Introduction
oooo

Preliminaries
oooooo

Models and Algorithms Lack a Simplicity Bias
oooooooooooooo

Learnability with a Simplicity Bias
oooooooooo

Conclusion
ooo

# Kolmogorov Complexity and How it Illuminates our Limitations to Let Machines Learn Simple Functions
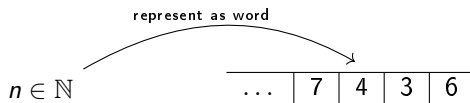
Lukas Rüttgers

IIS, Tsinghua University

July 3, 2024

1 Introduction

2 Preliminaries

3 Models and Algorithms Lack a Simplicity Bias
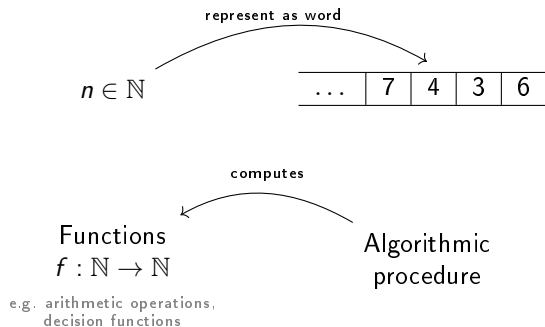
4 Learnability with a Simplicity Bias

5 Conclusion

Introduction
○●○○

Preliminaries
○○○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Motivation

How humans learn arithmetic:



represent as word

$$n \in \mathbb{N}$$

| ... | 7 | 4 | 3 | 6 |

## Motivation

How humans learn arithmetic:



represent as word

$n \in \mathbb{N}$

| ... | 7 | 4 | 3 | 6 |

computes

Functions
$f : \mathbb{N} \to \mathbb{N}$

e.g. arithmetic operations,
decision functions

Algorithmic
procedure

## Motivation

How humans learn arithmetic:



represent as word

$n \in \mathbb{N}$

| . . . | 7 | 4 | 3 | 6 |

computes

recursively applies
symbol transformations

Functions
$f : \mathbb{N} \to \mathbb{N}$

e.g. arithmetic operations,
decision functions
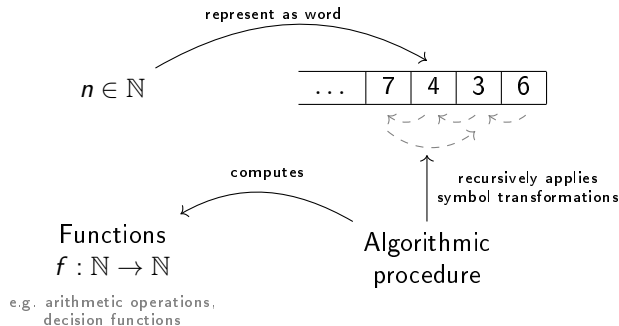
Algorithmic
procedure

Motivation

How humans learn arithmetic:



Expedient insights:

1. Recursive algorithmic descriptions generalize to unseen instances

represent as word

$n \in \mathbb{N}$

... | 7 | 4 | 3 | 6

computes

recursively applies
symbol transformations

Functions
$f : \mathbb{N} \to \mathbb{N}$

e.g. arithmetic operations,
decision functions

Algorithmic
procedure

## Motivation

How humans learn arithmetic:



Expedient insights:

1. Recursive algorithmic descriptions generalize to unseen instances

2. Generalization ability relies on some inductive simplicity bias

## Highlights

Kolmogorov Complexity naturally quantifies the *simplicity* of a function.

In light of this complexity measure we will see that

## Highlights

Kolmogorov Complexity naturally quantifies the *simplicity* of a function.

In light of this complexity measure we will see that

- feed-forward neural networks (and any another non-recursive models) cannot express some simple functions

## Highlights

Kolmogorov Complexity naturally quantifies the *simplicity* of a function.

In light of this complexity measure we will see that

- feed-forward neural networks (and any another non-recursive models) cannot express some simple functions
- incorporating Kolmogorov complexity as a simplicity bias into learning algorithms allows to
  - learn *any* computable function (e.g. prime numbers) with finite resources
  - learn some functions with *less* samples than usual (e.g. parity functions).

## Outline

**1** Introduction

**2** Preliminaries

**3** Models and Algorithms Lack a Simplicity Bias

**4** Learnability with a Simplicity Bias

**5** Conclusion

Introduction
○○○○

**Preliminaries**
●○○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

**1** Introduction

**2** Preliminaries

**3** Models and Algorithms Lack a Simplicity Bias

**4** Learnability with a Simplicity Bias

**5** Conclusion

## Supervised Learning

Objective: Learn $f : \mathcal{X} \to \mathcal{Y}$.

Given information: samples $(x_1, y_1), \ldots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, y_i = f(x_i)$

instances        labels

Introduction
oooo

**Preliminaries**
o●oooo

Models and Algorithms Lack a Simplicity Bias
ooooooooooooo

Learnability with a Simplicity Bias
ooooooooooo

Conclusion
ooo

## Supervised Learning

Objective: Learn $f : \mathcal{X} \to \mathcal{Y}$.

Given information: samples $(x_1, y_1), \ldots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, y_i = f(x_i)$

merely
minimize?

instances          labels

$\hat{R}(f') := \frac{1}{n} \sum_{i=1}^{n} \ell(f'(x_i), f(x_i)).$

Introduction
○○○○

**Preliminaries**
○●○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Supervised Learning

Objective: Learn $f : \mathcal{X} \to \mathcal{Y}$.

Given information: samples $(x_1, y_1), \ldots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, y_i = f(x_i)$ , $x_i \overset{i.i.d.}{\sim} P_{tr}$.

merely
minimize?

instances    labels

$\hat{R}(f') := \frac{1}{n} \sum_{i=1}^{n} \ell(f'(x_i), f(x_i)).$

fit $f$ inside
instance distribution

$R(f') := \mathbb{E}_{x \sim P_{tr}} \left[ \ell(f'(x), f(x)) \right].$

## Supervised Learning

Objective: Learn $f : \mathcal{X} \to \mathcal{Y}$.

Given information: samples $(x_1, y_1), \ldots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}, y_i = f(x_i)$ , $x_i \overset{i.i.d.}{\sim} P_{tr}$.

merely
minimize?

instances    labels

$\hat{R}(f') := \frac{1}{n} \sum_{i=1}^{n} \ell(f'(x_i), f(x_i)).$
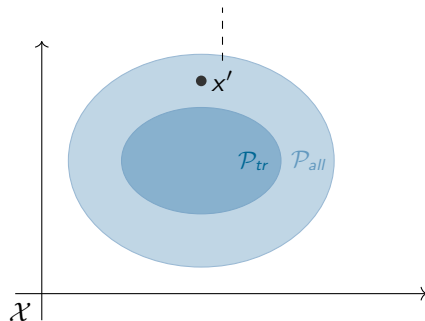
fit $f$ inside
instance distribution

$R(f') := \mathbb{E}_{x \sim P_{tr}} \left[ \ell(f'(x), f(x)) \right].$
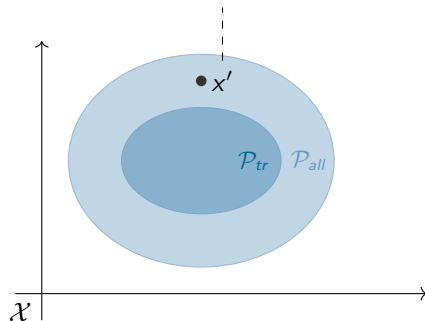
fit $f$ entirely

How does $f'$ generalize out-of-distribution?

## Domain Generalization

How shall we determine $f'(x')$?

