

WHAT DOES LEARNING MEAN?

Jeffrey Heinz



Stony Brook University

AI Institute
Stony Brook University
February 6, 2020

WHAT DOES LEARNING MEAN?

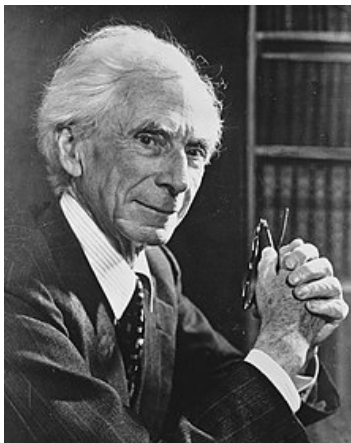
- Personal Motivation
- A Tour of Computational Learning Theory
- ...and consequences thereof
- Examples along this other path

Part I

Personal Motivation

WHAT DOES LEARNING MEAN?

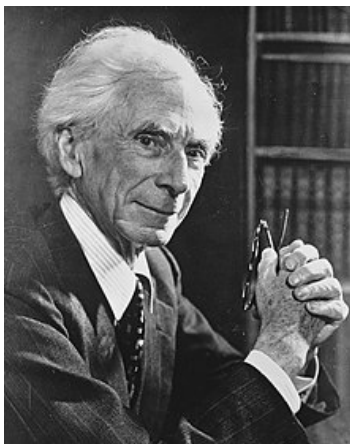
How comes it that human beings, whose contacts with the world are brief and personal and limited, are nevertheless able to know as much as they do know? (1935)



Bertrand Russell (circa 1957)

WHAT DOES LEARNING MEAN?

...if we are to be able
to draw inferences from
these data...we must know
...principles of some kind by
means of which such inferences
can be drawn. (1912)



Bertrand Russell (circa 1957)

ONE THING WE KNOW IS OUR OWN LANGUAGE

[Play video of Ella]

LEARNING LANGUAGE IS DECEPTIVELY HARD



THE HUMAN SPEECHOME PROJECT

“What would it take for the machines he made to think and talk? ‘I thought I could just read the literature on how kids do it, and that would give me a blueprint for building my language and learning robots,’ Roy told me.”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

THE HUMAN SPEECHOME PROJECT

“Over dinner one night, he boasted to Patel, who was then completing her PhD in human speech pathology, that he had already created a robot that was learning the same way kids learn. He was convinced that if it got the sort of input children get, the robot could learn from it.”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

THE HUMAN SPEECHOME PROJECT

“In all, they had captured 90,000 hours of video and 140,000 hours of audio. The 200 terabytes of data covered 85% of the first three years of their son’s life (and 18 months of his little sister’s).”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

THE HUMAN SPEECHOME PROJECT

“ ‘He said *fah*,’ Roy explained, ‘but he was actually clearly referring to a fish on the wall that we were both looking at. The way I knew it was not just coincidence was that right after he looked at it and said it, he turned to me. And he had this kind of look, like a cartoon lightbulb going off – an *Ah, now I get it* kind of look. He’s not even a year old, but there’s a conscious being, in the sense of being self-reflective.’ ”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

THE HUMAN SPEECHOME PROJECT

“ ‘I guess, putting on my AI hat, it was a humbling lesson,’ he continued. ‘A lesson of like, holy shit, there’s a lot more here.’ ”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

THE HUMAN SPEECHOME PROJECT

“Watching his son, Roy had been blown away by ‘the incredible sophistication of what a language learner in the flesh actually looks like and does’.”



Deb Roy



Rupal Patel

(as reported in the Guardian 2018)

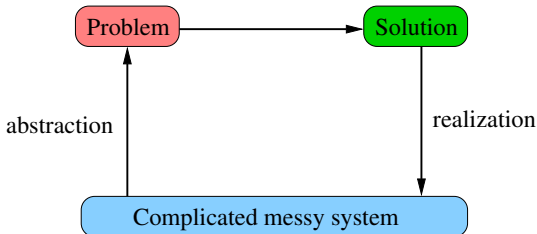
Part II

What is a learning problem?

DEFINING A LANGUAGE-LEARNING PROBLEM

- 1 What does it mean to know a language?
- 2 What does it mean to come by this language from experience?

One strategy is to identify *simple and clear* versions of these problems.



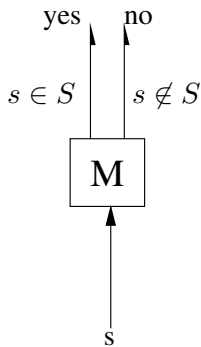
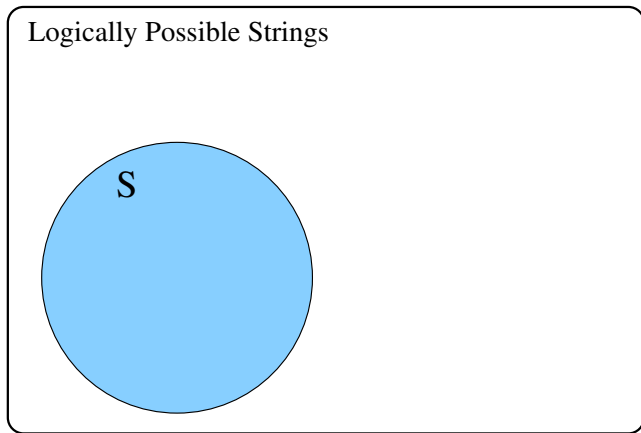
1. WHAT DOES IT MEAN TO KNOW A LANGUAGE?

Knowledge of language includes knowledge of which sequences are licit and which are not.

- John laughed and laughed. ✓
- John and laughed. ✗

1. WHAT DOES IT MEAN TO KNOW A LANGUAGE?

A Membership problem



VARIATIONS THEREOF

Functions on the string domain ...

Function Type	Output Type
$\Sigma^* \rightarrow \{T, F\}$	Booleans
$\Sigma^* \rightarrow \Sigma^*$	Strings
$\Sigma^* \rightarrow \mathbb{N}$	Natural Numbers
$\Sigma^* \rightarrow [0, 1]$	Reals in the Unit Interval
$\Sigma^* \rightarrow P(\Sigma^*)$	Stringsets
...	

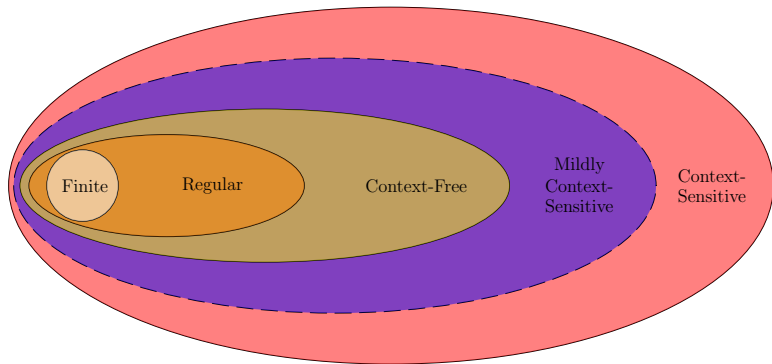
VARIATIONS THEREOF

Functions on the string domain ...

Function Type	Output Type
$\Sigma^* \rightarrow \{T, F\}$	Booleans
$\Sigma^* \rightarrow \Sigma^*$	Strings
$\Sigma^* \rightarrow \mathbb{N}$	Natural Numbers
$\Sigma^* \rightarrow [0, 1]$	Reals in the Unit Interval
$\Sigma^* \rightarrow P(\Sigma^*)$	Stringsets
...	

