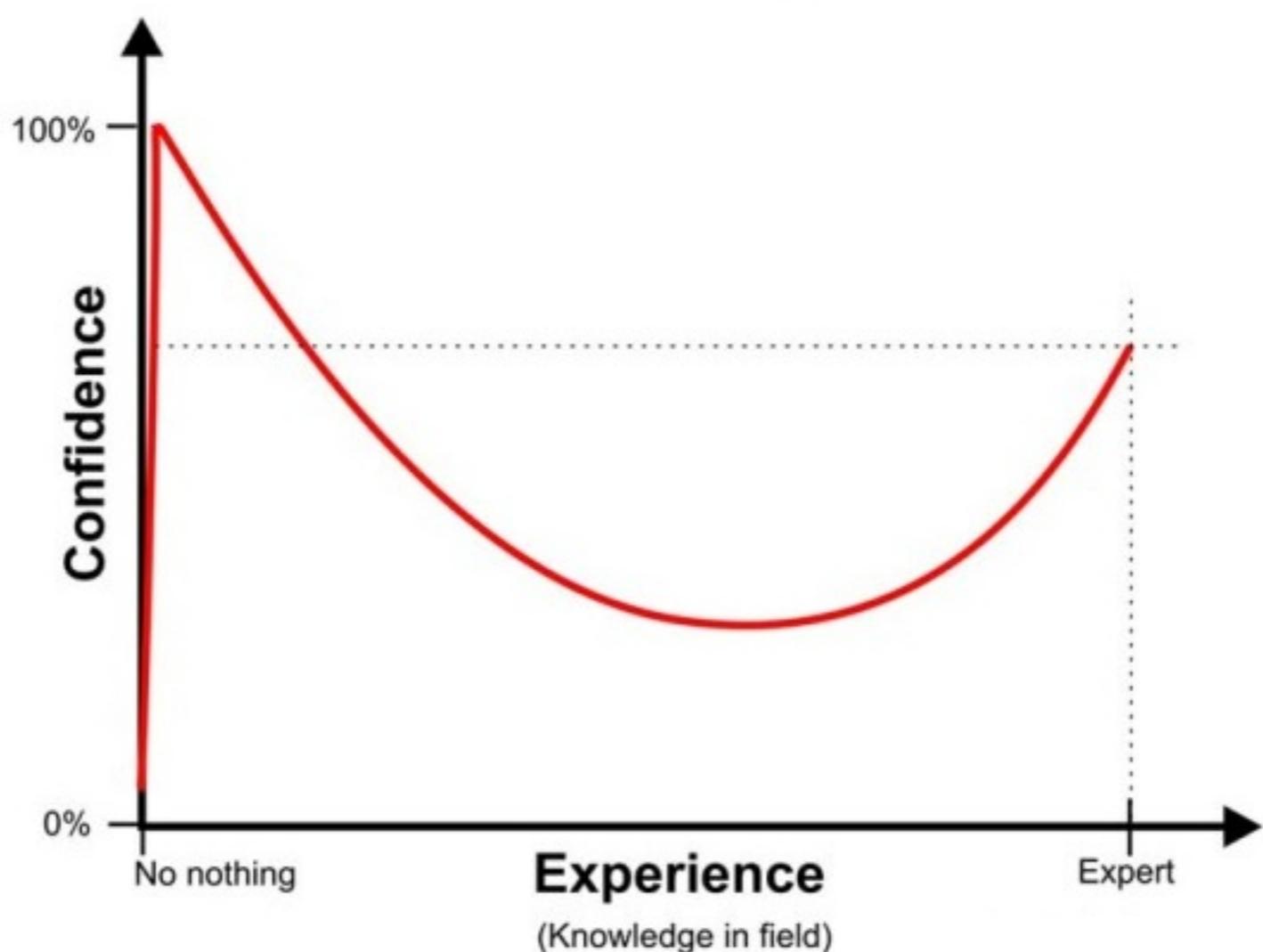


# Work to Machine!

DIAS ML course by Martin Topinka

<https://github.com/toastmaker/ml-dias>



# Introduce yourself and your goals



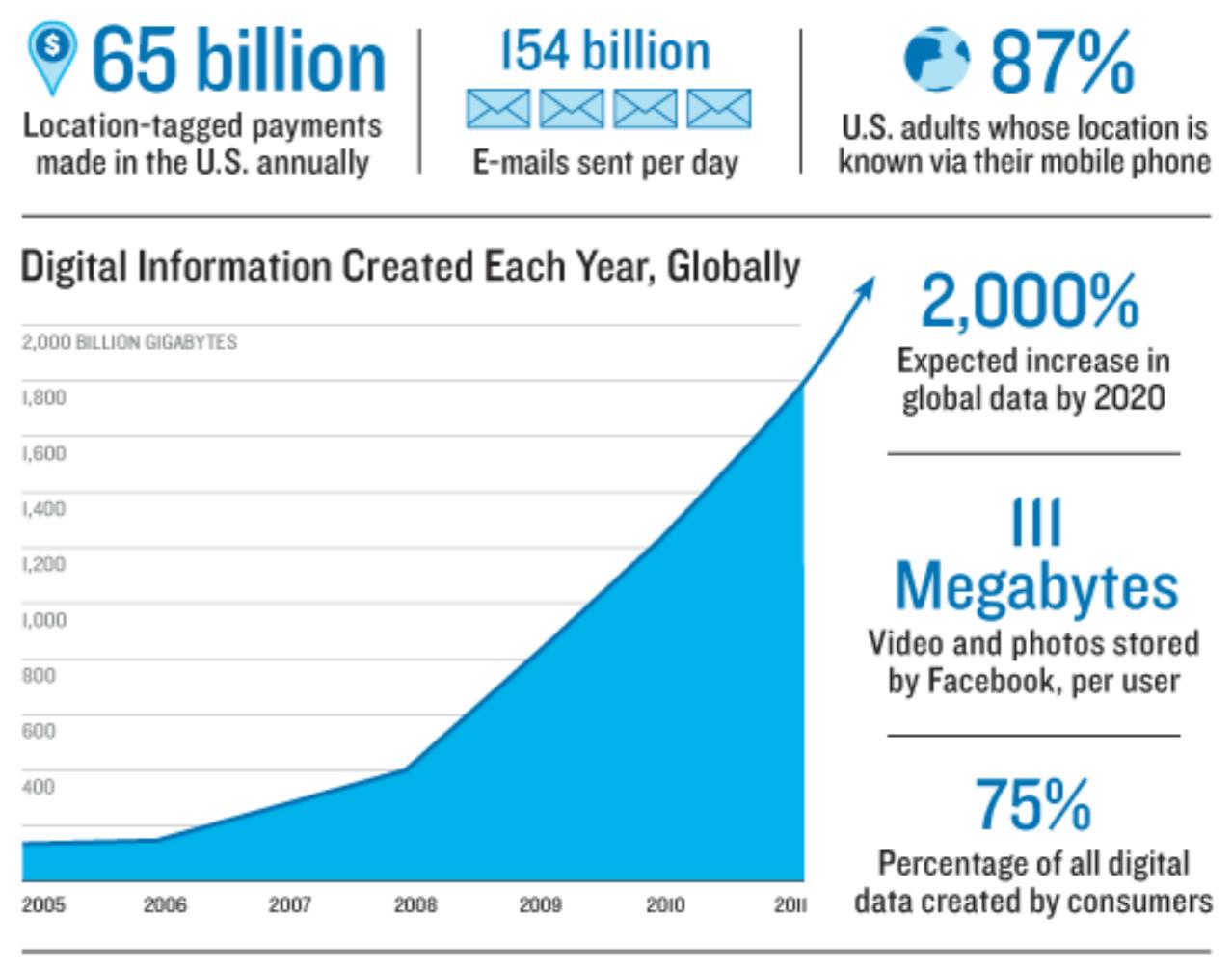
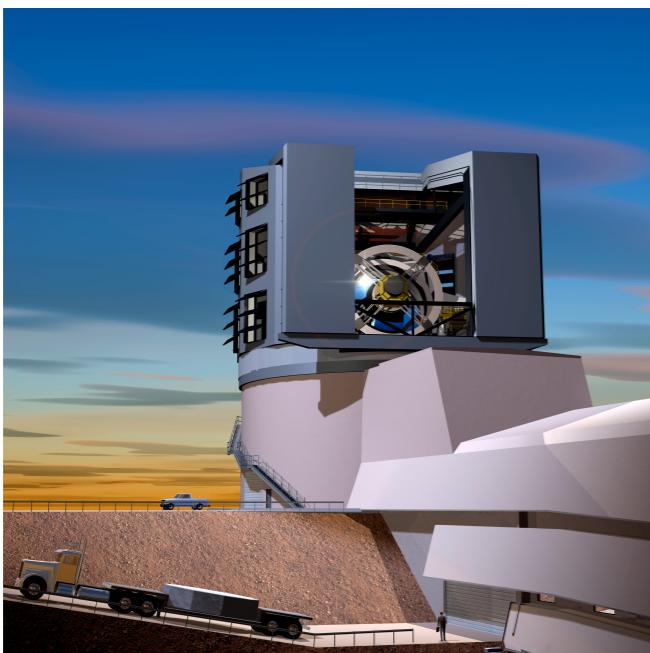
## 3Vs':

- **Volume (large archives)**
- **Velocity (continuous flow)**
- **Variety (complexity)**

## LSST

DR11 37  $10^9$  objects, 7  $10^{12}$  sources,  
5.5 million 3.2 Gigapixel images  
30 terabytes of data nightly

Final volume of raw image data = 60 PB  
Final image collection (DR11) = 0.5 EB  
Final catalog size (DR11) = 15 PB