Introduction
0000

**Preliminaries**
00●000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Domain Generalization

How shall we determine $f'(x')$? ← - - - - - - - -

Ahuja et al.: OOD generalization is impossible
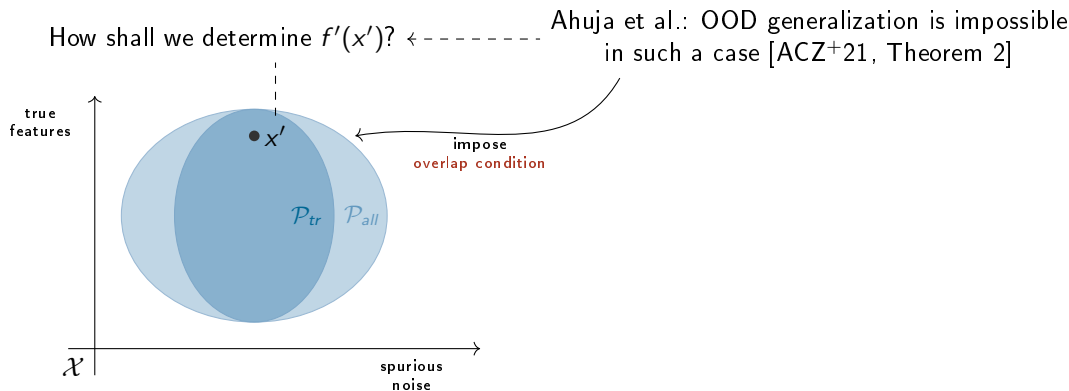in such a case [ACZ$^+$21, Theorem 2]

## Domain Generalization



How shall we determine $f'(x')$? ← - - - - - - - -    Ahuja et al.: OOD generalization is impossible
in such a case [ACZ+21, Theorem 2]

true
features

impose
overlap condition

$\mathcal{P}_{tr}$  $\mathcal{P}_{all}$

$\bullet\, x'$

$\mathcal{X}$

spurious
noise

## Domain Generalization

How shall we determine $f'(x')$? ← - - - - - - - -

Ahuja et al.: OOD generalization is impossible
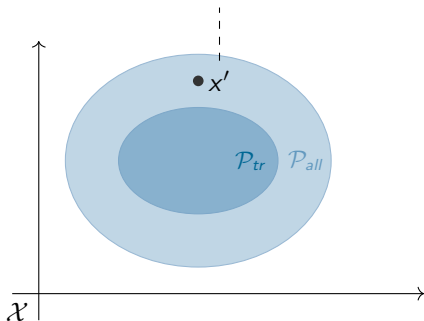in such a case [ACZ$^+$21, Theorem 2]



But aren't we posing too high
demands on "OOD generalization"?

## Inferring the simplest consistent function

Our maximum demand on OOD generalization should follow Ockham's razor:

*Infer the simplest function that remains consistent with the data.*

## Inferring the simplest consistent function

Our maximum demand on OOD generalization should follow Ockham's razor:

*Infer the simplest function that remains consistent with the data.*

How to define the simplicity of a function?

## Inferring the simplest consistent function

Our maximum demand on OOD generalization should follow Ockham's razor:

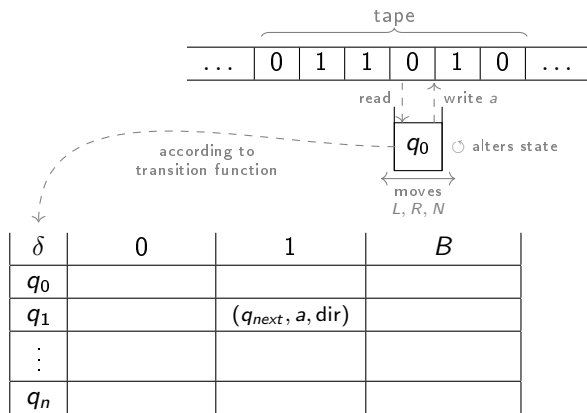*Infer the simplest function that remains consistent with the data.*

How to define the simplicity of a function?

Kolmogorov: the shortest description length of a program that produces this function.
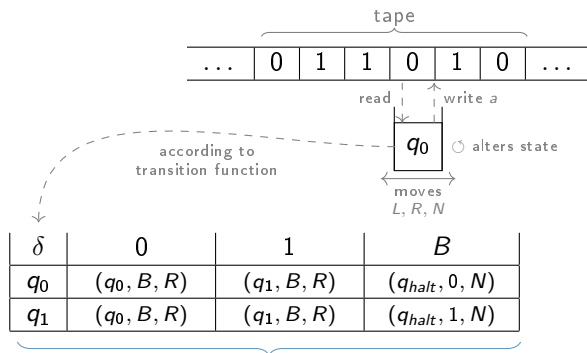
Turing Machine

## Turing Machines

A Turing Machine $\mathcal{T}$:



| $\delta$ | 0 | 1 | $B$ |
|---|---|---|---|
| $q_0$ | | | |
| $q_1$ | | $(q_{next}, a, \mathrm{dir})$ | |
| $\vdots$ | | | |
| $q_n$ | | | |

Introduction
○○○○

**Preliminaries**
○○○○●○

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Turing Machines

A Turing Machine $\mathcal{T}$:



| $\delta$ | 0 | 1 | $B$ |
|----------|---|---|-----|
| $q_0$ | $(q_0, B, R)$ | $(q_1, B, R)$ | $(q_{halt}, 0, N)$ |
| $q_1$ | $(q_0, B, R)$ | $(q_1, B, R)$ | $(q_{halt}, 1, N)$ |

**Example**: Modulo function $\text{mod}_2$

Introduction
0000

**Preliminaries**
000000●0

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Turing Machines

A Turing Machine $\mathcal{T}$ computes a *partial* computable function



$f : D_f \to \{0,1\}^*, D_f \subseteq \{0,1\}^*,$

$$f_{\mathcal{T}}(x) := \begin{cases} y, & \mathcal{T} \text{ halts and outputs } y, \\ \bot, & \mathcal{T} \text{ does not halt.} \end{cases}$$

$f$ is *total* computable if $D_f = \{0,1\}^*$.

| $\delta$ | 0 | 1 | $B$ |
|----------|---|---|-----|
| $q_0$ | | | |
| $q_1$ | | $(q_{next}, a, \mathrm{dir})$ | |
| $\vdots$ | | | |
| $q_n$ | | | |

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Turing Machines

A Turing Machine $\mathcal{T}$ computes a *partial* computable function   $f : D_f \to \{0,1\}^*, D_f \subseteq \{0,1\}^*$,



$$f_{\mathcal{T}}(x) := \begin{cases} y, & \mathcal{T} \text{ halts and outputs } y, \\ \bot, & \mathcal{T} \text{ does not halt.} \end{cases}$$
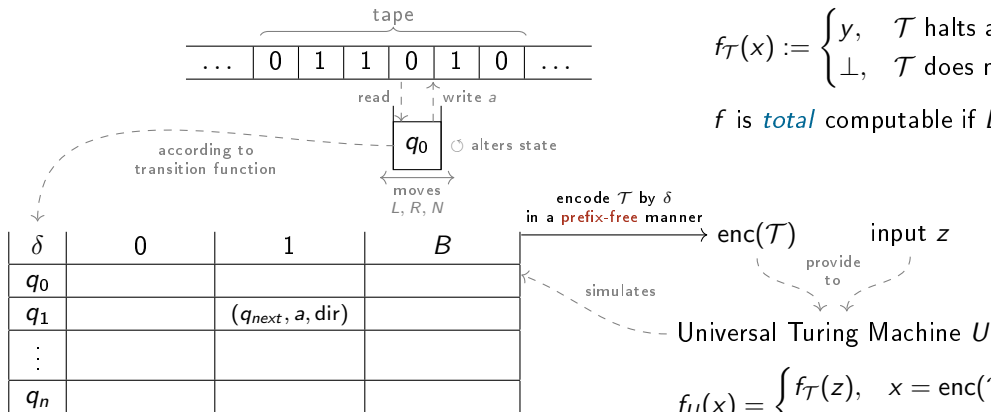
$f$ is *total* computable if $D_f = \{0,1\}^*$.

$$f_U(x) = \begin{cases} f_{\mathcal{T}}(z), & x = \mathrm{enc}(\mathcal{T})z \\ \bot, & \text{otherwise.} \end{cases}$$

## Kolmogorov Complexity

How many bits are needed to describe the encoding of a Turing Machine that computes $f$?

