# Causal Machine Learning

## Supervised Learning

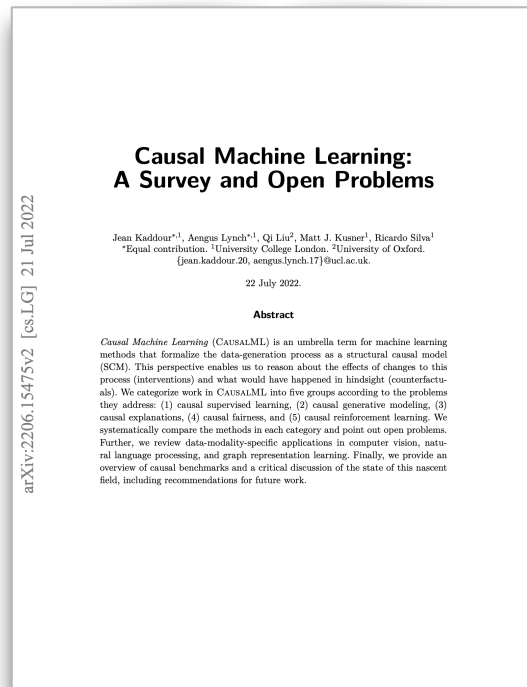Jordi Vitrià

jordi.vitria@ub.edu

UNIVERSITAT DE
BARCELONA

# Causal Machine Learning

Causal Machine Learning (CausalML) is an umbrella term for **machine learning methods** that are causally informed.
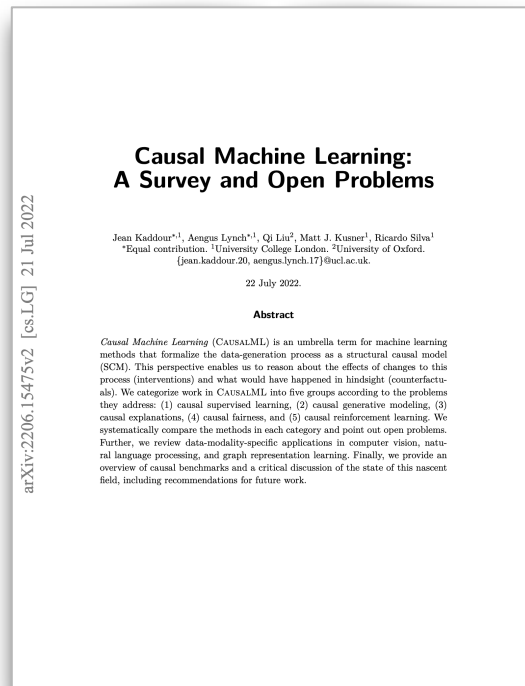
This perspective enables us to reason about the effects of changes in the data generation process (interventions) and what would have happened in hindsight (counterfactuals).

**Causal Machine Learning:
A Survey and Open Problems**

Jean Kaddour[*,1], Aengus Lynch[*,1], Qi Liu[2], Matt J. Kusner[1], Ricardo Silva[1]
[*]Equal contribution. [1]University College London. [2]University of Oxford.
{jean.kaddour.20, aengus.lynch.17}@ucl.ac.uk.

22 July 2022.

**Abstract**

*Causal Machine Learning* (CAUSALML) is an umbrella term for machine learning methods that formalize the data-generation process as a structural causal model (SCM). This perspective enables us to reason about the effects of changes to this process (interventions) and what would have happened in hindsight (counterfactuals). We categorize work in CAUSALML into five groups according to the problems they address: (1) causal supervised learning, (2) causal generative modeling, (3) causal explanations, (4) causal fairness, and (5) causal reinforcement learning. We systematically compare the methods in each category and point out open problems. Further, we review data-modality-specific applications in computer vision, natural language processing, and graph representation learning. Finally, we provide an overview of causal benchmarks and a critical discussion of the state of this nascent field, including recommendations for future work.

https://arxiv.org/pdf/2206.15475.pdf

# Causal Machine Learning

We can categorize work in CausalML into five groups according to the problems they address: (1) **causal supervised learning**, (2) causal generative modeling, (3) causal explanations, (4) causal fairness, and (5) causal reinforcement learning.

### Causal Machine Learning: A Survey and Open Problems

Jean Kaddour[*,1], Aengus Lynch[*,1], Qi Liu[2], Matt J. Kusner[1], Ricardo Silva[1]
[*]Equal contribution. [1]University College London. [2]University of Oxford.
{jean.kaddour.20, aengus.lynch.17}@ucl.ac.uk.

22 July 2022.

**Abstract**

*Causal Machine Learning* (CausalML) is an umbrella term for machine learning methods that formalize the data-generation process as a structural causal model (SCM). This perspective enables us to reason about the effects of changes to this process (interventions) and what would have happened in hindsight (counterfactuals). We categorize work in CausalML into five groups according to the problems they address: (1) causal supervised learning, (2) causal generative modeling, (3) causal explanations, (4) causal fairness, and (5) causal reinforcement learning. We systematically compare the methods in each category and point out open problems. Further, we review data-modality-specific applications in computer vision, natural language processing, and graph representation learning. Finally, we provide an overview of causal benchmarks and a critical discussion of the state of this nascent field, including recommendations for future work.

https://arxiv.org/pdf/2206.15475.pdf

# Causal Supervised Learning

The goal of supervised learning is to learn the conditional distribution $P(Y \mid X)$, or more generally $\mathbb{E}(Y \mid X)$, by training on data of the form $D = \{(x_i, y_i)\}_{i=1}^{N}$, where $X$ and $Y$ denote covariates and label, respectively.

