

S&DS 265 / 565

# Introductory Machine Learning

## Language Models

October 24

Yale

# Welcome back!

For today:

- Where have we been? Where are we going?
- Language models

## But first...

- Assignment 2 scores, midterm scores and mid-semester grades posted by Thursday
- Please email me (copy Joanne) with any concerns about grading, standing...
- Assignment 3 due this Thurs, October 26
- Assignment 4 out at same time (language models, word embeddings)

# Invitation: Panel discussion

Class on Tuesday, December 5:

We will have a panel discussion on

*Societal issues for AI and Machine Learning*

If you are interested in participating, please email me

**Subject: iML panel**

# Where we've been

Week	Dates	Topics	Demos & Tutorials	Lecture Slides	Readings and Notes	Assignments & Exams
1	Aug 31	Course overview		Thu: Course overview		
2	Sept 5, 7	Python and background concepts	<a href="#">Python elements</a> <a href="#">Covid trends</a>	Tue: Python elements Thu: Pandas and linear regression	Data8 Chapters 3, 4, 5	Quiz 1 <a href="#">Assn 1 out</a>
3	Sept 12, 14	Linear regression and classification	<a href="#">Covid trends (revisited)</a> <a href="#">Classification examples</a>	Tue: Regression concepts Thu: Classification	ISL Sections 3.1, 3.2, 3.5 Notes on regression ISL Sections 4.3, 4.4 Notes on classification	
4	Sept 19, 21	Stochastic gradient descent	<a href="#">SGD examples</a>	Tue: Classification (continued) Thu: Stochastic gradient descent	ISL Section 6.2.2 ISL Section 10.7.2	Assn 1 in <a href="#">Assn 2 out</a>
5	Sept 26, 28	Bias and variance, cross-validation	<a href="#">Bias-variance tradeoff</a> <a href="#">Covid trends (revisited)</a> <a href="#">California housing</a>	Tue: Bias and variance Thu: Cross-validation	ISL Section 2.2 ISL Section 5.1	Quiz 2
6	Oct 3, 5	Tree-based methods and principal components	<a href="#">Trees and forests</a> <a href="#">Visualizing trees</a> <a href="#">PCA examples</a>	Tue: Trees and Forests Thu: PCA	ISL Sections 8.1, 8.2 ISL Section 12.2	Assn 2 in <a href="#">Assn 3 out</a>
7	Oct 10, 12	PCA and dimension reduction	<a href="#">PCA revisited</a> <a href="#">Used for dimension reduction</a> <a href="#">Word embeddings</a>	Tue: PCA and word embeddings Thu: Embeddings and review	ISL Section 12.2	Quiz 3
8	Oct 17	Midterm exam (in class)			On Canvas: <a href="#">Practice midterms / Sample solns</a> <a href="#">Midterm / Sample soln</a>	

# Where we're going

			Topics			
8	Oct 17	Midterm exam (in class)		On Canvas: Practice midterms / Sample solns Midterm / Sample soln		
9	Oct 24, 26	Language models, topic models	GPT-3 demo Bayesian inference Topic models	Tue: Language models Thu: Topic models	OpenAI: Better language models Notes on Bayesian inference	Asn 3 in Asn 4 out
10	Oct 31, Nov 2	Introduction to neural networks	Sanity check Minimal neural network Regression examples	Tue: Neural networks Thu: Neural networks	ISL Sections 10.1, 10.2	Quiz 4
11	Nov 7, 9	Reinforcement learning	Q-learning	Tue: Reinforcement learning Thu: Deep reinforcement learning		Asn 4 in Asn 5 out
12	Nov 14, 16	Deep neural networks	Tensorflow playground Autoencoder examples	Tue: Reinforcement learning Thu: Deep reinforcement learning	ISL Section 10.7 Notes on backpropagation	Quiz 5
13	Nov 21, 23	No class, Thanksgiving break				
14	Nov 28, 30	Transformers and ChatGPT	ChatGPT demo	Tue: Transformers Thu: Human feedback and rewards		Asn 5 in
15	Dec 5, 7	Societal issues for machine learning		Tue: Panel discussion Thu: Course wrap up		Quiz 6
16	Fri, Dec 15, 2pm, Room TBA	Final exam		Registrar: Final exam schedule Practice final		

# For Today

- Language models
- Concepts...no new methods

# **When you text someone**

- Take out your cell phone...