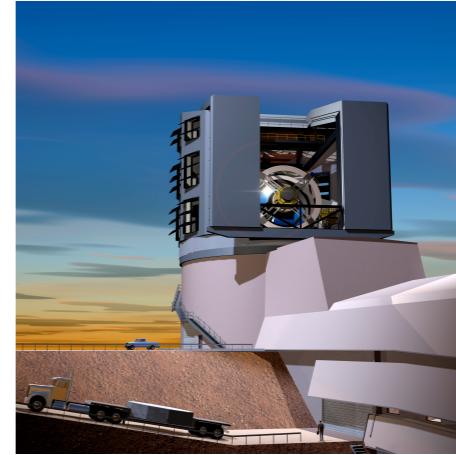


# Big Data in a Nutshell

- Data volume doubles every 18 months (Moore's law)

- LSST:

- 10 years long movie of the sky
- 30'000 GB/night (size of entire SDSS)
- 1'000'000 transient alerts/night in differential imaging (faster and more reliable than catalog cross-matching)
- 50% rubbish
- Human time/attention does **not** scale :-)

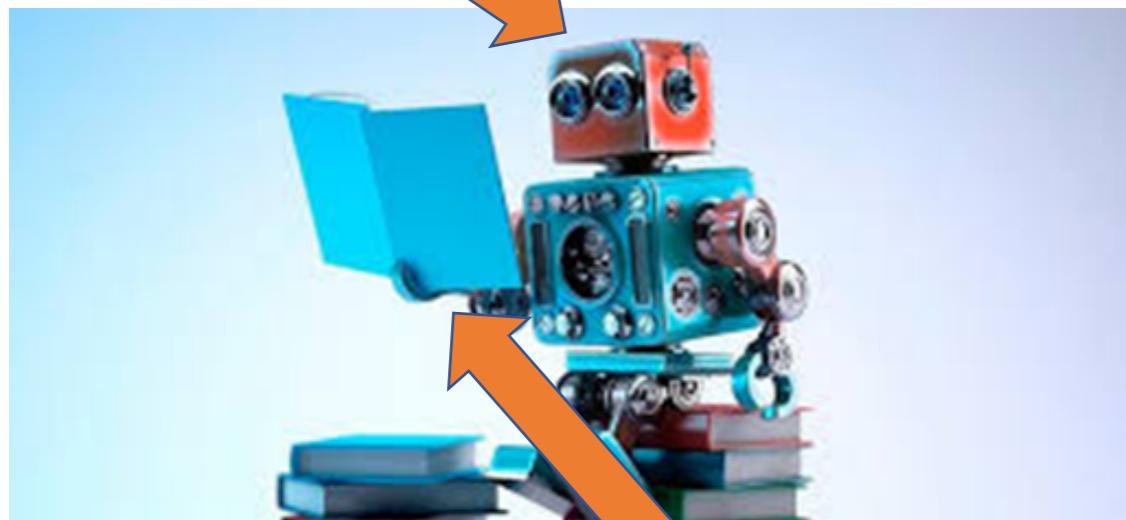


**Automated real-time classification is needed!**

The Fourth Paradigm: Data-Intensive Scientific Discovery

**Era of “Data-driven discoveries”**

Machine...



... Learning?!?

## 'AI IS THE NEW ELECTRICITY'



"Just as electricity transformed almost everything 100 years ago, today I actually have a hard time thinking of an industry that I don't think AI will transform in the next several years."

**Andrew Ng**

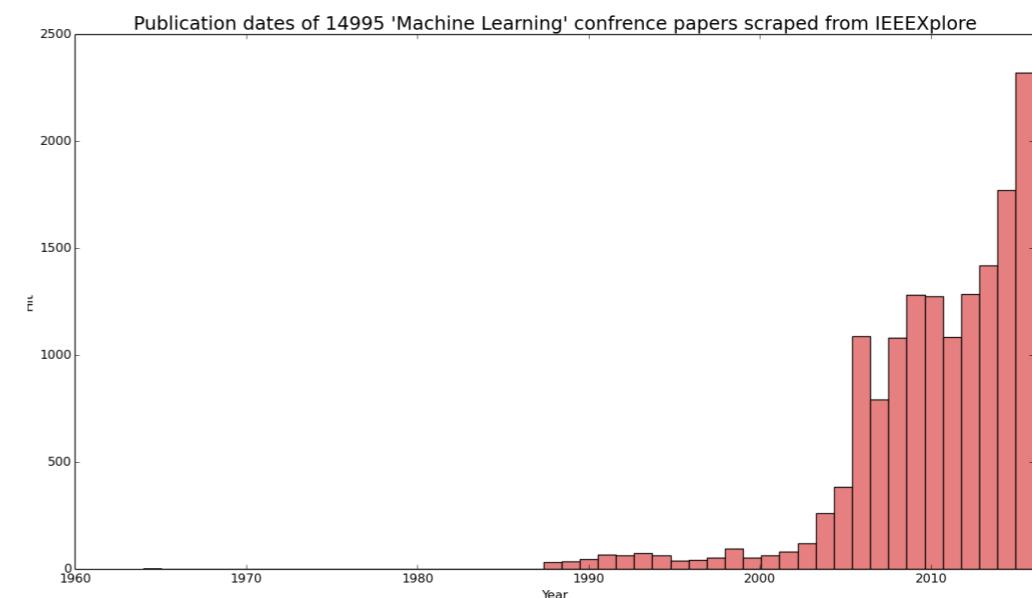
Former chief scientist at Baidu, Co-founder at Coursera

CBINSIGHTS

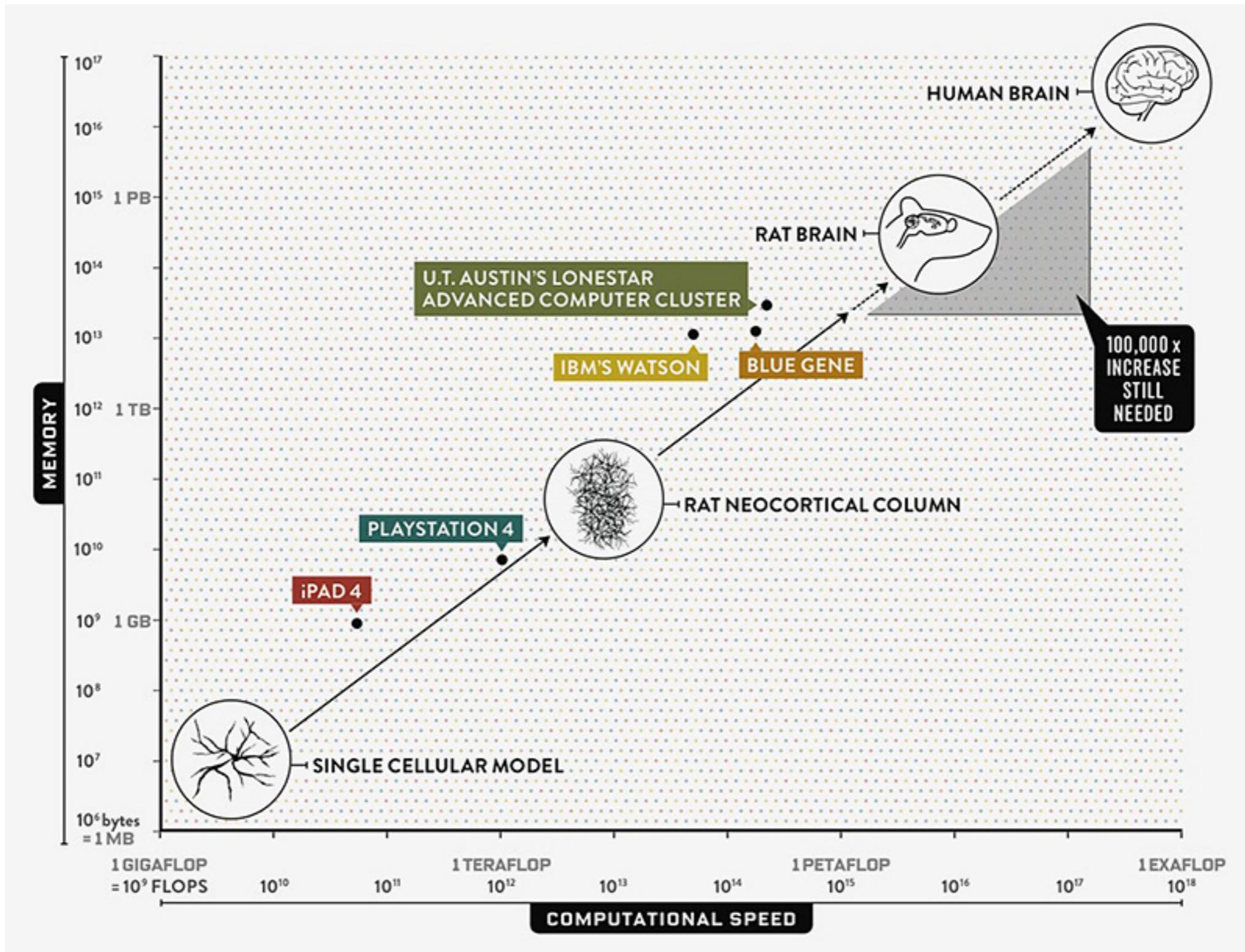
source: <https://www.gsb.stanford.edu/insights/andrew-ng-why-ai-new-electricity>

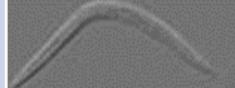
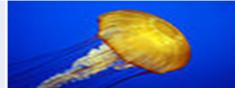
www.cbinsights.com

7



Machine learning algorithms can figure out how to perform tasks by generalising from examples (experience).



