



Superintelligence

Swiss Study Foundation

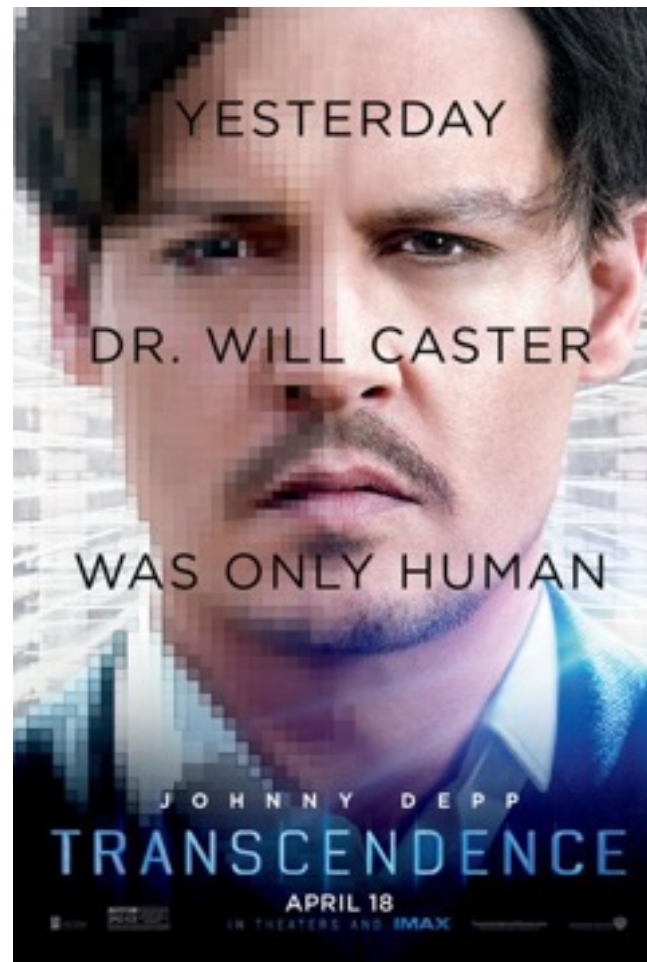
Kaspar Etter, kaspar.etter@gbs-schweiz.org
Adrian Hutter, adrian.hutter@gbs-schweiz.org

Solothurn, 12 June 2015
www.superintelligence.ch

Robin Li, Bill Gates, Elon Musk on AI



More Money for Entertainment



... than ensuring a good outcome!



Intelligence and Rationality

What are we talking about?

What is Intelligence?

«Innumerable tests are available for measuring intelligence, yet no one is quite certain of what intelligence is, or even just what it is that the available tests are measuring.» – R. L. Gregor

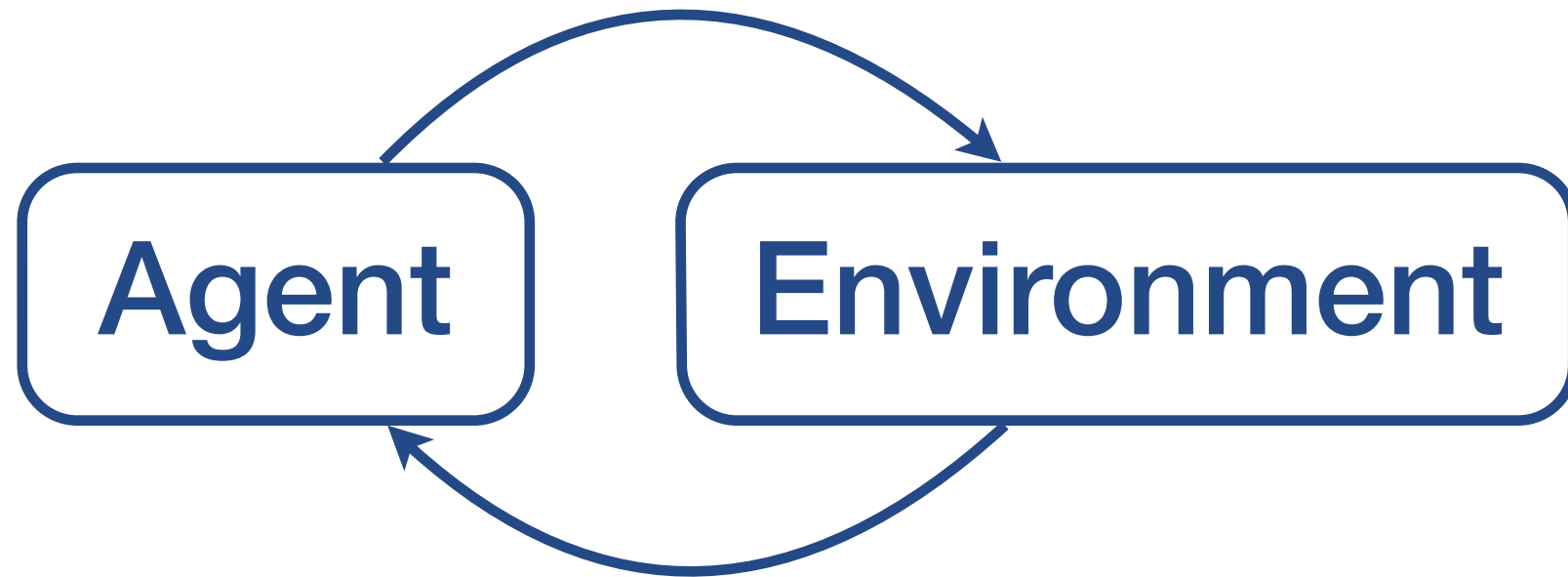
«The capacity to learn or to profit by experience.» – W. F. Dearborn

A Useful Definition

«Intelligence measures an agent's ability to achieve its goals in a wide range of unknown environments.»

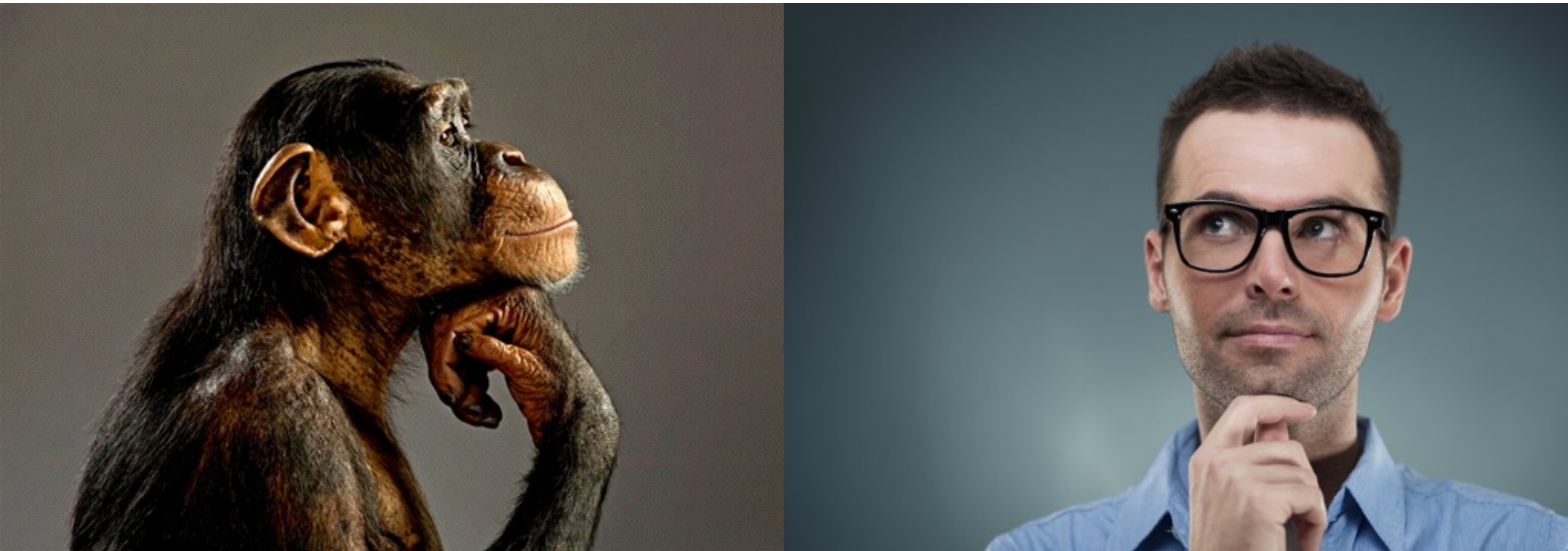
$$\text{Intelligence} = \frac{\text{Optimization Power}}{\text{Comp. Resources}}$$

Required Ingredients



Learn, predict, rate and plan!

Intelligence is a Big Deal



6 million years ago, 96% common DNA

Fast-Evolving Human DNA Leads to Bigger-Brained Mice **Superintelligence**
[phenomena.nationalgeographic.com/2015/\[...\]](http://phenomena.nationalgeographic.com/2015/[...]) Intelligence and Rationality

Technology = (Neutral) Lever

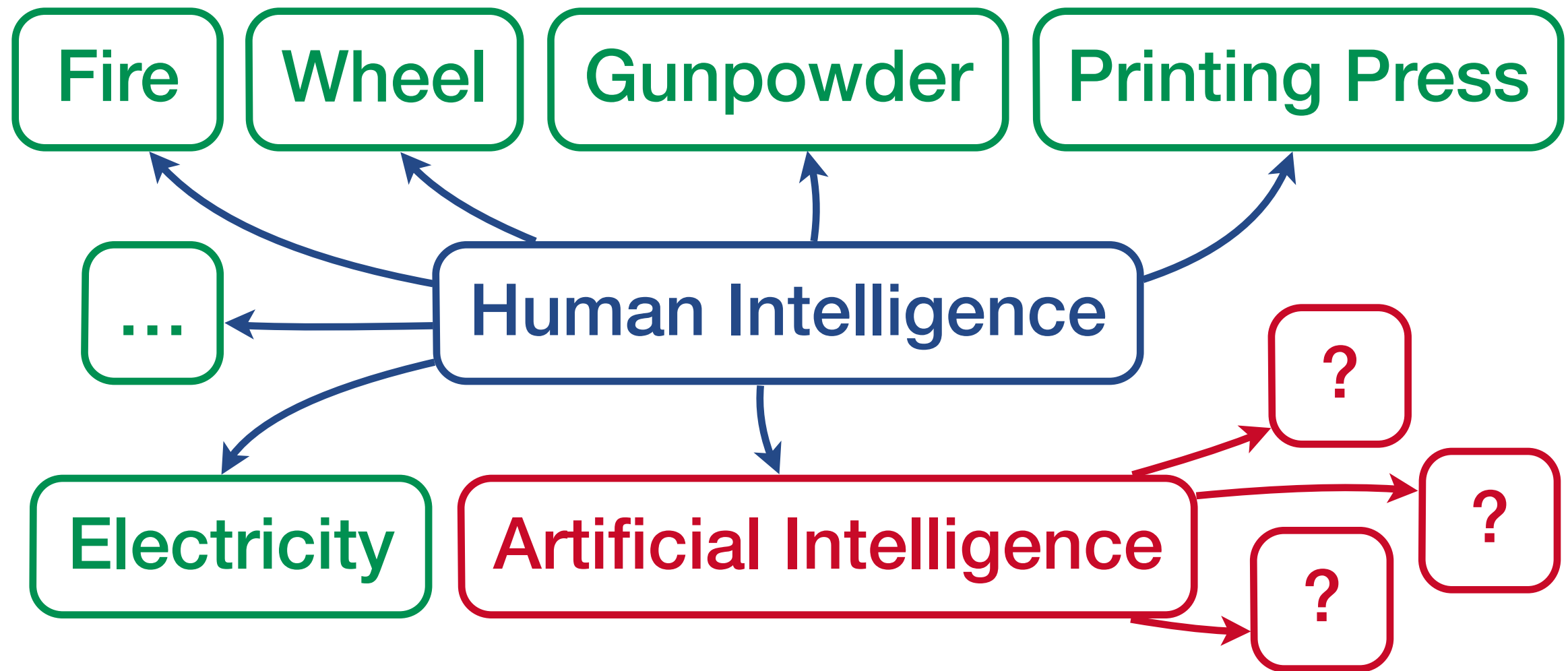
Greed and laziness drive us to

- increase our productivity
 - make our lives easier
- ... which is fine except:



**Our technological progress far
outperforms our moral progress!**

Artificial Intelligence



Intelligence is a technology like no other!