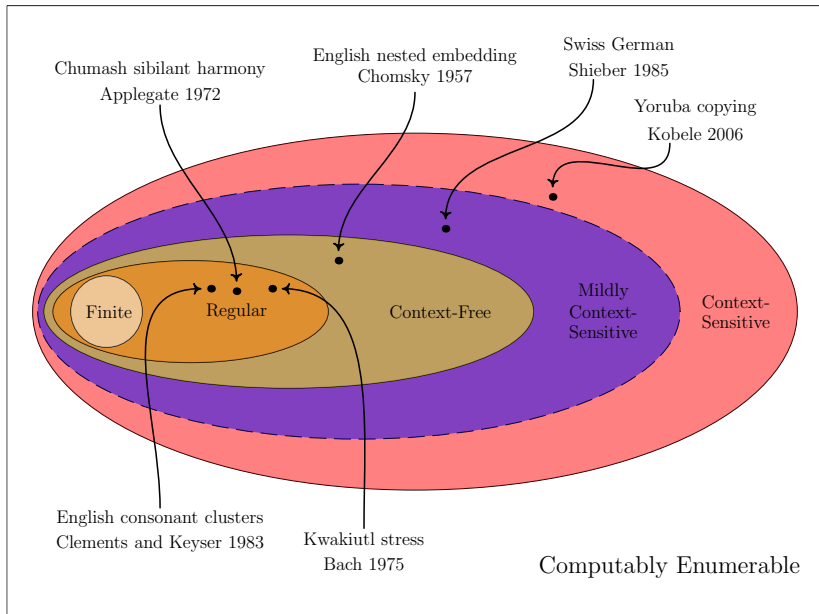
How are functions like those above classified?

CHOMSKY HIERARCHY



Computably Enumerable

CHOMSKY HIERARCHY

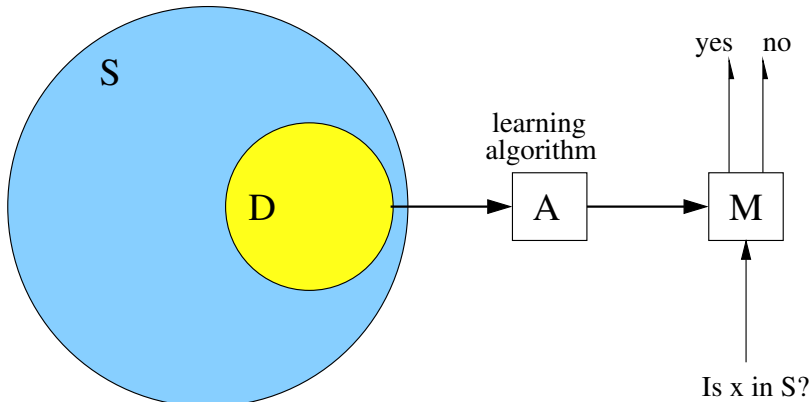


2. WHAT DOES IT MEAN TO COME BY THIS KNOWLEDGE FROM EXPERIENCE?

- 1 Which target functions?
- 2 What kinds of experience?
- 3 What counts as success?

Answering questions 1-2 are important for defining the **instance space** of the learning problem. Question 3 is about conditions on successful solutions.

IN PICTURES



for any S belonging to a class C and for any D from some
'legitimate' class of experience?

NOT ABOUT METHODS (MANY, MANY METHODS)

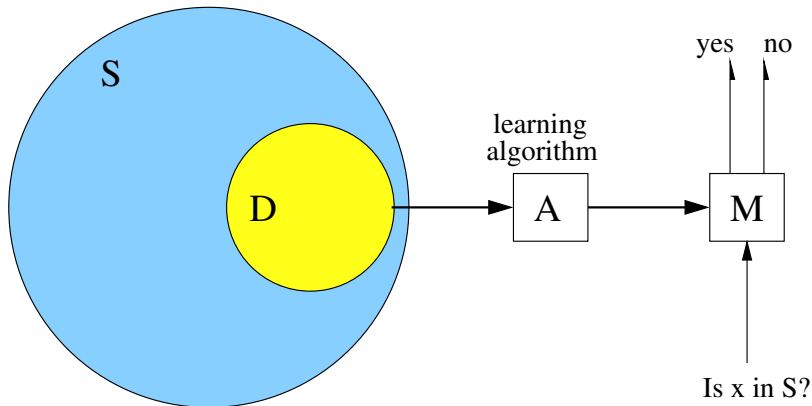
- ① Connectionism/Associative Learning (Rosenblatt 1959, McClelland and Rumelhart 1986, Kapatsinski 2018, a.o.)
- ② Bayesian methods (Bishop 2006, Kemp and Tenenbaum 2008, a.o.)
- ③ Probabilistic Graphical Models (Pearl 1988, Koller and Friedman 2010, a.o.)
- ④ State-merging (Feldman 1972, Angluin 1982, Oncina et al 1992, a.o.)
- ⑤ Statistical Relational Learning (De Raedt 2008, a.o.)
- ⑥ Minimum Description Length (Risannen 1978, Goldsmith a.o..)
- ⑦ Support Vector Machines (Vapnik 1995, 1998 a.o.)
- ⑧ ...

NOT ABOUT METHODS (MANY, MANY METHODS)

Newer methods

- ① Deep NNs (LeCun et al. 2015, Schmidhuber 2015, Goodfellow et al. 2016, a. MANY o.)
 - encoder-decoder networks
 - generative adversarial networks
 - ...
- ② Spectral Learning (Hsu et al 2009, Balle et al. 2012, 2014, a.o.)
- ③ Distributional Learning (Clark and Yoshinaka 2016, a.o.)
- ④ ...

IN PICTURES



CLT studies **conditions** on learning mechanisms/methods!

COMPUTATIONAL LEARNING THEORY

- 1 Identifications in the Limit (Gold 1967)
- 2 Identifications in the Limit with probability p with stochastic input (Angluin 1988a)
- 3 Active/Query Learning (Angluin 1988b)
- 4 Probably Approximately Correct (PAC) Learning (Valiant 1984)
- 5 Optimizing Objective Functions
- 6 ...

FILLING OUT THE INSTANCE SPACE: THE EXPERIENCE

- 1 It is a sequence.
- 2 It is finite.

w_0

w_1

w_2

\dots

w_n

↓ time

FILLING OUT THE INSTANCE SPACE: TYPES OF EXPERIENCE

1 Positive evidence

$$w_0 \in L$$

$$w_1 \in L$$

$$w_2 \in L$$

...

$$w_n \in L$$

↓ time

FILLING OUT THE INSTANCE SPACE: TYPES OF EXPERIENCE

2 Positive and negative evidence

$$w_0 \in L$$

$$w_1 \notin L$$

$$w_2 \notin L$$

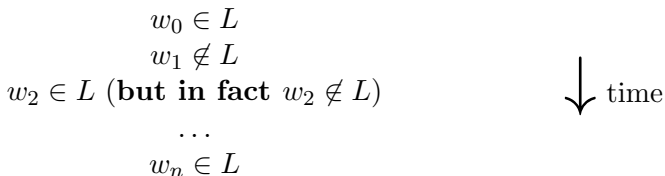
...

$$w_n \in L$$

↓ time

FILLING OUT THE INSTANCE SPACE: TYPES OF EXPERIENCE

③ Noisy evidence



FILLING OUT THE INSTANCE SPACE: TYPES OF EXPERIENCE

④ Queried Evidence

$$w_0 \in L$$

$$w_1 \notin L$$

$w_2 \in L$ (because learner
specifically asked about w_2)

...

$$w_n \in L$$

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$
w_2	$\varphi(\langle w_0, w_1, w_2 \rangle) = G_2$

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$
w_2	$\varphi(\langle w_0, w_1, w_2 \rangle) = G_2$
\dots	

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$
w_2	$\varphi(\langle w_0, w_1, w_2 \rangle) = G_2$
\dots	
w_n	$\varphi(\langle w_0, w_1, w_2, \dots, w_n \rangle) = G_n$

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$
w_2	$\varphi(\langle w_0, w_1, w_2 \rangle) = G_2$
\dots	
w_n	$\varphi(\langle w_0, w_1, w_2, \dots, w_n \rangle) = G_n$
\dots	

↓ time

WHAT COUNTS AS SUCCESS?