Name	# of neurons / # of synapses	Visuals
<i>Caenorhabditis elegans</i>	302	
<i>Hydra vulgaris</i>	5,600	
<i>Homarus americanus</i>	100,000	
<i>Blatta Orientalis</i>	1,000,000	
Nile Crocodile	80,500,000	
<b>Digital Reasoning NN (2015)</b>	<b>~86,000,000 (est.) / 1.6E11</b>	
<i>Rattus Rattatouillensis</i>	200,000,000	
Blue and yellow macaw	1,900,000,000	
Chimpanzee	28,000,000,000	
<b><i>Homo Sapiens Sapiens</i></b>	<b>86,000,000,000 / 1.5E14</b>	
African Elephant	257,000,000,000	

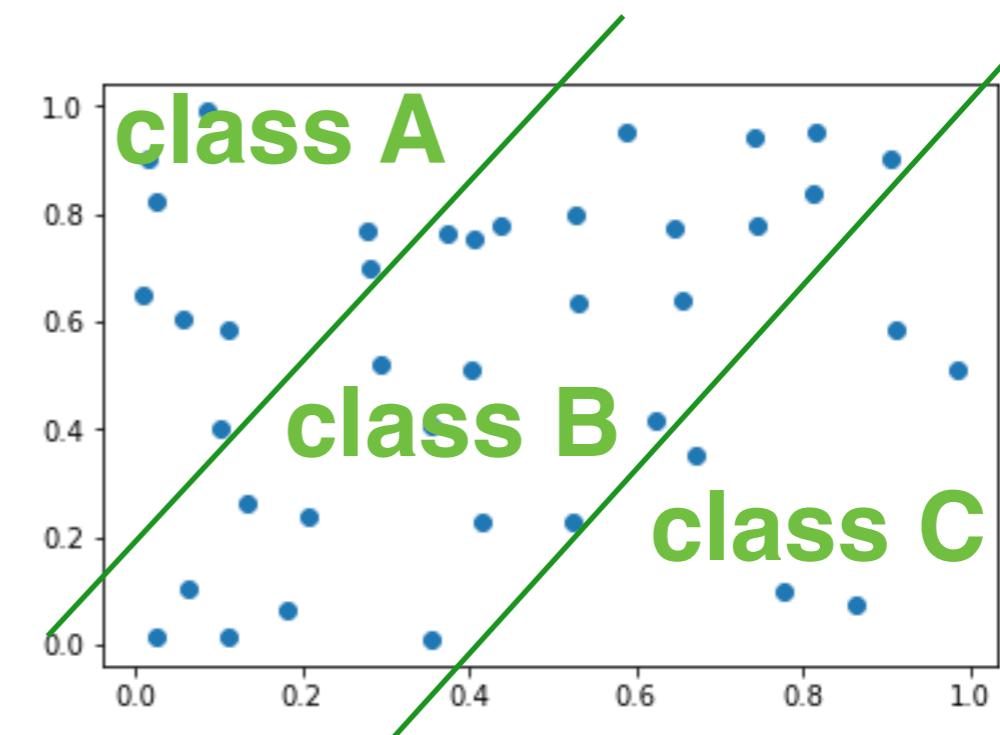
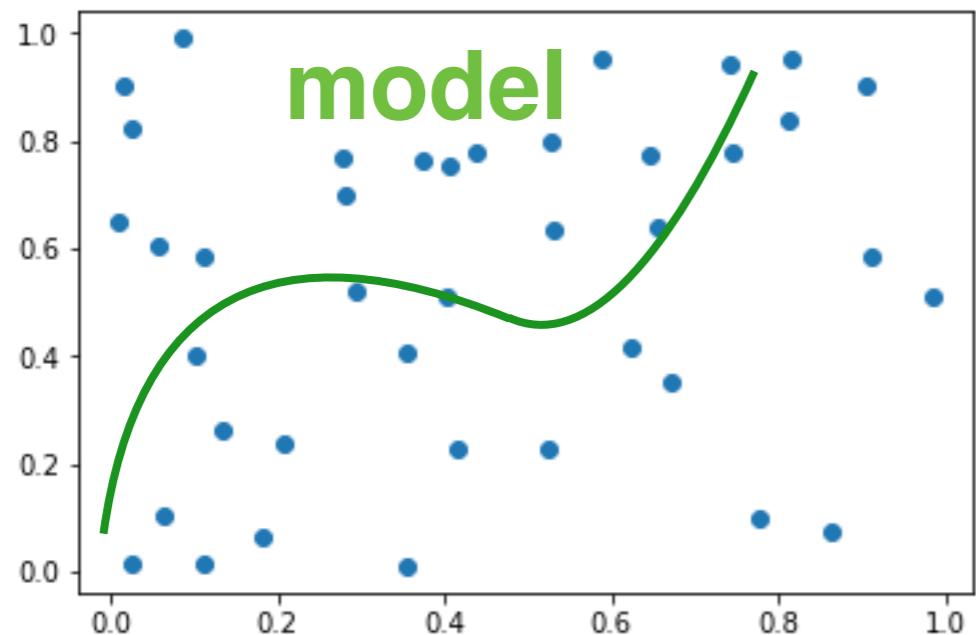
- **Supervised** - classification, regression (fitting)
- **Unsupervised** - clustering, graphs/trees, transformations, typically in multi-dim
- **Semi-supervised**, genetic algorithms, GANs...

*The difference between a physicist and an astronomer:*

*The physicist sees random 2D data and draws a curved line in it saying it's the model that describes the data.*

*The astronomer draws two parallel lines, saying these points belong to class A, these to class B and this is class C.*