# **When you text someone**

- Take out your cell phone...
- Text someone

# **When you text someone**

- Take out your cell phone...
- Text someone
- What did you notice?

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

- By the basic rules of conditional probability we can factor this as

$$p(w_1, \dots, w_n) = p(w_1)p(w_2 | w_1) \dots p(w_n | w_1, \dots, w_{n-1})$$

# Language models

- A language model is a way of *generating* any sequence of words

$$P(\text{"the whole forest had been anesthetized"}) =$$

$$\begin{aligned} & P(\text{"the"}) \times P(\text{"whole"} | \text{"the"}) \\ & \quad \times P(\text{"forest"} | \text{"the whole"}) \\ & \quad \times P(\text{"had"} | \text{"the whole forest"}) \\ & \quad \times P(\text{"been"} | \text{"the whole forest had"}) \\ & \quad \times P(\text{"anesthetized"} | \text{"the whole forest had been"}) \end{aligned}$$

# Remixing Noon

Text generated from Channel Skin by Jeff Noon

"The whole forest had been anesthetised, her temples wired, her senses stimulated, her eyes were not on Eva, not on Eva, not on Eva, not on anybody in that same realm, the land of dreams and nightmares."

Viability: 0.00000326%

[https://revdancatt.com/2017/03/01/markov\\_noon](https://revdancatt.com/2017/03/01/markov_noon)

# Text generation

- Words generated one-by-one
- A word is chosen by sampling from a probability distribution
- Result is purely synthetic text—generative AI

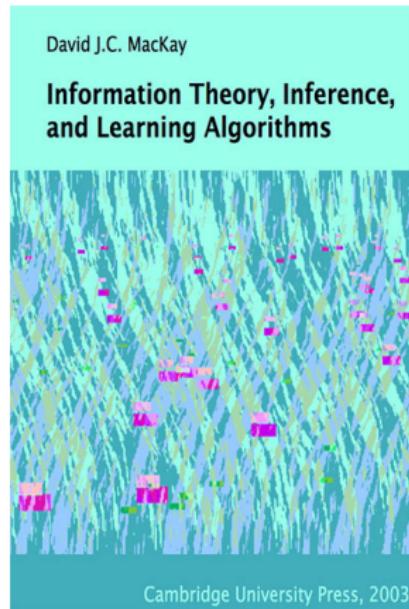
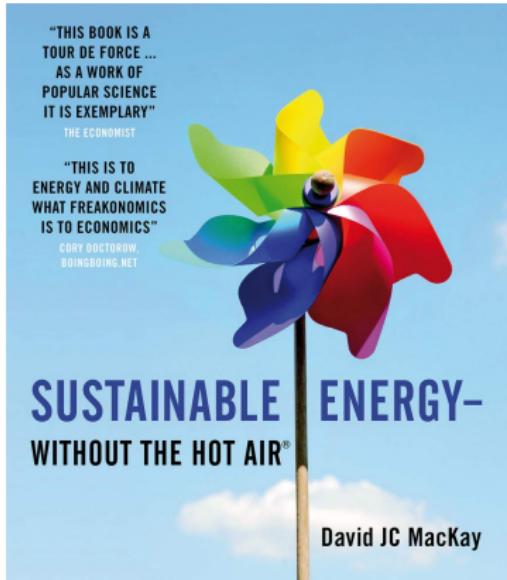
# Uses of language models

- Speech recognition
- Machine translation
- Text compression
- Texting
- Email completion
- Image captioning
- Mind reading from fMRI

---

Often built using Bayes' rule:  $P(\text{signal} \mid \text{words}) \propto P(\text{words} \mid \text{signal}) \cdot P(\text{words})$

# David MacKay



# Dasher: LMs for assistive devices

Language models enable new modes of text input:

[https://www.youtube.com/watch?v=quw\\_Kci4fUg](https://www.youtube.com/watch?v=quw_Kci4fUg)

<https://youtu.be/QxFEUk3J89Q?t=72>

Dasher poetry:

<https://www.youtube.com/watch?v=x-WLiY2p1LQ>

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

- By the basic rules of conditional probability we can factor this as

$$p(w_1, \dots, w_n) = p(w_1)p(w_2 | w_1) \dots p(w_n | w_1, \dots, w_{n-1})$$

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

- By the basic rules of conditional probability we can factor this as

$$p(w_1, \dots, w_n) = p(w_1)p(w_2 | w_1) \dots p(w_n | w_1, \dots, w_{n-1})$$

- The number of *histories* grows as  $V^{n-1}$ . Number of parameters in model grows as  $V^n$ , where  $V$  is number of words in vocabulary.

