

36 Causality

This chapter was written by Victor Veitch and Alex D'Amour.

36.1 Introduction

The bulk of machine learning considers relationships between observed variables with the goal of summarizing these relationships in a manner that allows predictions on similar data. However, for many problems, our main interest is to predict how system would change if it were observed under different conditions. For instance, in healthcare, we are interested in whether a patient will recover if given a certain treatment (as opposed to whether treatment and recovery are associated in the observed data). **Causal inference** addresses how to formalize such problems, determine whether they can be solved, and, if so, how to solve them. This chapter covers the fundamentals of this subject. Code examples for the discussed methods are available at <https://github.com/vveitch/causality-tutorials>.

To make the gap between observed data modeling and causal inference concrete, consider the relationships depicted in Figure 36.1a and Figure 36.1b. Figure 36.1a shows the relationship between deaths by drowning and ice cream production in the United States in 1931 (the pattern holds across most years). Figure 36.1b shows the relationship between smoking and lung cancer across various countries. In each case, there is a strong positive association. Faced with this association, we might ask: could we reduce drowning deaths by banning ice cream? Could we reduce lung cancer by banning cigarettes? We intuitively understand that these interventional questions have different answers, despite the fact that the observed associations are similar. Determining the causal effect of some intervention in the world requires some such causal hypothesis about the world.

For concreteness, consider three possible explanations for the association between ice cream and drowning. Perhaps eating ice cream does cause people to drown—due to stomach cramps or similar. Or, perhaps, drownings increase demand for ice cream—the survivors eat huge quantities of ice cream to handle their grief. Or, the association may be due (at least in part) to a common cause: warm weather makes people more likely to eat ice cream and more likely to go swimming (and, hence, to drown). Under all three scenarios, we can observe exactly the same data, but the implications for an ice cream ban are very different. Hence, answering questions about what will happen under an intervention requires us to incorporate some causal knowledge of the world—e.g., which of these scenarios is plausible?

Our goal in this chapter to introduce the essentials of estimating causal effects. The high-level approach has three steps.

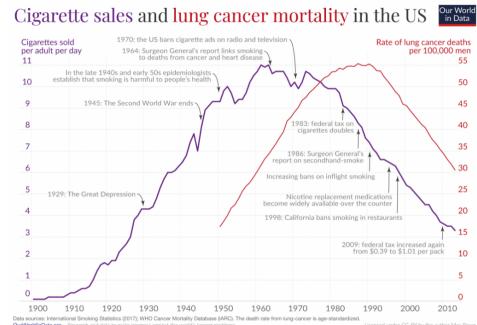
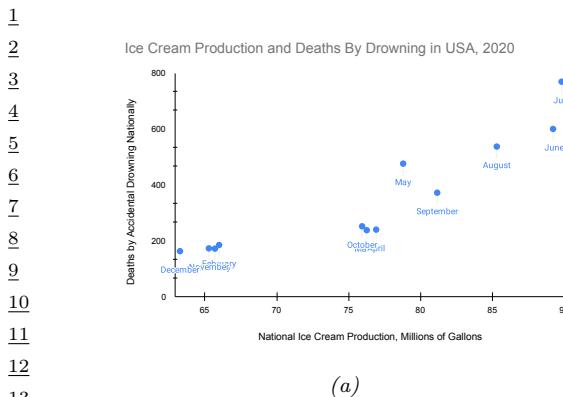


Figure 36.1: Correlation is not causation. (a) *Ice cream production is strongly associated with deaths by drowning.* Ice cream production data from the US Department of Agriculture National Agricultural Statistics Service. Drowning data from the National Center for Health Statistics at the United States Centers for Disease Control. (b) *Smoking is strongly associated with lung cancer.* From ourworldindata.org/smoking-big-problem-in-brief. Used with kind permission of Max Roser.

- **Causal Estimands:** The first step is to formally define the quantities we want to estimate. These are summaries of how the world would change under intervention, rather than summaries of the world as it has already been observed. E.g., we want to formalize “The expected number of drownings in the United States if we ban ice cream”.
- **Identification:** The next step is to identify the causal estimands with quantities that can, in principle, be estimated from observational data. This step involves codifying our causal knowledge of the world and translating this into a statement such as, “The causal effect is equal to the expected number of drownings after adjusting for month”. This step tells us what causal questions we could answer with perfect knowledge of the observed data distribution.
- **Estimation:** Finally, we must estimate the observable quantity using a finite data sample. The form of causal estimands favors certain efficient estimation procedures that allow us to exploit non-parametric (e.g., machine learning) predictive models.

