

PROBABILISTIC MACHINE LEARNING

LECTURE 01

INTRODUCTION

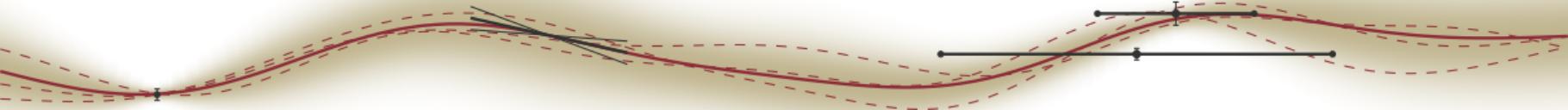
Philipp Hennig

20 April 2020

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



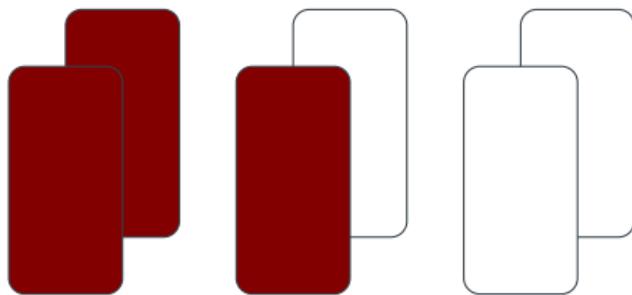
FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING





Which Card?

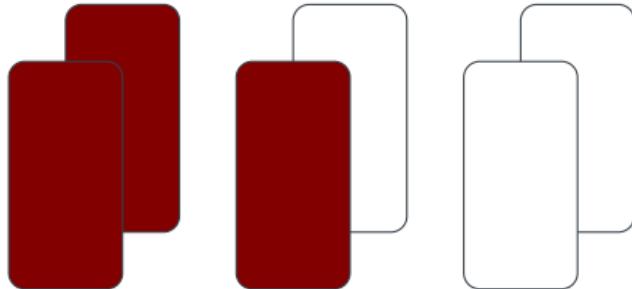
an opening experiment





Which Card?

an opening experiment



- ▶ $\frac{1}{2}$
- ▶ $\frac{2}{3}$
- ▶ something else
- ▶ I don't know yet



An *inference* problem requires statements about the value of an *unobserved* (latent) variable x based on observations y which are **related** to x , but may not be sufficient to fully determine x . This requires a notion of **uncertainty**.



Life is Uncertain

the ability to reason under uncertainty is the hallmark of intelligence



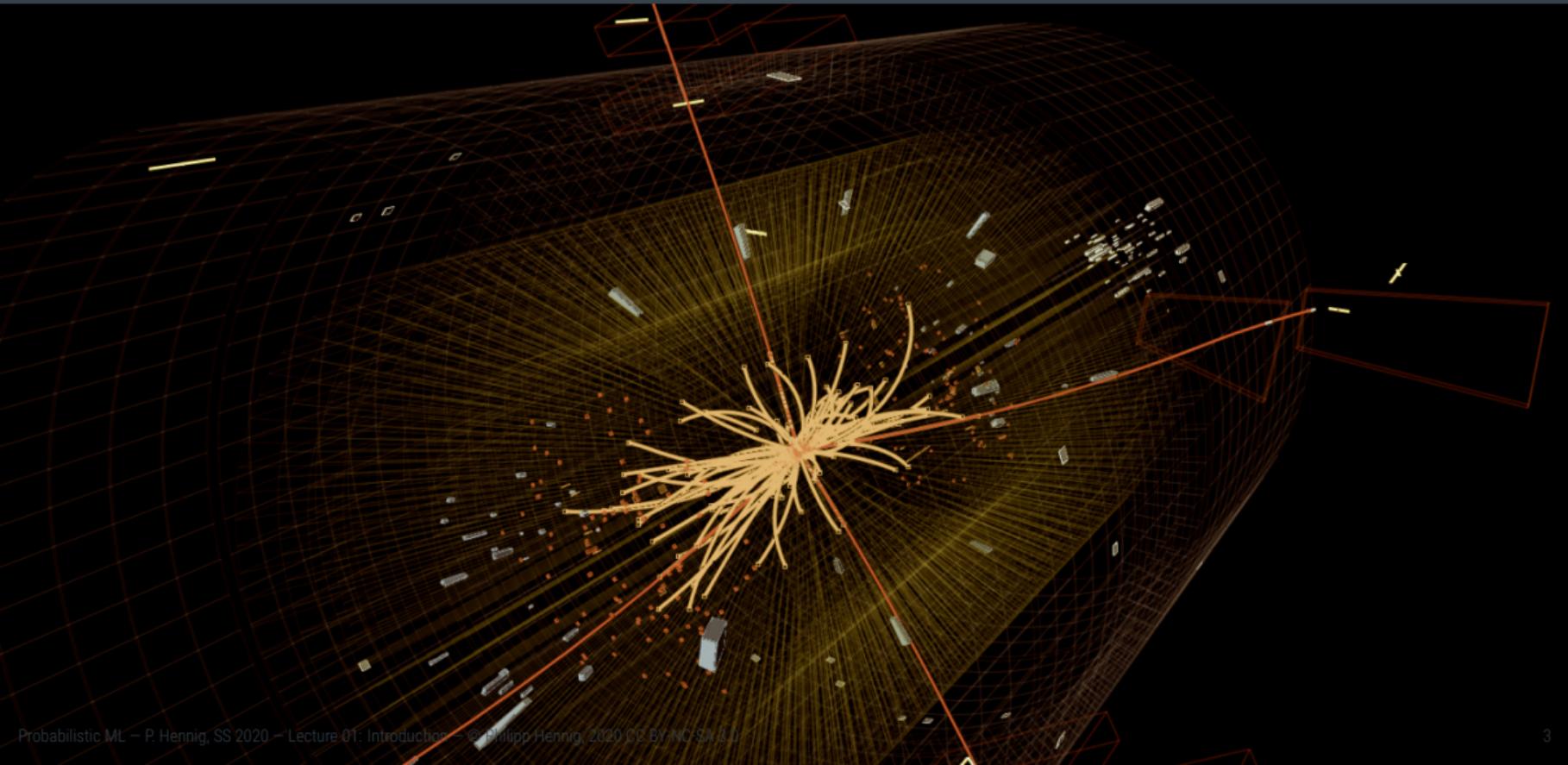
Life is Uncertain

the ability to reason under uncertainty is the hallmark of intelligence



Life is Uncertain

the ability to reason under uncertainty is the hallmark of intelligence



Life is Uncertain

the ability to reason under uncertainty is the hallmark of intelligence



[Luke Fildes, 1891 (Tate Modern)]