Rationality

**It's about
achieving
your goals!**



Rationality is the science of winning:

- **Epistemic rationality:** accurate beliefs
- **Instrumental rationality:** good strategy

What do we mean by “rationality”?

lesswrong.com/lw/31/what_do_we_mean_by_...

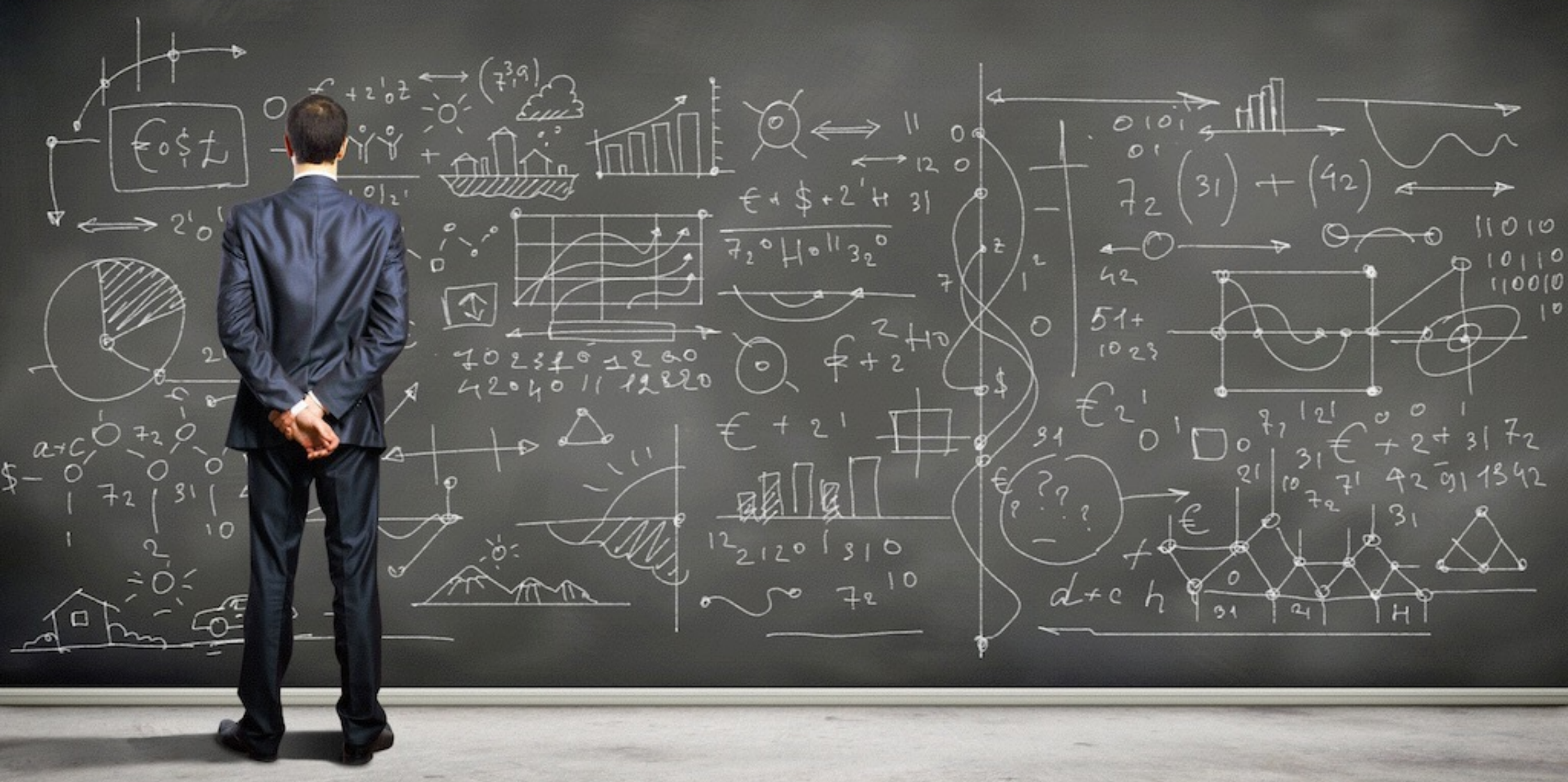
Superintelligence

Intelligence and Rationality

Rationality Continued

- **Normative Rationality:** probability theory, decision theory, game theory
- **Descriptive Rationality:** cognitive biases, system 1 and system 2
- **Prescriptive Rationality:** applied rationality; train:





Probability Theory

What can we know about the world?

What are Probabilities?

Probabilities can be interpreted as ...

- the frequency with which an event occurs if an experiment is repeated a large number of times (**Frequentist**)
- or as the credence in a statement (**Bayesian or subjective probability**).

What is the P that the 9th digit of π is 7?

Relevance of Probabilities

What bets should we take, given our uncertainty about the world we live in?

**Living is
acting,
acting is
betting!**



Properties of Probabilities

- For all propositions p : $0 \leq P(p) \leq 1$
- If p is certainly true: $P(p) = 1$
- If p and q are mutually exclusive:
$$P(p \cup q) = P(p) + P(q)$$
- Consequences: $P(\neg p) = 1 - P(p)$
- And: $P(p \cup \neg p) = P(p) + P(\neg p) = 1$

Conditional Probabilities

$P(A | B)$ is the probability of A given B.

Example:

Alice has two children. You learn either

a) that the older child is a girl – or

b) that at least one of them is a girl.

What is the P of two girls in a) and b)?

Joint Probabilities

$P(A, B)$: probability that A and B occur

$P(A, B)$ not determined by $P(A)$ & $P(B)$:

A = “she voted for SP”, B = “... for SVP”

$P(A) = 27\%$, $P(B) = 19\%$, $P(A, B) = 0\%$

A = “older than 55”, B = “older than 65”

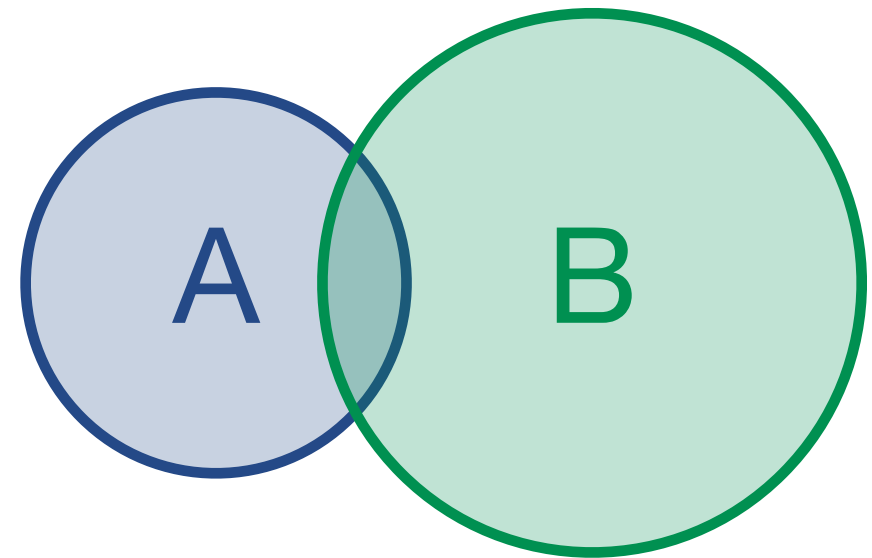
$P(A) = 27\%$, $P(B) = 19\%$, $P(A, B) = 19\%$

Joint and Conditional Probabil.

$$P(A, B) = P(A | B) * P(B)$$

$$P(B, A) = P(B | A) * P(A)$$

$$P(A, B) = P(B, A)$$



$$P(A | B) * P(B) = P(B | A) * P(A)$$

Bayes' Theorem

- Thomas Bayes
- 1701 – 1761
- English Statistician



$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Testing for a Disease

- 1% of people have a certain disease
- 80% with disease are tested positively
- 9.6% w/o disease are tested positively
- Given a positive test, how likely is “it”?
- Out of 100 people, 1 are ill, 99 aren't
- Given positive test, only with 7.8% ill

Bayesian Inference

How to update your beliefs:

$$P_{new}(h) = P_{old}(h) \cdot \frac{P_{old}(e|h)}{P_{old}(e)}$$

The confidence in a hypothesis **h** increases if the evidence **e** is more likely to happen given **h** than without.

Base Rate Fallacy

We tend to replace $P(h|e)$ with $P(e|h)$ instead of updating the prior belief!

$$P_{new}(h) = P_{old}(h) \cdot \frac{P_{old}(e|h)}{P_{old}(e)}$$

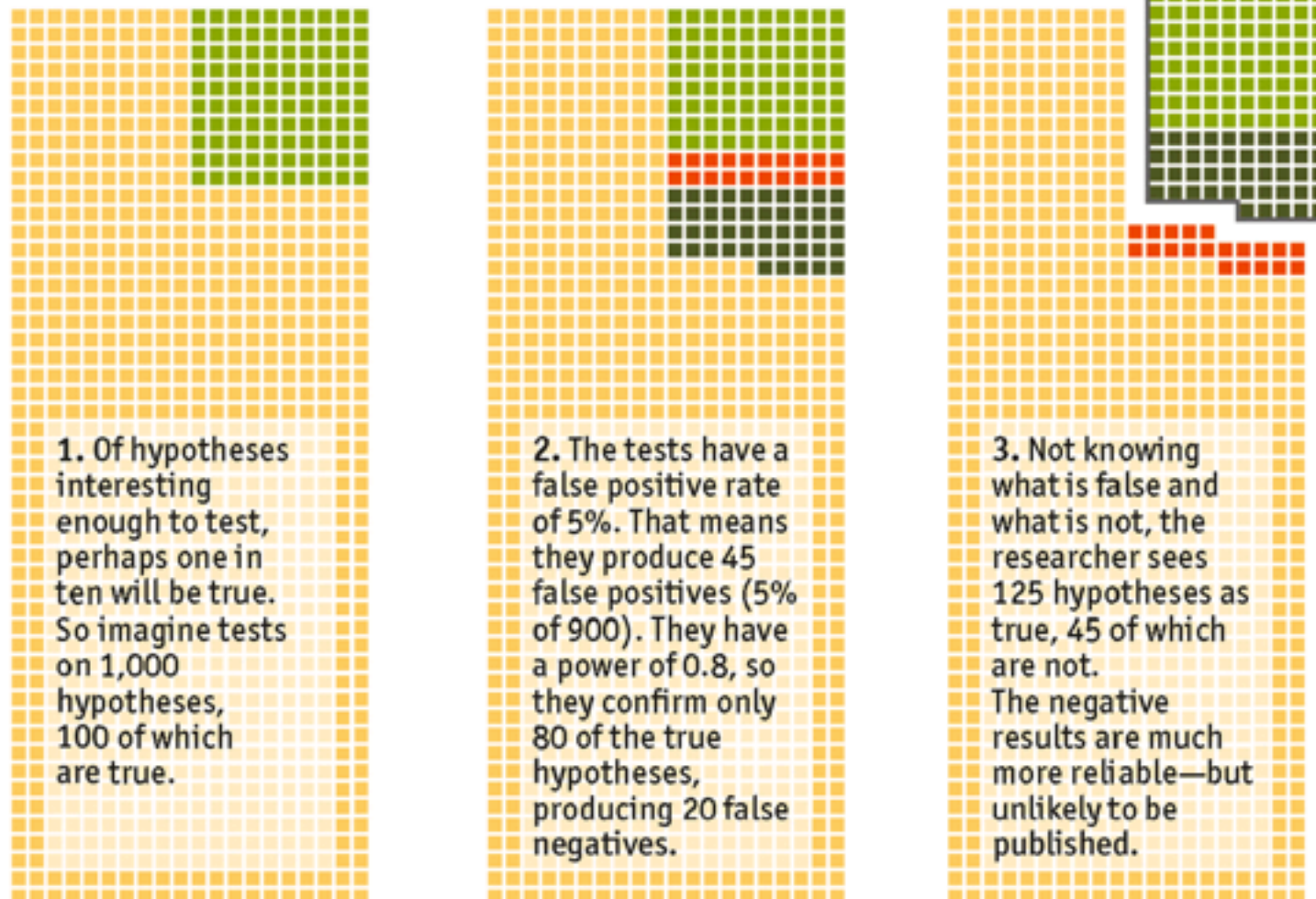
**Extraordinary claims require
extraordinary evidence!**