In this chapter, we'll mainly focus on the estimation of the causal effect of an intervention averaged over all members of a population, known as the **Average Treatment Effect** or **ATE**. This is the most common problem in applied causal inference work. It is in some sense the simplest problem, and will allow us to concretely explain the use and importance of the fundamental causal concepts. These causal concepts include structural causal models, causal graphical models, the do-calculus, and efficient estimation using influence function techniques. This problem is also useful for understanding the role that standard predictive modeling and machine learning play in estimating causal quantities.

36.2 Causal Formalism

In causal inference, the goal is to use data to learn about how the outcome in the world would change under intervention. In order to make such inferences, we must also make use of our causal knowledge

of the world. This requires a formalism that lets us make the notion of intervention precise and lets us encode our causal knowledge as assumptions.

36.2.1 Structural Causal Models

Consider a setting in which we observe four variables from a population of people: A_i , an indicator of whether or not person i smoked at a particular age, Y_i , an indicator of whether or not person i developed lung cancer at a later age, H_i , a “health consciousness” index that measures a person’s health-consciousness (perhaps constructed from a set of survey responses about attitudes toward health), and G_i , an indicator for whether the person has a genetic predisposition towards cancer. Suppose we observe a dataset of these variables drawn independently and identically from a population, $(A_i, Y_i, H_i) \stackrel{\text{iid}}{\sim} P^{\text{obs}}$, where “obs” stands for “observed”.

In standard practice, we model data like these using probabilistic models. Notably, there are many different ways to specify a probabilistic model for the same observed distribution. For example, we could write a probabilistic model for P^{obs} as

$$A \sim P^{\text{obs}}(A) \quad (36.1)$$

$$H|A \sim P^{\text{obs}}(H|A) \quad (36.2)$$

$$Y|A, H \sim P^{\text{obs}}(Y|H, A) \quad (36.3)$$

$$G|A, H, Y \sim P^{\text{obs}}(G|A, H, Y) \quad (36.4)$$

This is a valid factorization, and sampling variables in this order would yield valid samples from the joint distribution P^{obs} . However, this factorization does not map well to a mechanistic understanding of how these variables are causally related in the world. In particular, it is perhaps more plausible that health-consciousness H causally precedes smoking status A , since a person’s health-consciousness would influence their decision to smoke.

These intuitions about causal ordering are intimately tied to the notion of intervention. Here, we will focus on a notion of intervention that can be represented in terms of “structural” models that describe mechanistic relationships between variables. The fundamental objects that we will reason about are **structural causal models**, or SCM’s. SCM’s resemble probabilistic models, but they encode additional assumptions (see also Section 4.7). Specifically, SCM’s serve two purposes: they describe a probabilistic model *and* they provide semantics for transforming the data-generating process through intervention.

Formally, SCM’s describe a mechanistic data generating process with an ordered sequence of equations that resemble assignment operations in a program. Each variable in a system is determined by combining other modeled variables (the causes) with exogenous “noise” according to some (unknown) deterministic function. For instance, a plausible SCM for P^{obs} might be

$$G \leftarrow f_G(\xi_0) \quad (36.5)$$

$$H \leftarrow f_H(\xi_1) \quad (36.6)$$

$$A \leftarrow f_A(H, \xi_2) \quad (36.7)$$

$$Y \leftarrow f_Y(G, H, A, \xi_3) \quad (36.8)$$

where the (unknown) functions f are fixed, and the variables ξ are unmeasured causes, modeled as independent random “noise” variables. Conceptually, the functions f_G, f_H, f_A, f_Y describe deter-

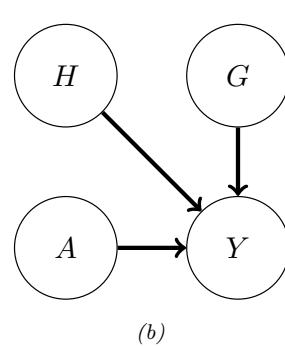
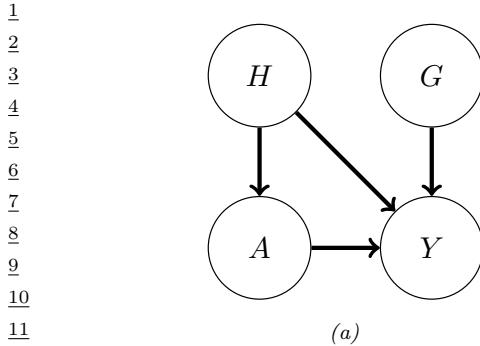


Figure 36.2: (a) Causal graph illustrating relationships between smoking A , cancer Y , health consciousness H , and genetic cancer pre-disposition G . (b) “Mutilated” causal graph illustrating relationships under an intervention on smoking A .

ministic physical relationships in the real world, while the variables ξ are hidden causes that are sufficient to distinguish each unit i in the population. Because we assume that each observed unit i is drawn at random from the population, we model ξ as random noise.