Deductive and Plausible Reasoning

Limits of propositional logic



[image: Todd Neville]

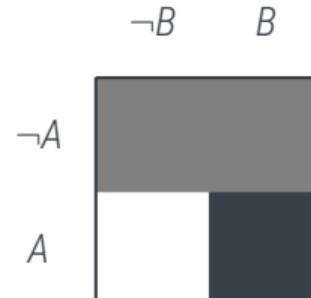


Deductive Reasoning:

A is true thus B is true
 B is false thus A is false

$$A \Rightarrow B$$

modus ponens
modus tollens



if A is true, then B is true

Plausible Reasoning:

$P(B | A) > P(B)$
 A is true thus B becomes more plausible
 B is false thus A becomes less plausible
 B is true thus A becomes more plausible
 A is false thus B becomes less plausible

if A is true, the B becomes more plausible



The actual science of logic is conversant at present only with things either certain, impossible, or entirely doubtful, none of which (fortunately) we have to reason on. Therefore the true logic for this world is the calculus of Probabilities, which takes account of the magnitude of the probability which is, or ought to be, in a reasonable man's mind.

James Clerk Maxwell (1831–1879)

[image source: BBC / public domain]



Goals of this course:

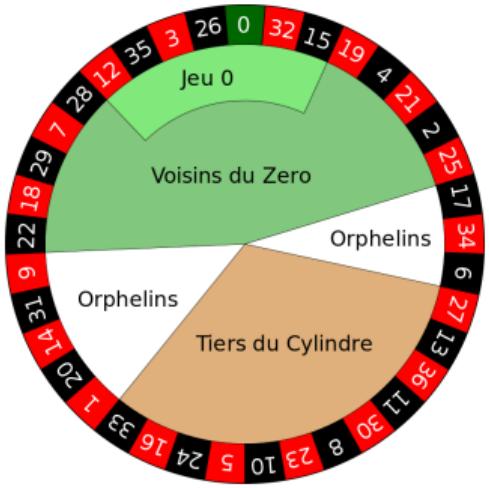
- ▶ Establish a *formal framework* for **probable reasoning**
- ▶ Use it to build *powerful* inference mechanisms for **real-world problems**
- ▶ Develop the *technical tools* necessary to implement **inference** in practice

Probabilities Distribute Truth

Towards a Theory of Probability



[image: wikipedia / Ann-Kathrin Schalkamp]



0			MANQUE
1	2	3	
4	5	6	
7	8	9	
10	11	12	
13	14	15	
16	17	18	
19	20	21	
22	23	24	
25	26	27	
PAIR			IMPARI
28	29	30	
31	32	33	
34	35	36	
12 ^P	12 ^M	12 ^D	

Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

§ 1. Axiome².

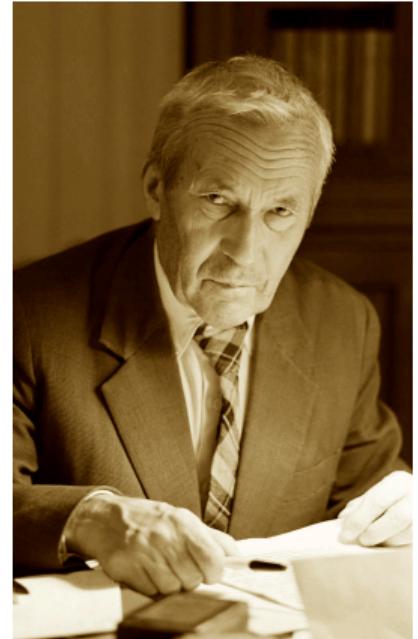
Es sei E eine Menge von Elementen ξ, η, ζ, \dots , welche man *elementare Ereignisse* nennt, und \mathfrak{F} eine Menge von Teilmengen aus E ; die Elemente der Menge \mathfrak{F} werden weiter *zufällige Ereignisse* genannt.

- I. \mathfrak{F} ist ein Mengenkörper³.
- II. \mathfrak{F} enthält die Menge E .
- III. Jeder Menge A aus \mathfrak{F} ist eine nichtnegative reelle Zahl $P(A)$ zugeordnet. Diese Zahl $P(A)$ nennt man die *Wahrscheinlichkeit* des Ereignisses A .
- IV. $P(E) = 1$.

V. Wenn A und B disjunkt sind, so gilt

$$P(A + B) = P(A) + P(B).$$

Ein Mengensystem \mathfrak{F} mit einer bestimmten Zuordnung der Zahlen $P(A)$, welche den Axiomen I–V genügt, nennt man ein *Wahrscheinlichkeitsfeld*.