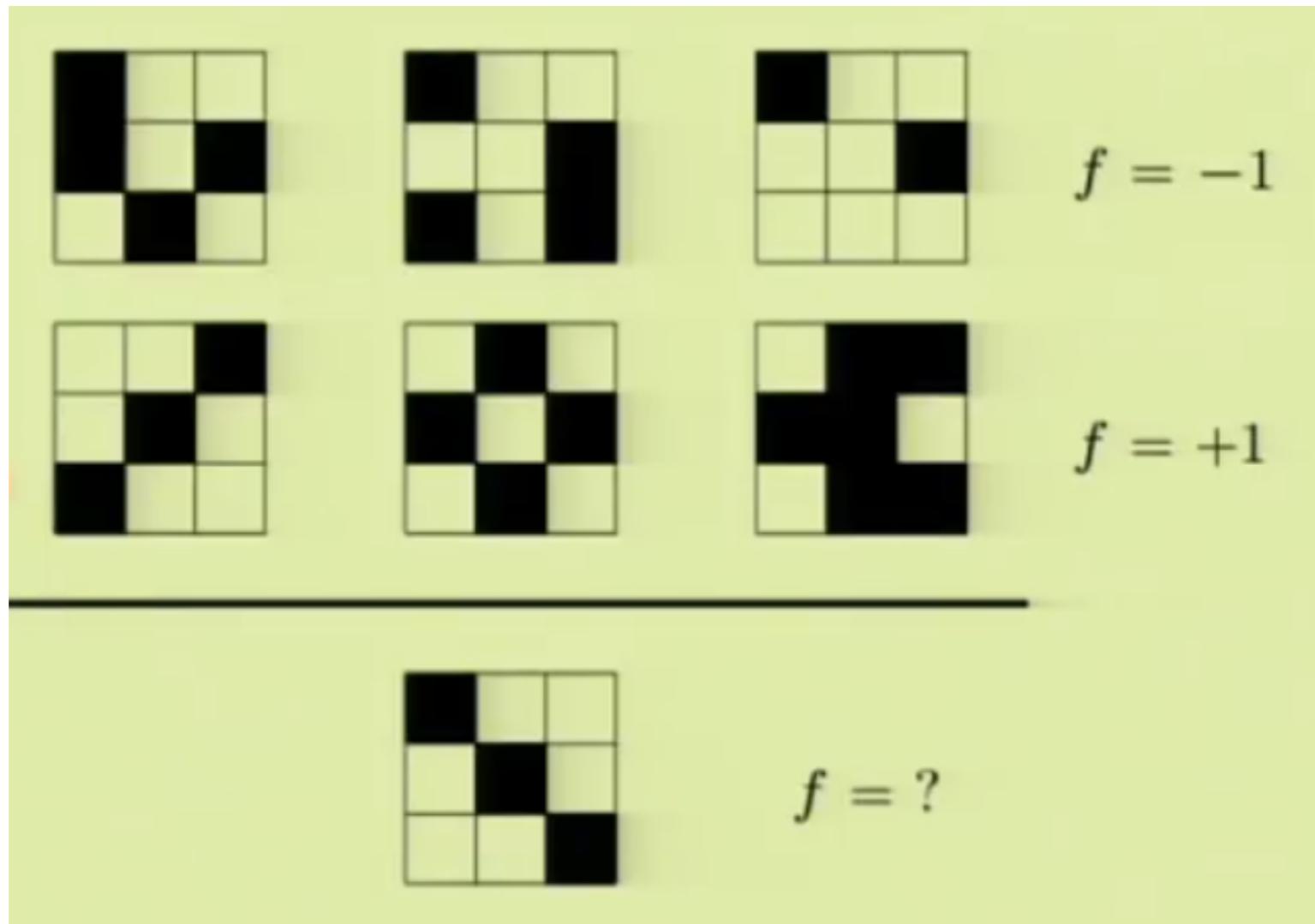
— Andy Lawrence, private communication



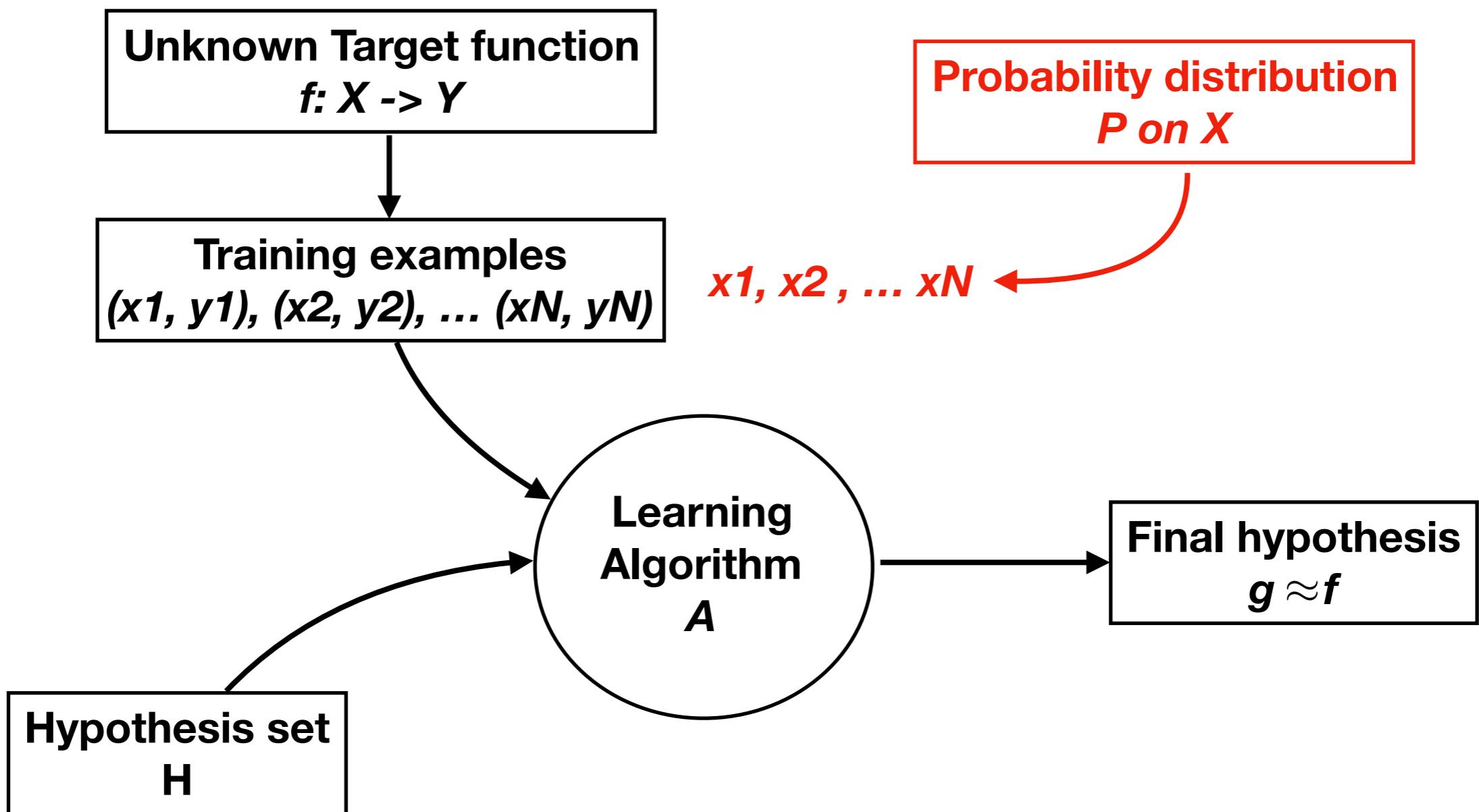
# Essentials of Learning

- Pattern must exist
- Mapping “target function” is unknown or expensive to calculate
- We have the data (and computing resources...)
- Data sample is representative

# Target function...

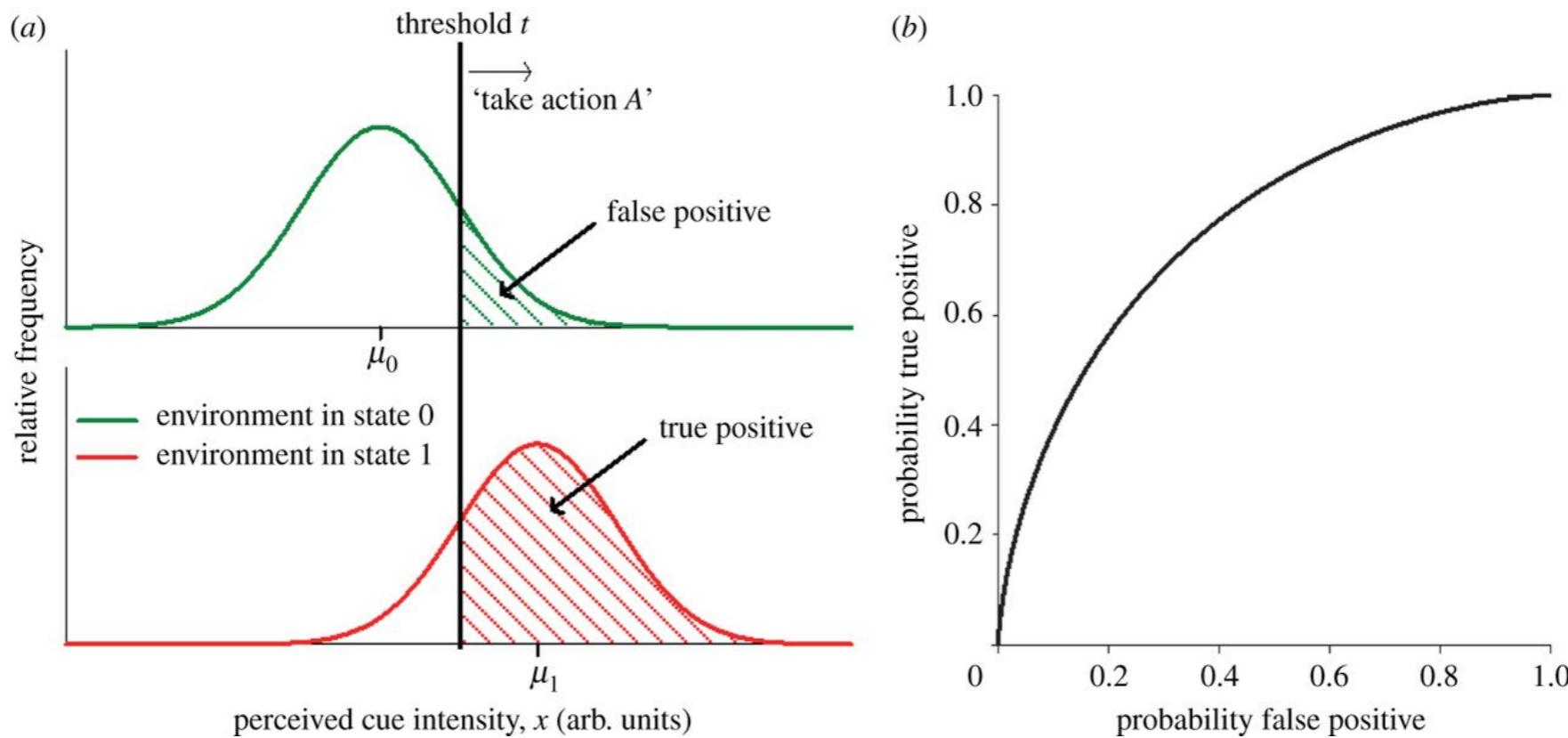


# Learning Diagram



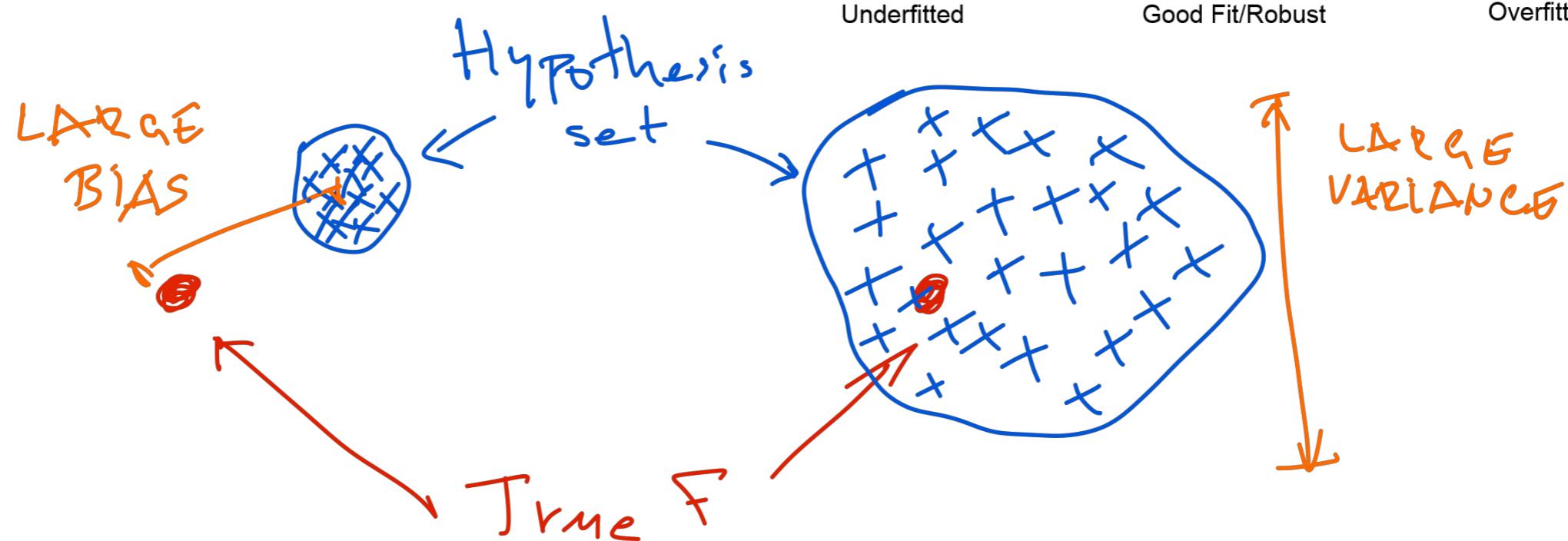
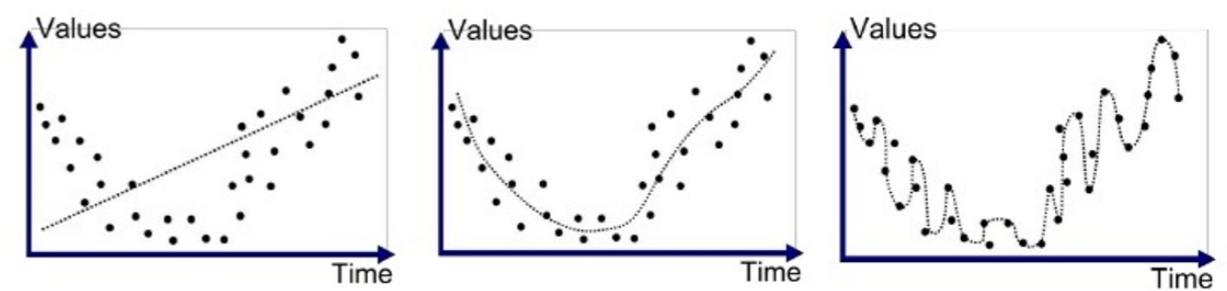
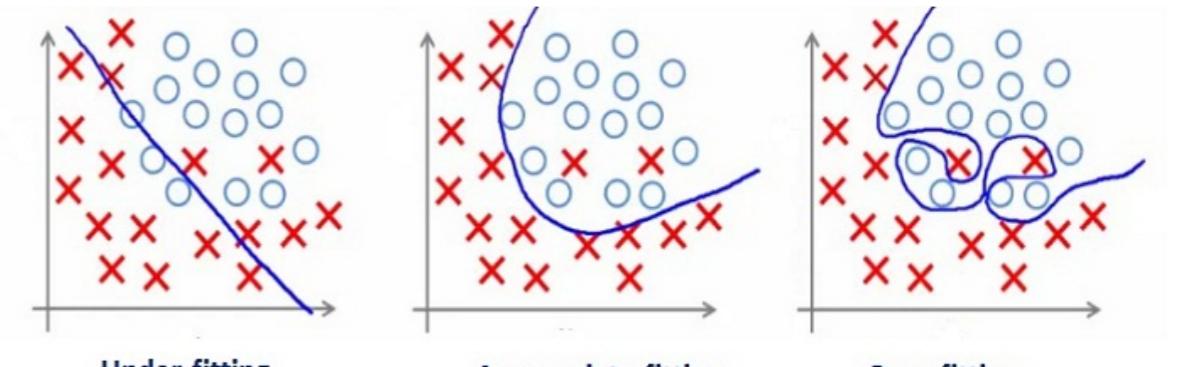
# Loss function