SCM’s imply probabilistic models, but not the other way around. For example, our example SCM implies probabilistic model for the observed data based on the factorization $P^{\text{obs}}(G, H, A, Y) = P^{\text{obs}}(G)P^{\text{obs}}(H)P^{\text{obs}}(A | H)P^{\text{obs}}(Y | A, H)$. Thus, we could sample from the SCM in the same way we would from a probabilistic model: draw a set of noise variables ξ and evaluate each assignment operation in the SCM in order.

Beyond the probabilistic model, an SCM encodes additional assumptions about the effects of interventions. This can be formalized using the **do calculus** (as in the verb “to do”), which we describe in Section 36.8; But in brief, interventions are represented by replacing assignment statements. For example, if we were interested in the distribution of Y in the hypothetical scenario that smoking were eliminated, we could set the second line of the SCM to be $A \leftarrow 0$. We would denote this by $P(Y|\text{do}(A=0), H)$. Because the f functions in the SCM are assumed to be invariant mechanistic relationships, the SCM encodes the assumption that this edited SCM generates data that we would see if we really applied this intervention in the world. Thus, the ordering of statements in an SCM are load-bearing: they imply substantive assumptions about how the world changes in response to interventions. This is in contrast to more standard probabilistic models where variables can be rearranged by applications of Bayes Rule without changing the substantive implications of the model. (See also Section 4.7.3.)

We note that structural causal model may not incorporate all possible notions of causality. For example, laws based on conserved quantities or equilibria—e.g., the ideal gas law—do not trivially map to SCMs, though these are fundamental in disciplines such as physics and economics. Nonetheless, we will confine our discussion to SCMs.

42

36.2.2 Causal DAGs

SCM’s encode many details about the assumed generative process of a system, but often it is useful to reason about causal problems at a higher level of abstraction. In particular, it is often useful

to separate the causal structure of a problem from the particular functional form of those causal relationships. **Causal graphs** provide this level of abstraction. A causal graph specifies which variables causally affect other variables, but leaves the parametric form of the structural equations f unspecified. Given an SCM, the corresponding causal graph can be drawn as follows: for each line of the SCM, draw arrows from the variables on the right hand side to variables on the left hand side. The causal DAG for our smoking-cancer example is shown in Figure 36.2. In this way, causal DAGs are related to SCMs in the same way that probabilistic graphical models (PGMs) are related to probabilistic models.

In fact, in the same way that SCMs imply a probabilistic model, causal DAGs imply a PGM. Functionally, causal graphs behave as probabilistic graphical models (Chapter 4). They imply conditional independence relationships between the variables in the observed data in same way. They obey the Markov property: If $X \leftarrow Y \rightarrow Z$ then $X \perp\!\!\!\perp Z|Y$; recall d-separation (Section 4.2.4.1). Additionally, if $X \rightarrow Y \leftarrow Z$ then, usually, $X \not\perp\!\!\!\perp Z|Y$ (even if X and Z are marginally independent). In this case, Y is called a **collider** for X and Z .

Conceptually, the difference between causal DAGs and PGMs is that probabilistic graphical models encode our assumptions about statistical relationships, whereas causal graphs encode our (stronger) assumptions about causal relationships. Such causal relationships can be used to derive how statistical relationships would change under intervention.

Causal graphs also allow us to reason about the causal and non-causal origins of statistical dependencies in observed data without specifying a full SCM. In a causal graph, two variables—say, A and D —can be statistically associated in different ways. First, there can be a directed path from (ancestor) A to (descendant) D . In this case, A is a causal ancestor of D and interventions on A will propagate through to change D ; $P(D|\text{do}(A = a)) \neq P(D|\text{do}(A = a'))$. For example, smoking is a causal ancestor of cancer in our example. Alternatively, A and D could share a common cause—there is some variable C such that there is a directed path from C to A and from C to D . If A and D are associated only through such a path then interventions on A will not change the distribution of D . However, it is still the case that $P(D|A = a) \neq P(D|A = a')$ —observing different values of A changes our guess for the value of D . The reason is that A carries information about C , which carries information about D . For example, suppose we lived in a world where there was no effect of smoking on developing cancer (e.g., everybody vapes), there would nevertheless be an association between smoking and cancer because of the path $A \leftarrow H \rightarrow Y$. The existence of such “backdoor paths” is one core reason that statistical and causal association are not the same. Of course, more complicated variants of these associations are possible—e.g., C is itself only associated with A through a backdoor path—but this already captures the key distinction between causal and non-causal paths.

