTreeLearning(examples, attributes, {})

function DecisionTreeLearning(examples, attributes, parent_examples) returns a tree
    if examples is empty then return PluralityValue(parent_examples)
    else if all examples have the same label then return the label
    else if attributes is empty then return PluralityValue(examples)
    else
        A ← argmaxa∈attributes Importance(a, examples)
        tree ← a new decision tree with root test A
        for each value  $v_k$  of A do #for categorical features; for numerical features, use e.g. binning
            exs ← {e: e∈examples and e.A= $v_k$ }
            subtree ← DecisionTreeLearning(exs, attributes-A, examples)
            add a branch to tree with label ( $A=v_k$ ) and subtree subtree
    return tree

```



- **PluralityValue(examples)** selects the **most common output** among examples
- **Importance(attribute, examples)** selects the **most important attribute**
- On ties, both functions choose randomly

Choosing an attribute

How to implement Importance (attribute, examples)

Idea: A good attribute splits examples into subsets that are (ideally) “all pos” or “all neg”

Example

- Question: “**Would I wait if the crowdedness is x ?**”

Answer: “ $x = \text{None}$: **no**; $x = \text{Some}$: **yes**; $x = \text{Full}$: **not clear**”

```

graph TD
    Root[Patrons?] -- None --> NodeNone[ ]
    Root -- Some --> NodeSome[ ]
    Root -- Full --> NodeFull[ ]
    NodeNone --- DotNone1(( ))
    NodeNone --- DotNone2(( ))
    NodeSome --- DotSome1(( ))
    NodeSome --- DotSome2(( ))
    NodeSome --- DotSome3(( ))
    NodeFull --- DotFull1(( ))
    NodeFull --- DotFull2(( ))
    NodeFull --- DotFull3(( ))
    NodeFull --- DotFull4(( ))
    NodeFull --- DotFull5(( ))
  
```
- Question: “**Would I wait if the restaurant’s type is x ?**”

Answer: “ $\forall x$: **fifty-fifty**”

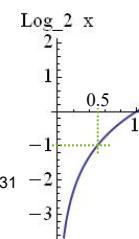
```

graph TD
    Root[Type?] -- French --> NodeFrench[ ]
    Root -- Italian --> NodeItalian[ ]
    Root -- Thai --> NodeThai[ ]
    Root -- Burger --> NodeBurger[ ]
    NodeFrench --- DotFrench1(( ))
    NodeFrench --- DotFrench2(( ))
    NodeItalian --- DotItalian1(( ))
    NodeItalian --- DotItalian2(( ))
    NodeThai --- DotThai1(( ))
    NodeThai --- DotThai2(( ))
    NodeThai --- DotThai3(( ))
    NodeThai --- DotThai4(( ))
    NodeBurger --- DotBurger1(( ))
    NodeBurger --- DotBurger2(( ))
    NodeBurger --- DotBurger3(( ))
    NodeBurger --- DotBurger4(( ))
  
```
- Patrons* is better choice: gives information about the classification

Recap: Information theory

- Information answers questions:** The more cluelessness an observation **removes**, the more information it contains
- Inversely proportional to **entropy** (**uncertainty** of a random variable)

 - A Boolean answer with **prior** $<0.5, 0.5>$ has entropy= **1 bit** (if we remove this uncertainty, we gain 1 bit of info.)