Andrey N. Kolmogorov
(1903–1987)

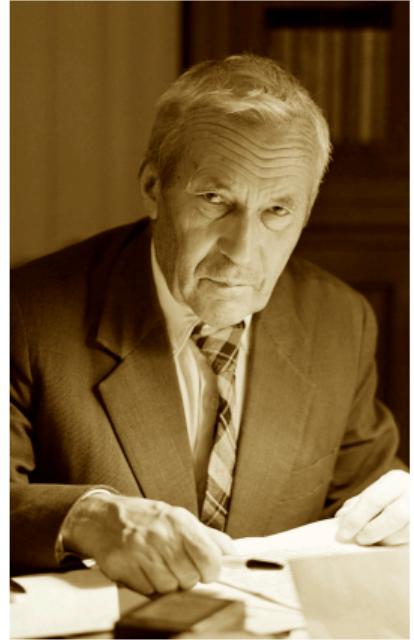


Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

* Vgl. HAUSDORFF: Mengenlehre 1927 S. 78. Ein Mengensystem heißt ein Körper, wenn Summe Durchschnitt und Differenz von zwei Mengen des Systems wieder dem System angehören. Jeder nicht leere Mengenkörper enthält die Nullmenge 0. Wir bezeichnen mit HAUSDORFF den Durchschnitt von A und B mit AB , die Vereinigungsmenge von A und B im Falle $AB = 0$ mit $A + B$, allgemein aber mit $A + B$, und die Differenz von A und B mit $A - B$. Das Komplement $E - A$ der Menge A wird durch A' bezeichnet. Die elementaren Rechengesetze für Mengen und ihre Durchschnitte, Summen und Differenzen werden weiter als bekannt vorausgesetzt. Mengen aus \mathfrak{F} werden weiter mit großen lateinischen Buchstaben bezeichnet.



Andrey N. Kolmogorov
(1903–1987)

Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

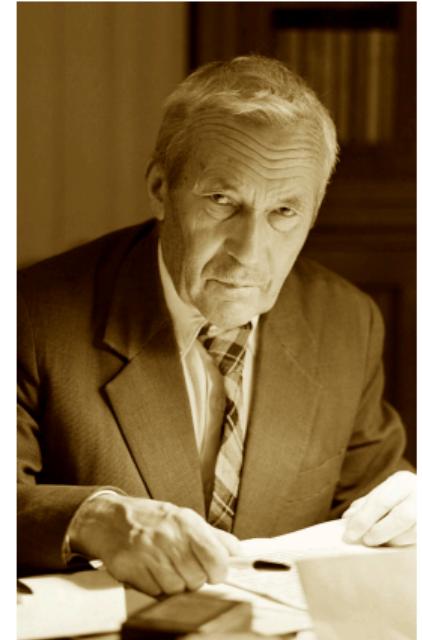
E is the roulette wheel
F is the table

Definition (σ -algebra, measurable sets & spaces)

Let E be a space of elementary events. Consider the power set 2^E , and let $\mathfrak{F} \subset 2^E$ be a set of subsets of E . Elements of \mathfrak{F} are called random events. If \mathfrak{F} satisfies the following properties, it is called a σ -algebra.

1. $E \in \mathfrak{F}$
2. $(A, B \in \mathfrak{F}) \Rightarrow (A - B \in \mathfrak{F})$ F is closed under the difference
3. $(A_1, A_2, \dots \in \mathfrak{F}) \Rightarrow (\bigcup_{i=1}^{\infty} A_i \in \mathfrak{F} \quad \wedge \quad \bigcap_{i=1}^{\infty} A_i \in \mathfrak{F})$

(this implies $\emptyset \in \mathfrak{F}$. If E is countable, then 2^E is a σ -algebra). If \mathfrak{F} is a σ -algebra, its elements are called measurable sets, and (E, \mathfrak{F}) is called a measurable space (or Borel space).



Andrey N. Kolmogorov
(1903–1987)

Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

Definition (Measure & Probability Measure)

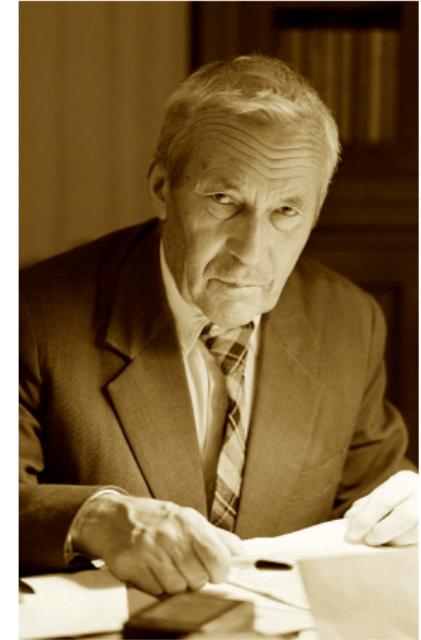
Let (E, \mathfrak{F}) be a **measurable space** (aka. Borel space). A nonnegative real function $P : \mathfrak{F} \rightarrow \mathbb{R}_{0,+}$ (III.) is called a **measure** if it satisfies the following properties:

1. $P(\emptyset) = 0$
2. For any countable sequence $\{A_i \in \mathfrak{F}\}_{i=1,\dots}$ of pairwise disjoint sets ($A_i \cap A_j = \emptyset$ if $i \neq j$), P satisfies **countable additivity** (aka. σ -additivity):

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i). \quad (\text{V.})$$

The measure P is called a **probability measure** if $P(E) = 1$.

(For probability measures, 1. is unnecessary). Then, (E, \mathfrak{F}, P) is called a **probability space**.



IV.

Andrey N. Kolmogorov
(1903–1987)



Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

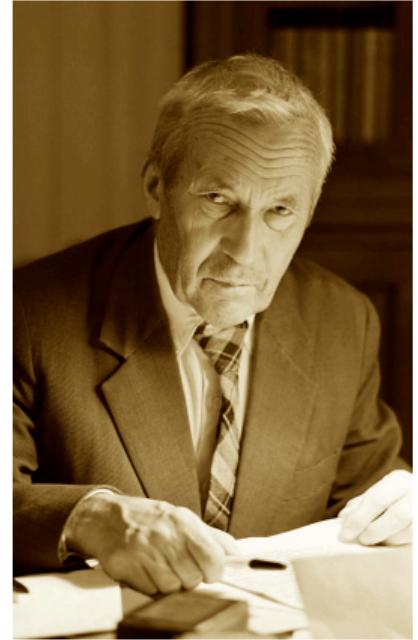
Theorem (Sum Rule)

From $A + \neg A = E$ we get

$$P(A) + P(\neg A) \stackrel{(V)}{=} P(E) \stackrel{(IV)}{=} 1, \quad \text{thus} \quad P(A) = 1 - P(\neg A).$$