One of the most fundamental principles in supervised learning is to assume that our data $D$ is independent and identically distributed (i.i.d.).

The validity of this assumption has been challenged; it has been famously called "*the big lie in machine learning*".

# Causal Supervised Learning

As an alternative to the i.i.d. assumption, we can assume that our data is sampled from interventional distributions governed by a causal model.

For a given dataset generated across a set of environments $\varepsilon$, $\{(x_i^e, y_i^e)_{i=1}^N\}_{e \in \varepsilon}$, we view each environment $e \in \varepsilon$ as being sampled from a separate interventional distribution.

How can we estimate $P(Y|X)$ in a principled, robust manner?

# Invariant Feature Learning

**Invariant feature learning** (IFL) is the task of identifying features of our data $X, X_c$, that are predictive of $Y$ across a range of environments $\mathcal{E}$.



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98

(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97

(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94

# Distribution Shifts

In this paper, authors provide a unifying framework for **specifying <span style="color:red">dataset shifts</span>** that can occur, analyzing model stability to these shifts, and determining conditions for achieving the lowest worst-case error across environments produced by these shifts.

This provides common ground so that we can begin to answer fundamental questions such as:

- **To what dataset shifts are the model's predictions stable vs unstable?** (Stability of the data generating model)

- **How will the model's performance be affected by these shifts?**



https://arxiv.org/pdf/1905.11374.pdf

# Distribution Shifts

**Example**: The goal is to diagnose pneumonia $Y$ from chest x-rays $Z$ and stylistic features of the image $X$ (i.e., orientation and coloring). The latent variable $W$ represents the hospital department the patient visited.



department

W

disease          style

Y                X

Z

x-ray image

In the pneumonia example, each department has its own protocols and equipment, so the style preferences $P(X \mid W)$ vary across departments.

# Distribution Shifts

Each environment is a different instantiation of that graph such that certain mechanisms differ.

Thus, the **factorization of the data distribution is the same in each environment**, but the **terms in the factorization corresponding to shifts will vary across environments**.

$$E = \{P(Z \,|\, Y, X)P(Y \,|\, W)P(X \,|\, W)P(W)\}$$

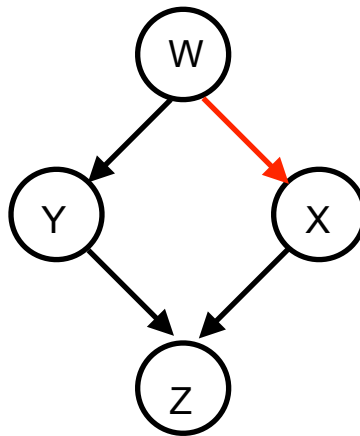# Distribution Shifts

**Key Result**: Distribution shifts can be expressed in terms of edges.

# Distribution Shifts

A graph and a set of edges which are marked as unstable defines an uncertainty set of environments whose distributions differ in the unstable factors.
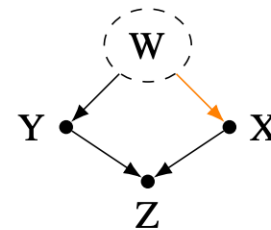


$$P(Z|Y,X)P(Y|W)P(X|W)P(W)$$

$$E = \{P(Z|Y,X)P(Y|W)P(X|W)P(W)\}$$
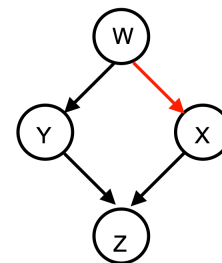
# Distribution Shifts

In this pneumonia example, if $W$ is unobserved, a model of $P(Y \,|\, X, Z)$ will learn an association between $Y$ and $X$ through $W$. Thus, $P(Y \,|\, X, Z)$ contains an **unstable** path, and this distribution is **unstable** to shifts in the style mechanism. This means that $P(Y \,|\, X, Z)$ is different in each environment.

By contrast, if $W$ were observed and we could condition on it, then $P(Y \,|\, X, Z, W)$ is **stable** to shifts in the style mechanism because all paths containing the unstable edge are blocked by $W$.

Thus, $P(Y \,|\, X, Z, W)$ is invariant across environments.

$P(Y \,|\, X, Z)$ is **unstable** because of the backdoor path.

# Distribution Shifts

In order to achieve stable distributions to shifts we can

- find the maximal set of features to **condition** on so that the resulting model is stable with respect to the foreseen shifts,

- **intervene** $(do(\ \cdot\ ))$ in variables with a shifted mechanism,

- compute **counterfactuals**.

$P(X \mid Y)$



$P(Y \mid V)$



$P(Y \mid V, Z, do(X)) = P(Z \mid X, Y)P(X \mid Y)P(Y \mid V)P(V)$



$P(Y* \mid X_{Y=0}^*, Z_{X=x}^*, V = v)$