 - A coin giving heads 99% of the time has entropy close to 0 ($\approx 0.08 \text{ bits}$ \rightarrow almost no **info.-gain** when observed)
 - Entropy in an observation** (having prior $<P_1, \dots, P_n>$): $H(\langle P_1, \dots, P_n \rangle) = -\sum_{i=1}^n P_i \log_2 P_i$



Information gain as splitting criterion

Suppose we have p positive and n negative examples at the root

- $H\left(\frac{p}{p+n}, \frac{n}{p+n}\right)$ bits needed to classify a new example
- E.g., for the 12 restaurant examples, $p = n = 6$, so we need overall 1 bit

An attribute A splits the examples E into **subsets** E_i (one per possible value)

- Each of which (we hope) **needs less information** to complete the classification
- Let E_i have p_i positive and n_i negative examples
 $\rightarrow H\left(\frac{p_i}{p_i+n_i}, \frac{n_i}{p_i+n_i}\right)$ bits needed to classify a new example
- **Expected** number of necessary bits per example over all branches i stemming from A is

$$\text{Remainder}(A) = \sum_i \frac{p_i + n_i}{p + n} H\left(\left(\frac{p_i}{p_i+n_i}, \frac{n_i}{p_i+n_i}\right)\right)$$

Entropy of branch i ,
weighted by branch's size

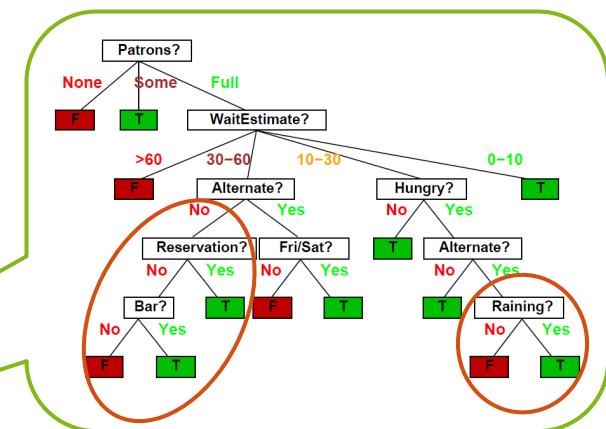
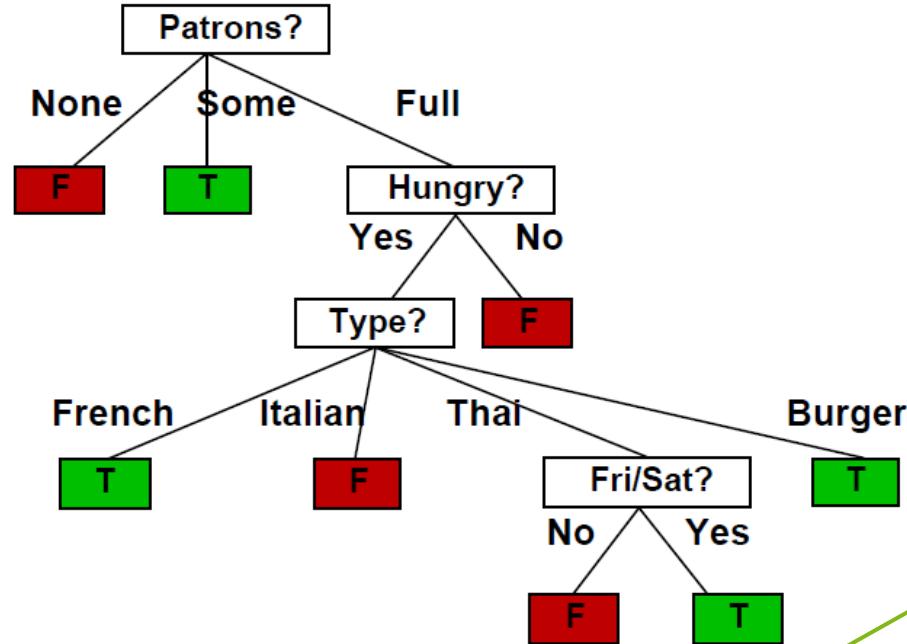
- For *Patrons* this is 0.459 bits, for *Type* this is (still) 1 bit
→ Choose the attribute that **minimizes** the **remaining information needed**, ...
→ i.e., maximizes **information gain**: $Gain(A) = H\left(\left(\frac{p}{p+n}, \frac{n}{p+n}\right)\right) - Remainder(A)$

Entropy of original problem

Entropy remaining after splitting on A

The learned decision tree

Based on our 12 examples

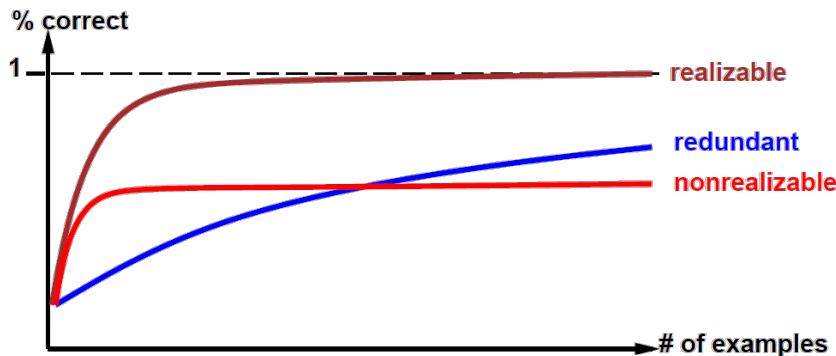
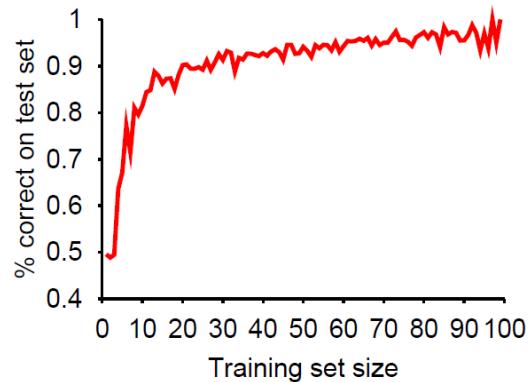


- Substantially simpler than “true” tree
→ E.g., *Reservation* and *Raining* are not needed (perfect classification possible without)
- A more complex hypothesis isn't justified by the small amount of data
→ But what makes one tree better than another?