- 1 Convergence.
- 2 Imagine an infinite sequence. Is there some point n after which the learner's hypothesis doesn't change (much)?

datum	Learner's Hypothesis
w_0	$\varphi(\langle w_0 \rangle) = G_0$
w_1	$\varphi(\langle w_0, w_1 \rangle) = G_1$
w_2	$\varphi(\langle w_0, w_1, w_2 \rangle) = G_2$
\dots	
w_n	$\varphi(\langle w_0, w_1, w_2, \dots, w_n \rangle) = G_n$
\dots	
w_m	$\varphi(\langle w_0, w_1, w_2, \dots, w_m \rangle) = G_m$

↓ time

Does
 $G_m \simeq G_n$?

FILLING OUT THE INSTANCE SPACE: WHICH EXPERIENCES?

Types of Experience

- 1 Positive-only or positive and negative evidence.
- 2 Noiseless or noisy evidence.
- 3 Queries allowed or not?

Which infinite sequences require convergence?

- 1 only **complete** ones? I.e. where every piece of information occurs at some finite point
- 2 only **computable** ones? I.e. the infinite sequence itself is describable by some grammar

SOME PARAMETERS FOR DEFINING LEARNING PROBLEMS

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

SOME PARAMETERS FOR DEFINING LEARNING PROBLEMS

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

1. Identification in the limit from positive data (Gold 1967)

SOME PARAMETERS FOR DEFINING LEARNING PROBLEMS

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

2. Identification in the limit from positive and negative data
(Gold 1967)

SOME PARAMETERS FOR DEFINING LEARNING PROBLEMS

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

3. Identification in the limit from positive data from c.e. texts
(Gold 1967)
4. Learning context-free and c.e. distributions
(Horning 1969, Angluin 1988)

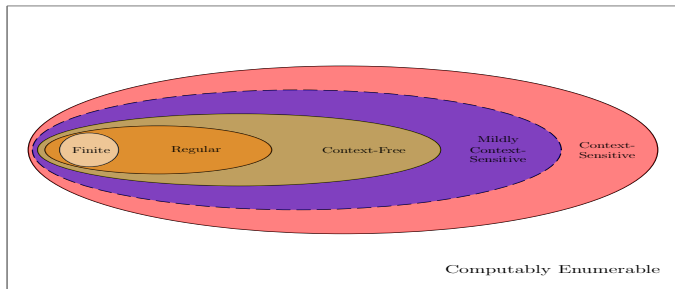
SOME PARAMETERS FOR DEFINING LEARNING PROBLEMS

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

5. Probably Approximately Correct learning
(Valiant 1984, Anthony and Biggs 1991, Kearns and Vazirani 1994)

RESULTS OF FORMAL LEARNING THEORIES: EXISTENCE

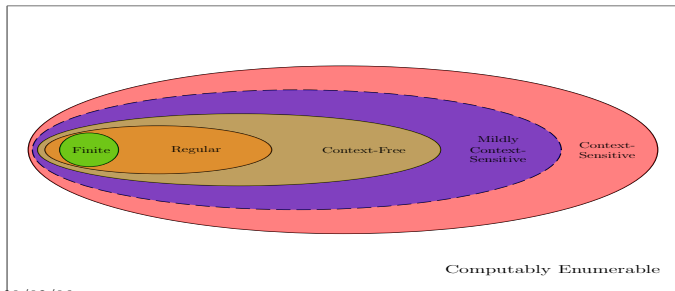
Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence



RESULTS OF FORMAL LEARNING THEORIES: EXISTENCE

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

1. Identification in the limit from positive data (Gold 1967)

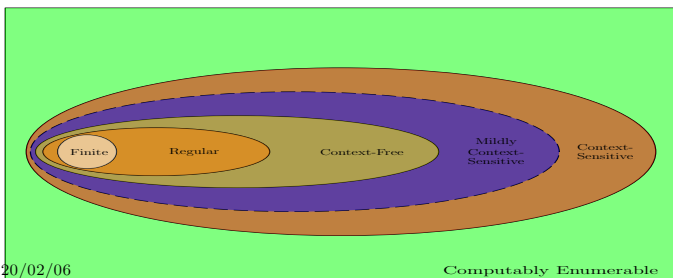


RESULTS OF FORMAL LEARNING THEORIES: EXISTENCE

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

2. Identification in the limit from positive and negative data

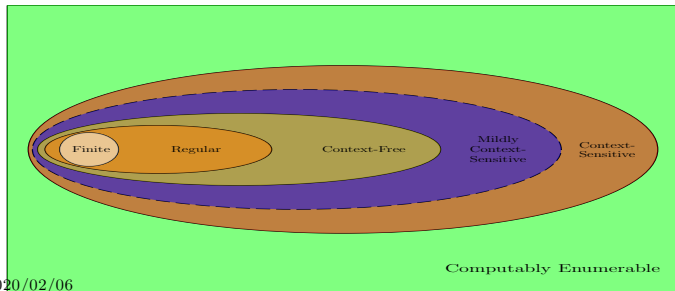
(Gold 1967)



RESULTS OF FORMAL LEARNING THEORIES: EXISTENCE

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

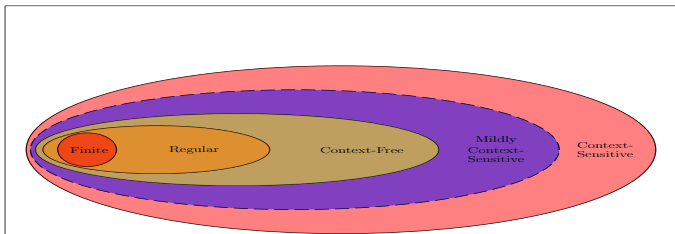
- Identification in the limit from positive data from c.e. texts (Gold 1967)
- Learning context-free and c.e. distributions (Horning 1969, Angluin 1988)



RESULTS OF FORMAL LEARNING THEORIES: EXISTENCE

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

5. Probably Approximately Correct learning
(Valiant 1984, Anthony and Biggs 1991, Kearns and Vazirani 1994)



RESULTS OF FORMAL LEARNING THEORY: FEASIBILITY

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

RESULTS OF FORMAL LEARNING THEORY: FEASIBILITY

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

1. Identification in the limit from positive data (Gold 1967)

No superfinite class is learnable.

The finite class is feasibly learnable.

RESULTS OF FORMAL LEARNING THEORY: FEASIBILITY

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

2. Identification in the limit from positive and negative data
(Gold 1967)

The c.e. class is learnable but NOT even the regular class is feasibly learnable (see appendix).

RESULTS OF FORMAL LEARNING THEORY: FEASIBILITY

Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

3. Identification in the limit from positive data from c.e. texts (Gold 1967)
4. Learning context-free and c.e. distributions (Horning 1969, Angluin 1988)

The c.e. class of languages and distributions is learnable but NOT even the regular class is feasibly learnable.