**Equivalence:** We write $U(p) \equiv f$ if

- $U(px) = f(x)$ for all $x \in D_f$, and

- $U$ does not halt on $px$ for all $x \in \{0,1\}^* \setminus D_f$.

## Kolmogorov Complexity

How many bits are needed to describe the encoding of a Turing Machine that computes $f$?

**Equivalence:** We write $U(p) \equiv f$ if

- $U(px) = f(x)$ for all $x \in D_f$, and

- $U$ does not halt on $px$ for all $x \in \{0,1\}^* \setminus D_f$.

We formalise this as the *Kolmogorov complexity* of a computable function $f$ (cf. [LV+08]):

$$K_U(f) := \min_{p \in \{0,1\}^*} \{ l(p) \mid U(p) \equiv f \}$$

Introduction
○○○○

**Preliminaries**
○○○○○●

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Kolmogorov Complexity

How many bits are needed to describe the encoding of a Turing Machine that computes $f$?

**Equivalence:** We write $U(p) \equiv f$ if

- $U(px) = f(x)$ for all $x \in D_f$, and

- $U$ does not halt on $px$ for all $x \in \{0, 1\}^* \setminus D_f$.

We formalise this as the *Kolmogorov complexity* of a computable function $f$ (cf. [LV$^+$08]):

$$K_U(f) := \min_{p \in \{0,1\}^*} \big\{ l(p) \mid U(p) \equiv f \big\}$$

Accordingly, the *conditional Kolmogorov complexity* given $z$ is defined as:

$$K_U(f \mid z) := \min_{p,p' \in \{0,1\}^*} \big\{ l(p) + l(p') \mid U(p[z]p') \equiv f \big\}$$

self-delimiting
encoding

**1** Introduction

**2** Preliminaries

**3** Models and Algorithms Lack a Simplicity Bias

**4** Learnability with a Simplicity Bias

**5** Conclusion

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0●0000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Our Limitations to Learn Simple Functions in Practice

### Models

have access to

function set
$$\tau = \{f_1, \ldots, f_j\}$$

### Simplicity

e.g. constants,
activation functions,
arithmetic operations,
logical expressions

### Algorithms

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0●00000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Our Limitations to Learn Simple Functions in Practice

### Models

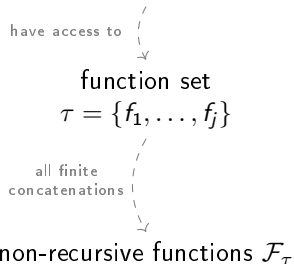have access to

function set
$$\tau = \{f_1, \ldots, f_j\}$$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

### Simplicity

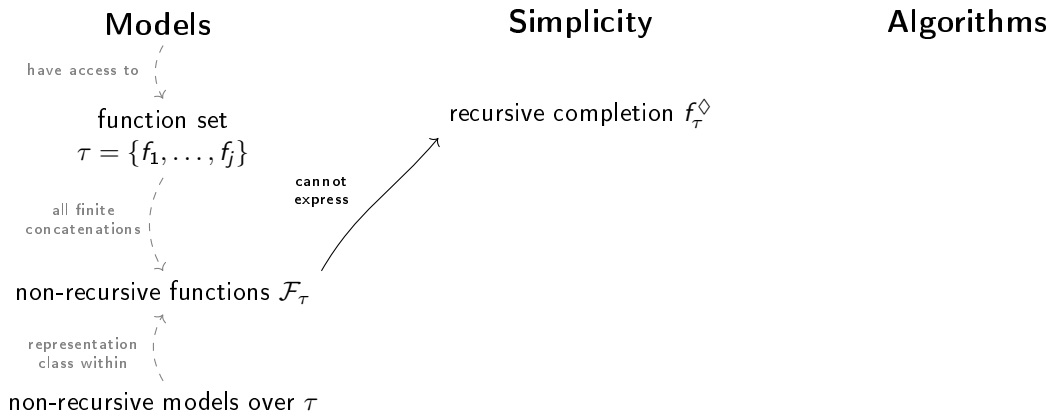### Algorithms

## Our Limitations to Learn Simple Functions in Practice

### Models

have access to
/
|
\
→

function set
$\tau = \{f_1, \ldots, f_j\}$
/
/
all finite    /
concatenations  |
\
\
→

non-recursive functions $\mathcal{F}_\tau$
→
/
representation  |
class within    \
\

non-recursive models over $\tau$

### Simplicity

### Algorithms

Introduction
○○○○

Preliminaries
○○○○○○

Models and Algorithms Lack a Simplicity Bias
○●○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Our Limitations to Learn Simple Functions in Practice

### Models

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

### Simplicity

recursive completion $f_\tau^\Diamond$

cannot
express

### Algorithms

## Our Limitations to Learn Simple Functions in Practice

### Models

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

cannot
express

### Simplicity

recursive completion $f_\tau^{\Diamond}$

constantly low
given $\tau$

Kolmogorov complexity

### Algorithms

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0●0000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Our Limitations to Learn Simple Functions in Practice



**Models**

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

**Simplicity**

recursive completion $f_\tau^\Diamond$

cannot
express

constantly low
given $\tau$

Kolmogorov complexity

**Algorithms**

ERM/IRM

for any finite dataset generated from $f_\tau^\Diamond$
infinitely many remain consistent and hence optimal

Introduction
oooo

Preliminaries
oooooo

Models and Algorithms Lack a Simplicity Bias
o●oooooooooo

Learnability with a Simplicity Bias
oooooooooo

Conclusion
ooo

## Our Limitations to Learn Simple Functions in Practice



**Models**

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

**Simplicity**

cannot
express

recursive completion $f_\tau^{\Diamond}$

constantly low
given $\tau$

Kolmogorov complexity

renders all consistent functions suboptimal
for sufficiently large datasets

for any finite dataset generated from $f_\tau^{\Diamond}$
infinitely many remain consistent and hence optimal

**Algorithms**

Simplicity bias
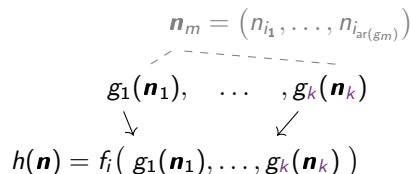+ ERM/IRM

ERM/IRM

## Inductive definition of non-recursive functions $\mathcal{F}_\tau$

**Given:** $\quad \tau := \{f_1, \ldots, f_j\}, f_i : \mathbb{N}^k \to \mathbb{N}$

**Base case:** $\quad$ Identity function $I \in \mathcal{F}_\tau$, $I(n_i) = n_i, n_i \in \mathbb{N}$

## Inductive definition of non-recursive functions $\mathcal{F}_\tau$

**Given:** $\tau := \{f_1, \ldots, f_j\}, f_i : \mathbb{N}^k \to \mathbb{N}$

**Base case:** Identity function $\mathrm{I} \in \mathcal{F}_\tau$, $\mathrm{I}(n_i) = n_i, n_i \in \mathbb{N}$

$$\boldsymbol{n}_m = \left(n_{i_1}, \ldots, n_{i_{\mathrm{ar}(g_m)}}\right)$$

**Induction step:** If $g_1, \ldots, g_k \in \mathcal{F}_\tau$, and $f_i \in \tau$ is $k$-ary,

$$g_1(\boldsymbol{n}_1), \quad \ldots \quad , g_k(\boldsymbol{n}_k)$$

then $h := f_i \circ (g_1, \ldots, g_k) \in \mathcal{F}_\tau$.

$$h(\boldsymbol{n}) = f_i\big( g_1(\boldsymbol{n}_1), \ldots, g_k(\boldsymbol{n}_k) \big)$$

Inductive definition of non-recursive functions $\mathcal{F}_\tau$

**Given:** $\tau := \{f_1, \ldots, f_j\}, f_i : \mathbb{N}^k \to \mathbb{N}$

**Base case:** Identity function $\mathrm{I} \in \mathcal{F}_\tau$, $\mathrm{I}(n_i) = n_i, n_i \in \mathbb{N}$

$$\boldsymbol{n}_m = \left(n_{i_1}, \ldots, n_{i_{\mathrm{ar}(g_m)}}\right)$$

**Induction step:** If $g_1, \ldots, g_k \in \mathcal{F}_\tau$, and $f_i \in \tau$ is $k$-ary,

$$g_1(\boldsymbol{n}_1), \quad \ldots \quad, g_k(\boldsymbol{n}_k)$$

then $h := f_i \circ (g_1, \ldots, g_k) \in \mathcal{F}_\tau$.

$$h(\boldsymbol{n}) = f_i\big(g_1(\boldsymbol{n}_1), \ldots, g_k(\boldsymbol{n}_k)\big)$$

**Example (Linear Functions):** $\tau = \{f_0, f_1, f_+\}$, $f_0(n_i) = 0$, $f_1(n_i) = 1$, $f_+(n_i, n_j) = n_i + n_j$.

Then, $\mathcal{F}_\tau = \{f : f(n) = an + b \mid a, b \in \mathbb{N}\}$.