Learning curves

Diagnosing learning problems

Learning curve, simplified: %correct on (full) test set only as a function of training set size



Accuracy shown in learning curve depends on

- **Realizability** (target function expressible in chosen hypothesis space?)
 - Non-realizability can be due to **missing attributes**
 - or **restricted hypothesis class** (e.g., a thresholded linear function might be overly simplistic)
- **Redundant** features
 - (e.g., loads of irrelevant attributes make learning difficult)



DEEP LEARNING SUCCESS STORY

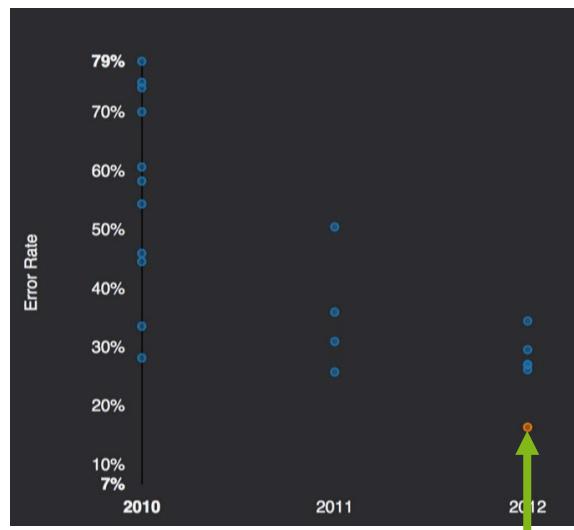
Why is this current hype about deep learning?

The ImageNet Competition



1000 categories

1 mio. training examples



2015: Computers learned to «see»

4.95% Microsoft (Feb 06)

→ super-human performance (human: 5.10%)

4.80% Google (Feb 11)

4.58% Baidu (May 11)

3.57% Microsoft (Dec 10)

2016: A summer of breakthroughs in ML

...enabled by deep learning

Impressive novelties within a summer's timespan

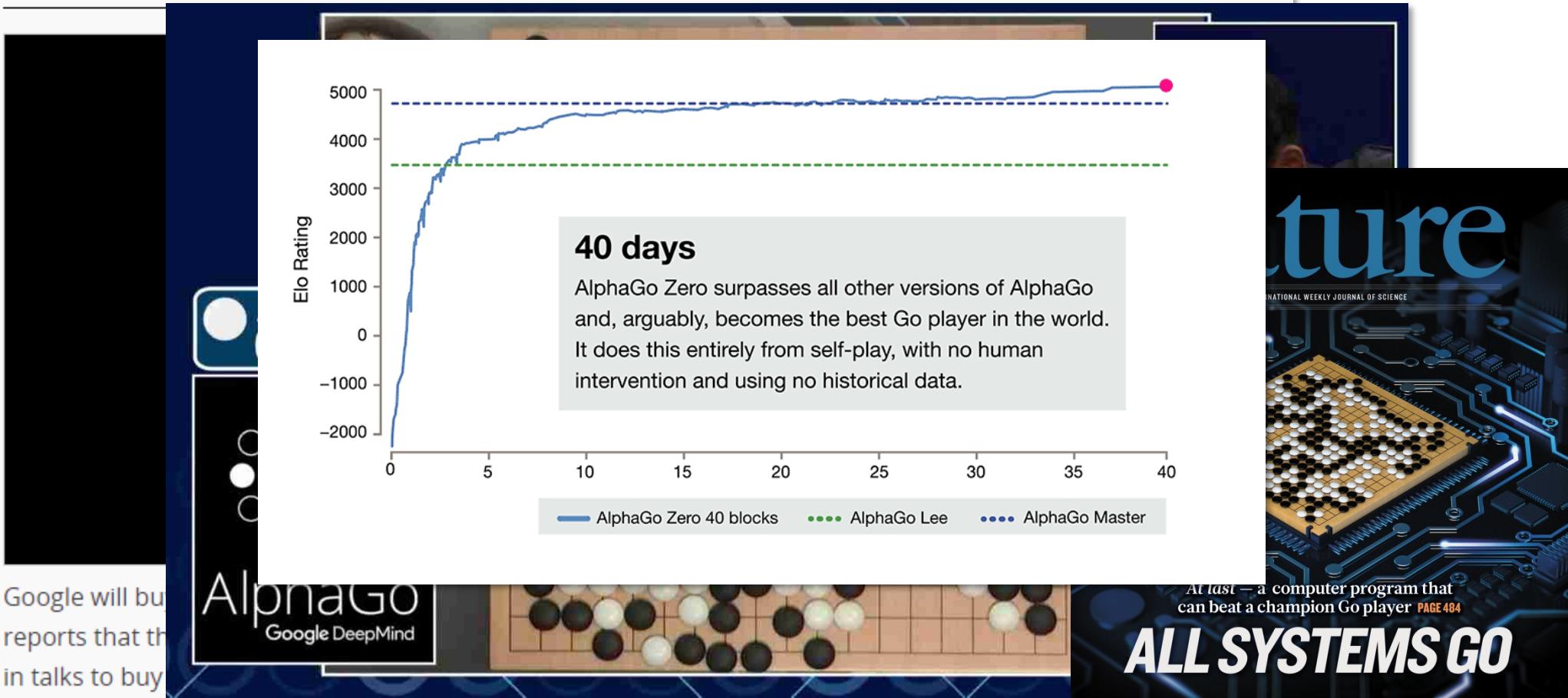
- Game playing: beating the human Go world champion
- Audio synthesis: Synthesizing speech & music sample by sample
- Art style transfer: Redraw the content of a picture in the style of any painting
- Image synthesis: Completion of missing parts in pictures
- Text synthesis: Generation of text in specific styles (e.g., Shakespeare, $L^A T_E X$, ...)
- Word vectors: Arithmetic with semantic meaning of text and images

➔ See next slides



Google Acquires Artificial Intelligence Startup DeepMind For More Than \$500M

Posted Jan 26, 2014 by Catherine Shu (@catherineshu)



The acquisition was originally confirmed by Google to Re/code.

CONSERVATION
