

# Invariance, Causality, and Applications

David M. Blei

Columbia University

Spring 2026

## Seminar overview

This seminar explores the idea of *invariance* and its role in causal inference and probabilistic modeling. We will study statistical methods for finding invariant relationships across multiple environments, how these ideas connect to causality and robustness, and how they extend to representation learning and empirical Bayes. Our goal is to explore research at the intersection of all of these ideas.

Our context for these subjects will be probabilistic modeling and machine learning. This is not a survey course or a gentle introduction. It is an exploratory seminar where we will think about and develop ideas together.

## Core themes

We will explore four core circles of ideas.

### 1. Causal inference

- Prediction under intervention
- Discovery of causal relationships
- Counterfactual reasoning
- Identification and estimation

### 2. Invariance

- Identifying relationships that remain stable across environments
- Connections to robustness and generalization
- Invariance as a criterion for scientific explanation

### 3. Causal representation learning

- High-dimensional observations with lower-dimensional causal structure
- Inferring latent causal variables (e.g., health states from medical images)
- The role and limits of “disentanglement”

### 4. Empirical Bayes

- Bayesian models that adapt to populations
- Learning priors from data
- A blend of Bayesian and frequentist thinking

## Prerequisites

I expect you to be fluent in probabilistic modeling and machine learning.

This typically means you have taken *Probabilistic Models and Machine Learning*. You should be able to design probabilistic models and implement approximate inference.

Ideas from this field—such as deep generative models, amortized variational inference, stochastic optimization, graphical models, probabilistic diffusion models, and others—will be in the *background* of our discussion. I expect you to know most of them well.

## Seminar logistics

- Enrollment cap: 25 students
- Devices: No laptops, tablets, or phones during seminar meetings
- Format: Discussion-based, instructor-led
- Expectations: Reading, participation, progress on a doctoral-level research project

## Readings and reader reports

We will read a mix of papers and books. There will usually be an assigned reading.

Every week, you must post a reader report to Slack. It should be prepared in LaTeX and under one page. It can include ideas, questions, confusions, opinions.

Before each class, please read the other students' reader reports.

## Final project

You must complete a research project by the end of the semester. This project should pursue original work around the themes of the class. Projects should aim to develop new methods for solving important applied problems.

We will start projects early. Every two weeks, all students will report on their progress.

In developing your project, you should be able to answer the following questions:

1. What problem am I solving? State it clearly.
2. Why is this problem important? To whom?
3. What is my strategy for solving it?
4. Why is this solution good? Where does it fall short?
5. How can I demonstrate that the solution works?

We hope that your project will use real data and pursue real questions. We prefer large, non-benchmark datasets. Examples:

- Electronic health records
- Neuroscience recordings
- Climate measurements

- Educational assessments
- Political polls
- Econometric data
- Genetic measurements

(We are less interested in finance applications.)

## Grading

- Final project: 75%
- Participation: 25%

Participation includes class discussion and reader reports.

## Schedule

The schedule is loose and will likely change as the semester evolves.

Week	Topic
1	Orientation: What this class is about.
2	Causality Foundations I
3	Causality Foundations II
4	Causality Foundations III
5	Invariance and Causality
6	Probabilistic Invariance
7	Invariant Risk Minimization
8	Causal Representation Learning I
9	Causal Representation Learning II
10	Causal Representation Learning III
11	Empirical Bayes I — Classical Foundations
12	Empirical Bayes II — Probabilistic Symmetries
13	Synthesis and discussion

## Reading list

Here are readings about each topic of the course. Starred readings are ones we will likely cover in some detail.

### Foundations of causal inference

- [Pearl et al. \(2016\)](#) [\*]
- [Peters et al. \(2018\)](#) [\*]
- [Pearl \(2009\)](#)
- [Pearl and MacKenzie \(2018\)](#)
- [Imbens and Rubin \(2015\)](#)
- [Morgan and Winship \(2015\)](#)
- [Hernan and Robins \(2020\)](#)

### Invariance and causality

- [Peters et al. \(2016\)](#) [\*]
- [Bühlmann \(2020\)](#) [\*]
- [Wu et al. \(2025b\)](#) [\*]
- [Wang et al. \(2024\)](#) [\*]
- [Pfister et al. \(2019\)](#)
- [Nguyen et al. \(2026\)](#)

### Causal representation learning and invariant risk minimization

- [Schölkopf et al. \(2021\)](#) [\*]
- [Arjovsky et al. \(2019\)](#) [\*]
- [Krueger et al. \(2021\)](#) [\*]
- [Uhler and Zhang \(2025\)](#) [\*]
- [Wang and Jordan \(2024\)](#)
- [Moran and Aragam \(2025\)](#)
- [von Kügelgen et al. \(2025\)](#)

### Synthetic controls

- [Abadie \(2021\)](#) [\*]
- [Athey et al. \(2021\)](#) [\*]
- [Shi et al. \(2022\)](#) [\*]
- [Abadie et al. \(2010\)](#)
- [Agarwal et al. \(2025\)](#)
- [Powell \(2026\)](#)
- [Rho et al. \(2026\)](#)

### Empirical Bayes

- [Efron \(2012\)](#) [\*]

- [Efron \(2019\) \[\\*\]](#)
- [Wu et al. \(2025a\) \[\\*\]](#)
- [Robbins \(1956\)](#)

## Datasets

Here are some datasets that might be useful.