Recall that our aim in introducing SCMs and causal graphs is to enable us to formalize our causal knowledge of the world and to make precise what interventional quantities we'd like to estimate. Writing down a causal graph gives a simple formal way to encode our knowledge of the causal structure of a problem. Usefully, this causal structure is sufficient to directly reason about the implications of interventions without fully specifying the underlying SCM. The key observation is that if a variable A is intervened on then, after intervention, none of the other variables are causes of A . That is, when we replace a line of an SCM with a statement directly assigning a variable a particular value, we cut off all dependencies that variable had on its causal parents. Accordingly, in the causal graph, the intervened on variable has no parents. This leads us to the **graph surgery** notion of intervention: an intervention that sets A to a is the operation that deletes all incoming edges to A in the graph, and then conditions on $A = a$ in the resulting probability distribution (which is defined

¹ by the conditional independence structure of the post-surgery graph). We'll use Pearl's do notation
² to denote this operation. $P(\mathbf{X}|\text{do}(A = a))$ is the distribution of \mathbf{X} given $A = a$ under the mutilated
³ graph that results from deleting all edges going into A . Similarly, $\mathbb{E}[\mathbf{X}|\text{do}(A = a)] \triangleq \mathbb{E}_{P(\mathbf{X}|\text{do}(A = a))}[\mathbf{X}]$.
⁴ Thus, we can formalize statements such as "The average effect of receiving drug A " as
⁵

$$\underline{6} \quad \text{ATE} = E[Y|\text{do}(A = 1)] - \mathbb{E}[Y|\text{do}(A = 0)], \quad (36.9)$$

⁷ where ATE stands for Average Treatment Effect.

⁸ For concreteness, consider our running example. We contrast the distribution that results by
⁹ conditioning on A with the distribution that results from intervening on A :

$$\underline{11} \quad P(Y, H, G|A = a) = P(Y|H, G, A = a)P(G)P(H|A = a) \quad (36.10)$$

$$\underline{12} \quad P(Y, H, G|\text{do}(A = a)) = P(Y|H, G, A = a)P(G)P(H) \quad (36.11)$$

¹³The key difference between these two distributions is that the standard conditional distribution
¹⁴describes a population where health consciousness H has the distribution that we observe among
¹⁵individuals with smoking status $A = a$, while the interventional distribution described a population
¹⁶where health consciousness H follows the marginal distribution among all individuals. For example,
¹⁷we would expect $P(H | A = \text{smoker})$ to put more mass on lower values of H than the marginal
¹⁸health consciousness distribution than the marginal distribution $P(H)$, which would also include
¹⁹non-smokers. The intervention distribution thus incorporates a hypothesis of how smoking would
²⁰affect the subpopulation individuals who tend to be too health conscious to smoke in the observed
²¹data.
²²

²³

²⁴ 36.2.3 Identification

²⁵

²⁶A central challenge in causal inference is that many different SCM's can produce identical distributions
²⁷of observed data. This means that, on the basis of observed data alone, we cannot uniquely identify
²⁸the SCM that generated it. This is true no matter how large of a data sample is available to us.

²⁹ For example, consider the setting where there is a treatment A that may or may not have an
³⁰effect on outcome Y , and where both the treatment and outcome are known to be affected by
³¹some *unobserved* common binary cause U . Now, we might be interested in the causal estimand
³² $E[Y|\text{do}(A = 1)]$. In general, we can't learn this quantity from the observed data. The problem
³³is that, we can't tell apart the case where the treatment has a strong effect from the case where
³⁴the treatment has no effect, but $U = 1$ both causes people to tend to be treated and increases the
³⁵probability of a positive outcome. The same observation shows we can't learn the (more complicated)
³⁶interventional distribution $P(Y|\text{do}(A = 1))$ (if we could learn this, then we'd get the average effect
³⁷automatically).

³⁸ Thus, an important part of causal inference is to augment the observed data with knowledge about
³⁹the underlying causal structure of the process under consideration. Often, these assumptions can
⁴⁰narrow the space of SCM's sufficiently so that there is only one value of the causal estimand that is
⁴¹compatible with the observed data. We say that the causal estimand is **identified** or **identifiable**
⁴²under a given set of assumptions if those assumptions are sufficient to provide a unique answer.
⁴³There are many different sets of sufficient conditions that yield identifiable causal effects; we call
⁴⁴each set of sufficient conditions an **identification strategy**.

⁴⁵ Given a set of assumptions about the underlying SCM, the most common way to show that a
⁴⁶causal estimand is identified is by construction. Specifically, if the causal estimand can be written
⁴⁷

entirely in terms of observable probability distributions, then it is identified. We call such a function of observed distributions a **statistical estimand**. Once such a statistical estimand has been recovered, we can then construct and analyze an estimator for that quantity using standard statistical tools. As an example of a statistical estimand, in the SCM above, it can be shown the ATE as defined in Equation (36.9), is equal to the following statistical estimand

$$\text{ATE} \stackrel{(*)}{=} \tau^{\text{ATE}} \triangleq \mathbb{E}[\mathbb{E}[Y|H, A = 1] - \mathbb{E}[Y|H, A = 0]], \quad (36.12)$$

where the equality (*) only holds because of some specific properties of the SCM. Note that the RHS above only involves conditional expectations between observed variables (there are no do operators), so τ^{ATE} is only a function of observable probability distributions.