SONGBIRDS À LA CARTE
Illegal harvest of millions of Mediterranean birds
PAGE 452

RESEARCH ETHICS
SAFEGUARD TRANSPARENCY
Don't let openness backfire on individuals
PAGE 453

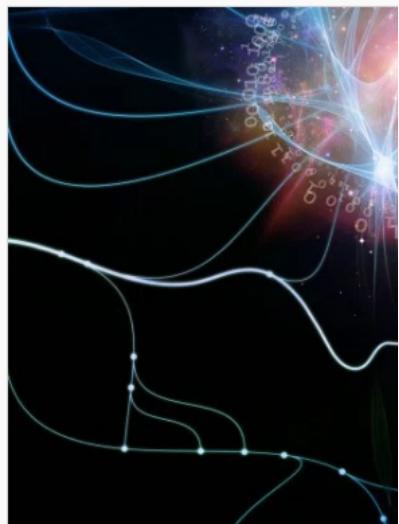
POPULAR SCIENCE
WHEN GENES GOT 'SELFISH'
Dawkins' calling card forty years on
PAGE 462

NATURE.COM/NATURE
26 January 2016 410
Vol. 529 No. 7587
047
9 770028053095

Google's WaveNet uses neural nets to generate eerily convincing speech and music

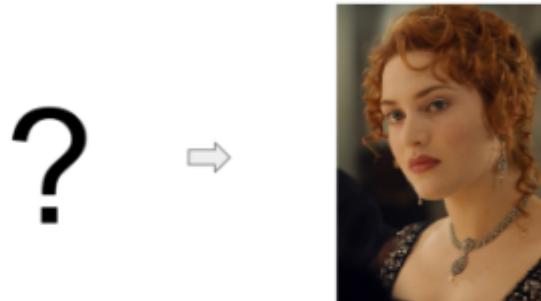
Posted Sep 9, 2016 by Devin Coldewey

Zurich University
of Applied Sciences



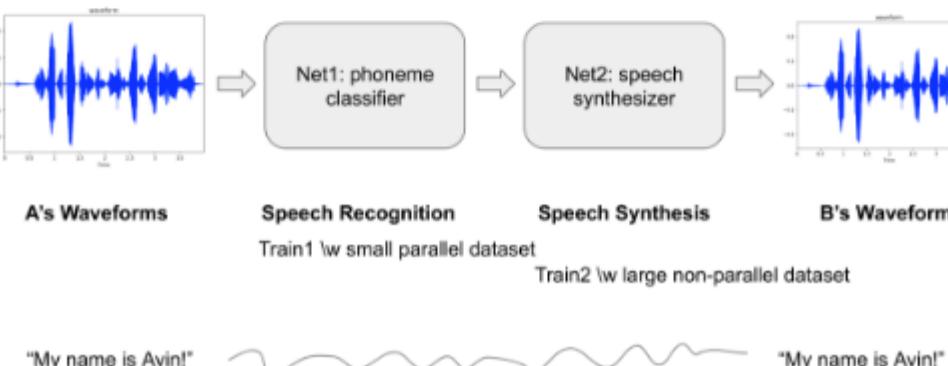
Intro

What if you could imitate a famous celebrity's voice or sing like a famous singer? This project started with a goal to convert someone's voice to a specific target voice. So called, it's voice style transfer. We worked on this project that aims to convert someone's voice to a famous English actress [Kate Winslet's voice](#). We implemented a deep neural networks to achieve that and more than 2 hours of audio book sentences read by Kate Winslet are used as a dataset.



Model Architecture

This is a many-to-one voice conversion system. The main significance of this work is that we could generate a target speaker's utterances without parallel data like <source's wav, target's wav>, <wav, text> or <wav, phone>, but only waveforms of the target speaker. (To make these parallel datasets needs a lot of effort.) All we need in this project is a number of waveforms of the target speaker's utterances and only a small set of <wav, phone> pairs from a number of anonymous speakers.



generated
speech from text

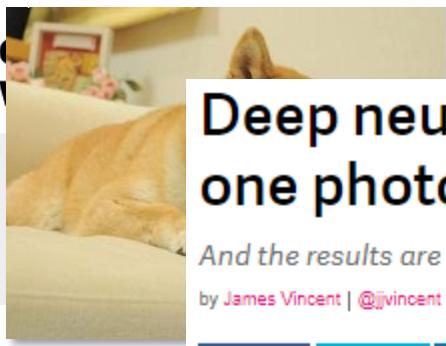
generated music
of creativity