## Roadmap



**Models**

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

**Simplicity**

recursive completion $f_\tau^\lozenge$

cannot
express

constantly low
given $\tau$

Kolmogorov complexity

renders all consistent functions suboptimal
for sufficiently large datasets

**Algorithms**

Simplicity bias
+ ERM/IRM

ERM/IRM

for any finite dataset generated from $f_\tau^\lozenge$
infinitely many remain consistent and hence optimal

Introduction
○○○○

Preliminaries
○○○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○●○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Recursive concatenation

Recursive concatenation of $k$-ary $f$:

$$f^{(0)}(n) := n, \qquad\qquad\qquad n \in \mathbb{N} \qquad (1)$$

$$f^{(m+1)}(n) := f\big(\underbrace{f^{(m)}(n), \ldots, f^{(m)}(n)}_{k \text{ times}}\big), \qquad n, m \in \mathbb{N}. \qquad (2)$$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000●00000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Recursive concatenation

Recursive concatenation of $k$-ary $f$:

$$f^{(0)}(n) := n, \qquad\qquad n \in \mathbb{N} \qquad\qquad (1)$$

$$f^{(m+1)}(n) := f\big(\underbrace{f^{(m)}(n), \ldots, f^{(m)}(n)}_{k \text{ times}}\big), \qquad n, m \in \mathbb{N}. \qquad (2)$$

**Example for** $2$-**ary** $f$:  $f^{(1)}(n)$ :

$f(\,\cdot\,,\,\cdot\,)$

$n \qquad\qquad n$

$f^{(2)}(n)$ :

$f(\,\cdot\,,\,\cdot\,)$

$f(\,\cdot\,,\,\cdot\,) \qquad f(\,\cdot\,,\,\cdot\,)$

$n \qquad n \quad n \qquad n$

## Recursive Completion

Function set $\tau = \{f_1, \ldots, f_j\}$

Recursive concatenation $f^{(m)}(n)$

## Recursive Completion

Function set $\tau = \{f_1, \ldots, f_j\}$

sum up

$$f_\tau := \sum_{i=1}^{j} f_i$$

Recursive concatenation $f^{(m)}(n)$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000●000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Recursive Completion

Function set $\tau = \{f_1, \ldots, f_j\}$

sum up

$f_\tau := \sum_{i=1}^{j} f_i$

Recursive concatenation $f^{(m)}(n)$

special case

**Recursive completion** $f^\diamond(n) := f^{(n)}(n)$

Introduction
○○○○

Preliminaries
○○○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○○●○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
○○○

## Recursive Completion

Function set $\tau = \{f_1, \ldots, f_j\}$

Recursive concatenation $f^{(m)}(n)$

sum up

special case

$f_\tau := \sum_{i=1}^{j} f_i$ — — — — — — — — — — — — — — Recursive completion $f^\Diamond(n) := f^{(n)}(n)$

combine

Recursive completion over $\tau$
$$f_\tau^\Diamond(n) := \sum_{i=1}^{j} f_i\big(f_\tau^{(n-1)}(n)\big)$$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000●00000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

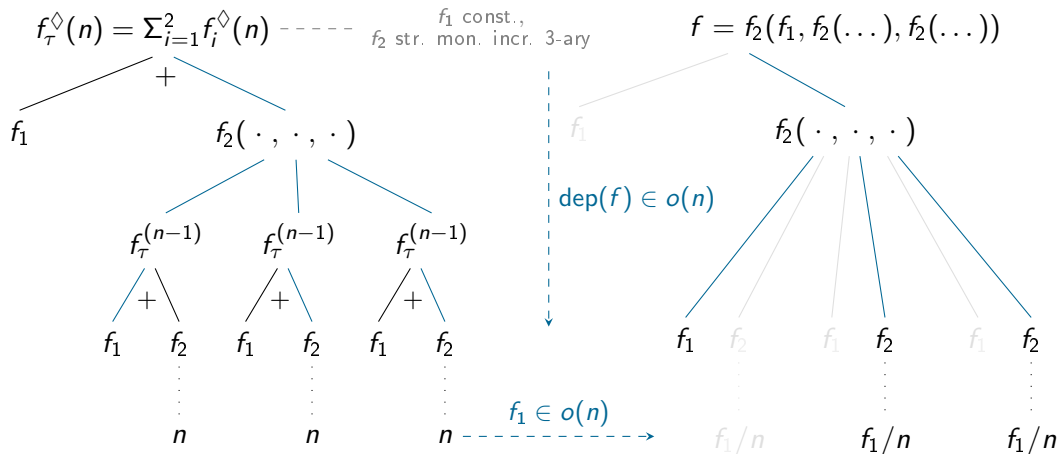## Why Recursive Completion Excels Every Non-Recursive Function

$$f_\tau^\diamond(n) = \Sigma_{i=1}^2 f_i^\diamond(n) \quad ----- \quad \begin{array}{l} f_1 \text{ const.,} \\ f_2 \text{ str. mon. incr. 3-ary} \end{array} \qquad f = f_2(f_1, f_2(\ldots), f_2(\ldots))$$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000●000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Why Recursive Completion Excels Every Non-Recursive Function

$$f_\tau^\diamond(n) = \Sigma_{i=1}^2 f_i^\diamond(n) \text{-----} \begin{array}{l} f_1 \text{ const.,} \\ f_2 \text{ str. mon. incr. 3-ary} \end{array}$$

$$f = f_2(f_1, f_2(\ldots), f_2(\ldots))$$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000●000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Why Recursive Completion Excels Every Non-Recursive Function

$$f_\tau^\diamond(n) = \Sigma_{i=1}^2 f_i^\diamond(n)$$

$f_1$ const.,
$f_2$ str. mon. incr. 3-ary

$$f = f_2(f_1, f_2(\dots), f_2(\dots))$$

$f_1$

$f_2(\cdot, \cdot, \cdot)$

$f_1$    $f_2$    $f_1$    $f_2$    $f_1$    $f_2$

$f_1/n$    $f_1/n$    $f_1/n$

## Why Recursive Completion Excels Every Non-Recursive Function

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000●00000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Why Recursive Completion Excels Every Non-Recursive Function



$$f_\tau^\diamond(n) = \Sigma_{i=1}^2 f_i^\diamond(n)$$

$f_1$ const.,
$f_2$ str. mon. incr. 3-ary

$$f = f_2(f_1, f_2(\dots), f_2(\dots))$$

$\mathrm{dep}(f) \in o(n)$

$f_1 \in o(n)$

## Non-recursion functions do not capture recursive completion

Fix an arbitrary function set $\tau = \{f_1, \ldots, f_j\}, f_i : \mathbb{N}^k \to \mathbb{N}$, where

- each $f_i$ is strictly monotonously increasing or bounded,
- some strictly monotonously increasing $f_i$ has arity $\mathrm{ar}(f_i) > 1$ (e.g. $f_+$).

For any non-recursive function $f \in \mathcal{F}_\tau$, there is an $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$, $f_\tau^\Diamond(n) > f(n)$.

## Non-recursion functions do not capture recursive completion

Fix an arbitrary function set $\tau = \{f_1, \ldots, f_j\}$, $f_i : \mathbb{N}^k \to \mathbb{N}$, where

- each $f_i$ is strictly monotonously increasing or bounded,
- some strictly monotonously increasing $f_i$ has arity $\mathrm{ar}(f_i) > 1$ (e.g. $f_+$).