# Language models

- A language model is a way of assigning a probability to any sequence of words (or string of text)

$$p(w_1, \dots, w_n)$$

- By the basic rules of conditional probability we can factor this as

$$p(w_1, \dots, w_n) = p(w_1)p(w_2 | w_1) \dots p(w_n | w_1, \dots, w_{n-1})$$

- The number of *histories* grows as  $V^{n-1}$ . Number of parameters in model grows as  $V^n$ , where  $V$  is number of words in vocabulary.
- What are some ways of reducing the number of parameters?

## One approach: Grouping histories

- Let  $g(w_1, \dots, w_n)$  be the group assigned to the history

# One approach: Grouping histories

- Let  $g(w_1, \dots, w_n)$  be the group assigned to the history
- Our model becomes

$$p(w_{n+1} | w_1, \dots, w_n) = p(w_{n+1} | g(w_1, \dots, w_n))$$

## One approach: Grouping histories

- Let  $g(w_1, \dots, w_n)$  be the group assigned to the history
- Our model becomes

$$p(w_{n+1} | w_1, \dots, w_n) = p(w_{n+1} | g(w_1, \dots, w_n))$$

- Number of parameters:  $O(V \cdot \text{number of groups})$

# One approach: Grouping histories

- Let  $g(w_1, \dots, w_n)$  be the group assigned to the history
- Our model becomes

$$p(w_{n+1} | w_1, \dots, w_n) = p(w_{n+1} | g(w_1, \dots, w_n))$$

- Number of parameters:  $O(V \cdot \text{number of groups})$
- What are some example groupings?

# Grouping histories

- Unigrams:  $g(w_1, \dots, w_n) = \emptyset$ .

---

The notation  $O(\cdot)$  is called "Big Oh" and means "no greater than a constant times", so that  $O(f(n))$  means a sequence that is bounded by  $C \cdot f(n)$  for large enough  $n$ .

# Grouping histories

- Unigrams:  $g(w_1, \dots, w_n) = \emptyset$ .
- Bigrams:  $g(w_1, \dots, w_n) = w_n$ .

---

The notation  $O(\cdot)$  is called "Big Oh" and means "no greater than a constant times", so that  $O(f(n))$  means a sequence that is bounded by  $C \cdot f(n)$  for large enough  $n$ .

# Grouping histories

- Unigrams:  $g(w_1, \dots, w_n) = \emptyset$ .
- Bigrams:  $g(w_1, \dots, w_n) = w_n$ .
- Trigrams:  $g(w_1, \dots, w_n) = (w_{n-1}, w_n)$ .

---

The notation  $O(\cdot)$  is called "Big Oh" and means "no greater than a constant times", so that  $O(f(n))$  means a sequence that is bounded by  $C \cdot f(n)$  for large enough  $n$ .

# Grouping histories

- Unigrams:  $g(w_1, \dots, w_n) = \emptyset$ .
- Bigrams:  $g(w_1, \dots, w_n) = w_n$ .
- Trigrams:  $g(w_1, \dots, w_n) = (w_{n-1}, w_n)$ .
- Number of parameters grows as  $O(V)$ ,  $O(V^2)$ , and  $O(V^3)$ , respectively.

---

The notation  $O(\cdot)$  is called "Big Oh" and means "no greater than a constant times", so that  $O(f(n))$  means a sequence that is bounded by  $C \cdot f(n)$  for large enough  $n$ .

# Estimating parameters

- The maximum likelihood estimate of a trigram model:

$$\hat{p}(w_3 | w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

# Estimating parameters

- The maximum likelihood estimate of a trigram model:

$$\hat{p}(w_3 | w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

- What are some problems with this model?