Generating speech from a piece of text computers, but it's pretty rare that the technique from researchers at Alphabet producing speech and even music that

Early systems used a large library of the large ruleset that described all the ways pieces were joined, or concatenated, creating most words, albeit with unconvincing generation of sound, making a library often less effective.

WaveNet, as the system is called, takes low a level as possible: one sample at a scratch — 16,000 samples per second.

1 Second



Deep neural networks can now transfer the style of one photo onto another

And the results are impressive

by James Vincent | @jvincent | Mar 30, 2017, 1:53pm EDT

f SHARE

tweet TWEET

in LINKEDIN

Computing

Algorithm
Artistic
Other In

A deep neural n
other images.

by Emerging Tech

The nature of art
of Vincent Van C
Edvard Munch's
humans recogni



Original photo

Reference photo

Result

You've probably heard of an AI technique known as "style transfer" — or, if you haven't heard of it, you've seen it. The process uses neural networks to apply the look and feel of one image to another, and appears in apps like [Prisma](#) and [Facebook](#). These style transfers, however, are stylistic, not photorealistic. They look good because they look like they've been painted. Now a group of researchers from Cornell University and Adobe have augmented

Ad closed by Google

Report this ad

AdChoices ▾



NOW TRENDING

...and the list could be continued

Brandon Amos About Blog

Image Completion with Deep Learning in TensorFlow

August 9, 2016



- Introduction
- Step 1: Interpreting images as samples from a probability distribution
 - How would you fill in the missing information?