There are many kinds of assumptions we might make about the SCM governing the process under consideration. For example, the following are assertions we might make about the system in our running example:

1. The probability of developing cancer is additive on the logit scale in A , G , and H (i.e., logistic regression is a well-specified model).
2. For each individual, smoking can never decrease the probability of developing cancer.
3. Whether someone smokes is influenced by their health consciousness H , but not by their genetic predisposition to cancer G .

These assumptions range from strong parametric assumptions fully specifying the form of the SCM equations, to non-parametric assumptions that only specify what the inputs to each equation are, leaving the form fully unspecified. Typically, assumptions that fully specify the parametric form are very strong, and would require far more detailed knowledge of the system under consideration than we actually have. The goal in identification arguments is to find a set of assumptions that are weak enough that they might be plausibly true for the system under consideration, but which are also strong enough to allow for identification of the causal effect.

If we are not willing to make any assumptions about the functional form of the SCM, then our assumptions are just about which variables affect (and do not affect) the other variables. In this sense, such which-affects-which assumptions are minimal. These assumptions are exactly the assumptions captured by writing down a (possibly incomplete) causal DAG, showing which variables are parents of each other variable. The graph may be incomplete because we may not know whether each possible edge is present in the physical system. For example, we might be unsure whether the gene G actually has a causal effect on health consciousness H . It is natural to ask to what extent we can identify causal effects only on the basis of partially specified causal DAGs. It turns out much progress can be made based on such non-parametric assumptions; we discuss this in detail in Section 36.8.

We will also discuss certain assumptions that cannot be encoded in a causal graph, but that are still weaker than assuming that full functional forms are known. For example, we might assume that the outcome is affected additively by the treatment and any confounders, with no interaction terms between them. These weaker assumptions can enable causal identification even when assuming the causal graph alone does not.

It is worth emphasizing that every causal identification strategy relies on assumptions that have some content that cannot be validated in the observed data. This follows directly from the ill-posedness of causal problems: if the assumptions used to identify causal quantities could be validated, that

1 would imply that the causal estimand was identifiable from the observed data alone. However, since
2 we know that there are many values of the causal estimand that are compatible with observed data,
3 it follows that the assumptions in our identification strategy must have unobservable implications.
4

5

6 36.2.4 Counterfactuals and the Causal Hierarchy

7 Structural causal models let us formalize and study a hierarchy of different kinds of query about the
8 system under consideration. The most familiar is observational queries: questions that are purely
9 about statistical associations (e.g., “Are smoking and lung cancer associated in the population this
10 sample was drawn from?”). Next is interventional queries: questions about causal relationships at
11 the population level (e.g., “How much does smoking increase the probability of cancer in a given
12 population?”). The rest of this chapter is focused on the definition, identification, and estimation of
13 interventional queries. Finally, there are counterfactual queries: questions about causal relationships
14 at the level of specific individuals, had something been different (e.g., “Would Alice have developed
15 cancer had she not smoked?”). This causal hierarchy was popularized by [Pea09a, Ch. 1].
16

17 **Interventional queries** concern the prospective effect of an intervention on an outcome; for
18 example, if we intervene and prevent a randomly sampled individual from smoking, what is the
19 probability they develop lung cancer? Ultimately, the probability statement here is about our
20 uncertainty about the “noise” variables ξ in the SCM. These are the unmeasured factors specific to
21 the randomly selected individual. The distribution is determined by the population from which that
22 individual is sampled. Thus, interventional queries are statements about populations. Interventional
23 queries can be written in terms of conditional distributions using do-notation, e.g. $P(Y|\text{do}(A = 0))$.
24 In our example, this represents the distribution of lung cancer outcomes for an individual selected at
25 random and prevented from smoking.

26 **Counterfactual queries** concern how an observed outcome might have been different had an
27 intervention been applied in the past. Counterfactual queries are often framed in terms of attributing
28 a given outcome to a particular cause. For example, would Alice have developed cancer had she not
29 smoked? Did most smokers with lung cancer develop cancer because they smoked? Counterfactual
30 queries are so called because they require a comparison of counterfactual outcomes within individuals.
31 In the formalism of SCM’s, counterfactual outcomes for an individual i are generated by running the
32 same values of ξ_i through differently intervened SCM’s. Counterfactual outcomes are often written
33 in terms of **potential outcomes** notation. In our running smoking example, this would look like:

$$\begin{aligned} \text{34} \quad Y_i(a) &\triangleq f_Y(G_i, H_i, a, \xi_{3,i}). \\ \text{35} \end{aligned} \tag{36.13}$$

36 That is, $Y_i(a)$ is the outcome we would have seen had A been set to a while all of $G_i, H_i, \xi_{3,i}$ were
37 kept fixed.