Broken Science

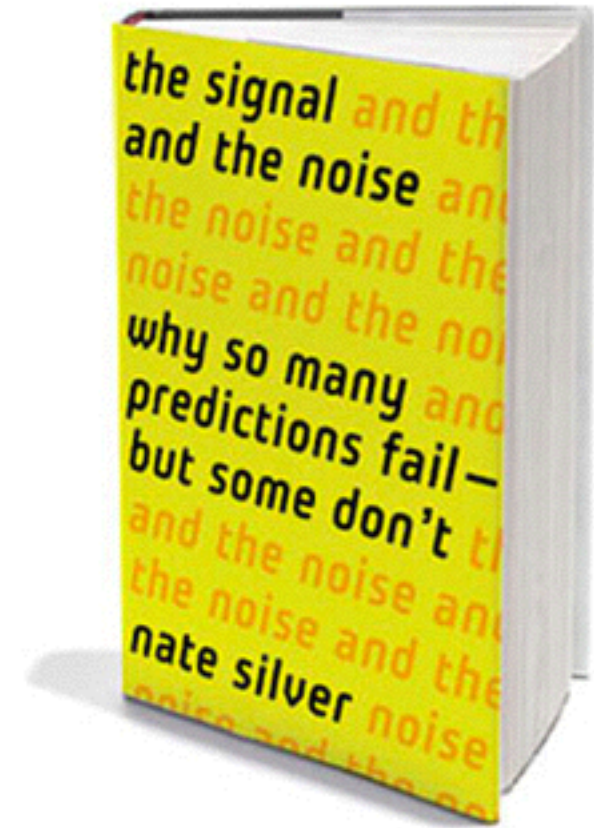
Unlikely results

How a small proportion of false positives can prove very misleading

False True False negatives False positives



Source: *The Economist*



**Null-Hypothesis
Testing: $P(e|h)$
instead of $P(h|e)$!**

Conservation of Exp. Evidence

$$P(h) = P(h \cap e) + P(h \cap \bar{e})$$

$$P(h) = P(h|e) \cdot P(e) + P(h|\bar{e}) \cdot P(\bar{e})$$

$$P(e) + P(\bar{e}) = 1$$

$$\Rightarrow P(h|\bar{e}) \leq P(h) \leq P(h|e)$$

**Absence of evidence is
evidence of absence!**

Prior Probability

Where do we get the priors from?

- Irrelevant given enough evidence
- Knowledge about the statistics
- Symmetries (even distribution)

But:

- What is the prior that Zeus exists?
- The prior that the universe is finite?

Solomonoff Induction

Ray Solomonoff, 1960: Universal Prior
Combines and formalizes ideas by

- **Epicurus:** “Keep all hypotheses that are consistent with the data.”
- **Ockham:** “Among all hypotheses consistent with the observations, choose the simplest.” (O.’s Razor)

Kolmogorov Complexity

Simplicity: Not how easy something is to understand for humans – but rather the length of the shortest program that is able to reproduce the observations.

The Solomonoff prior is exponentially small in this length & not computable.
(Formal def. involves *Turing machines*.)



Decision Theory (Part 1)

How shall we decide in the face of uncertainty?

Utility Function

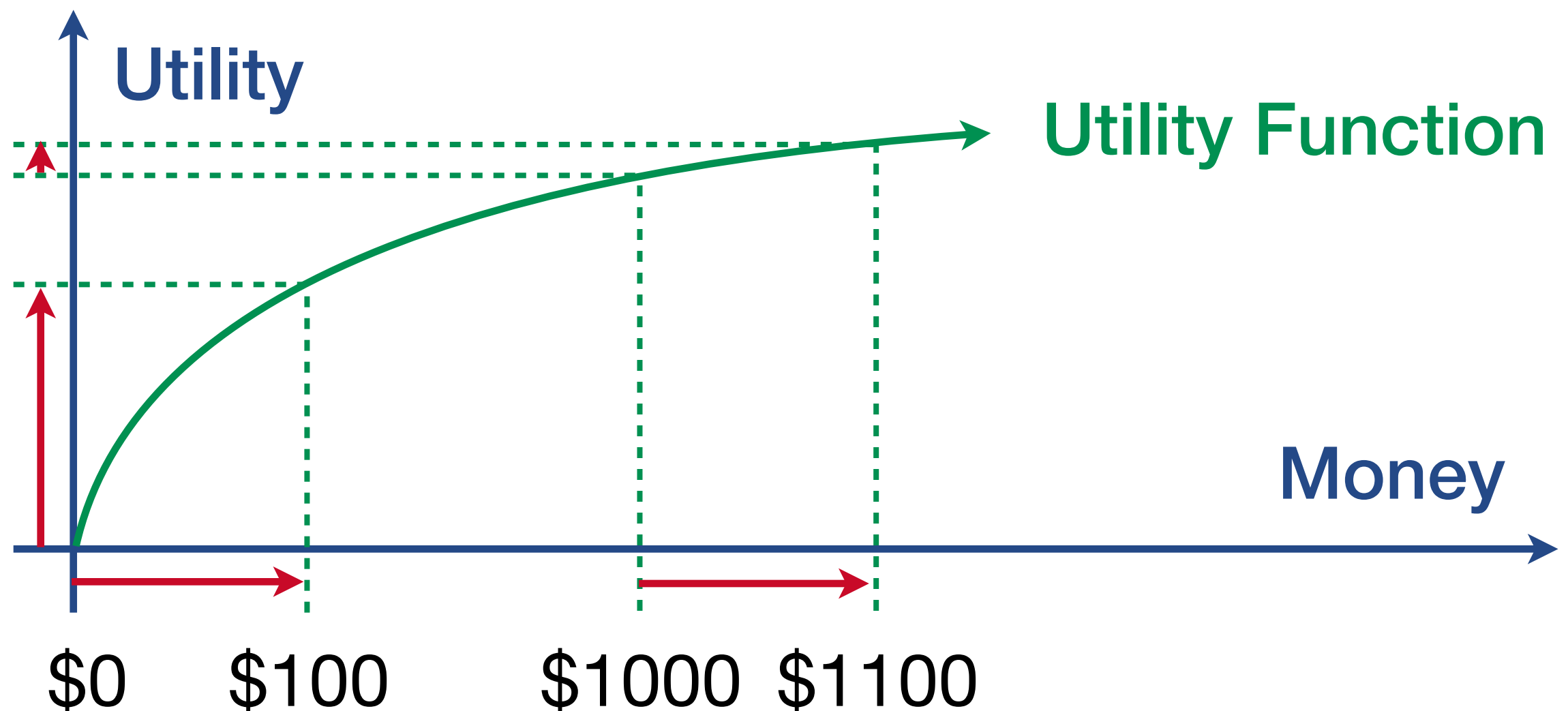
A utility function describes an agent's **preferences** over different outcomes.

Being nourished is a state of **higher utility** for animals than being starving.

A utility function doesn't have to be an **explicit or conscious goal** of an agent.

Marginal Utility

Money has diminishing marginal utility:



Expected Utility

Under **uncertainty** (i.e. in reality) we can only maximize expected utilities.

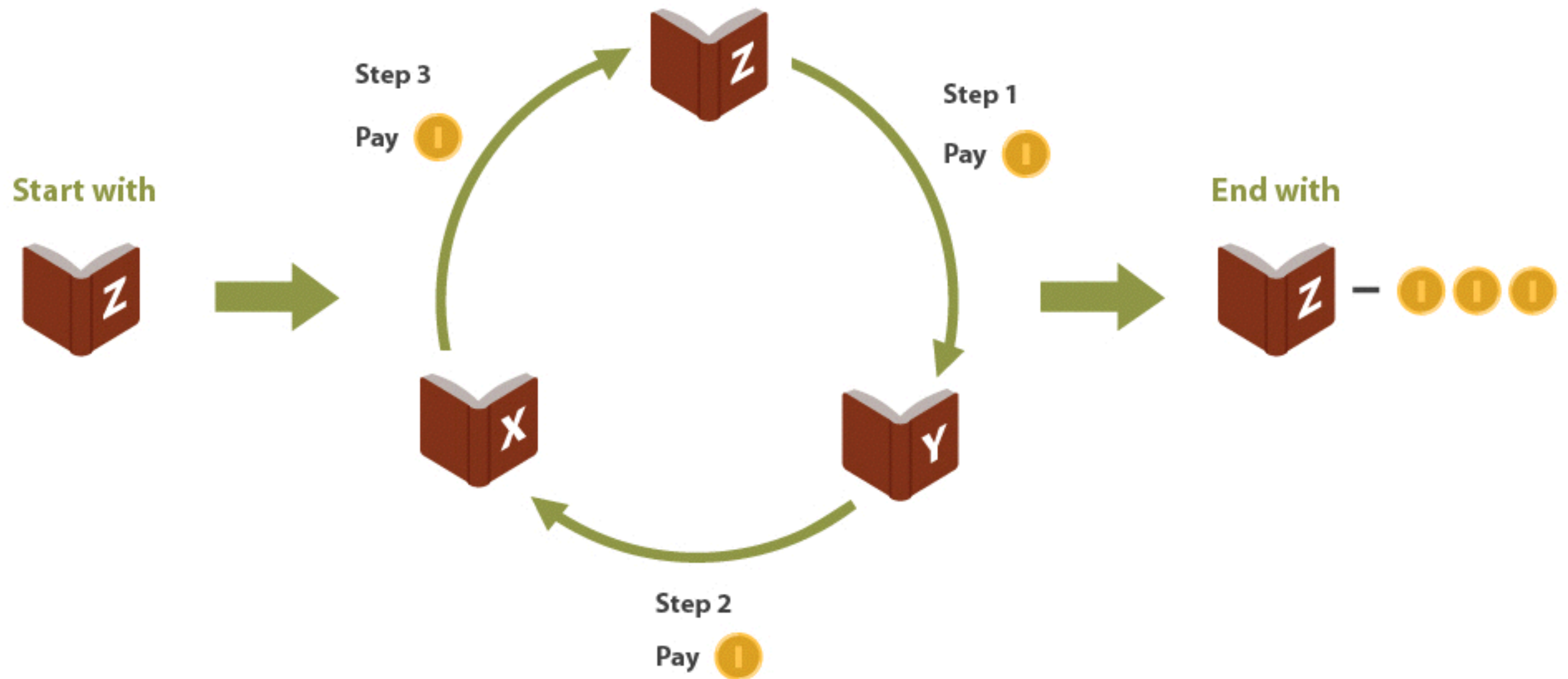
If there are 3 outcomes with **utilities** u_1, u_2, u_3 and **probabilities** p_1, p_2, p_3 , the expected utility is calculated as $p_1 * u_1 + p_2 * u_2 + p_3 * u_3$.

Von Neumann–Morgenstern

Von Neumann and Morgenstern, 1947:
Every *rational* agent (i.e. one satisfying four natural axioms) behaves *as if* it tries to maximize the expected utility of some utility function.

Axioms: Completeness, Transitivity, Continuity and Independence (VNM).