For any non-recursive function $f \in \mathcal{F}_\tau$, there is an $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$, $f_\tau^\Diamond(n) > f(n)$.



This result can be extended to $f_i : \mathbb{Z}^k \to \mathbb{Z}$ with realistic assumptions.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000●000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Roadmap



**Models**

have access to

function set
$\tau = \{f_1, \ldots, f_j\}$

cannot
express

all finite
concatenations

non-recursive functions $\mathcal{F}_\tau$

representation
class within

non-recursive models over $\tau$

**Simplicity**

recursive completion $f_\tau^\Diamond$

constantly low
given $\tau$

Kolmogorov complexity

renders all consistent functions suboptimal
for sufficiently large datasets

**Algorithms**

Simplicity bias
+ ERM/IRM

ERM/IRM

for any finite dataset generated from $f_\tau^\Diamond$
infinitely many remain consistent and hence optimal

Uniform Simplicity of Recursive Completion

Uniform TM $\mathcal{T}_\Diamond$:

input tape

| ... | $enc(\mathcal{T}_1)$ | $ar(f_1)$ | ... | $enc(\mathcal{T}_j)$ | $ar(f_j)$ | $n$ | ... |

loop tape

| ... | | ... |

computation tape

| ... | | ... |

accumulation tape

| ... | | ... |

output tape

| ... | | ... |

## Uniform Simplicity of Recursive Completion

Uniform TM $\mathcal{T}_\Diamond$:

| | input tape | | | | | | |
|---|---|---|---|---|---|---|---|
| ... | $enc(\mathcal{T}_1)$ | $ar(f_1)$ | ... | $enc(\mathcal{T}_j)$ | $ar(f_j)$ | $n$ | ... |

| loop tape | | |
|---|---|---|
| ... | $k$ | ... |

0. Initialize
$k := n$
$x := n$

| computation tape | | |
|---|---|---|
| ... | | ... |

| output tape | | |
|---|---|---|
| ... | $x$ | ... |

| accumulation tape | | |
|---|---|---|
| ... | | ... |

## Uniform Simplicity of Recursive Completion

Uniform TM $\mathcal{T}_\Diamond$:

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0000000000●00

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Uniform Simplicity of Recursive Completion

Uniform TM $\mathcal{T}_\Diamond$:

input tape

| ... | $enc(\mathcal{T}_1)$ | $ar(f_1)$ | ... | $enc(\mathcal{T}_j)$ | $ar(f_j)$ | $n$ | ... |

fetch description
and arity
of each $\mathcal{T}_i$

0. Initialize
$k := n$
$x := n$

loop tape

| ... | $k$ | ... |

$\circlearrowleft$ while $k > 0$

$k$-=1

computation tape

| ... | simulate $\mathcal{T}_i$ on $x$ | ... |

1. Provide $x$ as
input to each simulation

output tape

| ... | $x$ | ... |

2. Add outputs to

accumulation tape

| ... | $\sum_{i=1}^{j} \mathcal{T}_j(x)$ | ... |

Introduction
oooo

Preliminaries
oooooo

Models and Algorithms Lack a Simplicity Bias
ooooooooooo●oo

Learnability with a Simplicity Bias
oooooooooo

Conclusion
ooo

## Uniform Simplicity of Recursive Completion

Uniform TM $\mathcal{T}_\Diamond$:

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000000●0

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Roadmap

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0000000000●

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^\Diamond(n) > f'(n)$ for all $n \geq n_0$.
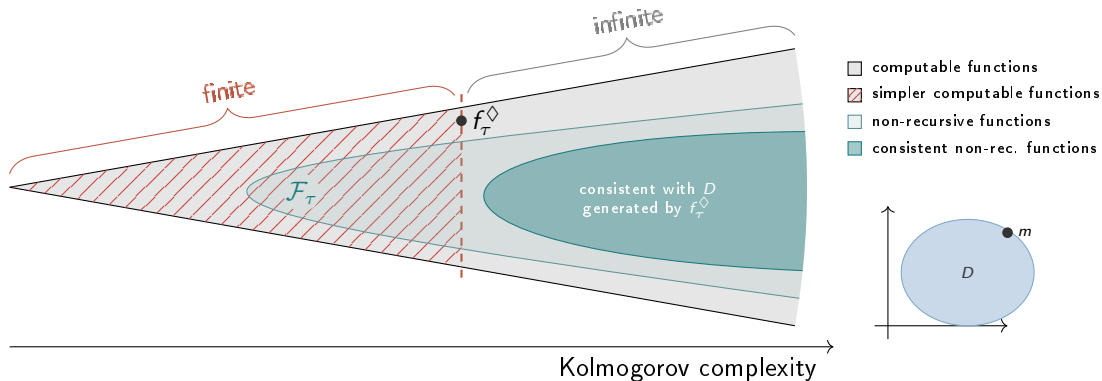
$\square$ computable functions

$\square$ non-recursive functions

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000000●

Learnability with a Simplicity Bias
0000000000
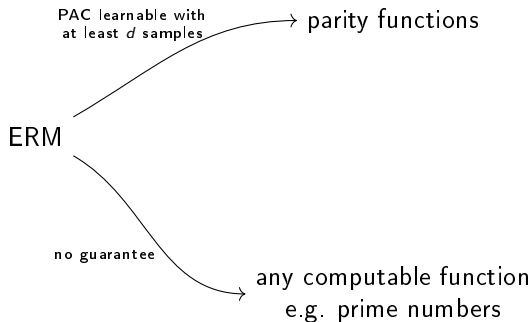
Conclusion
000

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^\diamond(n) > f'(n)$ for all $n \geq n_0$.
But non-recursive functions can still memorize the training data.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000000●

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^\Diamond(n) > f'(n)$ for all $n \geq n_0$.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0000000000●

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^\Diamond(n) > f'(n)$ for all $n \geq n_0$.
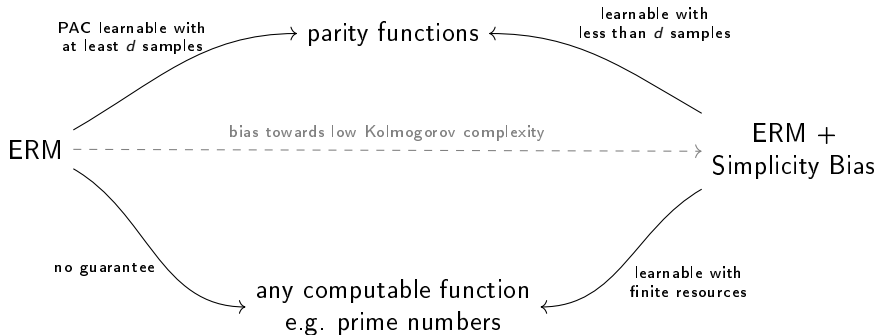→ Include large enough sample $(m, f_\tau^\Diamond(m))$ into dataset $D$.

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^\Diamond(n) > f'(n)$ for all $n \geq n_0$.
→ Include large enough sample $(m, f_\tau^\Diamond(m))$ into dataset $D$.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000000●

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Eliminating non-recursive functions with a simplicity bias

For any $f' \in \mathcal{F}_\tau$, there exists an $n_0$ such that $f_\tau^{\Diamond}(n) > f'(n)$ for all $n \geq n_0$.
→ Include large enough sample $(m, f_\tau^{\Diamond}(m))$ into dataset $D$.