# Sparse data problem

The next group of slides presentsk one way of quantifying the problem of sparse data in language modeling

# Half Earth

Aug 06, 2018

by Foundation Staff

0 Comment

Google Earth, Half-Earth Project, Map of Life

By Jeremy Malczyk, Michelle Duong, Ajay Ranipeta, Chris Heltne, Walter Jetz of Map of Life, Yale University, and the E.O. Wilson Biodiversity Foundation Half-Earth Project

This article originally in *Medium*, July 30, 2018

---

## The Significance of Biodiversity



Narrow-billed Tody, *Todus angustirostris*. Photo by Julie Hart.



Learn how you can be part of bringing  
Half-Earth to life.

<https://eowilsonfoundation.org/mapping-species-for-half-earth/>

<https://www.half-earthproject.org/>

# Specious species



- A naturalist (say, Edward O. Wilson) explores a region and observes animals (organisms).

# Specious species



- A naturalist (say, Edward O. Wilson) explores a region and observes animals (organisms).
- He finds 100,000 animals and 5,000 species.

# Specious species



- A naturalist (say, Edward O. Wilson) explores a region and observes animals (organisms).
- He finds 100,000 animals and 5,000 species.
- What is the chance I'll find a new species?

# Specious species



- A naturalist (say, Edward O. Wilson) explores a region and observes animals (organisms).
- He finds 100,000 animals and 5,000 species.
- What is the chance I'll find a new species?
- What if Wilson observes 100 unique species?

# Missing species: Good-Turing

- Wilson observes 100,000 animals and 5,000 species, 100 of them are unique (only one observation of that species).

# Missing species: Good-Turing

- Wilson observes 100,000 animals and 5,000 species, 100 of them are unique (only one observation of that species).
- Good-Turing: Estimate of the probability that the next animal is a new species?  $\hat{p}_{GT} = 100/100,000 = 10^{-3}$ .

# Missing species: Good-Turing

- Wilson observes 100,000 animals and 5,000 species, 100 of them are unique (only one observation of that species).
- Good-Turing: Estimate of the probability that the next animal is a new species?  $\hat{p}_{GT} = 100/100,000 = 10^{-3}$ .
- This is an estimate of the missing probability mass.

## Sparse data: Species $\approx$ words

- Suppose we have a corpus of 500 million words. I count trigrams and find that 50 million of them are unique.

## Sparse data: Species $\approx$ words

- Suppose we have a corpus of 500 million words. I count trigrams and find that 50 million of them are unique.
- When I see a new trigram, 10% of the time it won't have been seen before. (This is a typical number.)

## Sparse data: Species $\approx$ words

- Suppose we have a corpus of 500 million words. I count trigrams and find that 50 million of them are unique.
- When I see a new trigram, 10% of the time it won't have been seen before. (This is a typical number.)
- This means that the maximum likelihood estimate (MLE) is zero, and the probability that my model predicts the next word will be zero.

## Sparse data: Species $\approx$ words

- Suppose we have a corpus of 500 million words. I count trigrams and find that 50 million of them are unique.
- When I see a new trigram, 10% of the time it won't have been seen before. (This is a typical number.)
- This means that the maximum likelihood estimate (MLE) is zero, and the probability that my model predicts the next word will be zero.
- The MLE is supported on the observed data. We need to spread out the probability over unseen events.

# Estimating parameters

- The maximum likelihood estimate of a trigram model:

$$\hat{p}(w_3 | w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

# Estimating parameters

- The maximum likelihood estimate of a trigram model:

$$\hat{p}(w_3 | w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

- Some kind of “shrinkage” or smoothing needs to be done.

# Estimating parameters

- The maximum likelihood estimate of a trigram model:

$$\hat{p}(w_3 | w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\text{count}(w_1, w_2)}$$

- Some kind of “shrinkage” or smoothing needs to be done.
- How else can the model be strengthened?

# Interpolation

Linear interpolation:

$$p(w_3 | w_1, w_2) = \lambda_3 \hat{p}(w_3 | w_1, w_2) + \lambda_2 \hat{p}(w_3 | w_2) + \lambda_1 \hat{p}(w_3)$$

where  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ .

# Interpolation

Linear interpolation:

$$p(w_3 | w_1, w_2) = \lambda_3 \hat{p}(w_3 | w_1, w_2) + \lambda_2 \hat{p}(w_3 | w_2) + \lambda_1 \hat{p}(w_3)$$

where  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ .

This is a type of “mixture model”

# How good is a language model?

The next group of slides presents a useful way of quantifying how good a language model is

## Recall: Geometric mean

The *arithmetic mean* of 1/4, 4, 8 is

$$\frac{1}{3} \left( \frac{1}{4} + 4 + 8 \right) = 4.08\bar{3}$$

## Recall: Geometric mean

The *arithmetic mean* of 1/4, 4, 8 is

$$\frac{1}{3} \left( \frac{1}{4} + 4 + 8 \right) = 4.08\bar{3}$$

The *geometric mean* of 1/4, 4, 8 is

$$\sqrt[3]{\frac{1}{4} \cdot 4 \cdot 8} = 2$$

## Recall: Geometric mean

The *arithmetic mean* of 1/4, 4, 8 is

$$\frac{1}{3} \left( \frac{1}{4} + 4 + 8 \right) = 4.08\bar{3}$$

The *geometric mean* of 1/4, 4, 8 is

$$\sqrt[3]{\frac{1}{4} \cdot 4 \cdot 8} = 2$$

The geometric mean is no greater than the arithmetic mean

## Recall: Geometric mean

The *geometric mean* of  $x_1, \dots, x_n$  is

$$\sqrt[n]{x_1 x_2 \cdots x_n} = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

## Recall: Geometric mean

The *geometric mean* of  $x_1, \dots, x_n$  is

$$\sqrt[n]{x_1 x_2 \cdots x_n} = \left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}}$$

# How good is a language model? Perplexity

*Perplexity* is defined as

$$\text{Perplexity}(\theta) = \left( \prod_{i=1}^N p_\theta(w_i | w_{1:i-1}) \right)^{-\frac{1}{N}}$$

where  $w_1, w_2, \dots, w_N$  is a large chunk of text that wasn't used to train the language model.

# How good is a language model? Perplexity

- Perplexity is the inverse of the geometric mean of the word probabilities
- If the perplexity is 100, the model predicts, on average, as if there were 100 equally likely words to follow
- This is the (geometric) average “branching factor” for the model on real text

# Remixing Noon

Text generated from Channel Skin by Jeff Noon

"The whole forest had been anesthetised, her temples wired, her senses stimulated. her eyes were not on Eva, not on Eva, not on anybody in that same realm, the land of dreams and nightmares."

Viability: 0.000000326%

The diagram illustrates the flow of text from the beginning of the sentence to its end, forming a circle. The text is color-coded: purple for the first part, green for the middle, and red for the final part.

**Text Flow:**

- Start:** The whole forest had been anesthetised, her temples wired, her senses stimulated.
- Middle:** her eyes were not on Eva, not on Eva, not on anybody in that same realm, the land of dreams and nightmares.
- End:** signals her teeth even, lips cradles of spider alive with came to goddamn teeth even instrument room.

**Annotations:**

- Top Left:** The whole forest had been anesthetised, her temples wired, her senses stimulated.
- Bottom Left:** her eyes were not on Eva, not on Eva, not on anybody in that same realm, the land of dreams and nightmares.
- Bottom Right:** signals her teeth even, lips cradles of spider alive with came to goddamn teeth even instrument room.
- Top Right:** The whole forest had been anesthetised, her temples wired, her senses stimulated.

# Modern language models

Suppose a computer program assigns a “score” to possible next words:

$$s(v; \underbrace{w_1, \dots, w_n}_{\text{word history}})$$

# Modern language models

Suppose a computer program assigns a “score” to possible next words:

$$s(v; \underbrace{w_1, \dots, w_n}_{\text{word history}})$$

Can convert this to a language model by the “softmax” operation:

$$p(w | w_1, \dots, w_n) = \frac{\exp(s(w; w_1, \dots, w_n))}{\sum_{v \in V} \exp(s(v; w_1, \dots, w_n))}$$

# Modern language models

Suppose a computer program assigns a “score” to possible next words:

$$s(v; \underbrace{w_1, \dots, w_n}_{\text{word history}})$$

Can convert this to a language model by the “softmax” operation:

$$p(w | w_1, \dots, w_n) = \frac{\exp(s(w; w_1, \dots, w_n))}{\sum_{v \in V} \exp(s(v; w_1, \dots, w_n))}$$

In ChatGPT, the function  $s(v; w_{1:n})$  is learned on large amounts of text (unsupervised) using a type of deep neural network called a *transformer*.