Money-Pump Argument



Loss Aversion

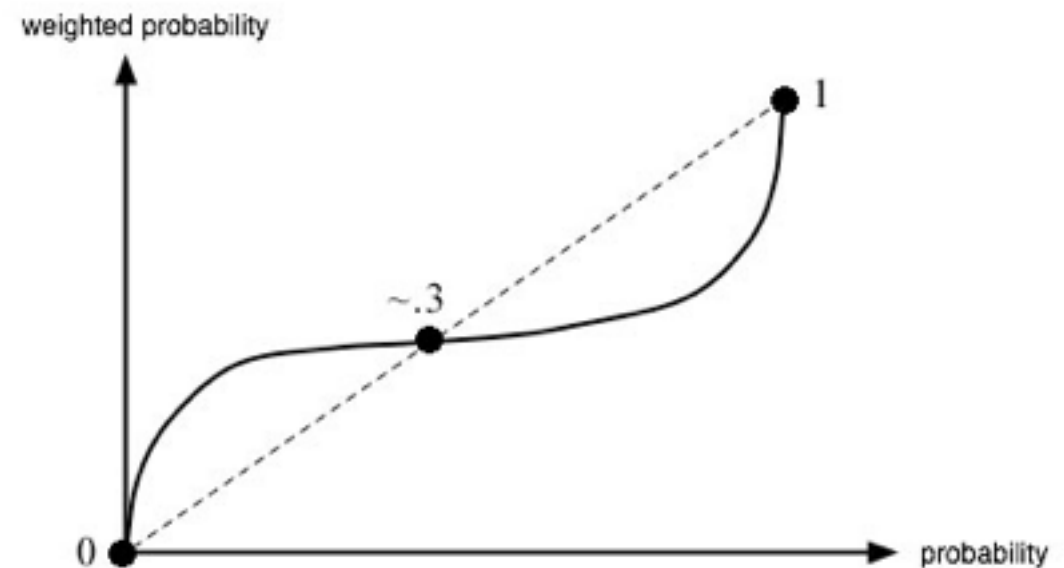
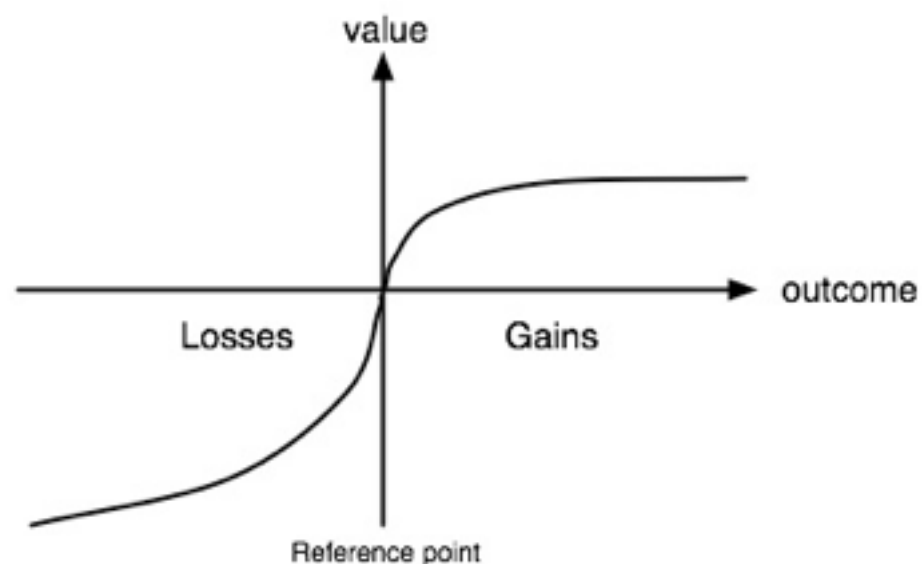
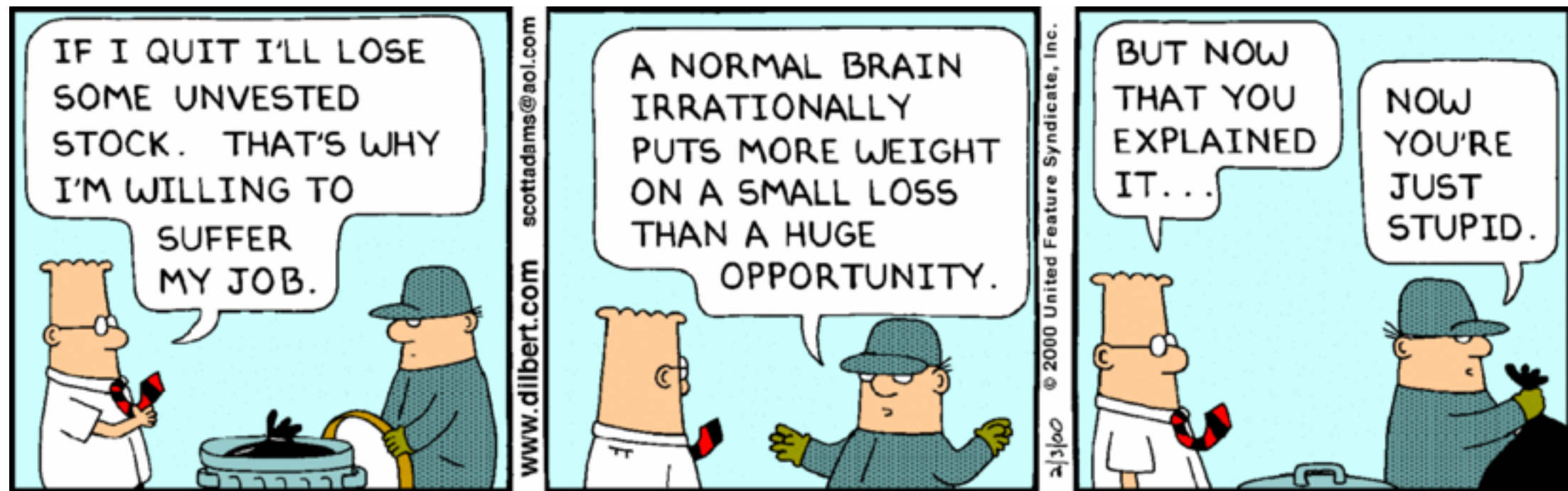
You have been given \$1000. Choose:

- Win \$1000 with 50% (risky)
- Win \$500 with certainty (safe)

You have been given \$2000. Choose:

- Lose \$1000 with 50% (risky)
- Lose \$500 with certainty (safe)

Prospect Theory



Tversky/Kahneman:

A: 72%, B: 28%

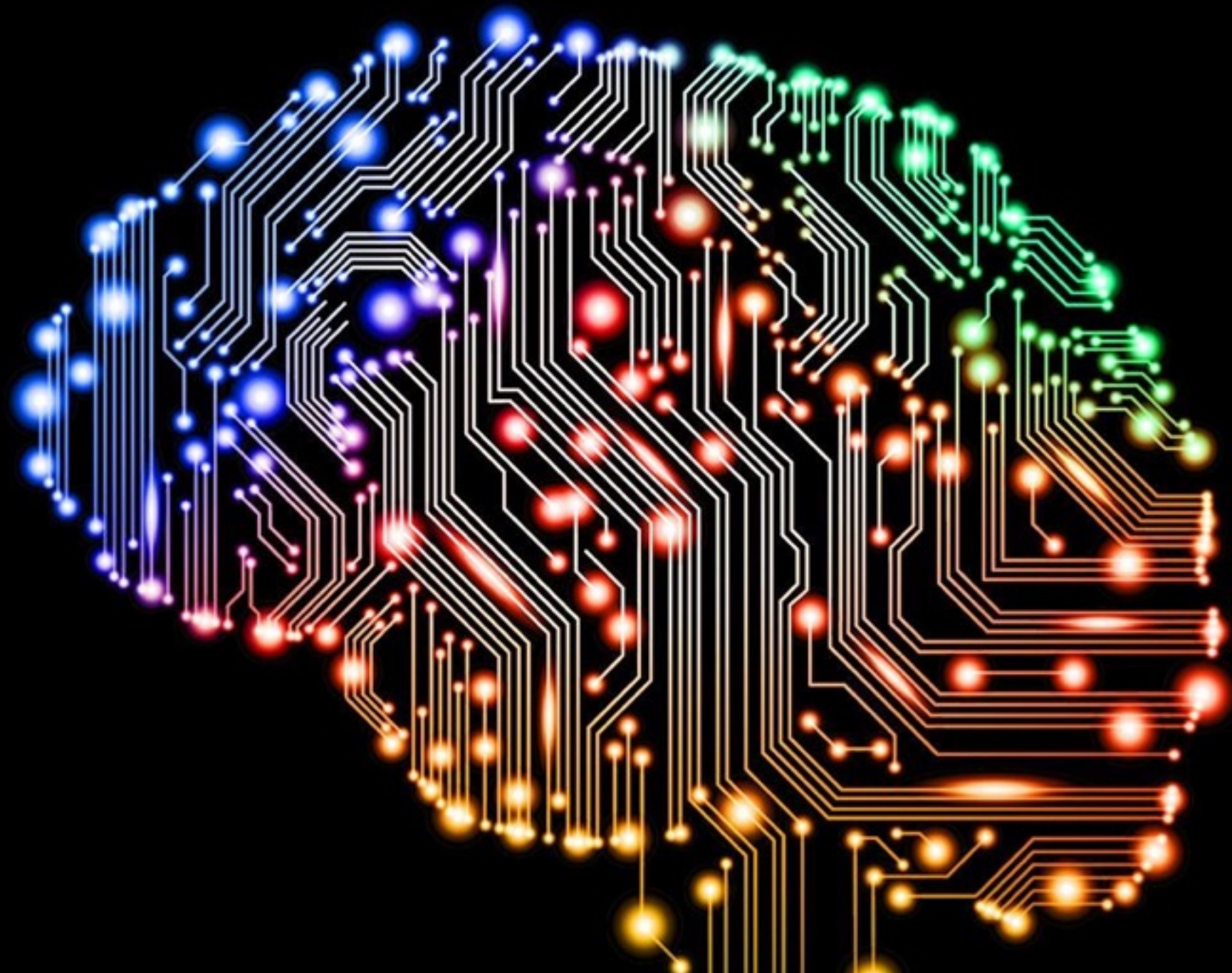
C: 22%, D: 78%

Framing Effect

Outbreak of disease, choose between:

- A: 200 out of 600 people will be saved
 - B: $\frac{1}{3}$ pr. of saving everyone, $\frac{2}{3}$ no one
-
- C: 400 out of 600 people will die
 - D: $\frac{1}{3}$ prob. nobody will die, $\frac{2}{3}$ 600 die

Don't trust your (moral) intuitions!



Universal Intelligence

Can we define an optimally intelligent agent?

Swiss AI Lab in Lugano

World-leading in pattern recognition with Artificial Neural Networks (ANNs).

Theory of optimally intelligent agents:

- AIXI (Marcus Hutter)
- Gödel Machine (Jürgen Schmidhuber)

Former PhD student co-founded DeepMind (sold to Google for \$500m, 2014).

An Optimal Agent

Active agents in a known environment maximize expected utility.

Passive agents in an unknown environment use Solomonoff induction to get probability distribution over possible environments.

What about active agents in unknown environments (like reality)?

Universal Intelligence: AIXI

Developed by Marcus Hutter in 2000

$$a_k := \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} [r_k + \dots + r_m] \sum_{q: U(q, a_1 \dots a_m) = o_1 r_1 \dots o_m r_m} 2^{-\ell(q)}$$

At each step, update your probability distribution over all possible worlds (Bayes/Solomonoff) and choose the action which maximizes expected utility over all remaining steps.

Monte Carlo AIXI

It is the most intelligent agent possible.

Like the Solomonoff prior, AIXI cannot be computed but only approximated.

AIXI is *not* intended as a proposal for building an AI, but as an upper bound on how intelligent an agent can be.

Shortcomings of AIXI

Cartesian model: The agent and the environment are modelled as separate Turing machines. In reality, the agent is part of the world in which it lives.

Wireheading: The most intelligent thing to do for a reinforcement learner is to get control over its reward-channel.

Gödel Machine

Named after Kurt Gödel (1906–1978),
developed by J. Schmidhuber (2003).

Mathematically rigorous, general, fully
self-referential, self-improving, optimally
efficient **problem solver**.

The GM solves all large enough problems
almost as quickly as if it already knew the
best (unknown) algorithm for solving them.



Decision Theory (Part 2)

What kind of problems can we run into?