And from $A = A \cap (B + \neg B)$, using the notation $P(A, B) = P(A \cap B)$ for the **joint probability** of A and B , we get the **Sum Rule**

$$P(A) \stackrel{(V)}{=} P(A, B) + P(A, \neg B).$$



Andrey N. Kolmogorov
(1903–1987)



Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

Definition (Conditional Probability)

If $P(A) > 0$, the quotient

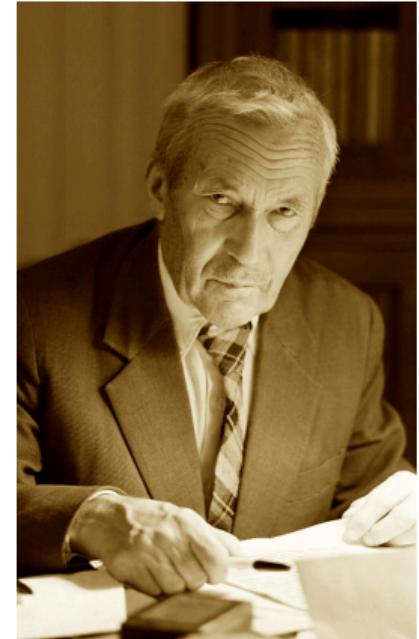
$$P(B | A) = \frac{P(A, B)}{P(A)}$$

is called the **conditional probability** of B given A . It immediately gives

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

It is easy to show that $P(B | A) \geq 0$ (III), $P(E | A) = 1$ (IV), and for $B \cap C = \emptyset$, we have $P(B + C | A) = P(B | A) + P(C | A)$ (V). Thus, for a fixed A , $(E, \mathfrak{F}, P(\cdot | A))$ is a probability space.

Note that $P(A | A) = \frac{P(A \cap A)}{P(A)} = 1$.



Andrey N. Kolmogorov
(1903–1987)

Kolmogorov's Axioms

Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

Theorem (Law of Total Probability)

Let $A_1 + A_2 + \dots + A_n = E$ and $A_i \cap A_j = \emptyset$ if $i \neq j$. Then, for any $X \in \mathfrak{F}$,

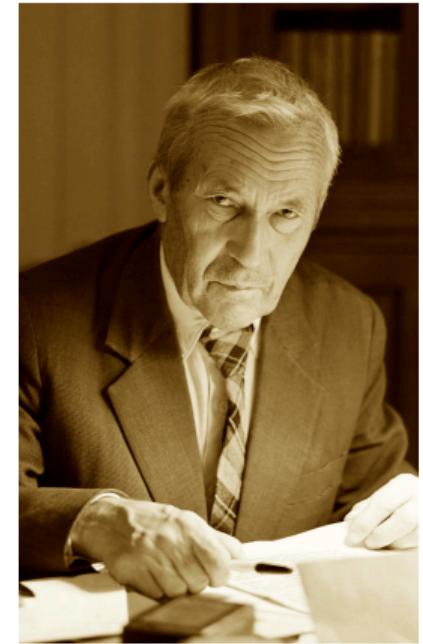
$$P(X) = \sum_{i=1}^n P(X | A_i)P(A_i).$$

Proof.

Because $X = E \cap X = \bigcup_{i=1}^n (A_i \cap X)$, we get from V that

$$P(X) = \sum_{i=1}^n P(A_i, X) \stackrel{\text{def.}}{=} \sum_{i=1}^n P(X | A_i)P(A_i).$$

□



Andrey N. Kolmogorov
(1903–1987)



Kolmogorov's Axioms

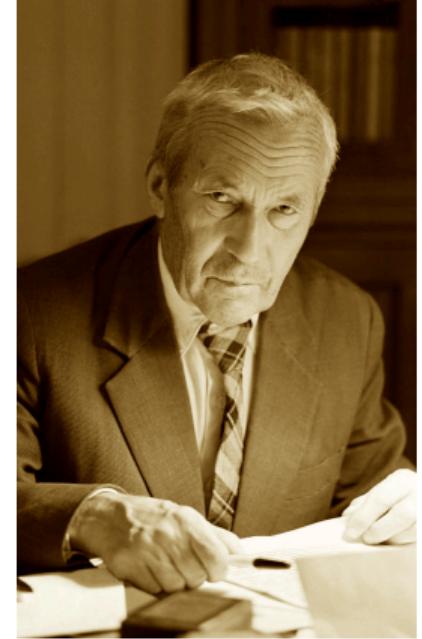
Plausibility as a Measure

[A.N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*, 1933]

Theorem (Bayes' Theorem)

Let $A_1 + A_2 + \dots + A_n = E$ and $A_i \cap A_j = \emptyset$ if $i \neq j$. Then, for any $X \in \mathfrak{F}$,

$$P(A_i | X) = \frac{P(A_i)P(X | A_i)}{\sum_{j=1}^n P(A_j)P(X | A_j)}.$$



Proof.

Apply the Sum Rule to the definition of the conditional probability. □

Andrey N. Kolmogorov
(1903–1987)



The Rules of Probability:

- ▶ the Sum Rule:

$$P(A) = P(A, B) + P(A, \neg B)$$

- ▶ the Product Rule:

$$P(A, B) = P(A \mid B) \cdot P(B) = P(B \mid A) \cdot P(A)$$

- ▶ Bayes' Theorem:

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)} = \frac{P(B \mid A)P(A)}{P(B, A) + P(B, \neg A)}$$

Bayes' Theorem

Inverting Probabilities

$$\underbrace{P(X | D)}_{\text{posterior for } X \text{ given } D} = \frac{\overbrace{P(X)}^{\text{prior for } X} \cdot \overbrace{P(D | X)}^{\text{likelihood for } X}}{\underbrace{P(D)}_{\text{evidence for the model}}} = \frac{P(X) \cdot P(D | X)}{\sum_{x \in \mathcal{X}} P(D | x) P(x)}$$

likelihood is when we have a distribution that is a function of a variable (X), which is being conditioned -> $P(D | X)$

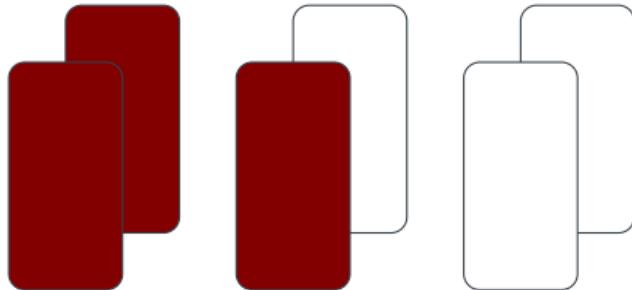
Bayes' Theorem tells us how to update the *belief* in a *hypothesis* X when observing $data D$.