- *CausalVerse*
  - Project page: <https://causal-verse.github.io/>
  - Hosted datasets and code: <https://huggingface.co/CausalVerse>
  - a benchmark with standardized tasks and evaluation protocols
- *All of Us Research Program*
  - Program overview: <https://allofus.nih.gov/>
  - Researcher Workbench: <https://www.researchallofus.org/>
  - a workbench with access to longitudinal EHR, survey, and genomics data.
- *MIMIC-IV*
  - Dataset page (PhysioNet): <https://physionet.org/content/mimiciv/>
  - EHR data used for treatment effect estimation, policy learning, and counterfactual evaluation
- *UK Biobank*
  - Main portal: <https://www.ukbiobank.ac.uk/>
  - Research access platform: <https://www.ukbiobank.ac.uk/enable-your-research>

## References

- Abadie, A. (2021). Using synthetic controls: Feasibility, data requirements, and methodological aspects. *Journal of Economic Literature*, 59(2):391–425.
- Abadie, A., Diamond, A., and Hainmueller, J. (2010). Synthetic control methods for comparative case studies: Estimating the effect of California’s tobacco control program. *Journal of the American Statistical Association*, 105(490):493–505.
- Agarwal, A., Shah, D., and Shen, D. (2025). Synthetic interventions: Extending synthetic controls to multiple treatments. *Operations Research*. Articles in Advance; published online December 8, 2025.
- Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. (2019). Invariant risk minimization. *arXiv:1907.02893*.
- Athey, S., Bayati, M., Doudchenko, N., Imbens, G., and Khosravi, K. (2021). Matrix completion methods for causal panel data models. *arXiv:1710.10251*.
- Bühlmann, P. (2020). Invariance, causality and robustness. *Statistical Science*, 35(3):404–426. 2018 Neyman Lecture, Special Issue on Causal Inference.
- Efron, B. (2012). *Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*. Cambridge University Press.
- Efron, B. (2019). Bayes, oracle Bayes, and empirical Bayes. *Statistical Science*, 34(2):177–201.
- Hernan, M. and Robins, J. (2020). *Causal Inference: What If?* Chapman & Hall/CRC.
- Imbens, G. and Rubin, D. (2015). *Causal Inference in Statistics, Social and Biomedical Sciences: An Introduction*. Cambridge University Press.
- Krueger, D., Caballero, E., Jacobsen, J., Zhang, A., Binas, J., Zhang, D., Priol, R. L., and Courville, A. (2021). Out-of-distribution generalization via risk extrapolation. *arXiv:2003.00688*.
- Moran, G. and Aragam, B. (2025). Towards interpretable deep generative models via causal representation learning. *arXiv preprint arXiv:2504.11609*.
- Morgan, S. and Winship, C. (2015). *Counterfactuals and Causal Inference*. Cambridge University Press, 2nd edition.
- Nguyen, N., Nguyen, P., Truong, T., Hoang, T., and Sugiyama, M. (2026). Causal graph learning via distributional invariance of cause-effect relationship. *Transactions on Machine Learning Research*. Published January 2026.
- Pearl, J. (2009). *Causality*. Cambridge University Press, 2nd edition.
- Pearl, J., Glymour, M., and Jewell, N. (2016). *Causal Inference in Statistics: A Primer*. John Wiley & Sons.
- Pearl, J. and MacKenzie, D. (2018). *The Book of Why*. Basic Books.

- Peters, J., Bühlmann, P., and Meinshausen, N. (2016). Causal inference by using invariant prediction: Identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012.
- Peters, J., Janzing, D., and Schoelkopf, B. (2018). *Elements of Causal Inference : Foundations and Learning Algorithms*. The MIT Press.
- Pfister, N., Bühlmann, P., and Peters, J. (2019). Invariant causal prediction for sequential data. *Journal of the American Statistical Association*, 114(527):1264–1276.
- Powell, D. (2026). Imperfect synthetic controls. *Journal of Applied Econometrics*.
- Rho, S., Illick, C., Narasipura, S., Abadie, A., Hsu, D., and Misra, V. (2026). Time-aware synthetic control. *arXiv:2601.03099*.
- Robbins, H. (1956). An empirical Bayes approach to statistics. In *Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability*, pages 131–148.
- Schölkopf, B., Locatello, F., Bauer, S., Ke, N., Kalchbrenner, N., Goyal, A., and Bengio, Y. (2021). Towards causal representation learning. *arXiv:2102.11107*.
- Shi, C., Sridhar, D., Misra, V., and Blei, D. (2022). On the assumptions of synthetic control methods. In *Artificial Intelligence and Statistics*.
- Uhler, C. and Zhang, J. (2025). Causal structure and representation learning with biomedical applications. *arXiv 2511.04790*.
- von Kügelgen, J., Ketterer, J., Shen, X., Meinshausen, N., and Peters, J. (2025). Representation learning for distributional perturbation extrapolation. *arXiv:2504.18522*.
- Wang, Y. and Jordan, M. (2024). Desiderata for representation learning: A causal perspective. *Journal of Machine Learning Research*, 25(275):1–65.
- Wang, Z., Hu, Y., Bühlmann, P., and Guo, Z. (2024). Causal invariance learning via efficient nonconvex optimization. *arXiv:2412.11850*.
- Wu, B., Weinstein, E. N., and Blei, D. M. (2025a). Bayesian empirical bayes: Simultaneous inference from probabilistic symmetries. *arXiv 2512.16239*.
- Wu, L., Yin, M., Wang, Y., Cunningham, J. P., and Blei, D. M. (2025b). Bayesian invariance modeling of multi-environment data. *arXiv 2506.22675*.