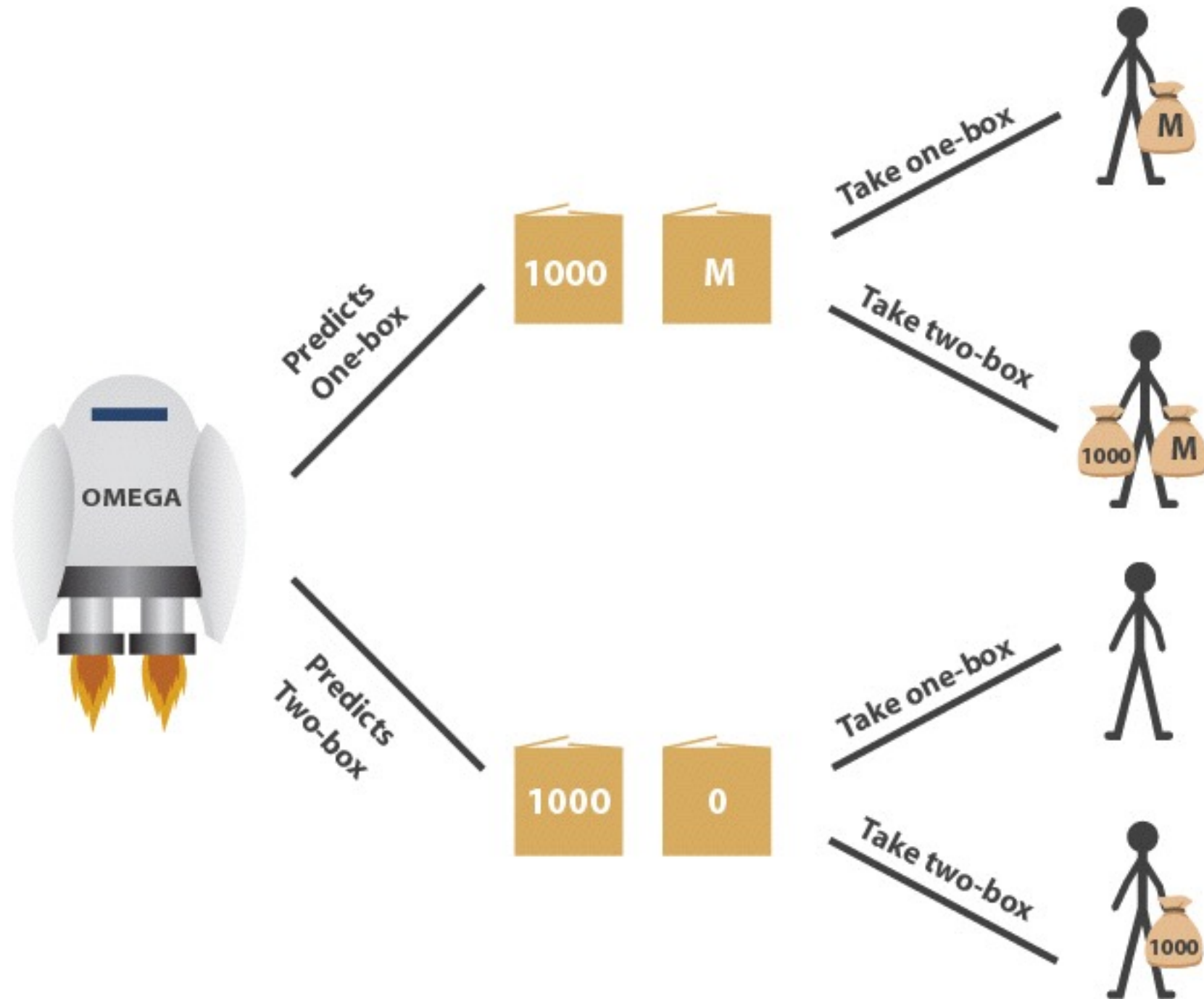
Complications

What if the payout-structure itself depends on the action you take?

Causal Decision Theory: Take the action that causes the best expected outcome.

Evidential Decision Theory: Choose the action which, *conditional on you having chosen it*, gives you the best e. outcome.

Newcomb's Paradox



Newcomb's Paradox
en.wikipedia.org/wiki/Newcomb%27s_paradox

Superintelligence
Decision Theory

Possible Answers

EDT > CDT?

Causal reasoning: Content of boxes is fixed and not affected by my decision. Taking both gives me more than one.

Evidential reasoning: If I take only one box, the predictor had predicted this and filled it with \$1m. If I take both boxes, I will walk away with only \$1k.

Where CDT (Arguably) > EDT

«Solomon is an ancient monarch vaguely reminiscent of the Israelite King. (Every part of this story is Biblically inaccurate.) He is pondering whether to summon Bathsheba, another man's wife. But Solomon is also fully informed as to the peculiar connection between his choice in this matter and the likelihood of his eventually suffering a successful revolt: "Kings have two basic personality types, charismatic and uncharismatic. A king's degree of charisma depends on his genetic make-up and early childhood experiences, and cannot be changed in adulthood. Now charismatic kings tend to act justly and uncharismatic kings unjustly. Successful revolts against charismatic kings are rare, whereas successful revolts against uncharismatic kings are frequent. Unjust acts themselves, though, do not cause successful revolts... Solomon does not know whether or not he is charismatic; he does know that it is unjust to send for another man's wife."»

Where Both (Arguably) Fail

- Blackmailing
- Parfit's Hitchhiker
- Counterfactual Mugging

Possible solution: Updateless DT (“Do what you would have precommitted.”)

EDT & CDT not reflectively consistent.

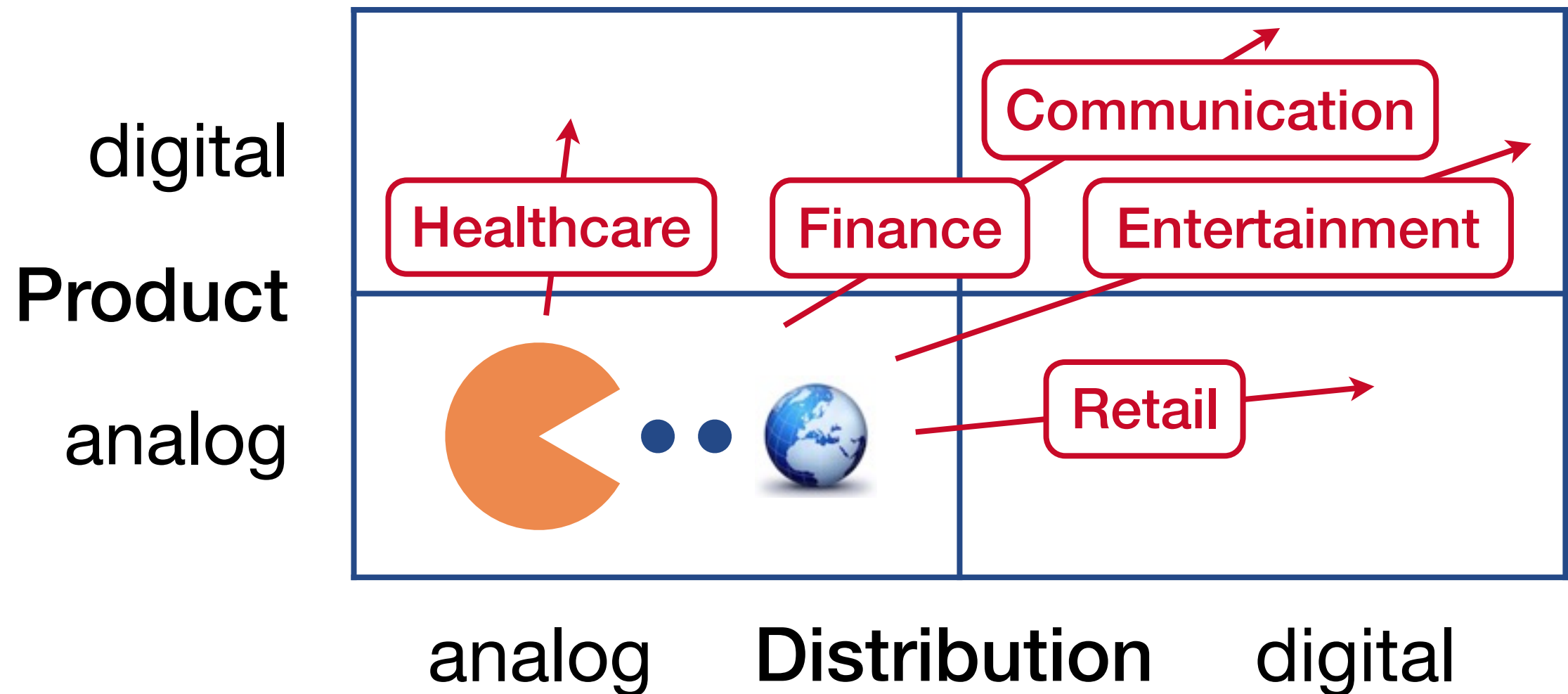


Machine Learning

How can algorithms separate the signal from the noise?

Software is eating the world

... because it is more productive!



Machine Intelligence LANDSCAPE

CORE TECHNOLOGIES

ARTIFICIAL INTELLIGENCE



DEEP LEARNING



MACHINE LEARNING



NLP PLATFORMS



PREDICTIVE APIS

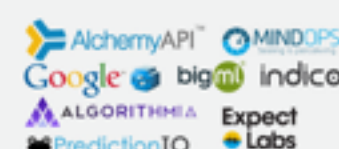


IMAGE RECOGNITION



SPEECH RECOGNITION



RETHINKING ENTERPRISE

SALES



SECURITY / AUTHENTICATION



FRAUD DETECTION



HR / RECRUITING



MARKETING



PERSONAL ASSISTANT



INTELLIGENCE TOOLS



RETHINKING INDUSTRIES

ADTECH



AGRICULTURE



EDUCATION



FINANCE



LEGAL



MANUFACTURING



MEDICAL



OIL AND GAS



MEDIA / CONTENT



CONSUMER FINANCE



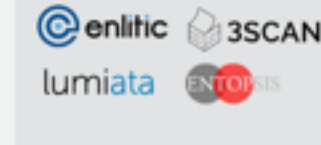
PHILANTHROPIES



AUTOMOTIVE



DIAGNOSTICS



RETAIL



RETHINKING HUMANS / HCI

AUGMENTED REALITY



GESTURAL COMPUTING



ROBOTICS



EMOTIONAL RECOGNITION



SUPPORTING TECHNOLOGIES

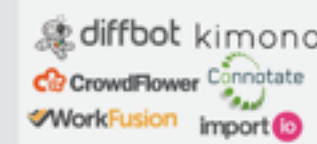
HARDWARE



DATA PREP



DATA COLLECTION



Computers have just learned ...

... how to see, read, write and classify.



Describes without errors

Describes with minor errors

Somewhat related to the image

Unrelated to the image



A person riding a motorcycle on a dirt road.



Two dogs play in the grass.



A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.



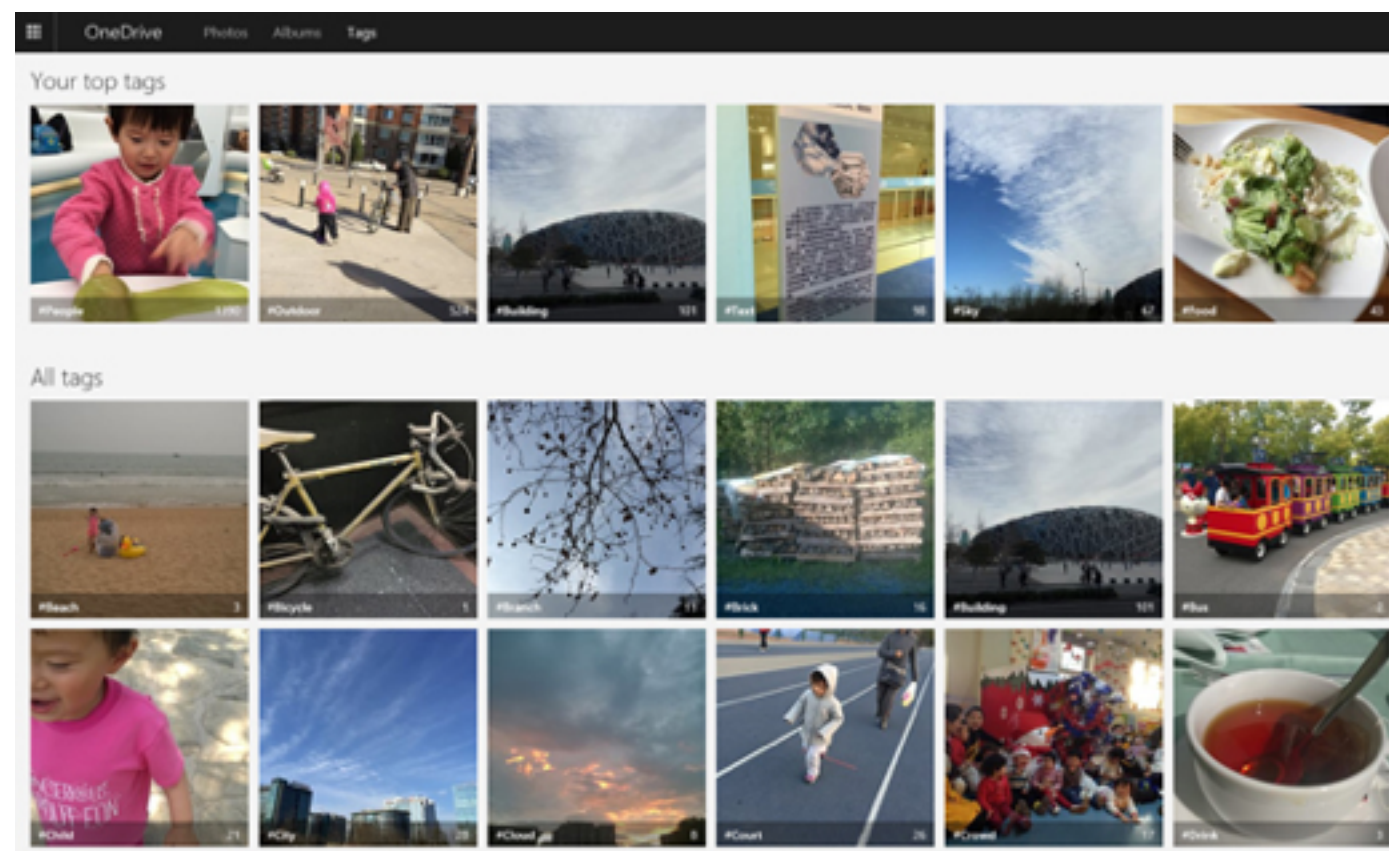
A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.