Kolmogorov complexity

□ computable functions
▨ simpler computable functions
□ non-recursive functions
▩ consistent non-rec. functions

1. Introduction

2. Preliminaries

3. Models and Algorithms Lack a Simplicity Bias

4. Learnability with a Simplicity Bias

5. Conclusion

## Teaser



PAC learnable with
at least $d$ samples

parity functions

ERM

no guarantee

any computable function
e.g. prime numbers

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
00000000000

Learnability with a Simplicity Bias
0●00000000

Conclusion
000

Teaser



PAC learnable with
at least $d$ samples

parity functions

learnable with
less than $d$ samples

ERM

bias towards low Kolmogorov complexity

ERM +
Simplicity Bias

no guarantee

any computable function
e.g. prime numbers

learnable with
finite resources
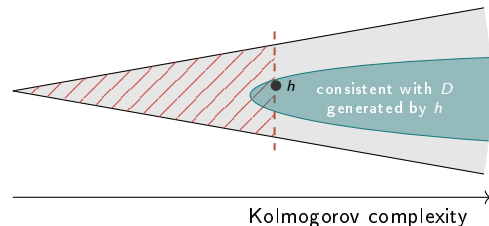
## PAC learning computable functions

### PAC learnable

$\mathcal{H}$ is PAC learnable if there is

- a learning algorithm $A$ *(ERM)*

- a sample number threshold $n_0(\varepsilon, \delta)$

such that for

- any error and failure probabilities $\varepsilon, \delta \in [0, 1)$,

- any hypothesis $h \in \mathcal{H}$ ,        $\Big\}$ hypothesis conditions

- any marginal distribution $P : \mathcal{X} \to [0, 1]$, and

- any dataset $D = \big\{ (X_i, h(X_i)) \mid i = 1, \ldots, n \big\}, X_i \overset{i.i.d}{\sim} P$, with $n \geq n_0(\varepsilon, \delta)$,

         $\Big\}$ data conditions

$$\Pr\big[ R(A(D)) \leq \varepsilon \big] \geq 1 - \delta.$$

## PAC learning computable functions

**PAC learnable** $\xleftarrow{\text{VC dimension}}$ $\xleftarrow{\text{Rademacher complexity}}$ ~~$\mathcal{H}$ finite?~~ ~~$\mathcal{H}$ bounded?~~

$\mathcal{H}$ is PAC learnable if there is

- a learning algorithm $A$ *(ERM)*

- a sample number threshold $n_0(\varepsilon, \delta)$

such that for

- any error and failure probabilities $\varepsilon, \delta \in [0, 1)$,

- any hypothesis $h \in \mathcal{H}$ , $\left.\right\}$ hypothesis conditions

- any marginal distribution $P : \mathcal{X} \to [0, 1]$, and

- any dataset $D = \big\{(X_i, h(X_i)) \mid i = 1, \ldots, n\big\}, X_i \overset{i.i.d}{\sim} P$, with $n \geq n_0(\varepsilon, \delta)$,

$\left.\right\}$ data conditions

$$\Pr\big[R(A(D)) \leq \varepsilon\big] \geq 1 - \delta.$$

$\bullet\, h$

consistent with $D$
generated by $h$

partial computable
functions $\mathcal{H}$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
00●0000000

Conclusion
000

## PAC learning computable functions

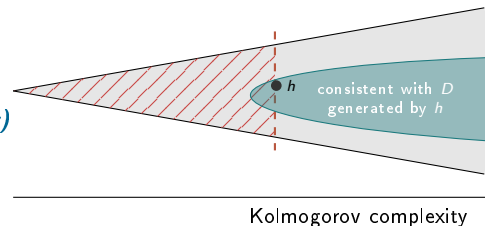### PAC learnable

$\mathcal{H}$ is PAC learnable if there is

- a learning algorithm $A$ *(ERM + Simplicity Bias)*

- a sample number threshold $n_0(\delta, k)$

such that for

- any failure probability $\delta \in (0, 1)$,

- any hypothesis $h \in \mathcal{H}$ **with Kolmogorov complexity $k = K(h)$**, $\Big\}$ hypothesis conditions

- any marginal distribution $P : \mathcal{X} \to [0, 1]$, and

- any dataset $D = \big\{(X_i, h(X_i)) \mid i = 1, \ldots, n\big\}, X_i \overset{i.i.d}{\sim} P$, with $n \geq n_0(\delta, k)$,

$$\Pr\big[\ \underbrace{A(D) = h}\ \big] \geq 1 - \delta.$$
$$\text{\footnotesize perfect learnability}$$



consistent with $D$
generated by $h$

Kolmogorov complexity

## Why conditioning learnability on the sample size is impossible in general

▷ **Example 1:** Unbounded Kolmogorov Complexity, but one sample suffices

▷ **Example 2:** Low Kolmogorov Complexity, but infinite dataset insufficient

## Why conditioning learnability on the sample size is impossible in general

▽ **Example 1:** Unbounded Kolmogorov Complexity, but one sample suffices

$D_y := \{(0, y)\}, y \in \{0, 1\}^*$.
For each $y$, there is a different simplest consistent function.
But any infinite function class is unbounded in terms of Kolmogorov complexity.

▷ **Example 2:** Low Kolmogorov Complexity, but infinite dataset insufficient

## Why conditioning learnability on the sample size is impossible in general

▷ **Example 1:** Unbounded Kolmogorov Complexity, but one sample suffices

▽ **Example 2:** Low Kolmogorov Complexity, but infinite dataset insufficient

We want to learn the modulo function $\mod_2(x) = x \mod 2$.
But the infinite dataset $D = \{(2n, 0) \mid n \in \mathbb{N}\}$ leaves the (simpler) constant function $f_0(x) = 0$ consistent.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Alternative conditions on the data

**Learnability**

$\mathcal{H}$ is learnable if there is

- a learning algorithm $A$ *(ERM + Simplicity Bias)*

- ~~a sample number threshold $n_0(\delta, k)$~~

such that for

- any failure probability $\delta \in (0, 1)$,

- any hypothesis $h \in \mathcal{H}$ with Kolmogorov complexity $k = K(h)$, $\quad \}$ hypothesis conditions



consistent with $D$ generated by $h$

Kolmogorov complexity

$\Big\}$ data conditions

**What conditions do $D$ and $P$ need to fulfil?**

$$\Pr\big[\; \underbrace{\boldsymbol{A(D) = h}}_{\substack{h \text{ is the simplest} \\ \text{consistent function}}} \;\big] \geq 1 - \delta.$$

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000●0000

Conclusion
000

## Conditioning learnability on functional information

Define the functional information in $D$ as
$$K_F(D) := \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}.$$

## Conditioning learnability on functional information

Define the functional information in $D$ as
$$K_F(D) := \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}.$$

This quantifies the information
that datasets convey about the functions
that could have generated them.

Look at the prior examples anew.



consistent with $D$
generated by $h$

$K_F(D)$    $K(h)$

Kolmogorov complexity

| | True function | Dataset | Sample Size | $K_F(D)$ |
|------|------|------|------|------|
| Ex. 1 | $f_y$ | $D_y$ | 1 | $K_F(D_y) = K(f_y)$ |
| Ex. 2 | $\mathrm{mod}_2$ | $D_0$ | $\infty$ | $K_F(D_0) \le K(f_0) < K(\mathrm{mod}_2)$ |

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
000000●0000

Conclusion
000

## Teaching prime numbers by enumerating them

Consider the *prime number decision function* $\mathbb{1}_{\mathbb{P}}(n) = \mathbb{1}\{n \in \mathbb{P}\}$.
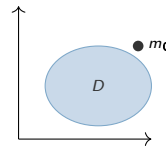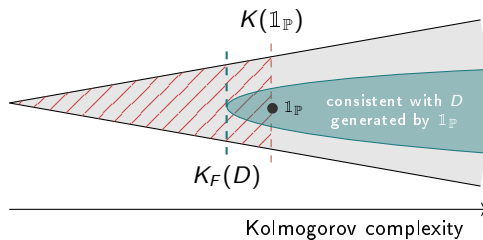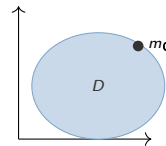
There exists an $m_0$ such that any dataset $D$ that contains $(n, \mathbb{1}_{\mathbb{P}}(n))$ for all $n \leq m_0$ renders $\mathbb{1}_{\mathbb{P}}$ the simplest consistent function with $D$ among all computable decision functions over $\mathbb{N}$.