- ▶ $P(D | X)$ is the *likelihood* of X , but the *(conditional) probability for D (given X)*
- ▶ the **model** is the entire thing – prior *and* likelihood
- ▶ despite the name, the prior is not necessarily what you know *before* seeing the data, but the marginal distribution $P(X) = \sum_{d \in \mathcal{D}} P(X, d)$ under *all* possible data.



Bayes' Theorem Appreciation Slides (1)

reasoning quantitatively under uncertainty



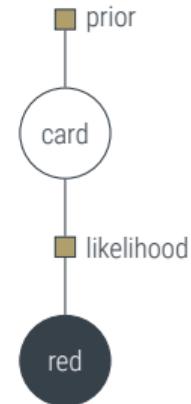
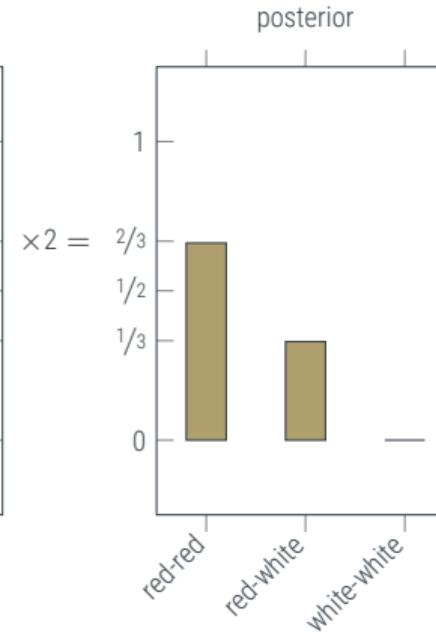
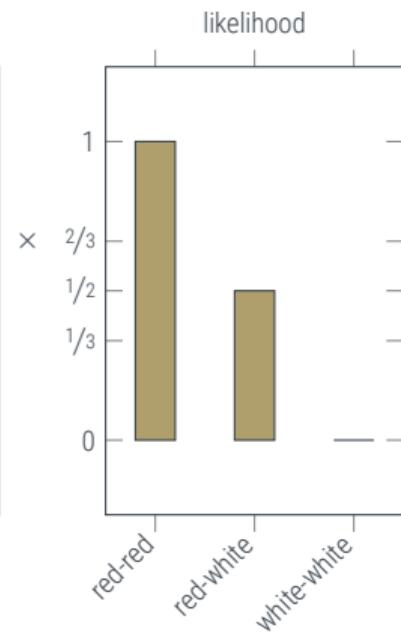
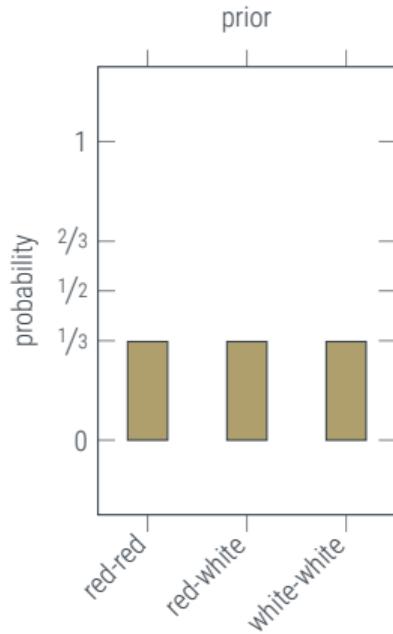
- ▶ $\frac{1}{2}$
- ▶ $\frac{2}{3}$
- ▶ something else
- ▶ I don't know yet

$$\begin{aligned} P(\text{card}|\text{color}) &= \frac{P(\text{card}) \cdot P(\text{color} | \text{card})}{P(\text{color})} = \frac{P(\text{card}) \cdot P(\text{color} | \text{card})}{\sum_{i=1}^3 P(\text{color} | \text{card} = i)P(\text{card} = i)} \\ &= \frac{\frac{1}{3} \cdot P(\text{color} | \text{card})}{\frac{1}{2}} = \frac{2}{3} \cdot \{1 \quad \frac{1}{2} \quad 0\} \end{aligned}$$

Note: It's all about the likelihood, not the prior!

Bayes' Theorem Appreciation Slides (1)

reasoning quantitatively under uncertainty

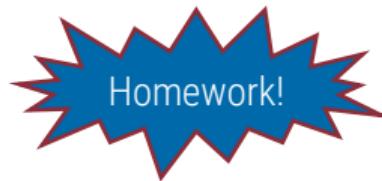


Plausible Reasoning, Revisited

have we succeeded in formalizing common sense?



$A = \text{"it will begin to rain by 6pm"}$
 $B = \text{"the sky will become cloudy before 6pm"}$



$$A \Rightarrow B$$

if A is true, the B is true

Assume: if A is true, then B is true ($A \Rightarrow B$)

if A is true, B becomes more plausible ($P(B | A) > P(B)$)

A is true thus B is true (**modus ponens**)

A is true thus B becomes more plausible

B is false thus A is false (**modus tollens**)

B is false thus A becomes less plausible

B is true thus A becomes more plausible

B is true thus A becomes more plausible

A is false thus B becomes less plausible

A is false thus B becomes less plausible



Bayes' Theorem Appreciation Slides (2)