Superhuman Image Recognition

- With convolutional neural networks
- 1.2 m training images, ~ 30 layers



Machine

A machine runs an algorithm which is

- a procedure for
- solving a specified problem
- in a finite number of steps
- (i.e. eventually producing an output).

Learning

- Building predictive models from data
- Algorithms improve with “experience”
- Useful if you don’t know the solution
- **Machine Learning: Learn and predict**
- **Data Mining: Discover new properties**

Feedback Mechanism

Categorization based on feedback:

- **Supervised Learning:** Given a set of inputs & outputs, learn a general rule
- **Unsupervised Learning:** Learn the structure in input without training set
- **Reinforcement Learning:** Achieve a certain goal in dynamic environment

Tasks and Models

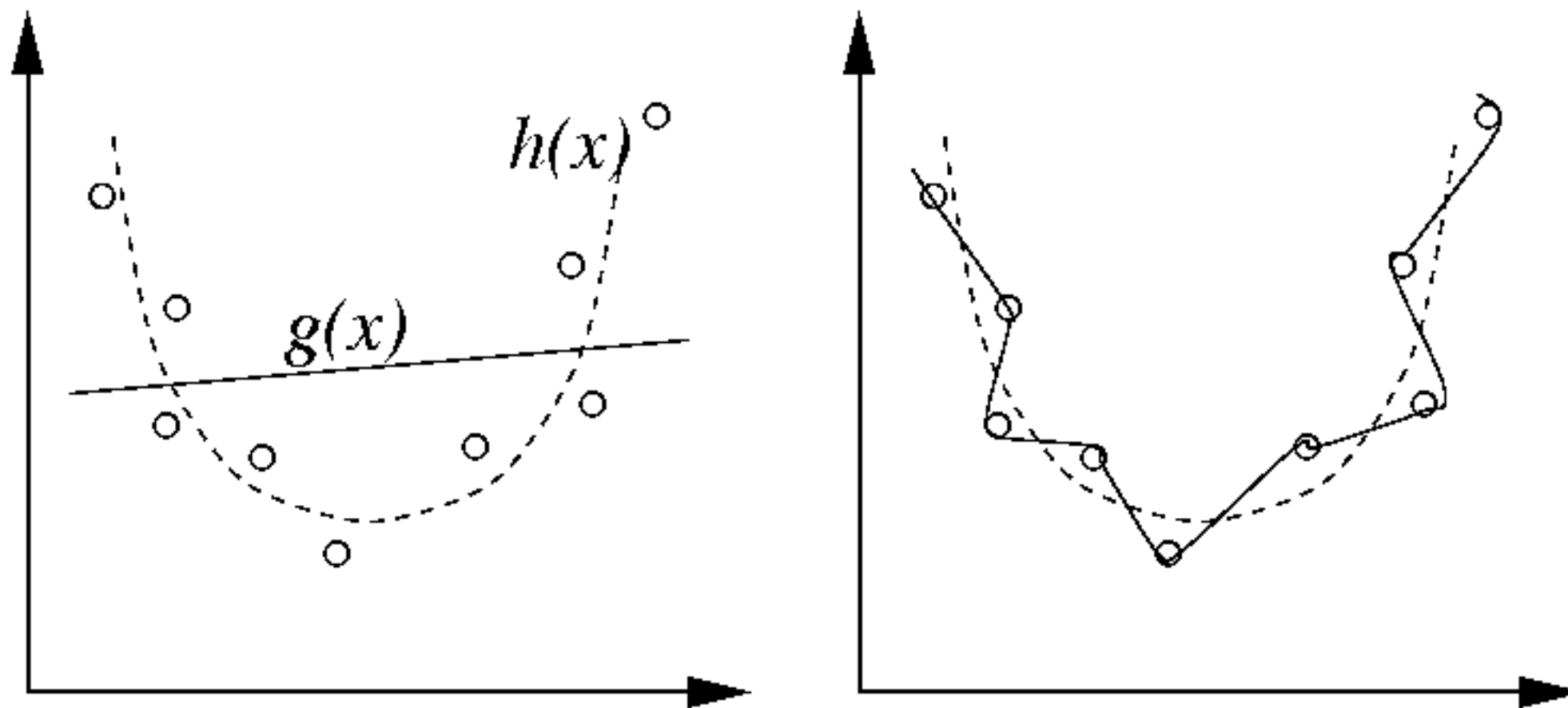
- **Classification:** Assign unseen input to discrete categories (supervised)
- **Regression:** Map the unseen input to continuous values (supervised)
- **Clustering:** Divide set of inputs into a number of groups (unsupervised)
- **Dimensionality Reduction:** Simplify

Applications

- Spam filtering
- Search engines
- Computer vision
- Medical diagnosis
- Stock market prediction
- Optical character recognition
- Recognizing credit card fraud
- Speech-recognition (e.g. SIRI)
- Autonomous systems (e.g. self-driving cars)



Problem: Overfitting



Cross-Validation

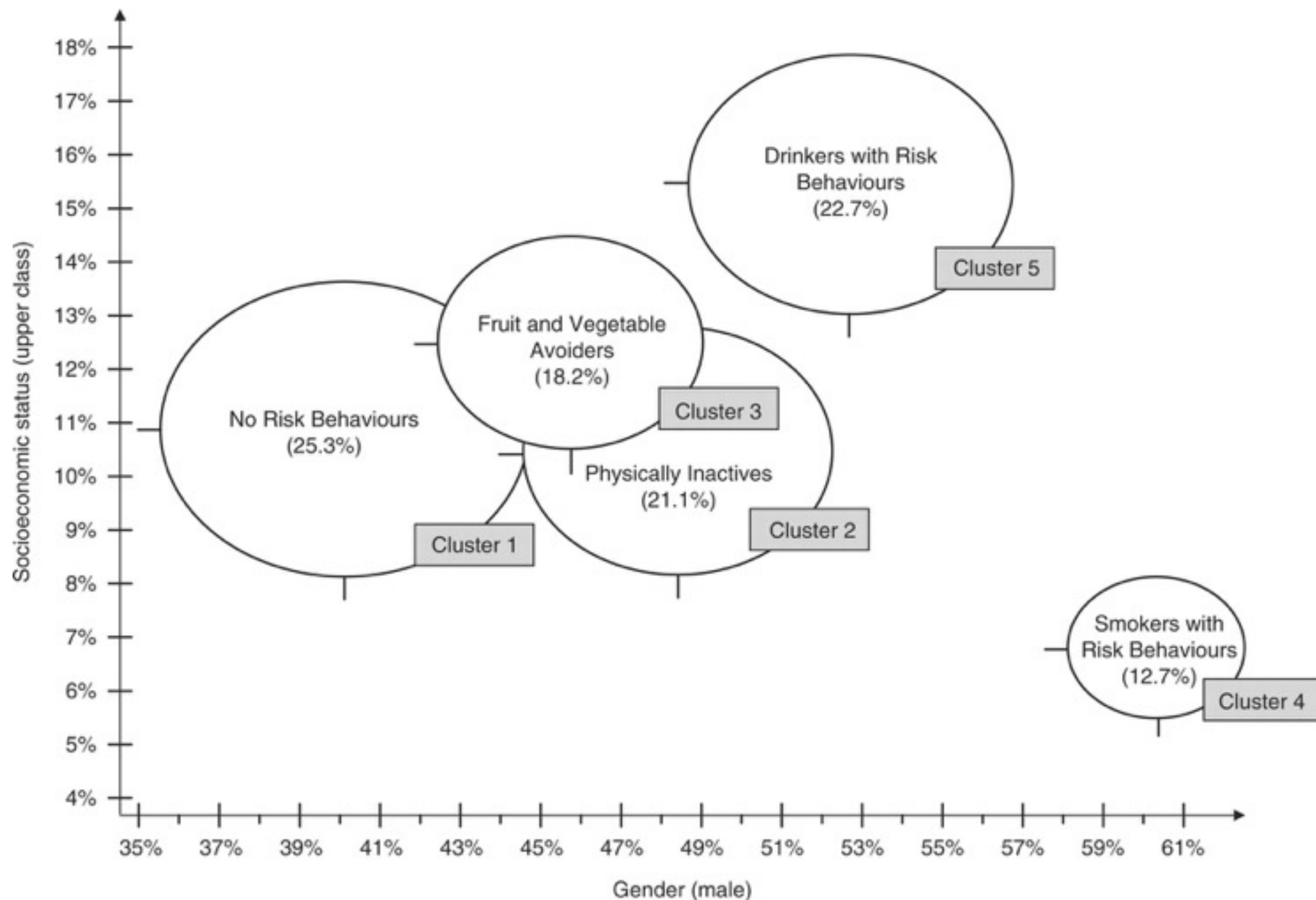
- Helps to prevent overfitting (i.e. find the optimal number of parameters).
- Helps to estimate how good a model will generalize to unseen data.

Partition data into a “training set” and a “validation set”.

Exploration-Exploitation Trade-Off

You have 1000 gambling machines in front of you, which give different pay-outs with different prob. distributions. You have 1000 trials for all machines. How many of your trials do you spend finding a good machine, how many to exploit the best machine found so far?

Clustering: Example



Clustering: k-means Algorithm

How to find **k clusters** in **N data points**?

Goal: Minimize the sum of the **distances** of each point to the **center** of its cluster.

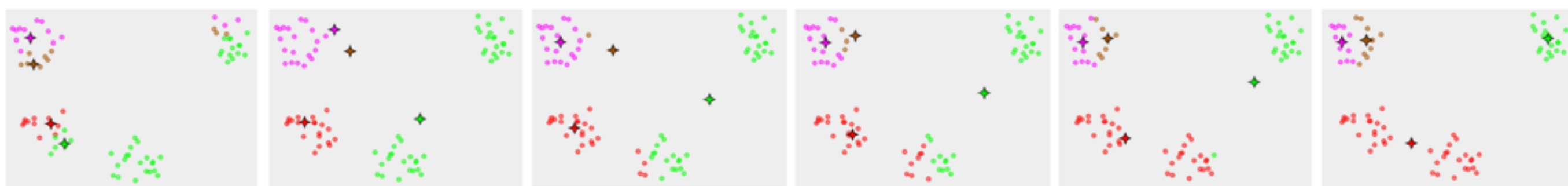
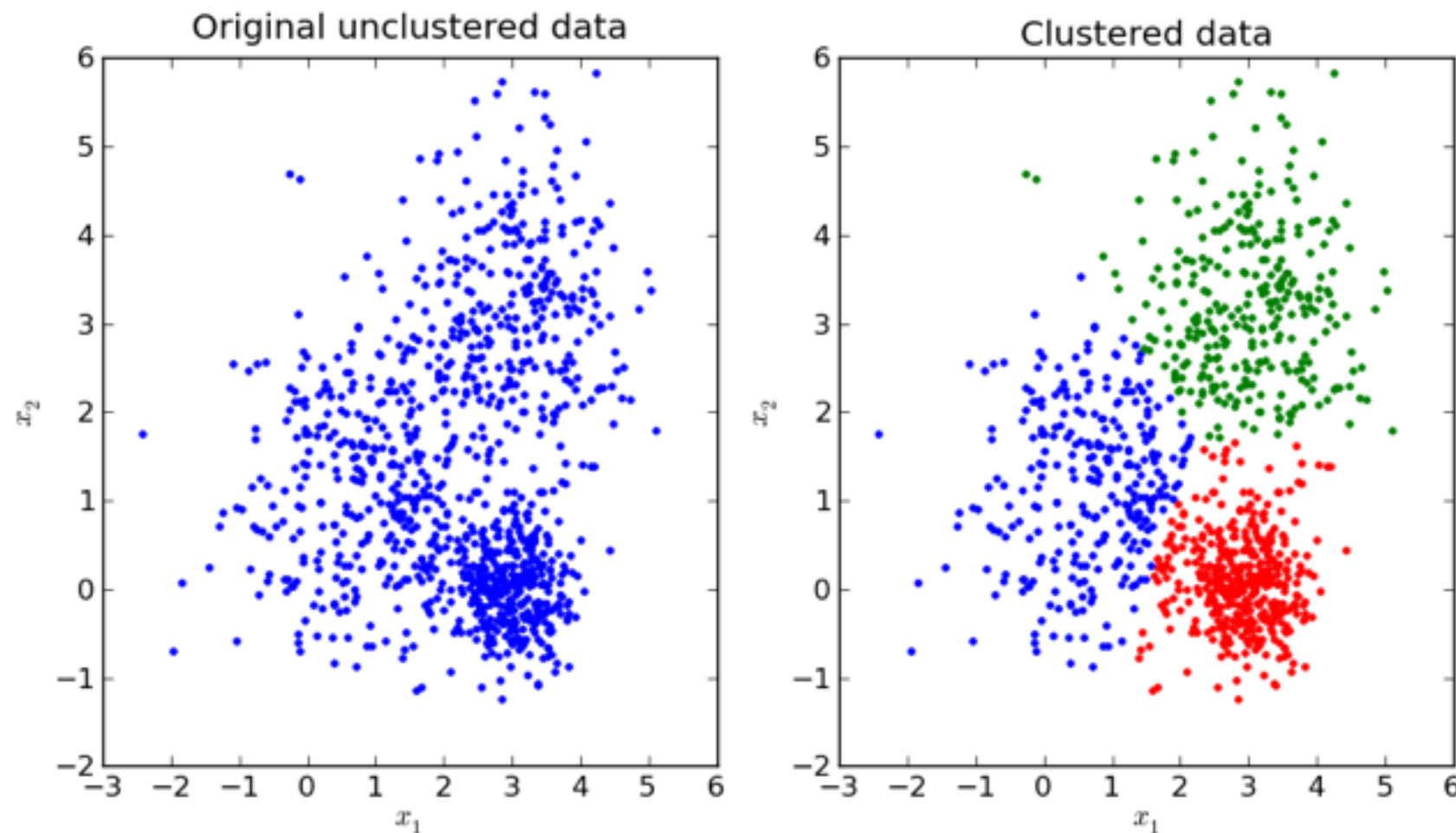
Solving this problem exactly is hard (NP).