- penalty for being wrong
- imbalanced classes

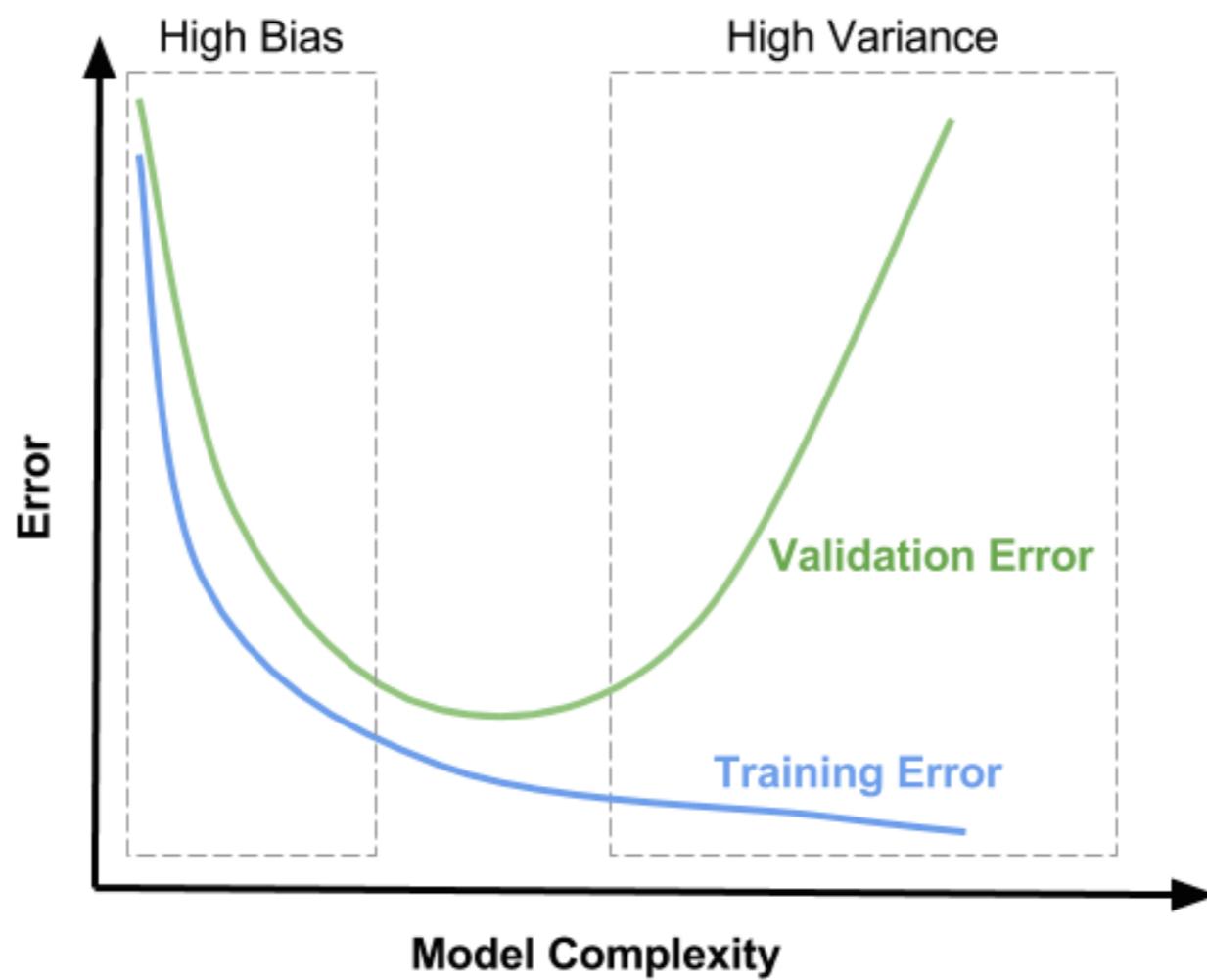


# Bias - Variance Trade-off

- Under-fitting/over-fitting
- in sample error vs out of sample error
- VC dimension

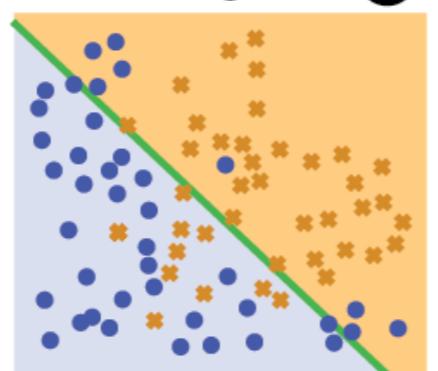


# (cross)-Validation



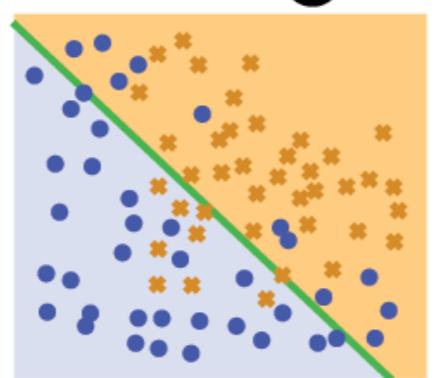
Model 1...

...on Training data. ①



\* 30 \* 10 error: 22.5%  
● 32 ● 8 acc.: 77.5%

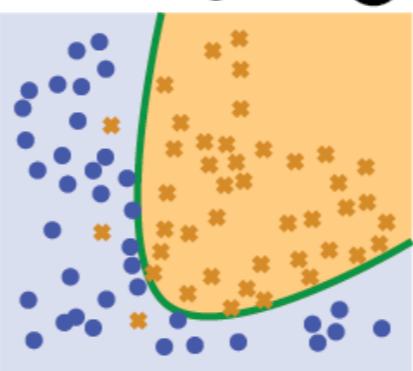
...on Test data. ④



\* 32 \* 8 error: 23.8%  
● 29 ● 11 acc.: 76.2%

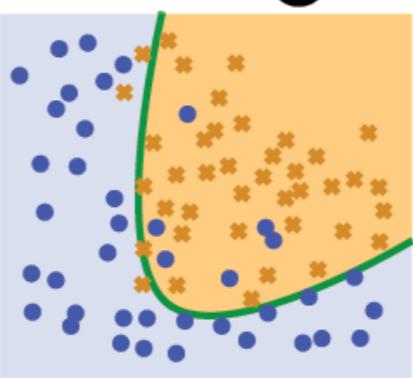
Model 2...

...on Training data. ②



\* 37 \* 3 error: 7.5%  
● 37 ● 3 acc.: 92.5%

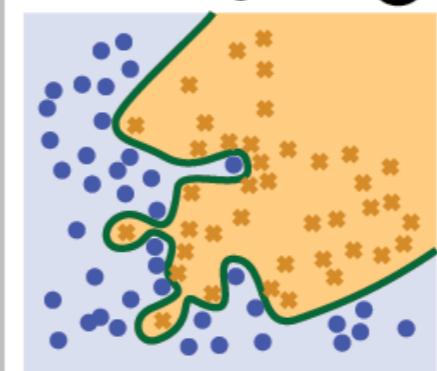
...on Test data. ⑤



\* 37 \* 3 error: 11.3%  
● 34 ● 6 acc.: 88.7%

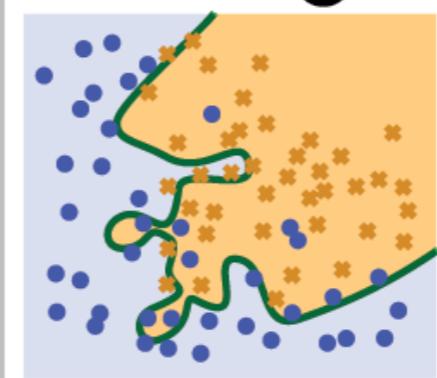
Model 3...