Bayesian reasoning can help in unintuitive situations

T - test result

C - coronavirus immunity

$$P(T | C) = 0.938 \quad P(\neg T | \neg C) = 0.96$$

$$\begin{aligned} P(C | T) &= \frac{P(T | C)P(C)}{P(T | C)P(C) + P(T | \neg C)P(\neg C)} \\ &= \frac{P(T | C)P(C)}{P(T | C)P(C) + (1 - P(\neg T | \neg C))(1 - P(C))} \\ &= \frac{0.938 \cdot P(C)}{0.938P(C) + 0.04(1 - P(C))} \\ &= \frac{0.938 \cdot P(C)}{0.898P(C) + 0.04} \end{aligned}$$

► $P(C) = 0.01 \rightarrow P(C | T) = 0.19$

Bayes' Theorem Appreciation Slides (3)

Bayesian reasoning matches human reasoning

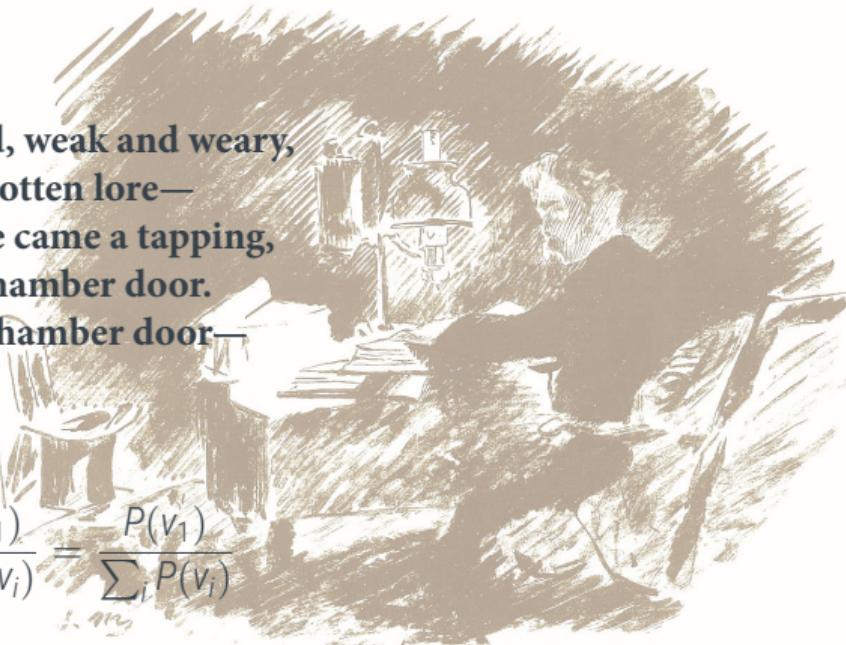


[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

Once upon a midnight dreary, while I pondered, weak and weary,
Over many a quaint and curious volume of forgotten lore—
While I nodded, nearly napping, suddenly there came a tapping,
As of some one gently rapping, rapping at my chamber door.
“Tis some visitor,” I muttered, “tapping at my chamber door—
Only this and nothing more.”

t - tapping the door
v₁ - visitor 1

$$P(v_1 | t) = \frac{P(t | v_1) \cdot P(v_1)}{P(t)} = \frac{P(t | v_1) \cdot P(v_1)}{P(t | v_1) \sum_i P(v_i)} = \frac{P(v_1)}{\sum_i P(v_i)}$$



If the likelihood is constant, data does not help.

Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

Ah, distinctly I remember it was in the bleak December;
And each separate dying ember wrought its ghost upon the floor.
Eagerly I wished the morrow;—vainly I had sought to borrow
From my books surcease of sorrow—sorrow for the lost Lenore—
For the rare and radiant maiden whom the angels name Lenore—
Nameless here for evermore.

$$P(\ell \mid t) = \frac{P(t \mid \ell) \cdot P(\ell)}{P(t)} = \frac{P(t \mid \ell) \cdot P(\ell)}{P(t \mid \ell)P(\ell) + \sum_{i \neq \ell} P(t \mid v_i)P(v_i)} \\ \lesssim 1$$



A very unlikely hypothesis can become dominant if it is the only one explaining the data well.

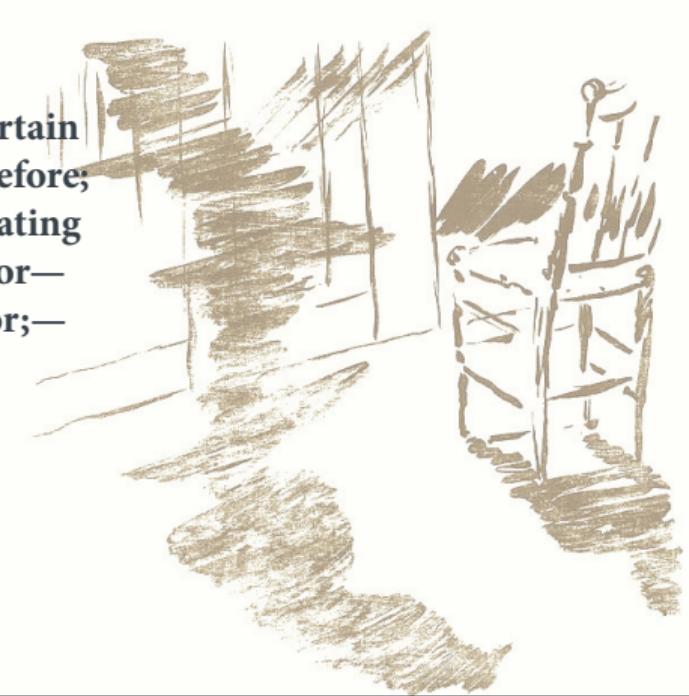
Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

And the silken, sad, uncertain rustling of each purple curtain
Thrilled me—filled me with fantastic terrors never felt before;
So that now, to still the beating of my heart, I stood repeating
“Tis some visitor entreating entrance at my chamber door—
Some late visitor entreating entrance at my chamber door;—
This it is and nothing more.”



$$\begin{aligned} P(\ell \mid t) &= \frac{P(t \mid \ell) \cdot P(\ell)}{P(t)} = \frac{P(t \mid \ell) \cdot 0}{P(t)} \\ &= 0 \end{aligned}$$