Choose **k** of the **N points** randomly.

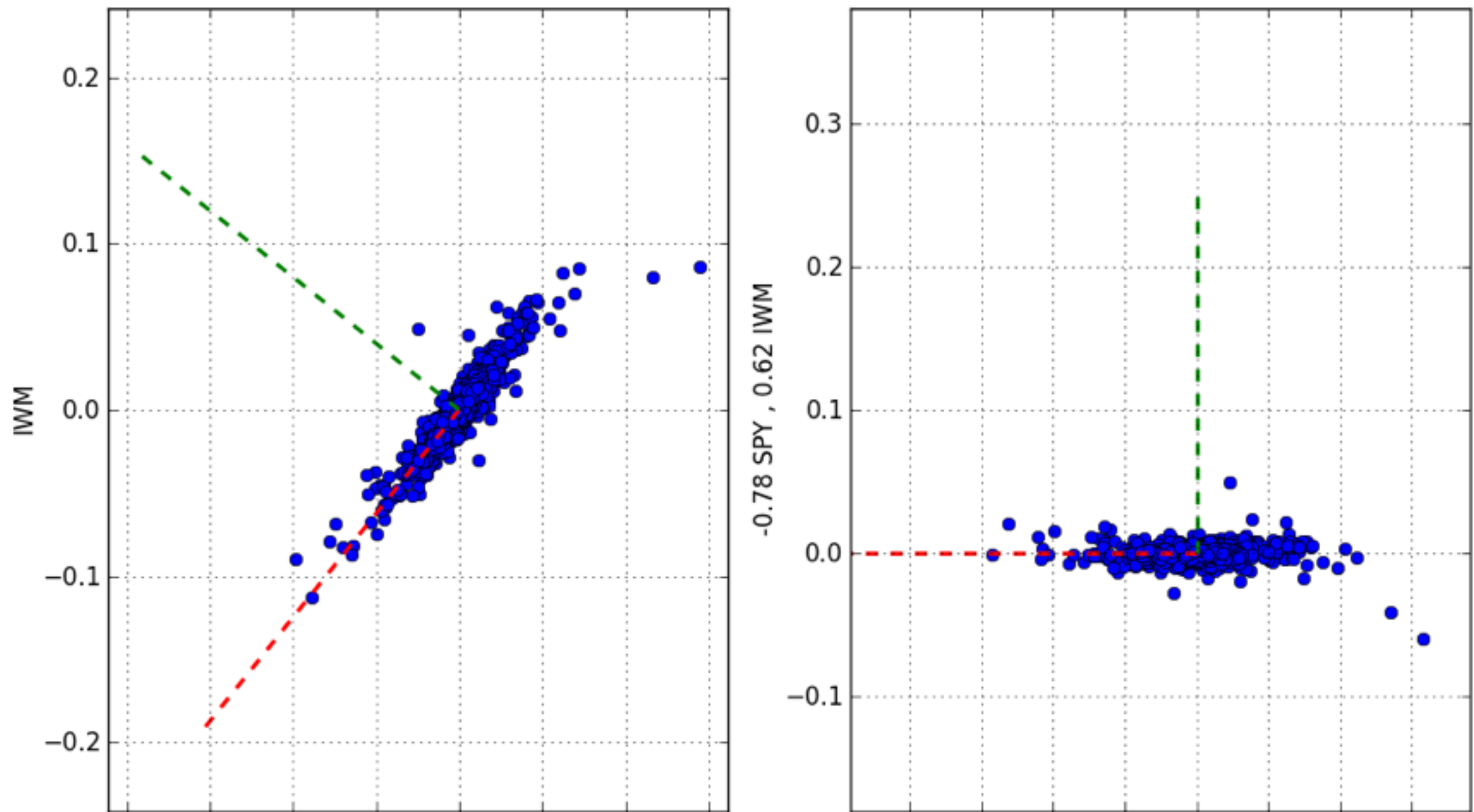
 Assign each point to the nearest mean.

Calculate the new means of each cluster.

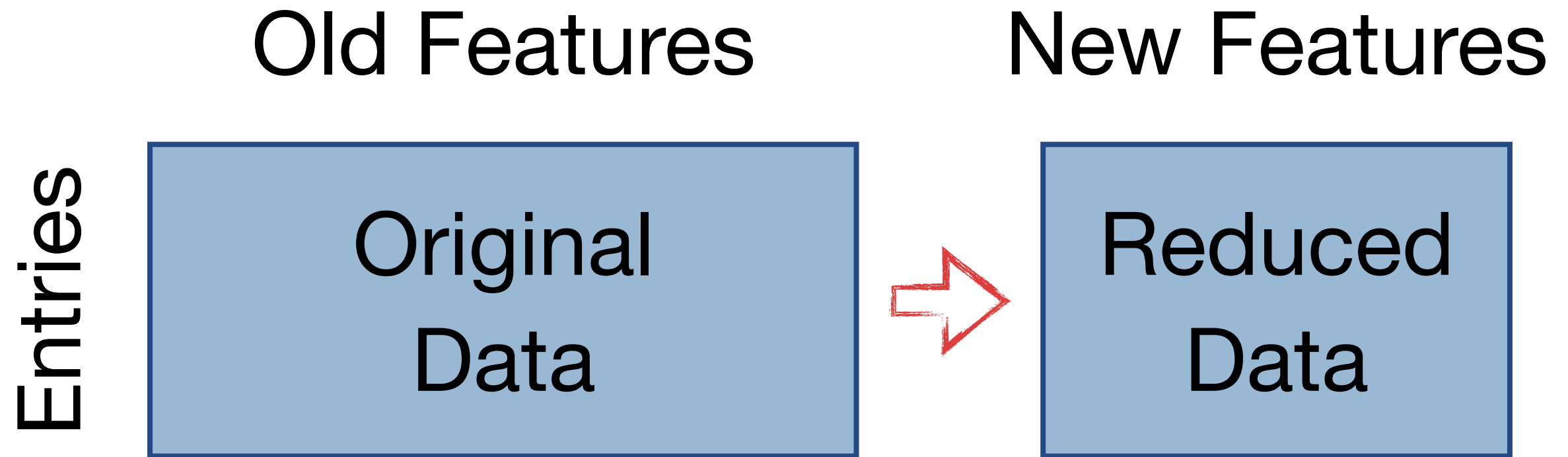
k-means Clustering Outcomes



Dimensionality Reduction



Goal: Reduce Dimensionality



Retain as much information as possible
with as little information as possible!

Recommendation Systems

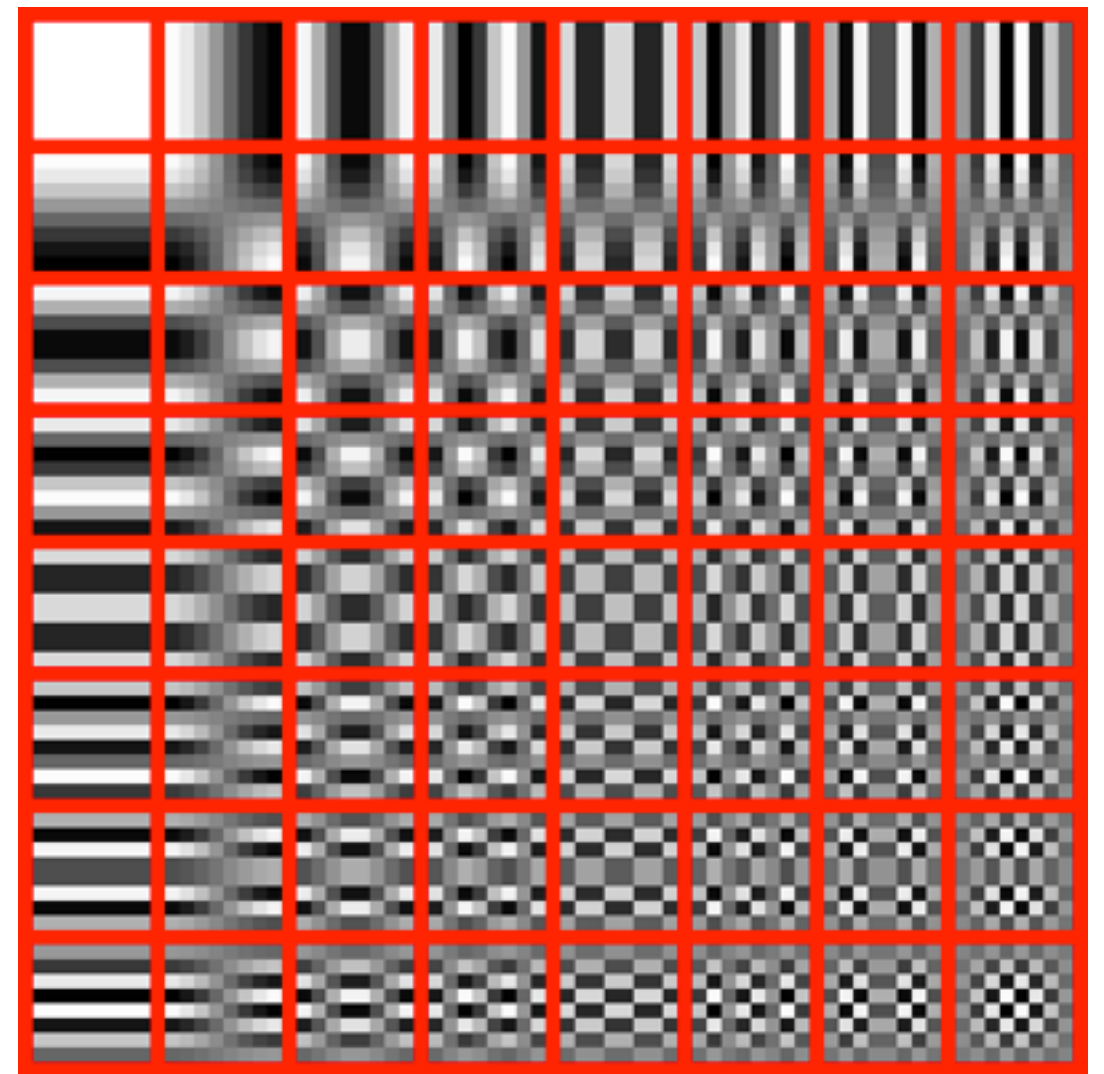
Example: If each song is a dimension, a user's taste can be viewed as vector.

Observation: As groups of users rate the same kind of songs high and low, try to distill the tastes as base vectors.

