

Terrance Boulton
Walter Scheirer *Editors*

A Unifying Framework for Formal Theories of Novelty

Discussions, Guidelines, and
Examples for Artificial Intelligence

Synthesis Lectures on Computer Vision

Series Editors

Gerard Medioni, University of Southern California, Los Angeles, CA, USA

Sven Dickinson, Department of Computer Science, University of Toronto, Toronto, ON,
Canada

This series publishes on topics pertaining to computer vision and pattern recognition. The scope follows the purview of premier computer science conferences, and includes the science of scene reconstruction, event detection, video tracking, object recognition, 3D pose estimation, learning, indexing, motion estimation, and image restoration. As a scientific discipline, computer vision is concerned with the theory behind artificial systems that extract information from images. The image data can take many forms, such as video sequences, views from multiple cameras, or multi-dimensional data from a medical scanner. As a technological discipline, computer vision seeks to apply its theories and models for the construction of computer vision systems, such as those in self-driving cars/navigation systems, medical image analysis, and industrial robots.

Terrance Boulton · Walter Scheirer
Editors

A Unifying Framework for Formal Theories of Novelty

Discussions, Guidelines, and Examples
for Artificial Intelligence

Editors

Terrance Boulton
Department of Computer Science
University of Colorado at Colorado Springs
Colorado Springs, CO, USA

Walter Scheirer
Department of Computer Science
and Engineering
University of Notre Dame
Notre Dame, IN, USA

ISSN 2153-1056

ISSN 2153-1064 (electronic)

Synthesis Lectures on Computer Vision

ISBN 978-3-031-33053-7

ISBN 978-3-031-33054-4 (eBook)

<https://doi.org/10.1007/978-3-031-33054-4>

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Nature Switzerland AG 2024

This work is subject to copyright. All rights are solely and exclusively licensed by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface: The Novelty Problem in AI

AI researchers these days are quick to tout the progress that has been made in the field over the past decade. From game playing to visual recognition, new capabilities are appearing all of the time for many different applications. And indeed, such achievements should be celebrated. However, some very useful AI capabilities remain out of reach. For instance, why aren't safe self-driving cars available in the market in 2023? A major limitation of today's AI systems has become apparent in the quest for autonomous systems that must operate in real environments: they cannot manage novelty in the environment they were designed for. That is to say, if something new appears, there is no capacity for an agent to detect, characterize, and learn how to handle it. Given the practically infinite number of ways an environment can configure itself, coupled with the routine appearance of new things within an environment, novelty can be a significant confound. This book is the first attempt to study novelty problems in a rigorous fashion through the use of a unifying framework for formal theories of novelty.

Ad hoc ways of addressing the novelty problem have proven to be insufficient. There exists a persistent belief that reinforcement learning is all that is needed for agents to manage novelty because any novelty can be learned over time. Similarly, there exists blind faith in the generalization properties of deep learning through invariant representation alone. Given the results found in this book for the simplest of AI domains, we can safely say that a more principled approach to novelty management is needed. Neither approach addresses the core detection problem at the classifier level, nor is there any capacity to characterize novelty, which can take on many different forms. On that latter point, what exactly does it mean for something to be novel? That isn't a question that can be answered using an off-the-shelf AI algorithm. The need for a theory matched to a specific domain can provide a better starting point for agent design.

A recent effort to address the novelty problem in AI has been the DARPA Science of Artificial Intelligence and Learning for Open-world Novelty (SAIL-ON) program. It established a research program to develop a set of engineering design principles for open-world learning in 2019 [1]. Throughout the four years of the program, a large consortium of academic, industry, and government researchers has collaborated on fundamental work looking into innovative strategies for effective open-world learning in both activity domains, e.g., interactive video games, and perceptual domains, e.g., datasets of images and videos. This book is the output of the program’s Novelty Working Group, which was charged with developing viable theories for the study of novelty in AI. Each chapter was contributed by different participants in that working group.

This book is organized in the following manner. Chapter 1 is the focal point of the book. It introduces a unifying framework for creating theories of novelty that are matched to specific domains. This includes definitions on different types of novelty, as well as constructs for building agents that can detect, characterize, and manage novelty. This framework is general and can apply to *any* domain in which novelty appears. To justify this claim, each subsequent chapter provides an example domain for which a theory is developed and evaluated. These chapters include: (1) a task overview, (2) definitions of dissimilarity and regret operators, (3) definitions of measurements and observations, (4) a description of novelty types and examples, (5) a set of experiments validating predictions made by the developed theory, and (6) concluding remarks.

The domain-specific chapters cover a broad range of activity and perceptual domains. Chapter 2 starts things off as simple as possible with a study of novelty in the 2D Cart-Pole activity domain. Chapter 3 extends the study of CartPole by examining a 3D version of the environment. Chapter 4 turns to the perceptual domain of image classification in computer vision. Chapter 5 discusses a related computer vision domain, handwriting recognition, which also contains elements of natural language processing. Chapter 6 pushes farther into the realm of natural language processing by studying contextual and semantic novelty in text. Chapter 7 comes back to activity domains with an examination of the game Monopoly. The book concludes in Chap. 8 by recapping what we have learned and suggesting the development of new theoretical directions that are interdisciplinary in nature.

Colorado Springs, USA
Notre Dame, USA

Terrance Boulton
Walter Scheirer

Acknowledgments This research was sponsored by the Defense Advanced Research Projects Agency (DARPA) and the Army Research Office (ARO) under multiple contracts/agreements including HR001120C0055, W911NF-20-2-0005, W911NF-20-2-0004, HQ0034-19-D-0001, and

W911NF2020009. The views contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of DARPA or ARO, or the US Government.