38 It is important to understand what distinguishes interventional and fundamentally counterfactual
39 queries. Just because a query can be written in terms of potential outcomes does not make it a
40 counterfactual query. For example, the average treatment effect, which is the canonical interventional
41 query, is easy to write in potential outcomes notation:
42

$$\begin{aligned} \text{43} \quad \text{ATE} &= \mathbb{E}[Y_i(1) - Y_i(0)]. \\ \text{44} \end{aligned} \tag{36.14}$$

45 Instead, the key dividing line between counterfactual and interventional queries is whether the query
46 requires knowing the joint distribution of potential outcomes within individuals, or whether marginal
47

1 distributions of potential outcomes across individuals will suffice. An important signature of a
2 counterfactual query is conditioning on the value of one potential outcome. For example, “the lung
3 cancer rate among smokers who developed cancer, had they not smoked” is a counterfactual query,
4 and can be written as:

$$\mathbb{E}[Y_i(0) \mid Y_i(1) = 1, A_i = 1] \quad (36.15)$$

10 Answering this query requires knowing how individual-level cancer outcomes are related (through
11 $\xi_{3,i}$) across the worlds where the each individual i did and did not smoke. Notably, this query cannot
12 be rewritten using do-notation, because it requires a distinction between $Y(0)$ and $Y(1)$ while the
13 ATE can: $\mathbb{E}[Y \mid \text{do}(A = 1)] - \mathbb{E}[Y \mid \text{do}(A = 0)]$.

14 Counterfactual queries require categorically more assumptions for identification than interventional
15 ones. For identifying interventional queries, knowing the DAG structure of an SCM is often sufficient,
16 while for counterfactual queries, some assumptions about the functional forms in the SCM are
17 necessary. This is because only one potential outcome is ever observed for each individual, so the
18 dependence between potential outcomes within individuals is not observable. For example, the data
19 in our running example provide no information on how individual-level smoking and non-smoking
20 cancer risk are related. Thus, answering a question like “Did smokers who developed cancer have lower
21 non-smoking cancer risk than smokers who did not develop cancer?”, requires additional assumptions
22 about how characteristics encoded in ξ_i are translated to cancer outcomes. To answer this question
23 without such assumptions, we would need to observe smokers who developed cancer in the alternate
24 world where they did not smoke. Because they compare how individuals would have turned out under
25 different generating processes, counterfactual queries are often referred to as “cross-world” quantities.
26 On the other hand, interventional queries only require understanding the marginal distributions of
27 potential outcomes $Y_i(0)$ and $Y_i(1)$ across individuals; thus, no cross-world information is necessary
28 at the individual level.

29 We conclude this section by noting that counterfactual outcomes and potential outcomes notation
30 are often conceptually useful, even if they are not used to explicitly answer counterfactual queries.
31 Many causal queries are more intuitive to formalize in terms of potential outcomes. E.g., “Would I
32 have smoked if I was more health conscious?” may be more intuitive than “Would a randomly sampled
33 individual from the same population have smoked had they been subject to an intervention that made
34 them more health concious?”. In fact, some schools of causal inference use potential outcomes, rather
35 than DAGs, as their primary conceptual building block [See IR15]. Causal graphs and potential
36 outcomes both provide ways to formulate interventional queries and causal assumptions. Ultimately,
37 these are mathematically equivalent. Nevertheless, practically, they have different strengths. The
38 main advantage of potential outcomes is that counterfactual statements often map more directly to
39 our mechanistic understanding of the world. This can make it easier to articulate causal desiderata
40 and causal assumptions we may wish to use. On the other hand, the potential outcomes notation
41 does not automatically distinguish between interventional and counterfactual queries. Additionally,
42 causal graphs often give an intuitive and easy way of articulating assumptions about structural
43 causal models involving many variables—potential outcomes get quickly unwieldy. In short: both
44 formalizations have distinct advantages, and those advantages are simply about how easy it is to
45 translate our causal understanding of the world into crisp mathematical assumptions.

36.3 Randomized Control Trials

We now turn to the business of estimating causal effects from data. We begin with **randomized control trials**, which are experiments designed to make the causal concerns as simple as possible.