## Teaching prime numbers by enumerating them

Consider the *prime number decision function* $\mathbb{1}_{\mathbb{P}}(n) = \mathbb{1}\{n \in \mathbb{P}\}$.
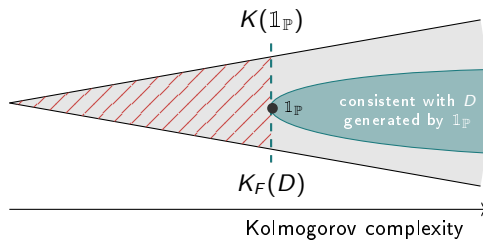
There exists an $m_0$ such that any dataset $D$ that contains $(n, \mathbb{1}_{\mathbb{P}}(n))$ for all $n \leq m_0$ renders $\mathbb{1}_{\mathbb{P}}$ the simplest consistent function with $D$ among all computable decision functions over $\mathbb{N}$.

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000●000

Conclusion
000

## Teaching prime numbers by enumerating them

Consider the *prime number decision function* $\mathbb{1}_\mathbb{P}(n) = \mathbb{1}\{n \in \mathbb{P}\}$.

There exists an $m_0$ such that any dataset $D$ that contains $(n, \mathbb{1}_\mathbb{P}(n))$ for all $n \leq m_0$ renders $\mathbb{1}_\mathbb{P}$ the simplest consistent function with $D$ among all computable decision functions over $\mathbb{N}$.
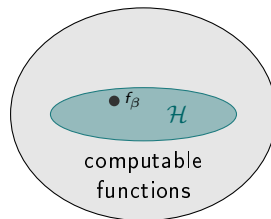
## Learning parity functions with less samples

Let $\mathcal{H} = \{f_\beta : \{0,1\}^d \to \{0,1\}, f(x) = \langle \beta, x \rangle \bmod 2 \mid \beta \in \{0,1\}^d \}$ be the class of parity functions over $d$-dimensional binary inputs.

Let $P = \mathrm{Ber}\left(\frac{1}{2}\right)^{\otimes d}$ be the uniform distribution over strings in $\{0,1\}^d$.

$$\Pr_{x \sim P}\left[f'_\beta(x) = f_\beta(x)\right] = \frac{1}{2}.$$



computable
functions

Introduction
oooo

Preliminaries
oooooo

Models and Algorithms Lack a Simplicity Bias
oooooooooooo

Learnability with a Simplicity Bias
oooooooo●oo

Conclusion
ooo

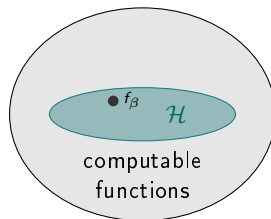## Learning parity functions with less samples

Let $\mathcal{H} = \left\{ f_\beta : \{0,1\}^d \to \{0,1\}, f(x) = \langle \beta, x \rangle \bmod 2 \mid \beta \in \{0,1\}^d \right\}$ be the class of parity functions over $d$-dimensional binary inputs.

Let $P = \mathrm{Ber}\left(\frac{1}{2}\right)^{\otimes d}$ be the uniform distribution over strings in $\{0,1\}^d$.

$$\Pr_{x \sim P}\left[ f'_\beta(x) = f_\beta(x) \right] = \frac{1}{2}. \xrightarrow{\hspace{3cm}}$$

At least $d$ samples necessary to render $f_\beta$ the only consistent function.
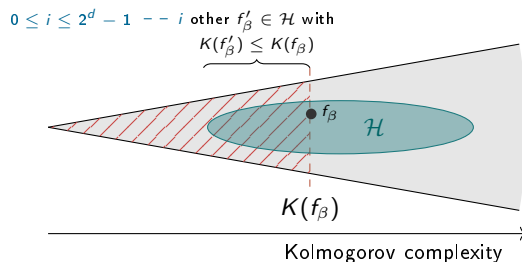


computable
functions

## Learning parity functions with less samples

Let $\mathcal{H} = \left\{ f_\beta : \{0,1\}^d \to \{0,1\}, f(x) = \langle \beta, x \rangle \bmod 2 \mid \beta \in \{0,1\}^d \right\}$ be the class of parity functions over $d$-dimensional binary inputs.

Let $P = \mathrm{Ber}\left(\frac{1}{2}\right)^{\otimes d}$ be the uniform distribution over strings in $\{0,1\}^d$.

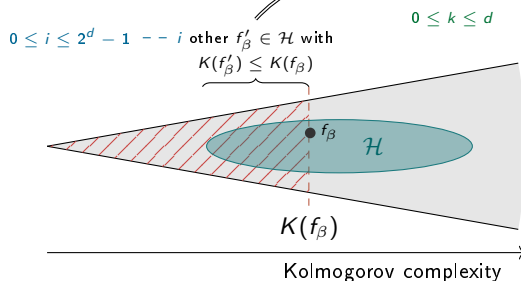$\Pr_{x \sim P}\left[ f'_\beta(x) = f_\beta(x) \right] = \frac{1}{2}$.

## Learning parity functions with less samples

Let $\mathcal{H} = \left\{ f_\beta : \{0,1\}^d \to \{0,1\}, f(x) = \langle \beta, x \rangle \mod 2 \mid \beta \in \{0,1\}^d \right\}$ be the class of parity functions over $d$-dimensional binary inputs.

Let $P = \text{Ber}\left(\frac{1}{2}\right)^{\otimes d}$ be the uniform distribution over strings in $\{0,1\}^d$.

$$\Pr_{x \sim P}\left[ f_\beta'(x) = f_\beta(x) \right] = \frac{1}{2}.$$

$k$ samples render $f_\beta$ the *simplest* consistent function with probability at least $1 - i \cdot \left(\frac{1}{2^k}\right)$.

$0 \leq k \leq d$

$0 \leq i \leq 2^d - 1$ — $i$ other $f_\beta' \in \mathcal{H}$ with $K(f_\beta') \leq K(f_\beta)$



$K(f_\beta)$

Kolmogorov complexity

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}$$

concatenate
samples

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound |  |  |
| ▷ Inconsistency implies lower bound |  |  |
| ▷ Monotonicity for supersets |  |  |
| ▷ Invariance under sample permutation |  |  |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}$$

concatenate samples

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| $\nabla$ Consistency implies upper bound | If $f$ is consistent with $D$, then $K_F(D) \leq K(f)$. | |
| $\triangleright$ Inconsistency implies lower bound | | |
| $\triangleright$ Monotonicity for supersets | | |
| $\triangleright$ Invariance under sample permutation | | |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \dashrightarrow \overset{\text{concatenate}}{\underset{\downarrow}{\text{samples}}}$$

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

| | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound | ✓ | |
| ▽ Inconsistency implies lower bound | If any $f$ with $K(f) < k$ is inconsistent with $D$, then $K_F(D) \geq k$. | |
| ▷ Monotonicity for supersets | | |
| ▷ Invariance under sample permutation | | |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \quad \text{--- concatenate samples}$$

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

| | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound | ✓ | |
| ▷ Inconsistency implies lower bound | ✓ | |
| ▽ Monotonicity for supersets | Any $D' \supset D$ adds constraints, hence $K_F(D) \leq K_F(D')$. | |
| ▷ Invariance under sample permutation | | |

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
000000000000

Learnability with a Simplicity Bias
0000000000

Conclusion
000

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \dashrightarrow \underset{\text{samples}}{\overset{\text{concatenate}}{\searrow}}$$

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound | ✓ | |
| ▷ Inconsistency implies lower bound | ✓ | |
| ▷ Monotonicity for supersets | ✓ | |
| ▽ Invariance under sample permutation | Constraints are unordered. | |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}$$

concatenate
samples

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|:---:|:---:|
| ▷ Consistency implies upper bound | ✓ | |
| ▷ Inconsistency implies lower bound | ✓ | |
| ▷ Monotonicity for supersets | ✓ | |
| ▷ Invariance under sample permutation | ✓ | |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \quad \text{---} \quad \begin{array}{c} \text{concatenate} \\ \text{samples} \end{array}$$

$$K_{JF}(D) = K([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n])$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|:---:|:---:|
| ▽ Consistency implies upper bound | ✓ | If $f$ is consistent with $D$, then $K_F(D) \leq K(f) + c$. |
| ▷ Inconsistency implies lower bound | ✓ | |
| ▷ Monotonicity for supersets | ✓ | |
| ▷ Invariance under sample permutation | ✓ | |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \dashrightarrow \substack{\text{concatenate} \\ \text{samples}}$$

$$K_{JF}(D) = K([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n])$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound | ✓ | $K_F(D) \leq K(f) + c$ |
| ▽ Inconsistency implies lower bound | ✓ | Notwithstanding $f(x_i) \neq y_i$, |
|  |  | potentially $f([x_1, \ldots, x_n]) = [y_1, \ldots, y_n]$. |
| ▷ Monotonicity for supersets | ✓ |  |
| ▷ Invariance under sample permutation | ✓ |  |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \dashrightarrow \begin{smallmatrix} \text{concatenate} \\ \text{samples} \end{smallmatrix}$$