Reference

1. DARPA. Teaching AI systems to adapt to dynamic environments, (2019)

Contents

1	A Unifying Framework for Novelty	1
	T. Boulton, D. S. Prijatelj, and W. Scheirer	
2	Novelty in 2D CartPole Domain	5
	P. A. Grabowicz, C. Pereyda, K. Clary, R. Stern, T. Boulton, D. Jensen, and L. B. Holder	
3	Novelty in 3D CartPole Domain	21
	T. Boulton, N. M. Windesheim, S. Zhou, C. Pereyda, and L. B. Holder	
4	Novelty in Image Classification	37
	A. Shrivastava, P. Kumar, Anubhav, C. Vondrick, W. Scheirer, D. S. Prijatelj, M. Jafarzadeh, T. Ahmad, S. Cruz, R. Rabinowitz, A. Al Shami, and T. Boulton	
5	Novelty in Handwriting Recognition	49
	D. S. Prijatelj, S. Grieggs, F. Yumoto, E. Robertson, and W. Scheirer	
6	Contextual and Semantic Novelty in Text	71
	N. Ma, B. Liu, and E. Robertson	
7	Multi-agent Game Domain: Monopoly	97
	T. Bonjour, M. Haliem, V. Aggarwal, M. Kejriwal, and B. Bhargava	
8	Concluding Thoughts	107
	T. Boulton and W. Scheirer	

List of Figures

Fig. 1.1	Main elements of the implicit theories of novelty; items with dashed outlines are outside of the task or agent but are critical to defining novelty. In the framework, a theory of novelty is obtained by specifying: world \mathcal{W} and the world dimensionality d' , observation space \mathcal{O} accessible to the agent, and agent space \mathcal{S} . The agent can only access world information indirectly through a perceptual operator \mathcal{P} with associated \mathcal{V}_t processed world regions. The agent α in state $s_t \in \mathcal{S}$ at time t , using state recognition function $f_t(x, s)$ to determine the action $a_t \in \mathcal{A}$ to be taken, which is used by an oracle's world state transition function T . Critical to defining novelty are task-dependent world dissimilarity functions $\mathcal{D}_{w, \mathcal{T}; E_t}$ with associated threshold δ_w , task-dependent observation space dissimilarity functions $\mathcal{D}_{o, \mathcal{T}; E_t}$ with threshold δ_o , world regret function $\mathcal{R}_{w, \mathcal{T}}$, observation space regret function $\mathcal{R}_{o, \mathcal{T}}$, and agent space regret function $\mathcal{R}_{a, \mathcal{T}}$. The framework also defines a task-dependent agent space dissimilarity function $\mathcal{D}_{\alpha, \mathcal{T}; E_t}$ which allows one to consider the difference between agent models for different worlds, <i>e.g.</i> , to measure learning. While not explicit in the figure, the world state and observed states may collapse, if perceptual operator is the identity. In addition, the world and observed spaces may include full or partial copies of the agent's memory/state/state recognition function. Every set of these operators/functions/values defines a different theory of novelty for its associated task	2
Fig. 2.1	Average dissimilarity, $\mathbb{E}_{\tilde{w}} \mathcal{D}_{o, \mathcal{T}}(w_t, \tilde{w}_t)$, between the future states expected and observed by agents that were tuned to the world w with incorrect value of the magnitude of pushing force, F_p (left panel), or a horizontal force acting on the cart, F_h (right panel), while tested in the world \tilde{w} . The expectation is computed over 20 samples of the initial world state w_0	8

- Fig. 2.2 Expected loss and regret of simulation-based agent when varying: **a** *push force* from 0 to 360 N, **b** a constant *horizontal force* from -10 to $+10$ N, **c** *gravity* from 0 to 360 m/s², **d** *pole length* from near 0 to 180 m. The heatmap (left) shows the difference in performance between an agent trained in agent's world, w , and evaluated on a different test world, \check{w} . The scatter plots show the expected loss (middle) and regret (right) in the novel environment as a function of difference-based dissimilarity, $\tilde{\mathcal{D}}_{w,\mathcal{T}}^d(w, \check{w})$ 12
- Fig. 2.3 Expected loss and regret of the DQN agent when varying: **a** *push force* from 0 to 360 N, **b** a constant *horizontal force* from -10 to $+10$ N, **c** *gravity* from 0 to 360 m/s², **d** *pole length* from near 0 to 180 m. The heatmap (left) shows the difference in performance between an agent trained in agent's world, w , and evaluated on a different test world, \check{w} . The scatter plots show the expected loss (middle) and regret (right) in the novel environment as a function of difference-based dissimilarity, $\tilde{\mathcal{D}}_{w,\mathcal{T}}^d(w, \check{w})$ 15
- Fig. 3.1 Our multi-type novelty experimental environment uses CartPole3D, where the goal is to keep the pole balanced. The environment has a controllable cart (green) balancing a pole (blue). It also has added independently moving environmental agents (red). The environment can be changed to provide for multiple subtypes of novelty. The Weibull Open World control-agent (WOW-agent) only sees the observational vector of 37–73 numeric values of the position/velocity of the cart, pole, the environmental agents, and the walls that define the world boundaries. In each episode, it receives the observations of each step, makes a control decision (left, right, front, back, none), and tries to keep the pole balanced for 200 timesteps. It reports the probability the world is novel in each episode. The image also shows episode (E=) and step number (S=), the WOW-agent's probability that the world has changed (WC=) to a novel state and the WOW-agent score (S=) 22