...on Training data. ③

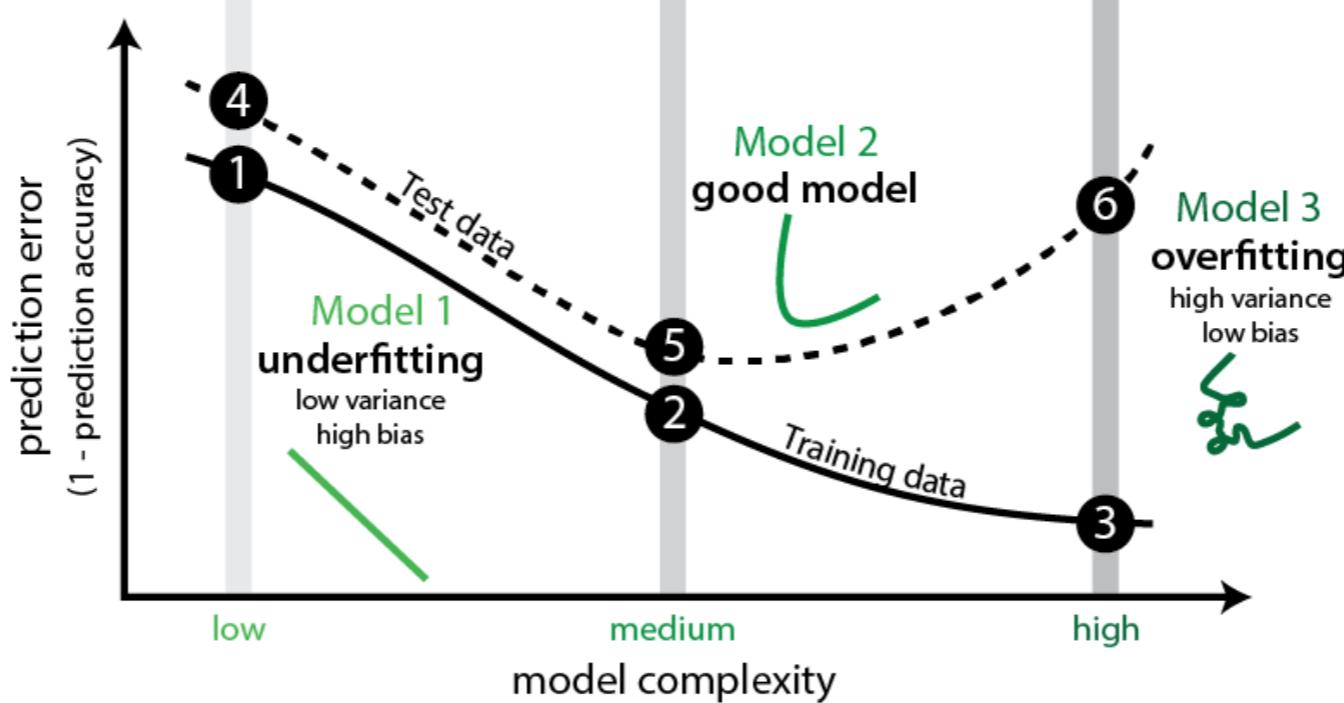


\* 37 \* 0 error: 0%  
● 37 ● 0 acc.: 100%

...on Test data. ⑥



\* 34 \* 6 error: 21.3%  
● 29 ● 11 acc.: 78.7%



	Remedies
High Bias	<ul style="list-style-type: none"> <li>• Train longer</li> <li>• Increase model complexity <ul style="list-style-type: none"> <li>• more features</li> <li>• more parameters,</li> <li>• richer architecture</li> </ul> </li> </ul>
High Variance	<ul style="list-style-type: none"> <li>• Get more data</li> <li>• Decrease model complexity <ul style="list-style-type: none"> <li>• less features</li> <li>• less parameters,</li> <li>• simpler architecture</li> </ul> </li> <li>• Regularization</li> <li>• Early stopping</li> <li>• Drop-out</li> </ul>

#### Data Augmentation:

- When more data are needed, make up new ones! (The way of the god.)
- Translate, rotate, flip, crop, lighten/darken, add noise, dephase, etc.



Cats



Dogs



Machines are lazy and love shortcuts



Dog?

50% dog, 50% cat?



# Correlations != Reasoning



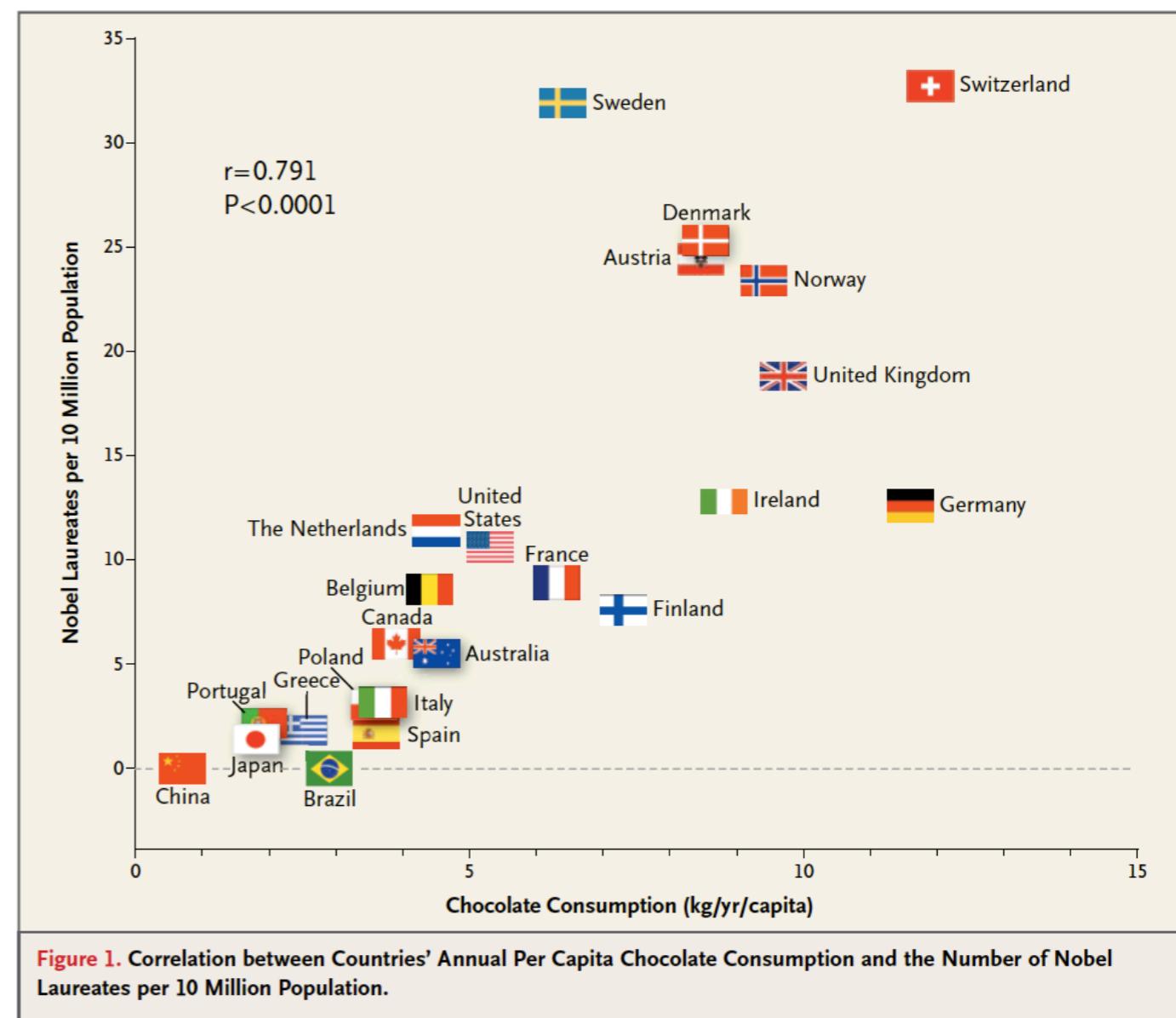
Justin King

@Justinkingnews

Headline: "Women who own horses live longer"

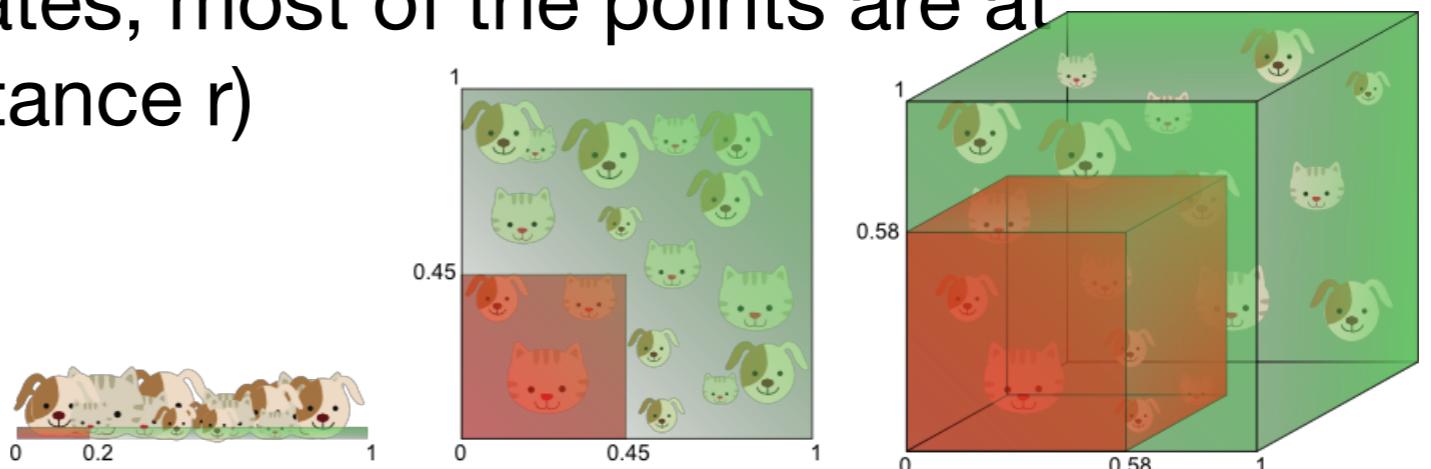
Implied correlation: horses make you live longer.