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|---|---|
| ▷ Consistency implies upper bound | ✓ | $K_F(D) \leq K(f) + c$ |
| ▷ Inconsistency implies lower bound | ✓ | ✗ |
| ▽ Monotonicity for supersets | ✓ | add label as another instance $D = \{(x,y)\} \quad D' = \{(x,y),(y,0)\}$ |
| ▷ Invariance under sample permutation | ✓ |  |

Introduction
Preliminaries
Models and Algorithms Lack a Simplicity Bias
**Learnability with a Simplicity Bias**
Conclusion

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\}$$

concatenate samples

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

| | $K_F(D)$ | $K_{JF}(D)$ |
|---|:---:|:---:|
| ▷ Consistency implies upper bound | ✓ | $K_F(D) \leq K(f) + c$ |
| ▷ Inconsistency implies lower bound | ✓ | ✗ |
| ▷ Monotonicity for supersets | ✓ | ✗ |
| ▽ Invariance under sample permutation | ✓ | $\underbrace{01101\ldots111011}_{\text{incompressible}} \xrightarrow{\pi} \underbrace{00\ldots011\ldots1}_{\text{compressible}}$ |

## Weakening functional information spoils desirable properties

Could we weaken the constraint-based formulation of $K_F(D)$?

$$K_F(D) = \min_{p \in \{0,1\}^*} \{l(p) \mid U(px_i) = y_i \text{ for all } (x_i, y_i) \in D\} \dashrightarrow \substack{\text{concatenate} \\ \text{samples}}$$

$$K_{JF}(D) = K\left([y_1, \ldots, y_n] \mid [x_1, \ldots, x_n]\right)$$

|  | $K_F(D)$ | $K_{JF}(D)$ |
|---|:---:|:---:|
| ▷ Consistency implies upper bound | ✓ | $K_F(D) \leq K(f) + c$ |
| ▷ Inconsistency implies lower bound | ✓ | ✗ |
| ▷ Monotonicity for supersets | ✓ | ✗ |
| ▷ Invariance under sample permutation | ✓ | ✗ |

Introduction
0000

Preliminaries
000000

Models and Algorithms Lack a Simplicity Bias
0000000000000

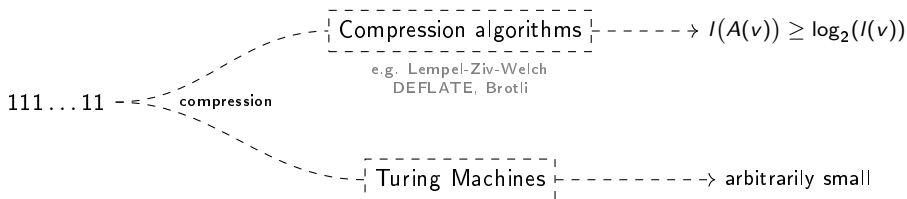Learnability with a Simplicity Bias
000000000●

Conclusion
000

## Compression algorithms cannot approximate Kolmogorov complexity

Kolmogorov Complexity is incomputable. But is there at least a viable approximation $A$ that satisfies

$$A(v) \geq \exp_2^{(k)}(a \cdot A(w) + b) \Rightarrow K(v) \geq K(w) \qquad \text{for some } a, b, k?$$

## Compression algorithms cannot approximate Kolmogorov complexity

Kolmogorov Complexity is incomputable. But is there at least a viable approximation $A$ that satisfies

$$A(v) \geq \exp_2^{(k)}(a \cdot A(w) + b) \Rightarrow K(v) \geq K(w) \qquad \text{for some } a, b, k?$$

Compression algorithms were employed in practice [LV+08, p. 696].
But their compression ratio is limited.



$111\ldots11$ --- compression

Compression algorithms ----→ $I(A(v)) \geq \log_2(I(v))$

e.g. Lempel-Ziv-Welch
DEFLATE, Brotli

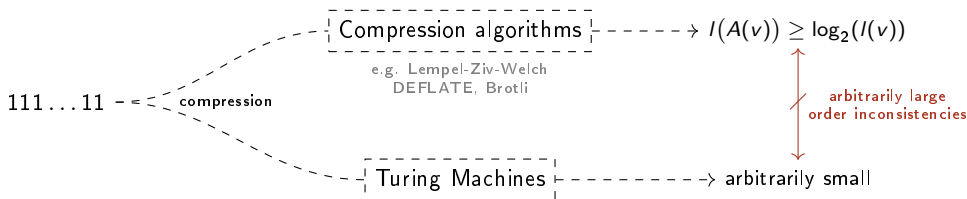Turing Machines --------→ arbitrarily small

## Compression algorithms cannot approximate Kolmogorov complexity

Kolmogorov Complexity is incomputable. But is there at least a viable approximation $A$ that satisfies

$$A(v) \geq \exp_2^{(k)}(a \cdot A(w) + b) \Rightarrow K(v) \geq K(w) \qquad \text{for some } a, b, k?$$

Compression algorithms were employed in practice [LV+08, p. 696].
But their compression ratio is limited.



$111 \ldots 11$  — compression

Compression algorithms
e.g. Lempel-Ziv-Welch
DEFLATE, Brotli

$I(A(v)) \geq \log_2(I(v))$

arbitrarily large
order inconsistencies

Turing Machines — arbitrarily small

Introduction
○○○○

Preliminaries
○○○○○○

Models and Algorithms Lack a Simplicity Bias
○○○○○○○○○○○○

Learnability with a Simplicity Bias
○○○○○○○○○○

Conclusion
●○○

1. Introduction

2. Preliminaries

3. Models and Algorithms Lack a Simplicity Bias

4. Learnability with a Simplicity Bias

5. Conclusion

## Key takeaways and future research

Key takeaways:

- Recursion is a powerful yet simple mechanism that feed-forward neural networks alone cannot express.

## Key takeaways and future research

Key takeaways:

- Recursion is a powerful yet simple mechanism that feed-forward neural networks alone cannot express.
- Out-of-distribution generalization guarantees could draw upon a simplicity bias.

## Key takeaways and future research

Key takeaways:

- Recursion is a powerful yet simple mechanism that feed-forward neural networks alone cannot express.

- Out-of-distribution generalization guarantees could draw upon a simplicity bias.

- Yet compression algorithms can not yield approximate guarantees on Kolmogorov complexity.

## Key takeaways and future research

Key takeaways:

- Recursion is a powerful yet simple mechanism that feed-forward neural networks alone cannot express.
- Out-of-distribution generalization guarantees could draw upon a simplicity bias.
- Yet compression algorithms can not yield approximate guarantees on Kolmogorov complexity.

Future research:

- Bestow learning algorithms with viable simplicity heuristics.

## Key takeaways and future research

Key takeaways:

- Recursion is a powerful yet simple mechanism that feed-forward neural networks alone cannot express.
- Out-of-distribution generalization guarantees could draw upon a simplicity bias.
- Yet compression algorithms can not yield approximate guarantees on Kolmogorov complexity.

Future research:

- Bestow learning algorithms with viable simplicity heuristics.
- How to efficiently learn recursive algorithms over discrete inputs?

## References I

[ACZ+21] Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-distribution generalization. Advances in Neural Information Processing Systems, 34:3438–3450, 2021.

[LV+08] Ming Li, Paul Vitányi, et al. An introduction to Kolmogorov complexity and its applications, volume 3. Springer, 2008.