RESULTS OF FORMAL LEARNING THEORY: FEASIBILITY

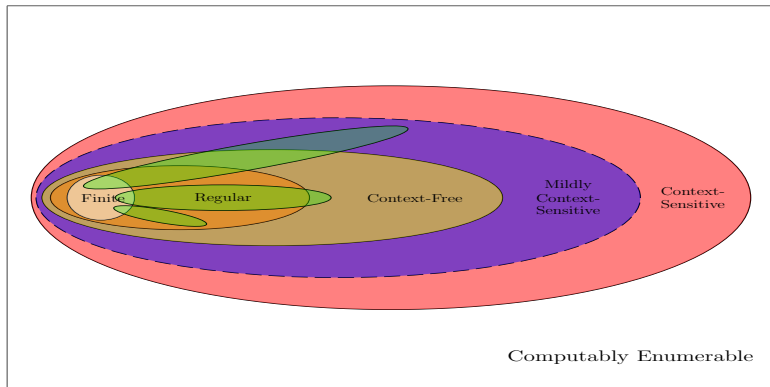
Makes learning easier	Makes learning harder
positive and negative evidence	positive evidence only
noiseless evidence	noisy evidence
queries permitted	queries not permitted
approximate convergence	exact convergence
complete infinite sequences	any infinite sequence
computable infinite sequences	any infinite sequence

5. Probably Approximately Correct learning
(Valiant 1984, Anthony and Biggs 1991, Kearns and Vazirani 1994)

Not even the finite class of languages is learnable.

FORMAL LEARNING THEORY: POSITIVE RESULTS

Many classes which cross-cut the Chomsky hierarchy and exclude some finite languages are feasibly learnable in the senses discussed.



(Angluin 1980, 1982, Garcia et al. 1990, Muggleton 1990, Denis et al. 2002, Fernau 2003, Yokomori 2003, Oates et al. 2006, Niyogi 2006, Clark and Eryaud 2007, Heinz 2008, to appear, Yoshinaka 2008, Case et al. 2009, de la Higuera 2010)

SUMMARY

- 1 Structured, restricted hypothesis spaces can be feasibly learned.
- 2 The positive learning results are proven results, and the proofs are often constructive.
- 3 The claim that “statistical learning” is more powerful than “symbolic learning” mischaracterizes the learning issues.
- 4 The real issue is how to reduce the size of the instance space.

GOLD 1967

Gold provides three ways to interpret his three main results:

- ① “The class of natural languages is much smaller than one would expect from our present models of syntax. That is, even if English is context-sensitive, it is not true that any context-sensitive language can occur naturally ...”
- ② “The child receives negative instances by being corrected in a way that we do not recognize ...”
- ③ “There is an a priori restriction on the class of texts [presentations of data; i.e. infinite sequences of experience] which can occur ...”

(Gold 1967: 453-4)

Part IV

Consequences of CLT (or The Perceptual Misconception Machine)

THE MAIN LESSON FROM CLT

There is no free lunch.

- 1 There is **no** algorithm that can feasibly learn **any** pattern P , even with lots of data from P .

THE MAIN LESSON FROM CLT

There is no free lunch.

- 1 There is **no** algorithm that can feasibly learn **any** pattern P , even with lots of data from P .
- 2 But—There are algorithms that can feasibly learn patterns **which belong to a suitably structured class C** .

Gold 1967, Angluin 1980, Valiant 1984,
Wolpert and McReady 1997, a.o.

THE PERPETUAL MOTION MACHINE



October 1920 issue of Popular Science magazine, on perpetual motion.

“Although scientists have established them to be impossible under the laws of physics, perpetual motion continues to capture the imagination of inventors.”

https://en.wikipedia.org/wiki/Perpetual_motion

THE PERPETUAL MISCONCEPTION MACHINE

\exists machine-learning algorithm A , \forall patterns P with enough data D from P : $A(D) \approx P$.

- 1 It's just not true.

THE PERPETUAL MISCONCEPTION MACHINE

\exists machine-learning algorithm A , \forall patterns P with enough data D from P : $A(D) \approx P$.

- 1 It's just not true.
- 2 What is true is this:
 \forall patterns P , \exists data D and ML A : $A(D) \approx P$.

THE PERPETUAL MISCONCEPTION MACHINE

\exists machine-learning algorithm A , \forall patterns P with enough data D from P : $A(D) \approx P$.

- 1 It's just not true.
- 2 What is true is this:
 \forall patterns P , \exists data D and ML A : $A(D) \approx P$.
- 3 In practice, the misconception means searching for A and D so that your approximation is better than everyone else's.

THE PERPETUAL MISCONCEPTION MACHINE

\exists machine-learning algorithm A , \forall patterns P with enough data D from P : $A(D) \approx P$.

- 1 It's just not true.
- 2 What is true is this:
 \forall patterns P , \exists data D and ML A : $A(D) \approx P$.
- 3 In practice, the misconception means searching for A and D so that your approximation is better than everyone else's.
- 4 With next pattern P' , we will have no guarantee A will work, we will have to search again.

COMPUTATIONAL LAWS OF LEARNING

Feasibly solving a learning problem requires defining a target class C of patterns.

- 1 The class C cannot be all patterns, or even all computable patterns.

COMPUTATIONAL LAWS OF LEARNING

Feasibly solving a learning problem requires defining a target class C of patterns.

- 1 The class C cannot be all patterns, or even all computable patterns.
- 2 Class C must have *more structure*, and many logically possible patterns *must be outside* of C .

COMPUTATIONAL LAWS OF LEARNING

Feasibly solving a learning problem requires defining a target class C of patterns.

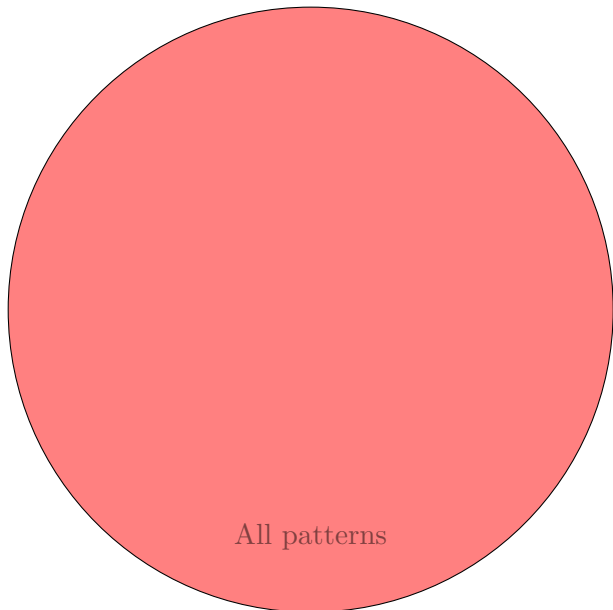
- 1 The class C cannot be all patterns, or even all computable patterns.
- 2 Class C must have *more structure*, and many logically possible patterns *must be outside* of C .
- 3 There is no avoiding prior knowledge.

COMPUTATIONAL LAWS OF LEARNING

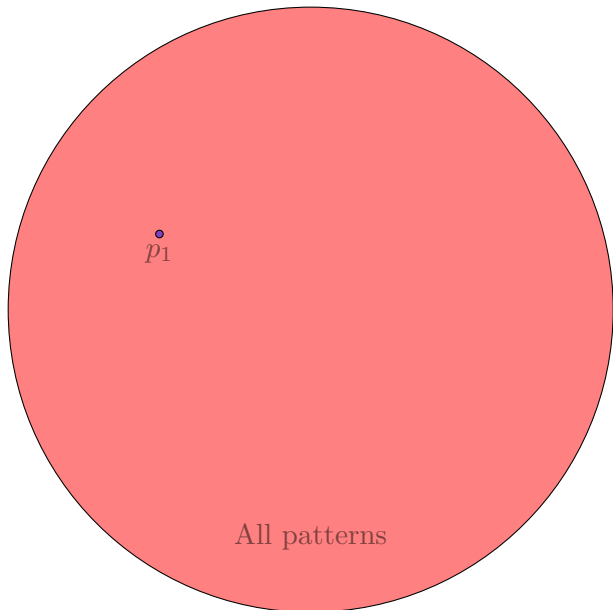
Feasibly solving a learning problem requires defining a target class C of patterns.