Reality: if you own a horse, you can probably afford health insurance.



# Curse of Dimensionality

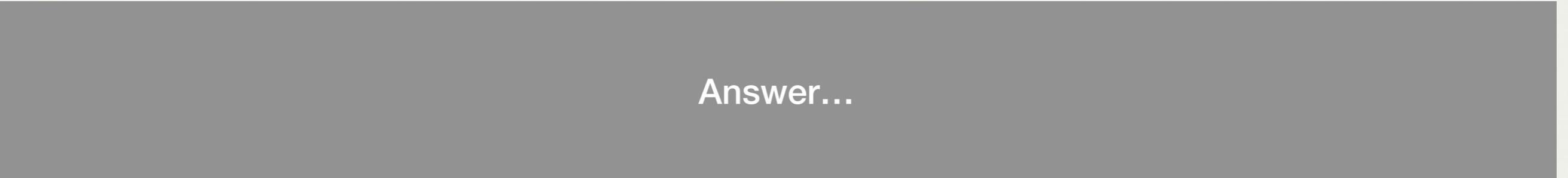
- Too many features
    - Expensive to store
    - Slowing down computation
    - Subject to *Dimensionality curse*
  - Sample space gets harder and harder to fill as dimensions grow
  - A reason why too many features lead to overfitting as data become **sparse**
- 
- “*If people can see in multi-dimensions we would not need machine learning*”
  - More and more data needed to fill the same % of space
  - Distance measure degenerates, most of the points are at the surface of a sphere (distance  $r$ )



# How Many Shades of Gray Can you Distinguish?



Answer...

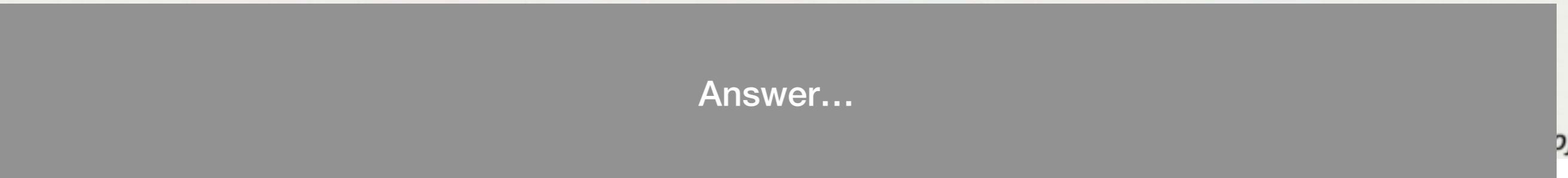


Value encodes continuous variables (less well)

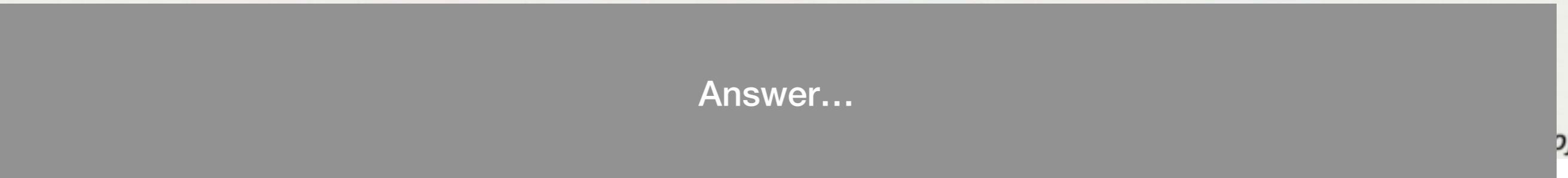
# How Many Colors?



Hue encodes nominal variables

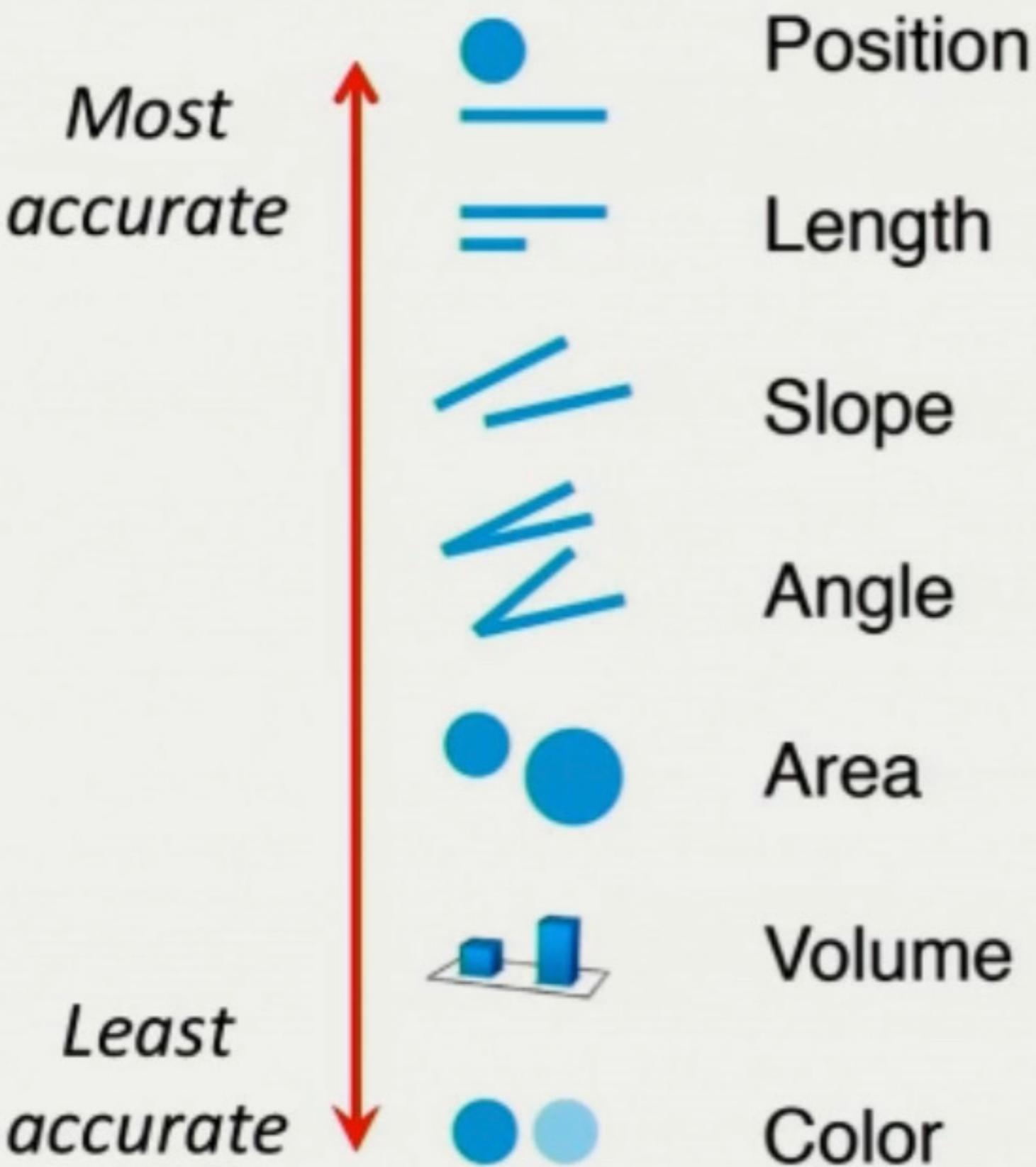


Answer...



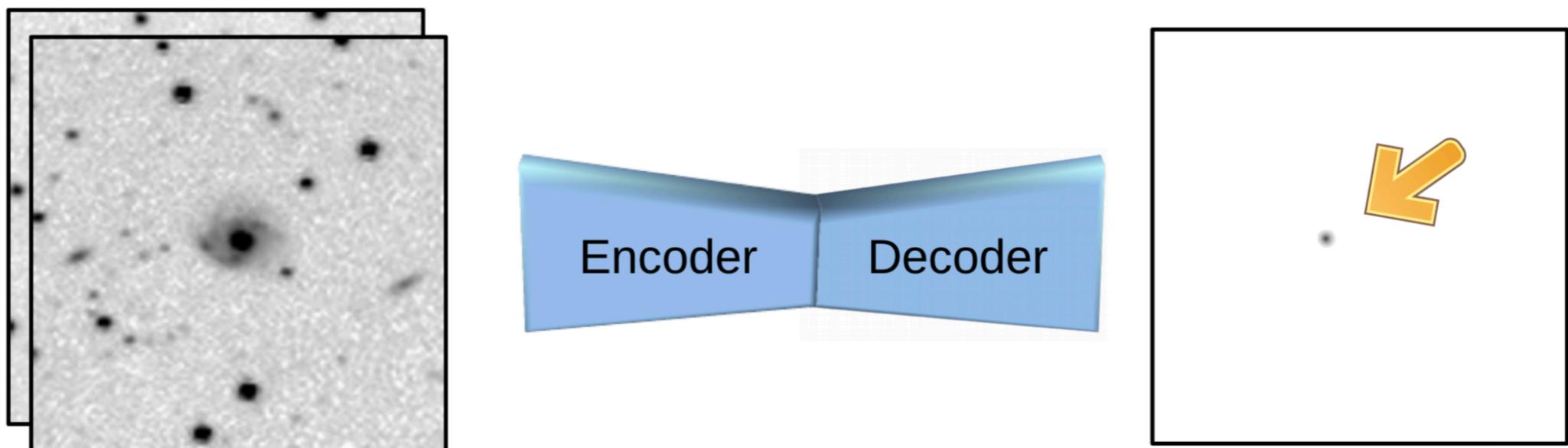
(off)