No data can revive an a priori impossible hypothesis

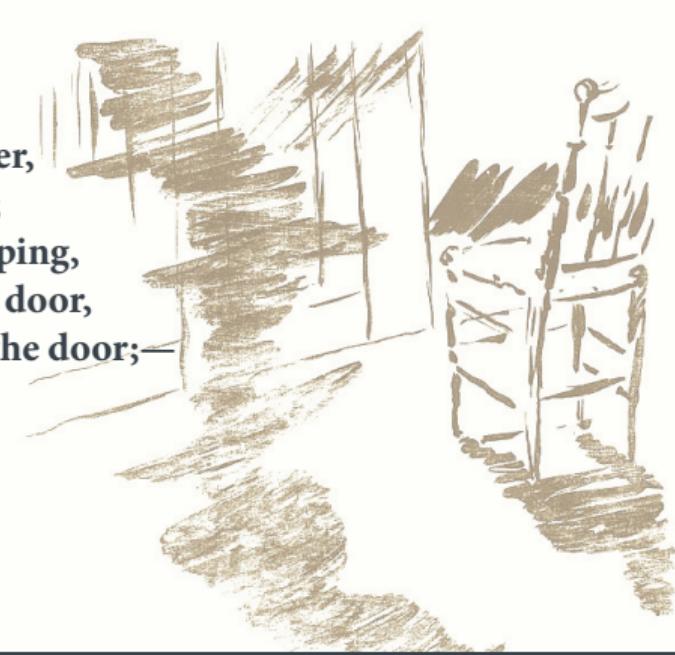
Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

Presently my soul grew stronger; hesitating then no longer,
"Sir," said I, "or Madam, truly your forgiveness I implore;
But the fact is I was napping, and so gently you came rapping,
And so faintly you came tapping, tapping at my chamber door,
That I scarce was sure I heard you"—here I opened wide the door;—
Darkness there and nothing more.



$$P(v) = 0$$

Additional evidence may force you to reconsider your prior.

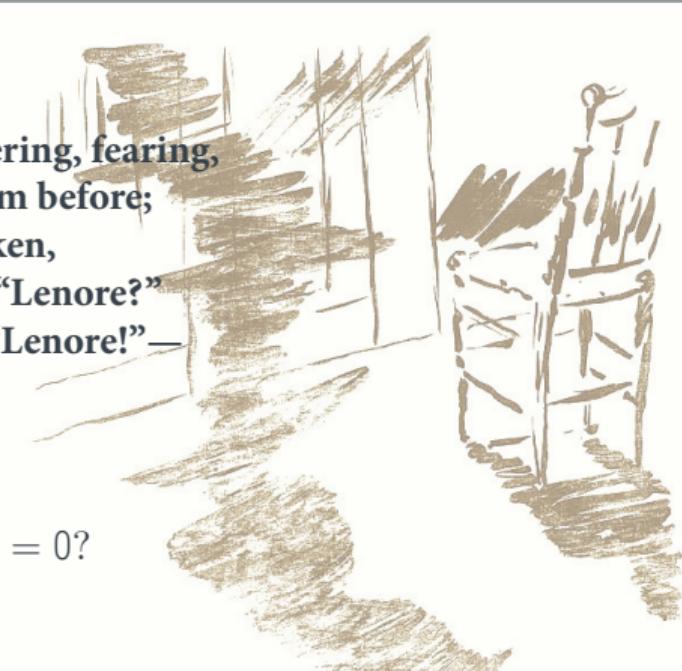
Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

Deep into that darkness peering, long I stood there wondering, fearing,
Doubting, dreaming dreams no mortal ever dared to dream before;
But the silence was unbroken, and the stillness gave no token,
And the only word there spoken was the whispered word, “Lenore?”
This I whispered, and an echo murmured back the word, “Lenore!”—
Merely this and nothing more.



$$P(v) = 0 \quad P(\ell) = 0 \quad P(t) = \sum_i P(t | v_i)P(v_i) = 0?$$

$P(E) = 1$. The hypothesis space has to contain *some* explanation for the data.

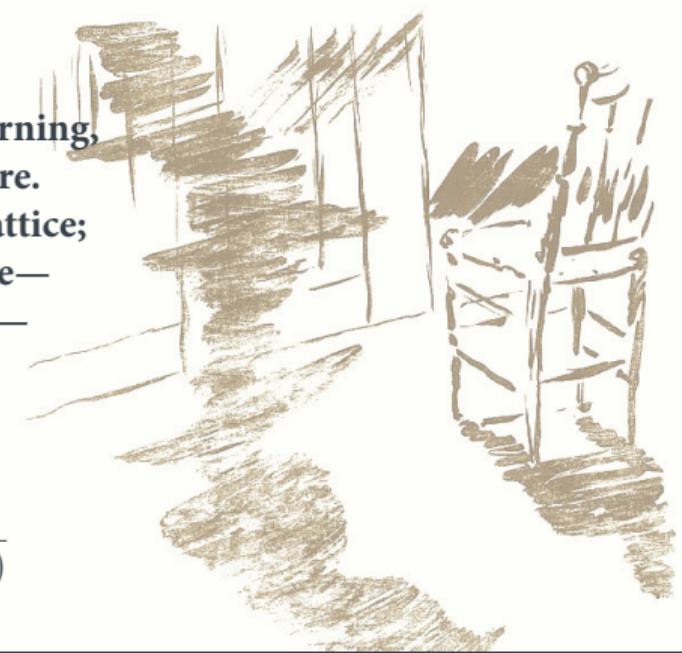
Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

Back into the chamber turning, all my soul within me burning,
Soon again I heard a tapping somewhat louder than before.
“Surely,” said I, “surely that is something at my window lattice;
Let me see, then, what thereat is, and this mystery explore—
Let my heart be still a moment and this mystery explore;—
’Tis the wind and nothing more!”



$$P(w | t) = \frac{P(t | w) \cdot P(w)}{P(t)} = \frac{P(t | w) \cdot P(w)}{P(t, w) + \sum_i P(t, v_i)}$$

The σ -algebra itself, rather than P , is often the most important prior assumption.