- 1 The class C cannot be all patterns, or even all computable patterns.
- 2 Class C must have *more structure*, and many logically possible patterns *must be outside* of C .
- 3 There is no avoiding prior knowledge.
- 4 Do not “confuse ignorance of biases with absence of biases.” (Rawski and Heinz 2019)

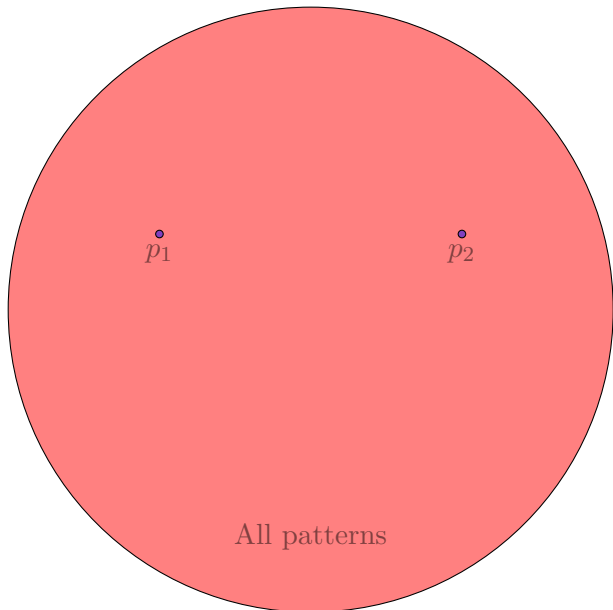
IN PICTURES: GIVEN ML ALGORITHM A



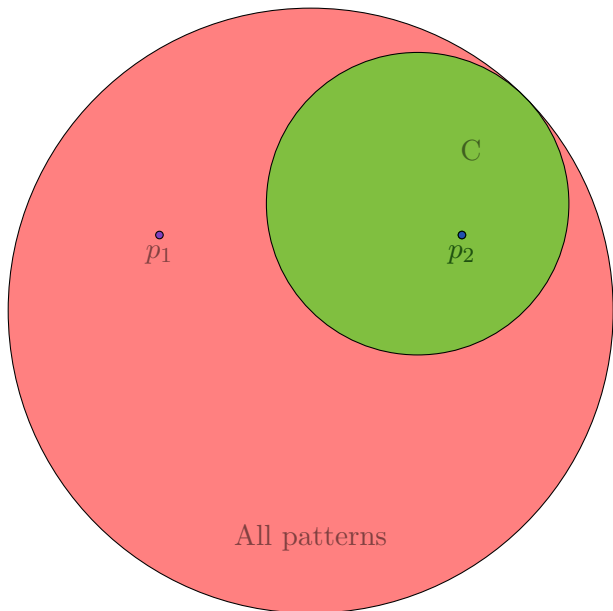
IN PICTURES: GIVEN ML ALGORITHM A



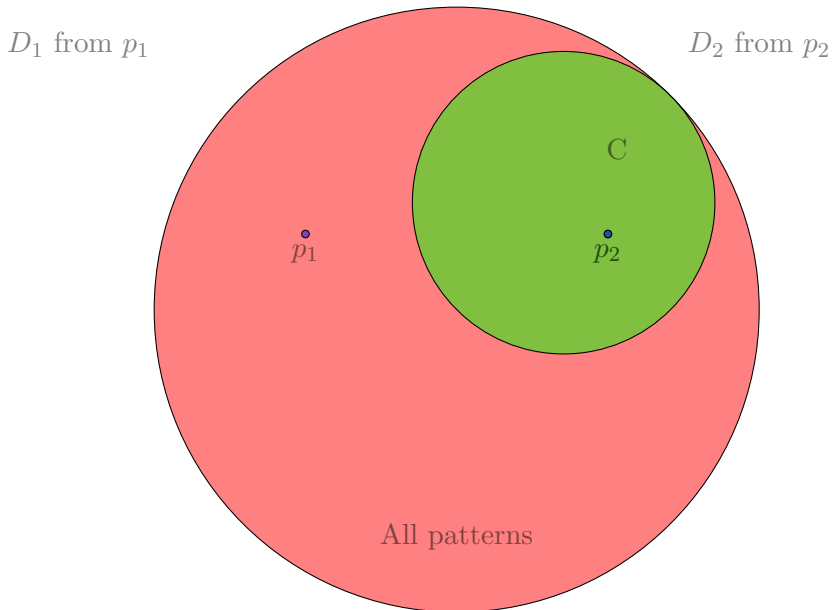
IN PICTURES: GIVEN ML ALGORITHM A



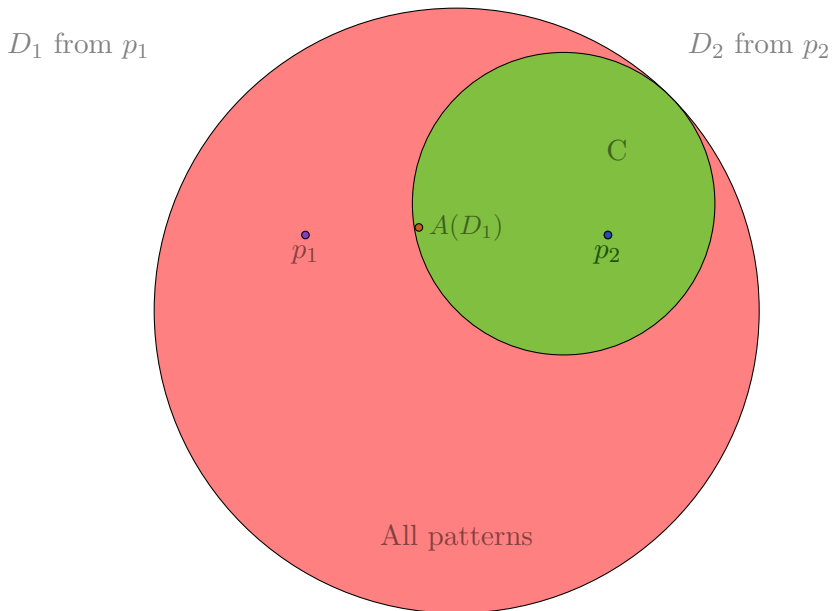
IN PICTURES: GIVEN ML ALGORITHM A



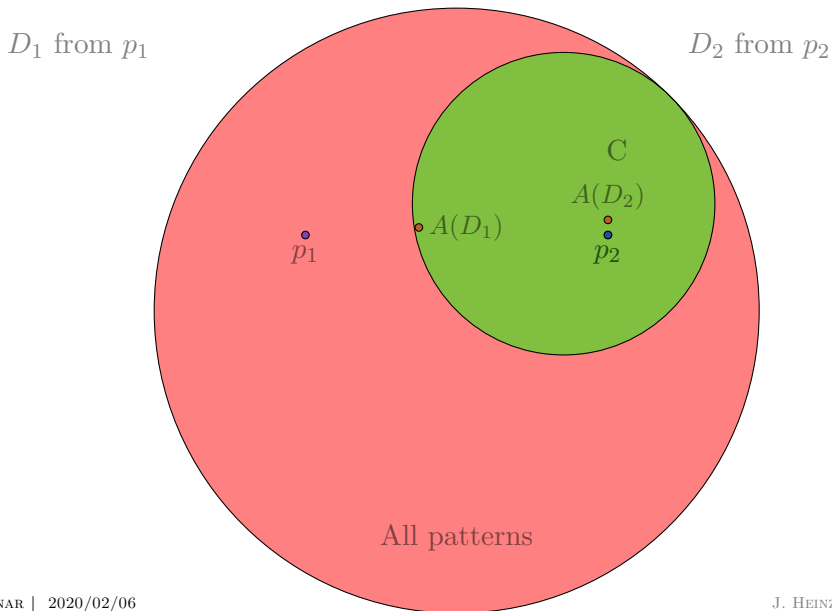
IN PICTURES: GIVEN ML ALGORITHM A



IN PICTURES: GIVEN ML ALGORITHM A



IN PICTURES: GIVEN ML ALGORITHM A