The simplest situation for causal estimation is when there are no common causes of A and Y . The world is rarely so obliging as to make this the case. However, sometimes we can design an experiment to enforce the no-common-causes structure. In randomized control trials we assign each participant to either the treatment or control group at random. Because random assignment does not depend on any property of the units in the study, there are no causes of treatment assignment, and hence also no common causes of Y and A .

In this case, it's straightforward to see that $P(Y|do(A = a)) = P(Y|a)$. This is essentially by definition of the graph surgery: since A has no parents, the mutilated graph is the same as the original graph. Indeed, the graph surgery definition is chosen to make this true: any sensible formalization of causality should have this identification result.

It is common to use RCTs to study the average treatment effect,

$$\text{ATE} = E[Y|do(A = 1)] - \mathbb{E}[Y|do(A = 0)]. \quad (36.16)$$

This is the expected difference between being assigned treatment and assigned no treatment for a randomly chosen member of the population. It's easy to see that in an RCT this causal quantity is identified as a parameter τ^{RCT} of the observational distribution:

$$\tau^{\text{RCT}} = \mathbb{E}[Y|A = 1] - \mathbb{E}[Y|A = 0].$$

Then, a natural estimator is:

$$\hat{\tau}^{\text{RCT}} \triangleq \frac{1}{n_A} \sum_{i:A_i=1} Y_i - \frac{1}{n - n_A} \sum_{i:A_i=0} Y_i, \quad (36.17)$$

where n_A is the number of units who received treatment. That is, we estimate the average treatment effect as the difference between the average outcome of the treated group and the average outcome of the untreated (control) group. ¹

Randomized control trials are the gold standard for estimating causal effects. This is because we know *by design* that there are no confounders that can produce alternative causal explanations of the data. In particular, the assumption of the triangle DAG—there are no unobserved confounders—is enforced by design. However, there are limitations. Most obviously, randomized control trials are sometimes infeasible to conduct. This could be due to expense, regulatory restrictions, or more fundamental difficulties (e.g., in developmental economics, the response of interest is sometimes collected decades after treatment). Additionally, it may be difficult to ensure that the participants in an RCT are representative of the population where the treatment will be deployed. For instance, participants in drug trials may skew younger and poorer than the population of patients who will ultimately take the drug.

⁴⁴ 1. There is a literature on efficient estimation of causal effects in RCT's going back to Fisher [Fis25] that employ more sophisticated estimators. See also Lin [Lin13a] and Bloniarz et al. [Blo+16] for more modern treatments.

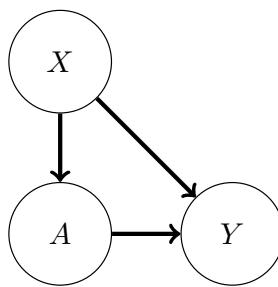


Figure 36.3: A causal DAG illustrating a situation where treatment A and outcome Y are both influenced by observed confounders X .

36.4 Confounder Adjustment

We now turn to the problem of estimating causal effects using observational (i.e., not experimental) data. The most common application of causal inference is estimating the average treatment effect (ATE) of an intervention. The ATE is also commonly called the **average causal effect**, or ACE. Here, we focus on the important special case where the treatment A is binary, and we observe the outcome Y as well as a set of common causes X that influence both A and Y .

36.4.1 Causal Estimand, Statistical Estimand, and Identification

Consider a problem where we observe treatment A , outcome Y , and covariates X , which are drawn i.i.d. from some unknown distribution P . We wish to learn the average treatment effect: the expected difference between being assigned treatment and assigned no treatment for a randomly chosen member of the population. Following the discussion in the introduction, there are three steps to learning this quantity: mathematically formalize the causal estimand, give conditions for the causal estimand to be identified as a statistical estimand, and, finally, estimate this statistical estimand from data. We now turn to the first two steps.

The average treatment effect is defined to be the difference between the average outcome if we intervened and set A to be 0, versus the average outcome if we intervened and set A to be 1. Using the do notation, we can write this formally as

$$\text{ATE} = \mathbb{E}[Y|\text{do}(A = 1)] - \mathbb{E}[Y|\text{do}(A = 0)]. \quad (36.18)$$

The next step is to articulate sufficient conditions for the ATE to be identified as a statistical estimand (a parameter of distribution P). The key issue is the possible presence of **confounders**. Confounders are “common cause” variables that affect both the treatment and outcome. When there are confounding variables in observed data, the sub-population of people who are observed to have received one level of the treatment A will differ from the rest of the population in ways that are relevant to their observed Y . For example, there is a strong positive association between horseback riding in childhood (treatment) and healthiness as an adult (outcome) [RB16]. However, both of these quantities are influenced by wealth X . The population of people who rode horses as children ($A = 1$) is wealthier than the population of people who did not. Accordingly, horseback-riding population

¹
² will have better health outcomes even if there is no actual causal benefit of horseback riding for adult
³ health.