 - But where does statistics fit in? These are images.
 - So how can we complete images?
- Step 2: Quickly generating fake images
 - Learning to generate new samples from an unknown probability distribution
 - [ML-Heavy] Generative Adversarial Net (GAN) building blocks
 - Using $G(z)$ to produce fake images
 - [ML-Heavy] Training DCGANs
 - Existing GANs
 - [ML-Heavy] Generating faces
 - Running DCGANs
- Step 3: Finding the right samples
 - Image completion
 - [ML-Heavy] Content-aware fill
 - [ML-Heavy] Semantic Image Inpainting
 - [ML-Heavy] Generating fake images
 - Completing your images
- Conclusion
- Partial bibliography
- Bonus: Incomplete list of papers

Introduction

Content-aware fill is a powerful technique for image completion and inpainting. It can do content-aware fill, image completion, and semantic image inpainting. "Semantic Image Inpainting" shows how to use deep learning to fill in some deeper portions of images. This section can be skipped if you're not interested in completion from images of faces. I have a blog post about image completion: tensorflow-image-completion.tensorflow.

We'll approach image completion in three steps:

1. We'll first interpret what we want to do.
2. This interpretation will lead us to a solution.
3. Then we'll find the solution.



...and the list could be continued

Image Completion with Deep Learning in TensorFlow

August 9, 2016



- Introduction
- Step 1: Interpreting images as samples from a probability distribution
 - How would you fill in the missing information?
 - But where does statistics fit in? These are images.
 - So how can we complete images?
- Step 2: Quickly generating fake images
 - Learning to generate new samples from an unknown probability distribution
 - [ML-Heavy] Generative Adversarial Net (GAN) building blocks
 - Using $G(z)$ to produce fake images
 - [ML-Heavy] Training DCGANs
 - Existing GANs
 - [ML-Heavy] Generating faces
 - Running DCGANs
- Step 3: Finding the right samples
 - Image completion
 - [ML-Heavy] Content-aware fill
 - [ML-Heavy] Semantic Image Inpainting
 - [ML-Heavy] Generating fake images
 - Completing your images
- Conclusion
- Partial bibliography
- Bonus: Incomplete list of papers

Introduction

Content-aware fill is a powerful technique for image completion and inpainting. It can do content-aware fill, image completion, and semantic image inpainting. "Semantic Image Inpainting" shows how to use deep learning to fill in some deeper portions of images. This section can be skipped if you're not interested in completion from images of faces. I have a blog post about image completion: tensorflow-image-completion.tensorflow.

We'll approach image completion in three steps:

1. We'll first interpret what we want to do.
2. This interpretation will lead us to a solution.
3. Then we'll find the solution.

...and the list could be continued

Andrej Karpathy blog

About Hacker's guide to Neural Networks

The Unreasonable Effectiveness of Recurrent Neural Networks

TECH

Nvidia AI Generates Fake Faces Based On Real Celebs

BY STEPHANIE MLOT 10.21.2017 :: 10:00AM EST

32 SHARES



STAY ON TARGET

AI Shelley Pens Truly Creepy Horror Stories—And You Can Help

Neural Network Serves Up Truly Frightening Halloween Costume Ideas

Celebrity scandals are about to get a lot more complicated.

Nvidia has developed a way of producing photo-quality, AI-generated human profiles—by using famous faces.