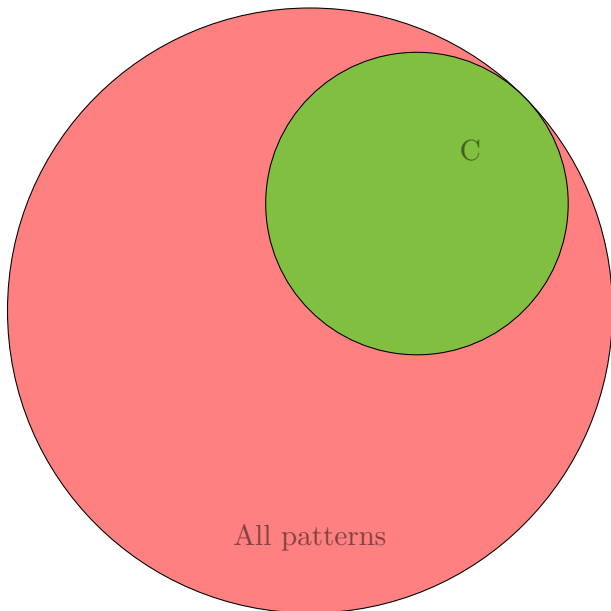
THE PERPETUAL MISCONCEPTION MACHINE

When you believe in things that you don't understand
then you suffer.

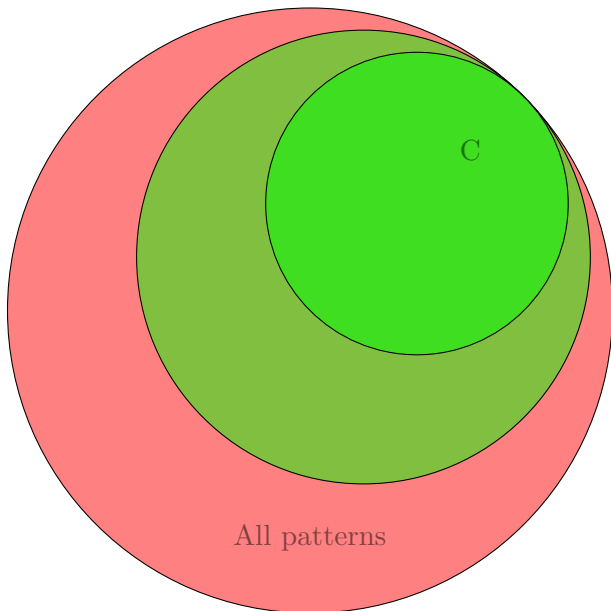
– Stevie Wonder



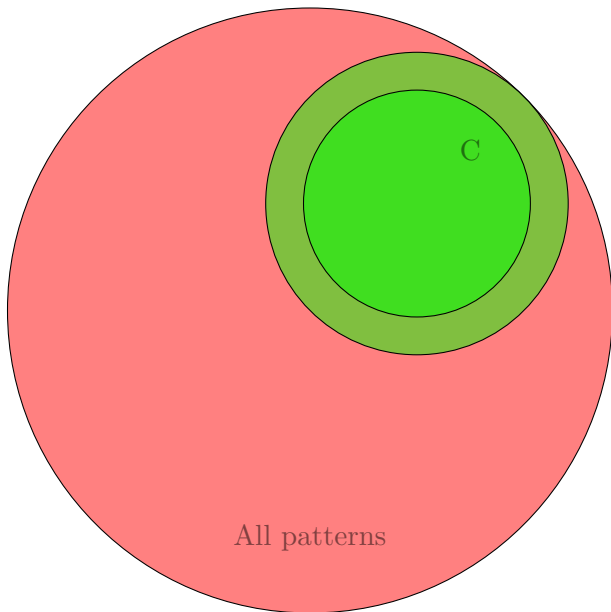
GO SMALLER, NOT BIGGER!



GO SMALLER, NOT BIGGER!



GO SMALLER, NOT BIGGER!



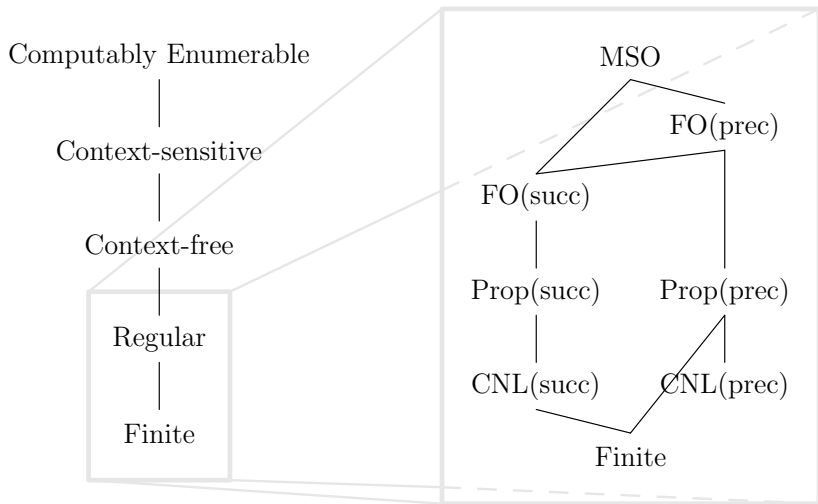
All patterns

C

Part III

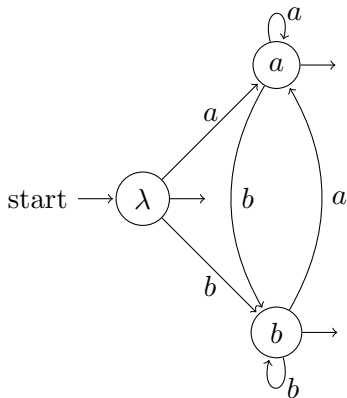
What does ‘going smaller’ look like?

REGULAR AND SUBREGULAR LANGUAGES



CLASSES DEFINED BY SINGLE DFAS

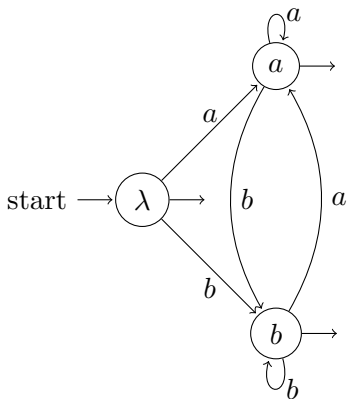
Example: Strictly 2-Local Languages



Parameters
$\theta_{\times a}$
$\theta_{\times b}$
$\theta_{\times \times}$
θ_{aa}
θ_{ab}
$\theta_{a \times}$
θ_{ba}
θ_{bb}
$\theta_{b \times}$