# Modern language models

In fact, the score is computed as

$$s(v; \underbrace{w_1, \dots, w_n}_{\text{word history}}) = \beta_v^T g(w_1, \dots, w_n)$$

- the “grouping function”  $g(w_1, \dots, w_n)$  is like an embedding, that maps the word history to a high dimensional vector
- the weights  $\beta_v$  are essentially parameters in a big logistic regression

We'll see how this is done when we discuss neural networks



# Better Language Models and Their Implications

We've trained a large-scale unsupervised language model which generates coherent paragraphs of text, achieves state-of-the-art performance on many language modeling benchmarks, and performs rudimentary reading comprehension, machine translation, question answering, and summarization—all without task-specific training.



February 14, 2019  
24 minute read

# Summary of today: Language models

- A language model is used to predict or generate the next word
- Used many different applications
- Probabilities need to be “smoothed” to avoid zeros
- Perplexity is a measure of a language model’s predictive power
- ChatGPT and descendants represent frontier of AI

**By the way...**

### Yale researchers create map of undiscovered life

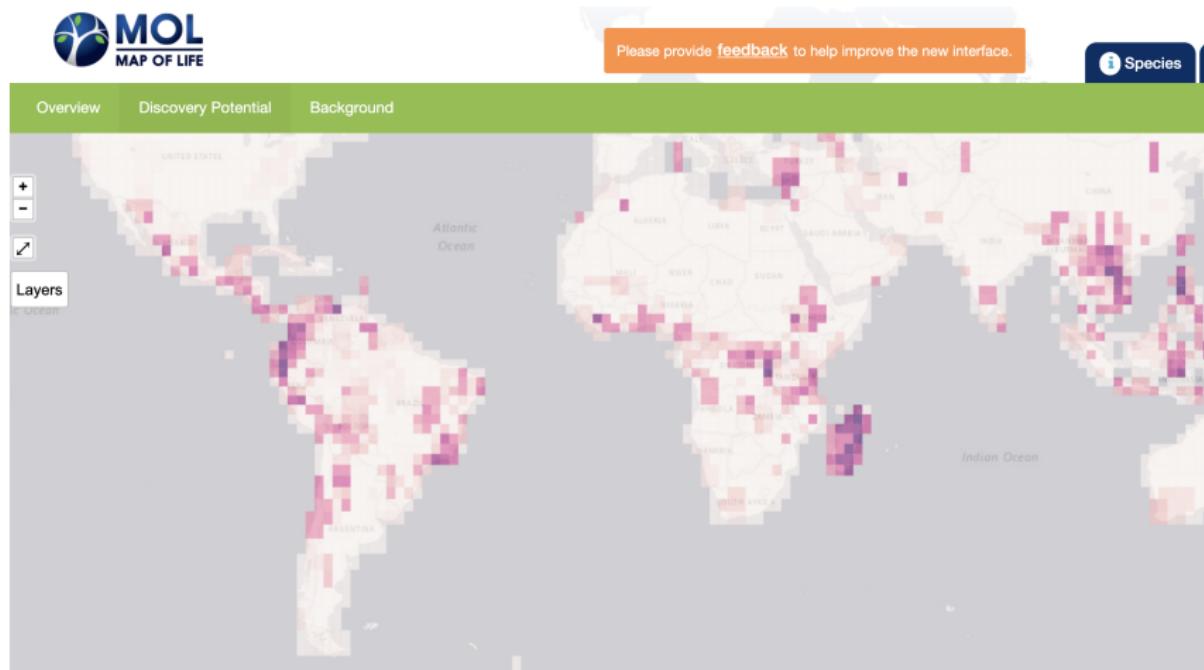
By Bill Hathaway | MARCH 22, 2021



Less than a decade after unveiling the "[Map of Life](#)," a global database that marks the distribution of known species across the planet, Yale researchers have launched an ambitious and perhaps even more important project — creating a map of where life has yet to be discovered.

For [Walter Jetz](#), a professor of ecology and evolutionary biology at Yale who spearheaded the Map of Life project, the new effort is a moral imperative that can help support biodiversity discovery and preservation around the world.

# Map of Life



# Map of Life



[contact us](#)    [login](#)    [register](#)

## *Putting biodiversity on the map*