the morning paper

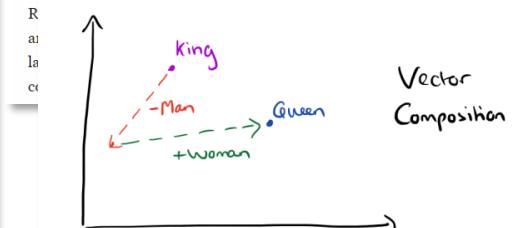
The amazing power of word vectors

APRIL 21, 2016

For today's post, I've drawn material not just from one paper, but from five! The subject matter is 'word2vec' – the work of Mikolov et al. at Google on efficient vector representations of words (and what you can do with them). The papers are:

- ★ Efficient Estimation of Word Representations in Vector Space – Mikolov et al. 2013
- ★ Distributed Representations of Words and Phrases and their Compositionality – Mikolov et al. 2013
- ★ Linguistic Regularities in Continuous Space Word Representations – Mikolov et al. 2013
- ★ word2vec Parameter Learning Explained – Rong 2014
- ★ word2vec Explained: Deriving Mikolov et al's Negative Sampling Word-Embedding Method – Goldberg and Levy 2014

From the first of these papers ('Efficient estimation...') we get a description of the *Continuous Bag-of-Words* and *Continuous Skip-gram* models for learning word vectors (we'll talk about what a word vector is in a moment...). From the second paper we get more illustrations of the power of word vectors, some additional information on optimisations for the skip-gram model (hierarchical softmax and negative sampling), and a discussion of



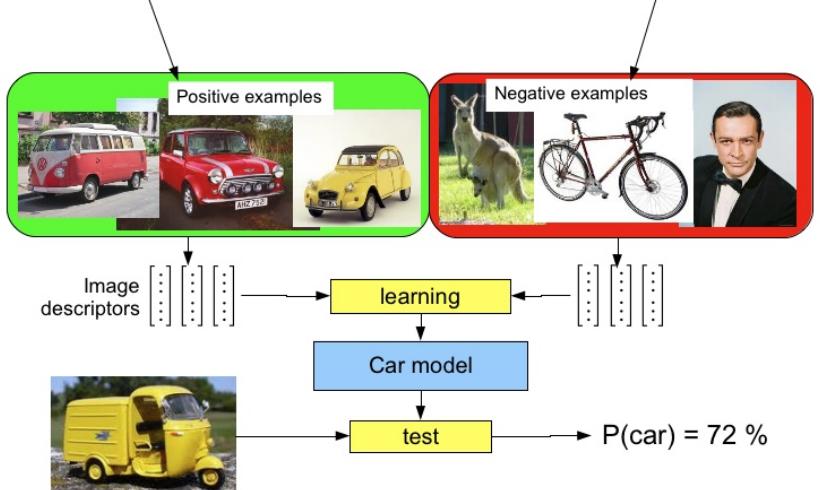
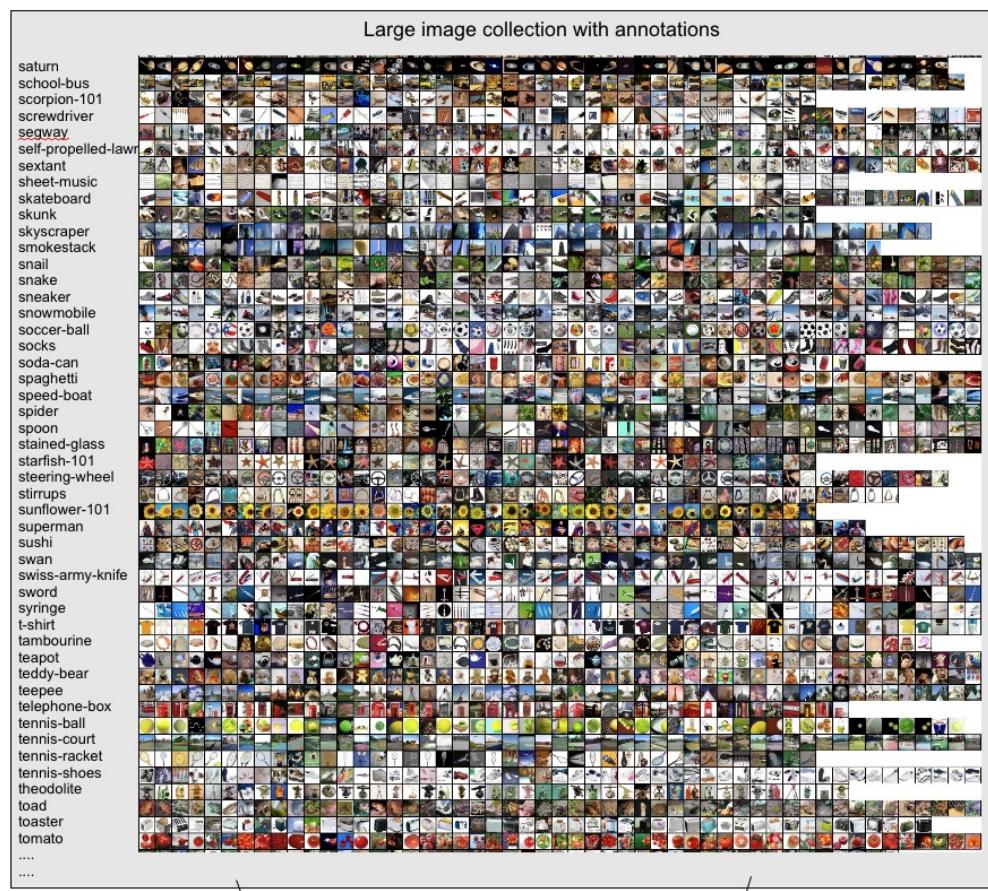
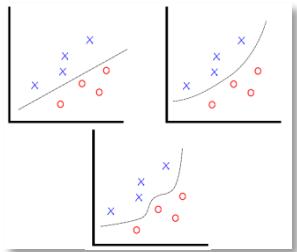
Inductive supervised learning

Assumption

- A **model** fit to *enough training examples*...