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



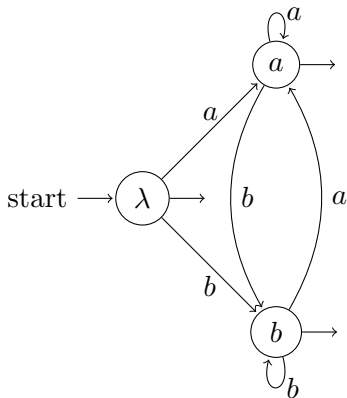
Parameters

 $\theta_{\times a}$ $\theta_{\times b}$ $\theta_{\times \times}$ θ_{aa} θ_{ab} $\theta_{a \times}$ θ_{ba} θ_{bb} $\theta_{b \times}$

$$D = \langle ab, aabb \rangle$$

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



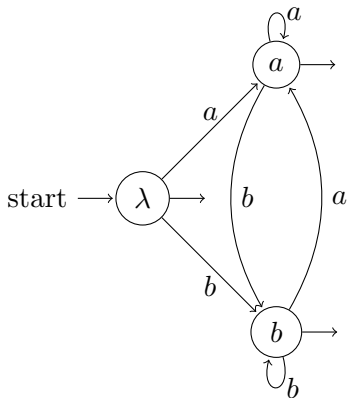
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	
$\theta_{b \times}$	1

$$D = \langle \textcolor{blue}{ab}, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



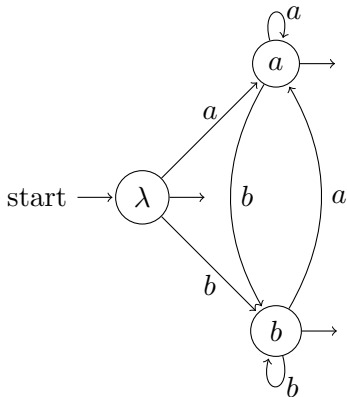
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	1
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	1
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

CLASSES DEFINED BY SINGLE DFAS

Example: Strictly 2-Local Languages



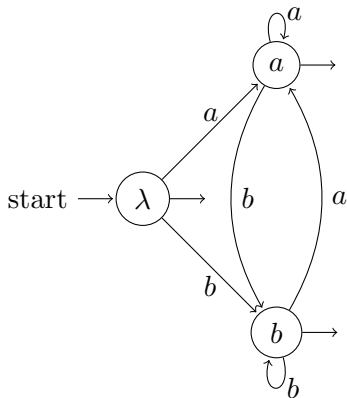
Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	0
$\theta_{\times \times}$	0
θ_{aa}	1
θ_{ab}	1
$\theta_{a \times}$	0
θ_{ba}	0
θ_{bb}	1
$\theta_{b \times}$	1

$$D = \langle ab, aabb \rangle$$

Smallest language consistent with D in C is obtained by passing D through DFA and ‘activating’ parsed transitions. (Heinz and Rogers 2013)

CLASSES DEFINED BY SINGLE DFAS

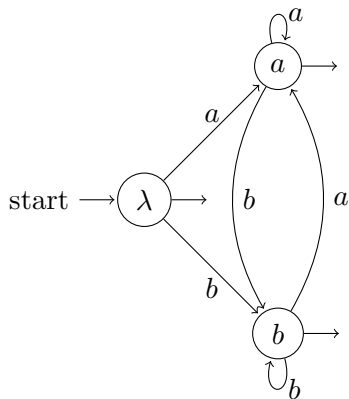
Example: Bigram model



Parameters
$\theta_{\times a}$
$\theta_{\times b}$
$\theta_{\times \times}$
θ_{aa}
θ_{ab}
$\theta_{a \times}$
θ_{ba}
θ_{bb}
$\theta_{b \times}$

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



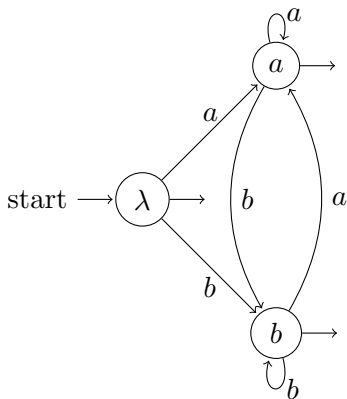
Parameters

 $\theta_{\times a}$ $\theta_{\times b}$ $\theta_{\times \times}$ θ_{aa} θ_{ab} $\theta_{a \times}$ θ_{ba} θ_{bb} $\theta_{b \times}$

$$D = \langle ab, aabb \rangle$$

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	
θ_{ab}	1
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	
$\theta_{b \times}$	1

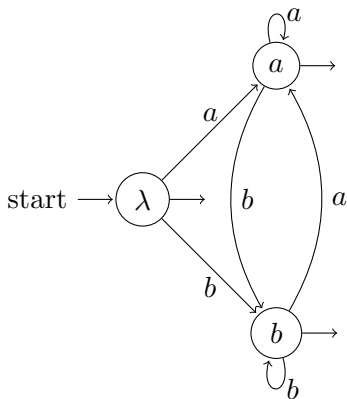
$$D = \langle \textcolor{blue}{ab}, aabb \rangle$$

Maximum Likelihood Estimate (MLE) is obtained by passing D through DFA and normalizing.

(Vidal et al. 2005)

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



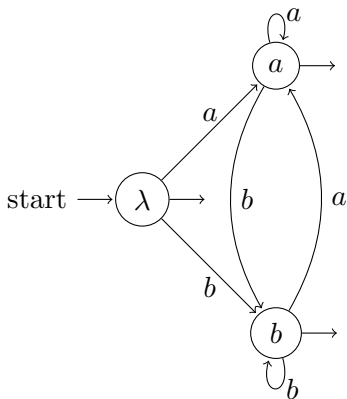
Parameters	
$\theta_{\times a}$	2
$\theta_{\times b}$	
$\theta_{\times \times}$	
θ_{aa}	1
θ_{ab}	2
$\theta_{a \times}$	
θ_{ba}	
θ_{bb}	1
$\theta_{b \times}$	2

$$D = \langle ab, aabb \rangle$$

Maximum Likelihood Estimate (MLE) is obtained by passing D through DFA and normalizing.
(Vidal et al. 2005)

CLASSES DEFINED BY SINGLE DFAS

Example: Bigram model



Parameters	
$\theta_{\times a}$	1
$\theta_{\times b}$	0
$\theta_{\times \times}$	0
θ_{aa}	1/3
θ_{ab}	2/3
$\theta_{a \times}$	0
θ_{ba}	0
θ_{bb}	1/3
$\theta_{b \times}$	2/3

$$D = \langle ab, aabb \rangle$$

Maximum Likelihood Estimate (MLE) is obtained by passing D through DFA and normalizing.

(Vidal et al. 2005)

LEARNING RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$		
language	$f : \Sigma^* \rightarrow [0, 1]$		

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013, Shibata and Heinz 2019)

LEARNING RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	
language	$f : \Sigma^* \rightarrow [0, 1]$		

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013, Shibata and Heinz 2019)

LEARNING RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013, Shibata and Heinz 2019)

LEARNING RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	✓
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013, Shibata and Heinz 2019)

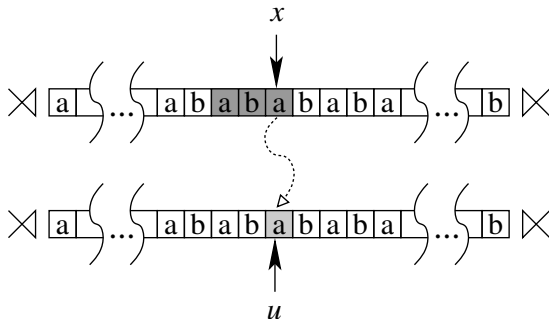
LEARNING RESULTS

		Class C defined with	
		single DFA	finitely many DFA
type of	$f : \Sigma^* \rightarrow \{0, 1\}$	✓	✓
language	$f : \Sigma^* \rightarrow [0, 1]$	✓	✓

- 1 For Boolean languages, the learning algorithms return the smallest language in C which includes D .
- 2 For Stochastic languages, the MLE returns the language in C which maximizes likelihood of D .