Apply Principal Component Analysis!

Lossy Compression

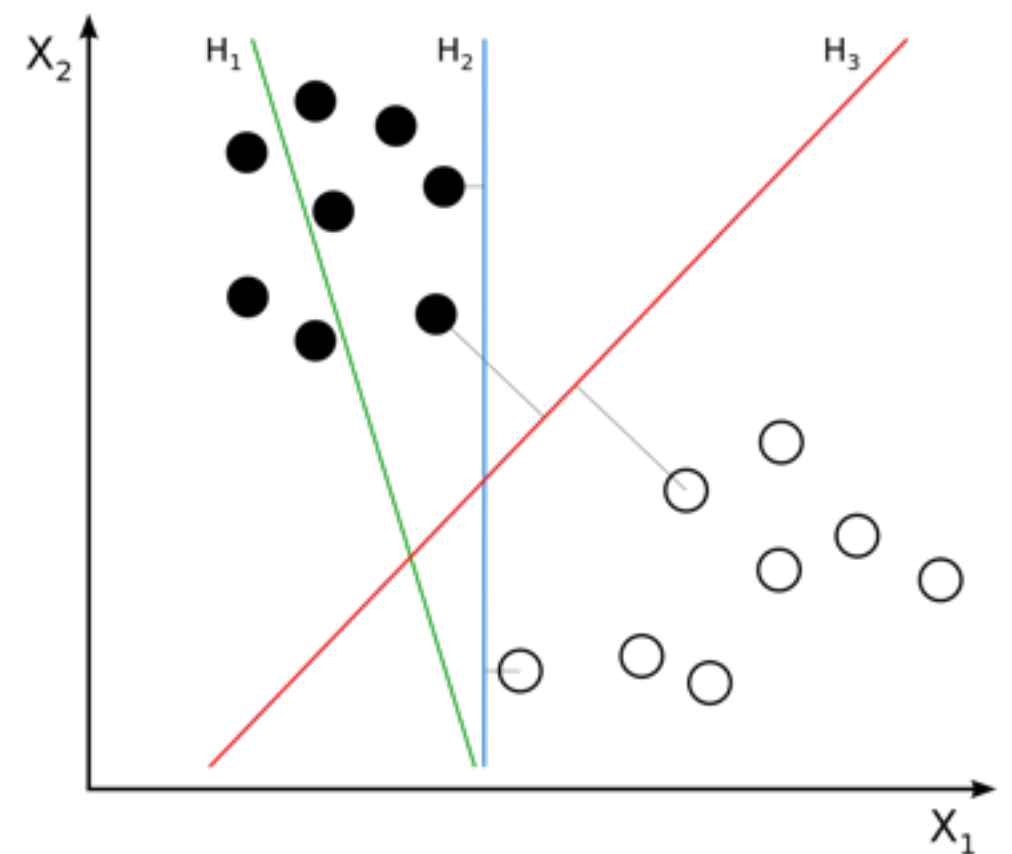
- Downside: New base vectors need to be stored to recover the data.
- Solution: Choose appropriate base vectors (in case of JPEG: DCT)!



Support Vector Machines

State-of-the-art supervised learning models used for classification.

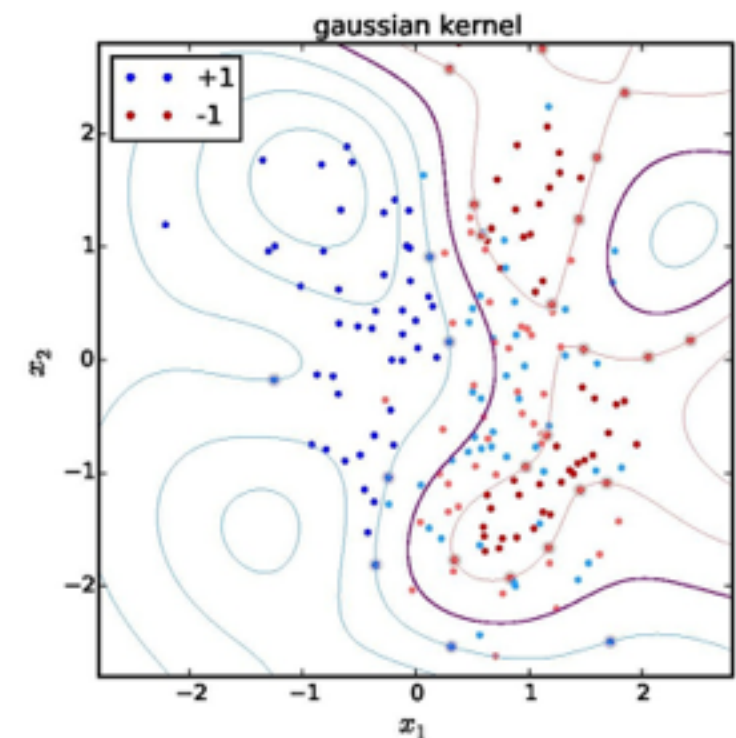
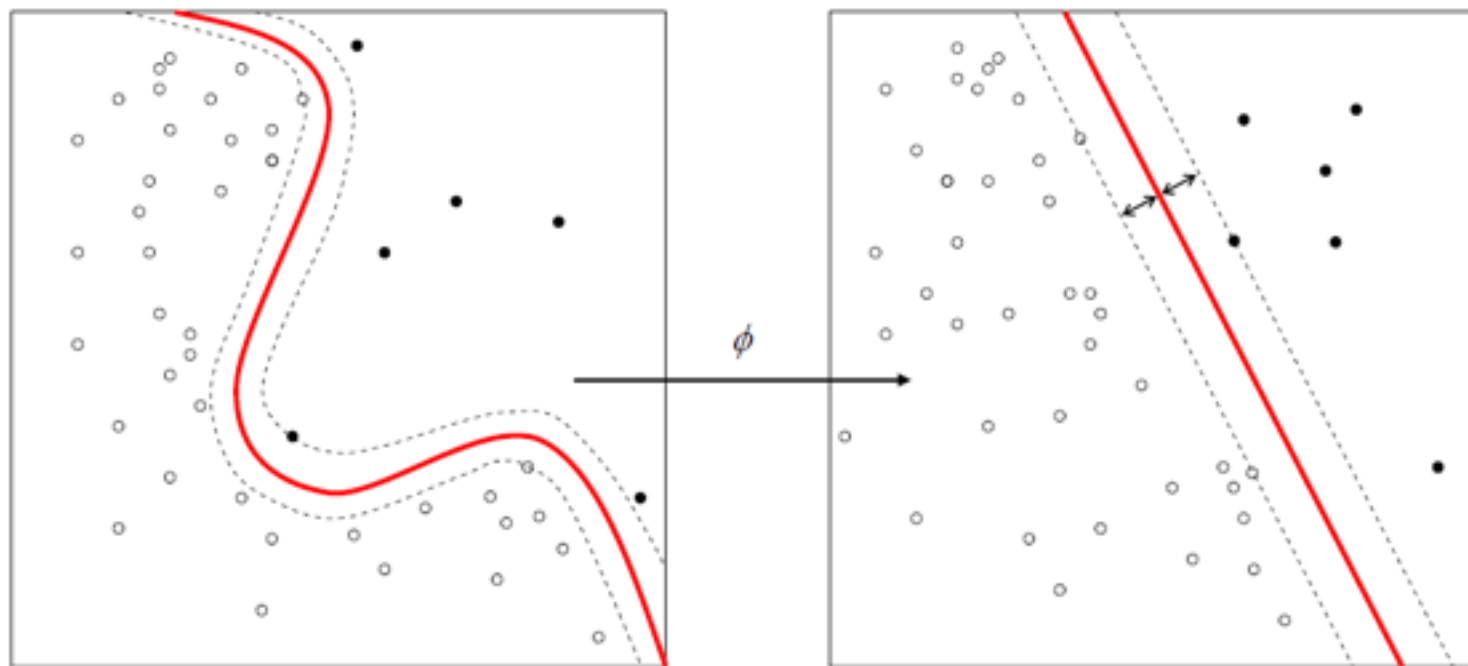
Binary, linear classifier:
Find a plane that max. distance to the closest data points from both classes.



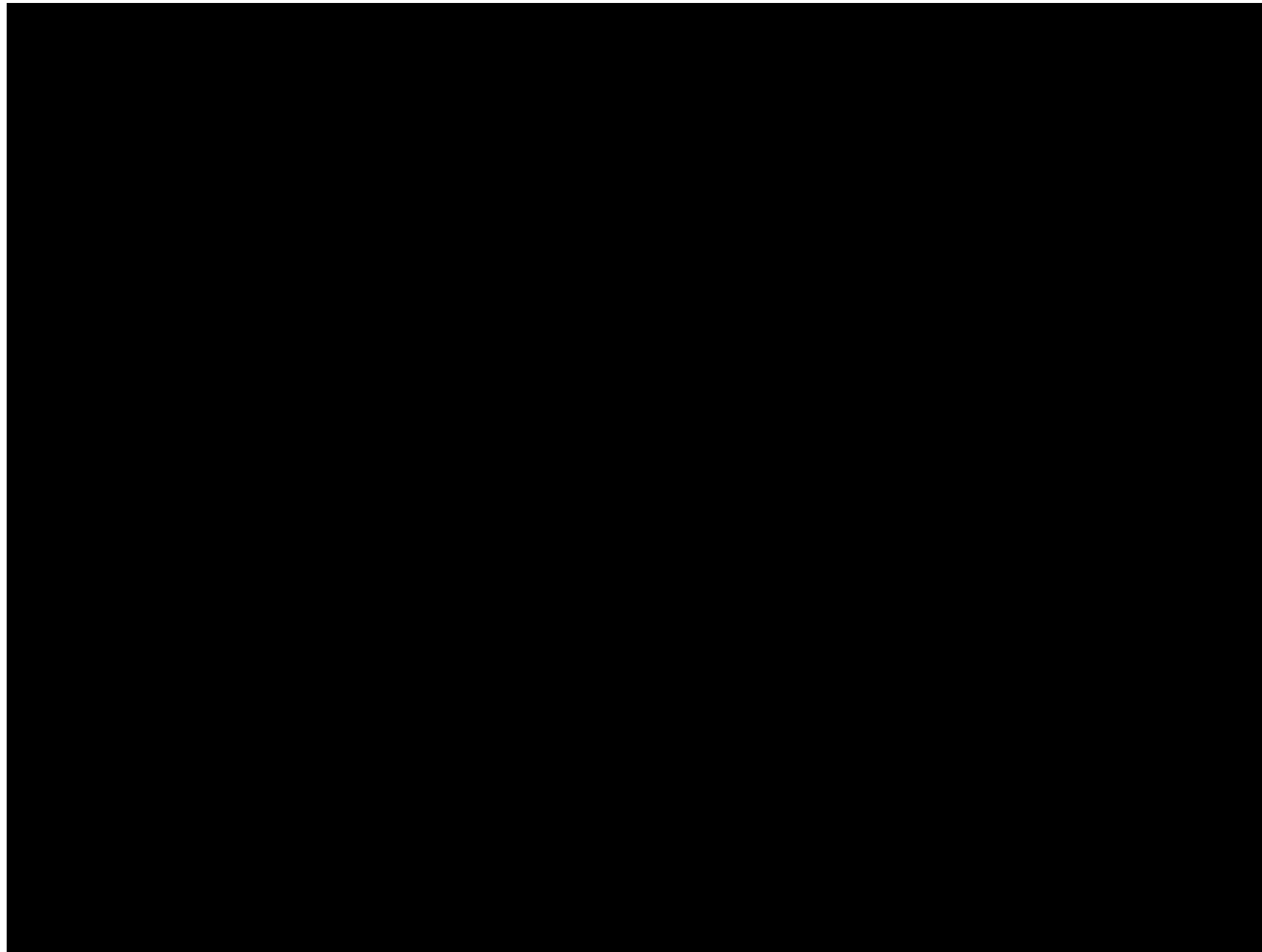
SVM's Extensions

In practice, there will often be no plane that perfectly separates the two classes.

Soft margin: Allow violations w. penalty.



Reinforcement Learning (CNN)



Google develops self-learning computer program
[www.theguardian.com/technology/2015/feb/25/\[...\]](http://www.theguardian.com/technology/2015/feb/25/[...])

Superintelligence
Machine Learning

Outlook: Huge Responsibility

... in the near-term. For the long-term, see www.superintelligence.ch.

Train CNN



Deploy CNN



Campaign to Stop Killer Robots
www.stopkillerrobots.org

Superintelligence
Machine Learning

75