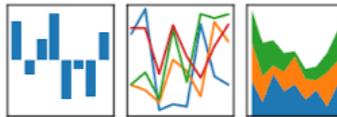
## Relative accuracy of the visualisation space axes:

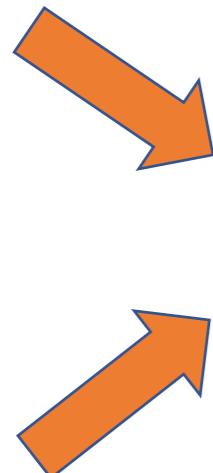


# Hands on Sessions

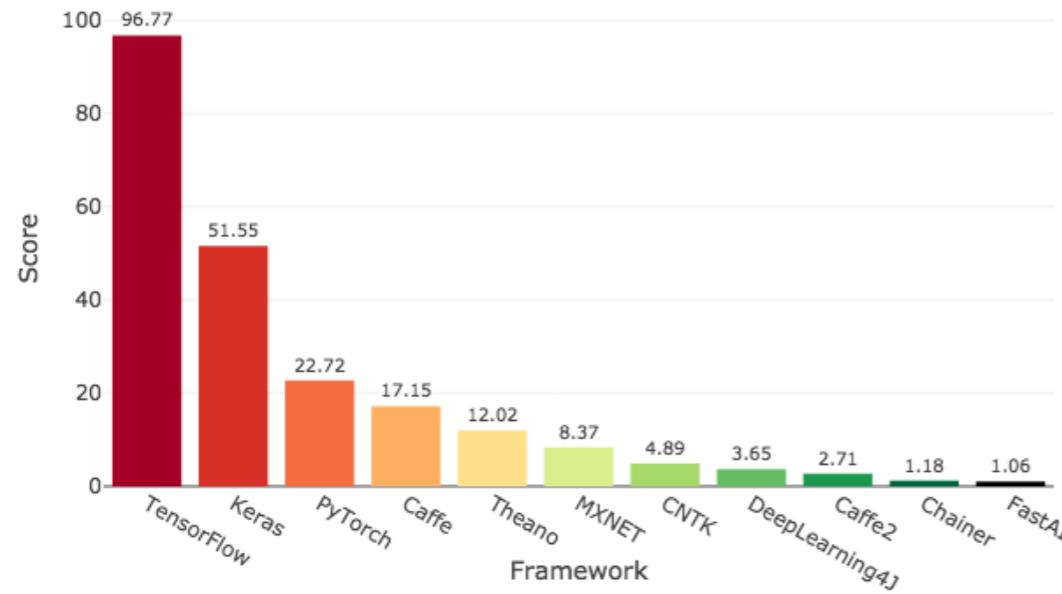
- Galaxy/Star classification (with PCA)
- Regression Photometric redshift with Random Forests
- “voodoo lecture” on Deep Learning with Convolutional Neural Networks on denoising/transient finding



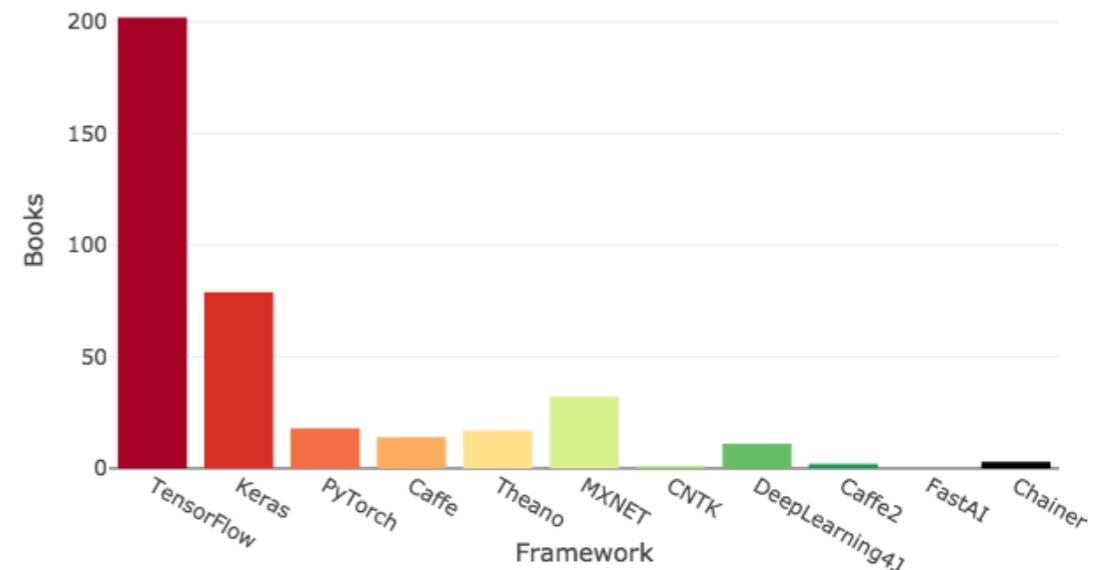
Name	Use	Logo
Pandas	Data Analysis $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$	pandas 
Spark	Distributed Computing	
Scikit-learn	Machine Learning Toolbox	
Keras	Deep Learning	
TensorFlow	Deep Learning	
Open-cv	Computer Vision	



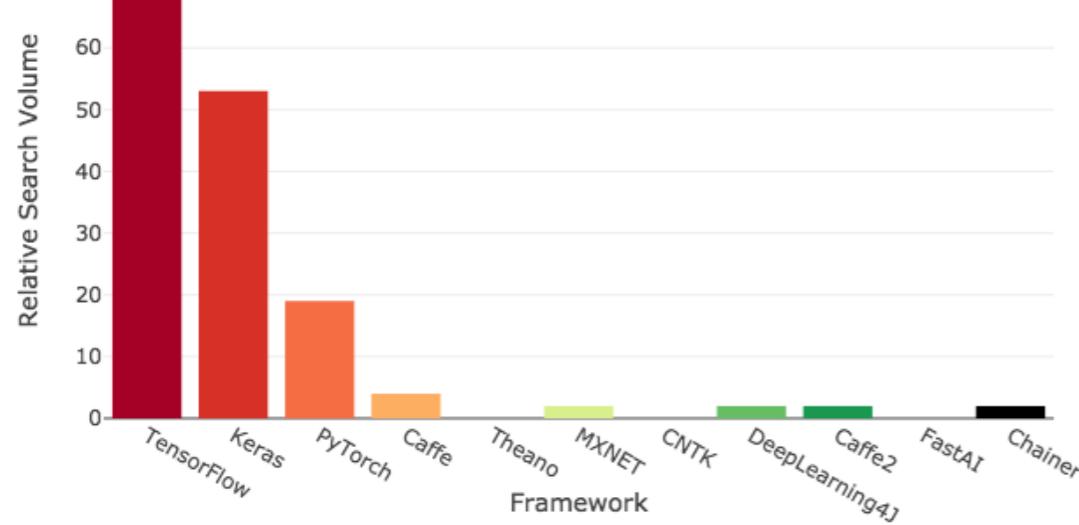
Deep Learning Framework Power Scores 2018



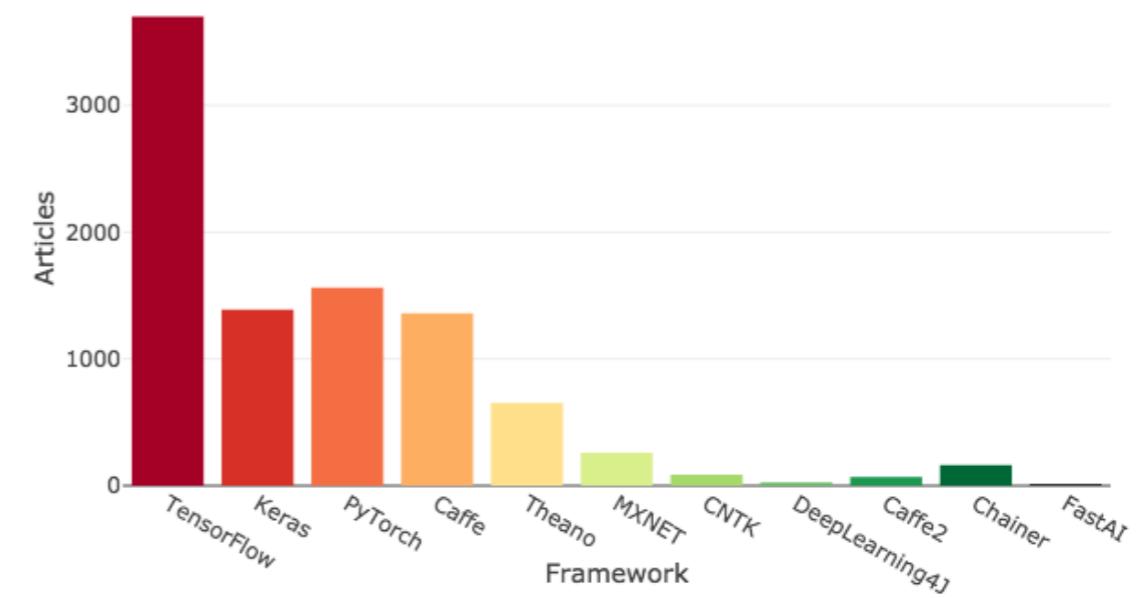
Amazon Books



Google Search Volume



ArXiv Articles



# Software prerequisites

- Python 3
- Anaconda [anaconda.org](https://anaconda.org) with Python 3.7
- You may create an environment  
`conda create --name=ml python=3.6`
- *scikit-learn, pandas, scikit-image, seaborn, tensorflow, keras* (may induce a downgrade to Python 3.6, but it's ok)  
`conda install xxx`
- Cheat sheets...

# ML seminar

- Aka astro-ph
- an algorithm overview, usage
- publications together?
- 1x every 2 weeks or month?