- ...will **generalize** well to unseen **test data**

Method

- **Search for parameters** of a given class of functions...
- ...such that every training input (e.g. an image) is mapped to the correct output **label** (e.g. «car»)



What does a neural network «see»? A hierarchy of progressively complex features

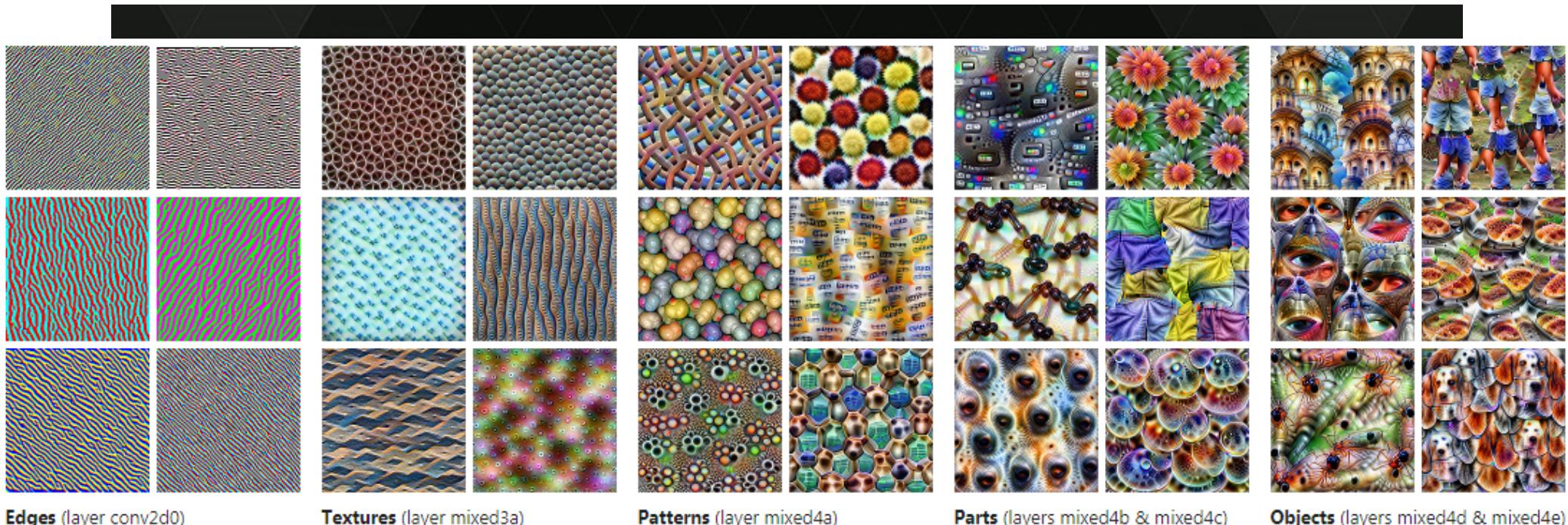


Image source: "Unsupervised Learning of Hierarchical Representations with Convolutional Deep Belief Networks" ICML 2009 & Comm. ACM 2011.
Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Ng.

Sources: <https://www.pinterest.com/explore/artificial-neural-network/>

Olah, et al., "Feature Visualization", Distill, 2017, <https://distill.pub/2017/feature-visualization/>.