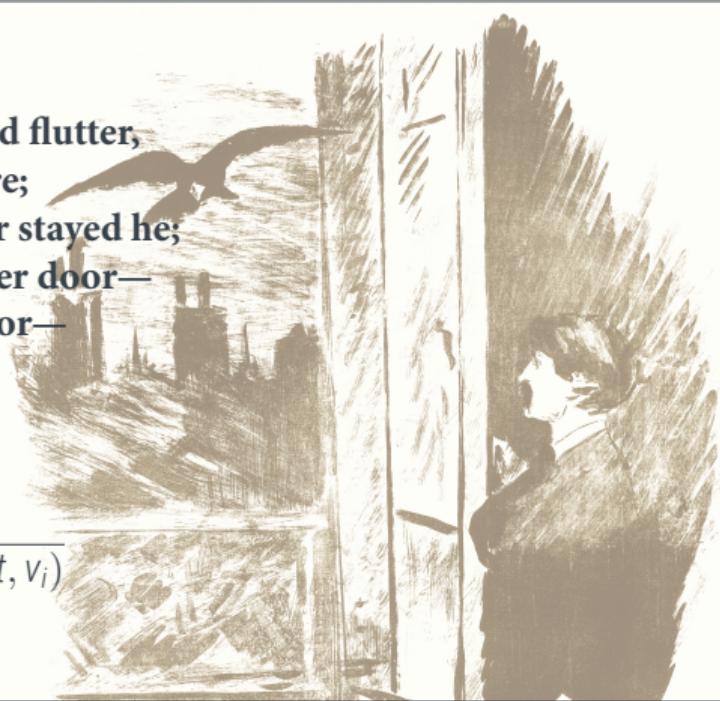
Bayes' Theorem Appreciation Slides (3)

rationality and nothing more



[Edgar Allan Poe, 1845 / Édouard Manet, 1875]

**Open here I flung the shutter, when, with many a flirt and flutter,
In there stepped a stately Raven of the saintly days of yore;
Not the least obeisance made he; not a minute stopped or stayed he;
But, with mien of lord or lady, perched above my chamber door—
Perched upon a bust of Pallas just above my chamber door—
Perched, and sat, and nothing more.**



$$P(r | t) = \frac{P(t | r) \cdot P(r)}{P(t)} = \frac{P(t | r) \cdot P(r)}{P(t, r) + P(t, w) + \sum_i P(t, v_i)}$$

Probabilistic reasoning is a mechanism, it does not replace creativity.



Probability theory is nothing but common sense reduced to calculation.

Pierre-Simon, marquis de Laplace (1749-1827)



Administrative Stuff

Let's meet!

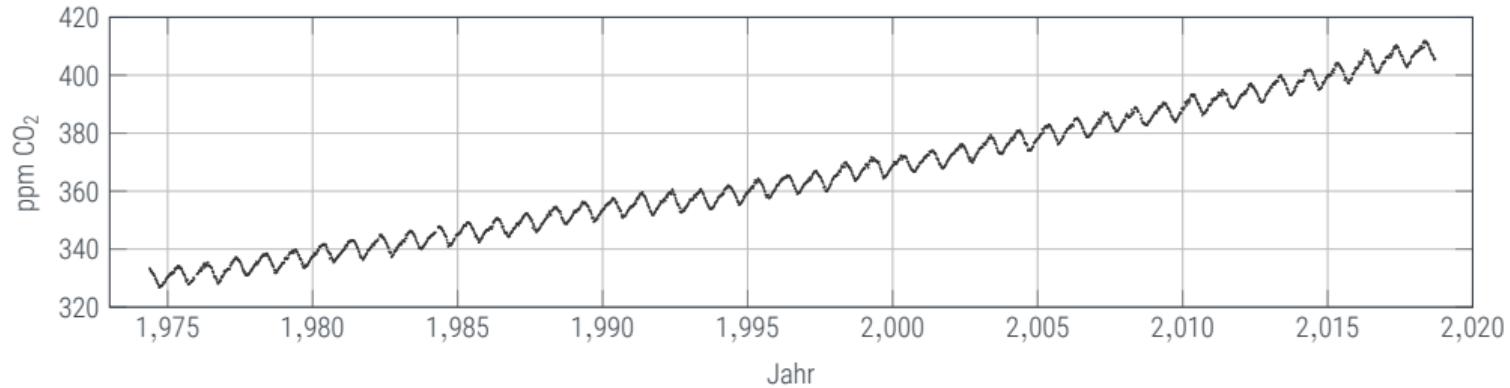
- ▶ *Flipped Classroom on Tuesday, 21 April, 10ct-12*
- ▶ sign-in will be distributed via Ilias
- ▶ sign up on Ilias!
- ▶ please provide feedback on Ilias

<https://uni-tuebingen.de/de/134452>



Exercises

Setting the Scene



- ▶ Goal: Reflect on the utility and power of deep learning (and what you know about it)
- ▶ Things to think about:
 - ▶ How do you include external structural knowledge?
 - ▶ Do you trust the prediction? Why, or why not?
 - ▶ Can your model do anything other than predict future data?

To be able to reason in an uncertain world,
whether to build intelligent machines or find scientific insight,
you have to understand probabilities.

- ▶ both philosophical (Cox) and formal (Kolmogorov) arguments lead to the same two rules:
 sum rule: $P(X) = P(X, Y) + P(X, \neg Y)$
 product rule: $P(X, Y) = P(X) \cdot P(Y | X)$
- ▶ their corollary, **Bayes' Theorem** provides the mechanism for inference:

$$\underbrace{P(X | D)}_{\text{posterior of } X \text{ given } D} = \frac{\overbrace{P(D | X)}^{\text{likelihood of } X \text{ under } D} \cdot \overbrace{P(X)}^{\text{prior of } X}}{\underbrace{P(D)}_{\text{evidence for the model}}}$$

- ▶ the fundamentality of probabilities has been debated at length. Probabilities are not the only inference system, but they are uniquely **general, expressive, and powerful**.
- ▶ **Machine learning and AI** can be approached in various ways. The probabilistic viewpoint is the closest we have to a **theory of everything** for ML.