(Vidal et al. 2005, Heinz and Rogers 2013, Shibata and Heinz 2019)

INPUT STRICTLY LOCAL FUNCTIONS

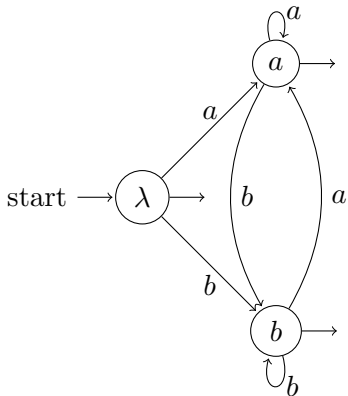


For every Input Strictly 3-Local function, the output string u of each input element x depends only on x and the two input elements previous to x . In other words, the contents of the lightly shaded cell only depends on the contents of the darkly shaded cells.

(Chandlee et al. 2014, 2015, Chandlee and Heinz 2018,
Chandlee et al. 2018)

CLASSES DEFINED BY SINGLE DFAS

Example: 2-ISL Model (Parameter values are strings!)



Parameters	
$\theta_{\times a}$	a
$\theta_{\times b}$	b
$\theta_{\times \times}$	\times
θ_{aa}	ba
θ_{ab}	b
$\theta_{a \times}$	\times
θ_{ba}	a
θ_{bb}	ab
$\theta_{b \times}$	\times

$aabb \mapsto ababab$

ISL FUNCTION LEARNING RESULTS

- Particular finite-state transducers can be used to represent ISL functions.
- Automata-inference techniques (de la Higuera 2010) are used to learn these transducers.

Theorems

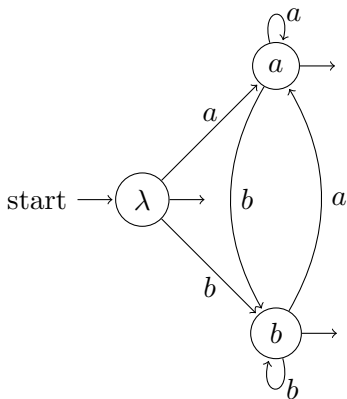
Given k and a sufficient sample of (u, s) pairs any k -ISL function can be *exactly* learned in *polynomial* time and data.

- ISLFLA (quadratic time and data)
- SOSFIA (linear time and data)
- OSLFLA (quadratic time and data)

(Chandlee et al. 2014, 2015, Jardine et al. 2014)

CLASSES DEFINED BY SINGLE DFAS

Example: 2-ISL Model (Parameters are stringsets!)

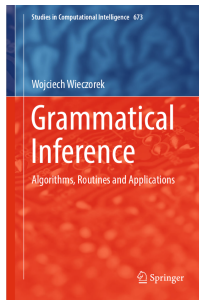
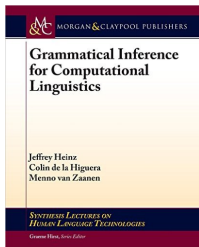
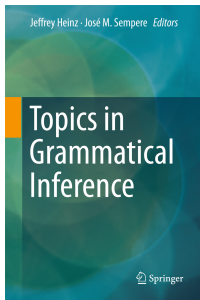
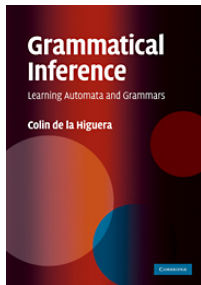


Parameters	
$\theta_{\times a}$	$\{a\}$
$\theta_{\times b}$	$\{b\}$
$\theta_{\times \times}$	$\{\times\}$
θ_{aa}	$\{a, ba\}$
θ_{ab}	$\{b\}$
$\theta_{a \times}$	$\{\times\}$
θ_{ba}	$\{a\}$
θ_{bb}	$\{b, ab\}$
$\theta_{b \times}$	$\{\times\}$

$$aabb \mapsto \{aabb, ababb, aabab, ababab\}$$

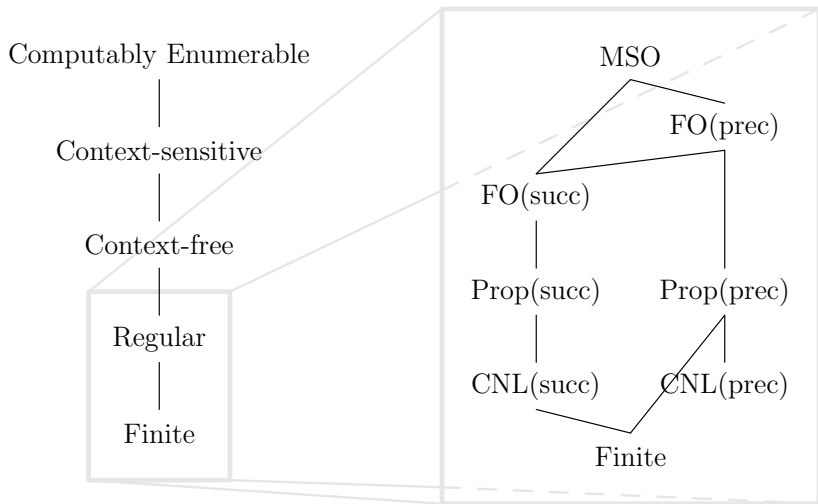
(work in progress with Kiran Eiden and Eric Schieferstein)

GRAMMATICAL INFERENCE

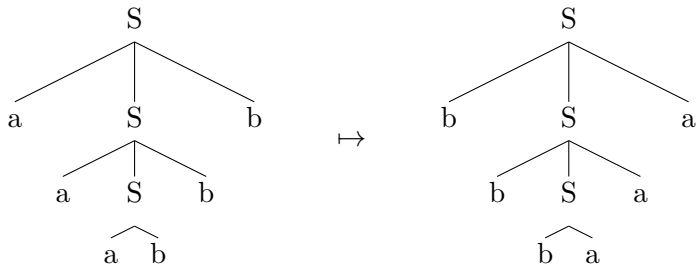


ICGI 2020 in NYC August 26-28!!
<https://grammarlearning.org/>

REGULAR AND SUBREGULAR LANGUAGES



ISL BOTTOM-UP TREE TRANSDUCERS



(Graf 2020, Ji and Heinz 2020)

Part IV

Summing Up

PERSONAL VIEW

- It is a fascinating question how computers—human or machine—can learn anything, and
- How we come by the knowledge of our own language is especially interesting.

Two key ideas

- 1 Mathematics and theoretical computer science can be developed to provide stronger/tighter characterizations of natural language patterns.
- 2 Computational Learning Theory stresses the importance and necessity of structured hypothesis spaces. Don't treat them cavalierly!

Appendix

Extra Slides

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).
- 5 Clark and Thollard (2004) present an algorithm which learns the class of PDFAs in a modified PAC setting. (See also Parekh and Hanover 2001.)

LEARNING RESULTS FOR REGULAR LANGUAGES

- 1 The class of all DFAs is not identifiable in the limit from positive data (Gold 1967).
- 2 It is NP-hard to find the minimal DFA consistent with a finite sample of positive and negative examples (Gold 1978).
- 3 Each DFA admits a characteristic sample D of positive and negative examples such that RPNI identifies the DFA from any superset of D in cubic time (Oncina and Garcia 1992, DuPont 1996).
- 4 ALEGRIA/RLIPS (based on RPNI) (Carrasco and Oncina 1994, 1999) learns the class of PDFAs in polynomial time with probability one (de la Higuera and Thollard 2001).
- 5 Clark and Thollard (2004) present an algorithm which learns the class of PDFAs in a modified PAC setting. (See also Parekh and Hanover 2001.)
- 6 Maximization-Expectation techniques are used to learn the class of PNFAs, but there is no guarantee to find a global optimum (Rabiner 1989).