⁴ We'll express the assumptions required for causal identification in the form of a causal DAG.
⁵ Namely, we consider the simple triangle DAG in Figure 36.3, where the treatment and outcome
⁶ are influenced by *observed* confounders X . It turns out that the assumption encoded by this DAG
⁷ suffices for identification. To understand why this is so, recall that the target causal effect is defined
⁸ according to the distribution we would see if the edge from X to A was removed (that's the meaning
⁹ of do). The key insight is that because the intervention only modifies the relationship between X
¹⁰ and A , the structural equation that generates outcomes Y given X and A , illustrated in Figure 36.3
¹¹ as the $A \rightarrow Y \leftarrow X$, is the same even after the $X \rightarrow Y$ edge is removed. For example, we might
¹² believe that the physiological processes by which smoking status A and confounders X produce
¹³ lung cancer Y remain the same, regardless of how the decision to smoke or not smoke was made.
¹⁴ Secondly, because the intervention does not change the composition of the population, we would also
¹⁵ expect the distribution of background characteristics X to be the same between the observational
¹⁶ and intervened processes.

¹⁷ With these insights about invariances between observed and interventional data, we can derive a
¹⁸ statistical estimand for the ATE as follows.

¹⁹ **Theorem 2** (Adjustment with No Unobserved Confounders). *We observe $A, Y, X \sim P$. Suppose
²⁰ that
²¹*

²² 1. (*Confounders observed*) *The data obeys the causal structure in Figure 36.3. In particular, X
²³ contains all common causes of A and Y and no variable in X is caused by A or Y .*

²⁴ 2. (*Overlap*) $0 < P(A = 1|X = x) < 1$ for all values of x . That is, there are no individuals for whom
²⁵ *treatment is always or never assigned*.

²⁷ Then, the average treatment effect is identified as $\text{ATE} = \tau$, where
²⁸

$$\tau = \mathbb{E}[\mathbb{E}[Y|A = 1, X]] - \mathbb{E}[\mathbb{E}[Y|A = 0, X]]. \quad (36.19)$$

³¹ *Proof.* First, we expand the ATE using the tower property of expectation, conditioning on X . Then,
³² we apply the invariances discussed above:

$$\text{ATE} = \mathbb{E}[Y|\text{do}(A = 1)] - \mathbb{E}[Y|\text{do}(A = 0)] \quad (36.20)$$

$$= \mathbb{E}[\mathbb{E}[Y|\text{do}(A = 1), X]] - \mathbb{E}[\mathbb{E}[Y|\text{do}(A = 0), X]] \quad (36.21)$$

$$= \mathbb{E}[\mathbb{E}[Y|A = 1, X]] - \mathbb{E}[\mathbb{E}[Y|A = 0, X]] \quad (36.22)$$

³⁸ The final equality is the key to passing from a causal to observational quantity. This follows because,
³⁹ from the causal graph, the conditional distribution of Y given A, X is the same in both the original
⁴⁰ graph and in the mutilated graph created by removing the edge from X to A . This mutilated graph
⁴¹ defines $P(Y|\text{do}(A = 1), X)$, so the equality holds.

⁴² The condition that $0 < P(A = 1|X = x) < 1$ is required for the first equality (the tower property)
⁴³ to be well defined. \square

⁴⁴ Note that Equation (36.19) is a function of only conditional expectations and distributions that
⁴⁵ appear in the observed data distribution (in particular, it contains no “do” operators). Thus, if we
⁴⁶

can fully characterize the observed data distribution P , we can map that distribution to a unique ATE.

It is useful to note how τ differs from the naive estimand $\mathbb{E}[Y|A = 1] - \mathbb{E}[Y|A = 0]$ that just reports the treatment-outcome association without adjusting for confounding. The comparison is especially clear when we write out the outer expectation in τ explicitly as an integral over X :

$$\tau = \int \mathbb{E}[Y | A = 1, X] P(X) dX - \int \mathbb{E}[Y | A = 0, X] P(X) dX \quad (36.23)$$

We can write the naive estimand in a similar form by applying the tower property of expectation:

$$\mathbb{E}[Y | A = 1] - \mathbb{E}[Y | A = 0] = \int \mathbb{E}[Y | A = 1, X] P(X | A = 1) dX - \int \mathbb{E}[Y | A = 0, X] P(X | A = 0) dX \quad (36.24)$$

The key difference is the probability distribution over X that is being integrated over. The observational difference in means integrates over the distinct conditional distributions of confounders X , depending on the value of A . On the other hand, in the ATE estimand τ , we integrate over the same distribution $P(X)$ for both levels of the treatment.

Overlap In addition to the assumption on the causal structure, identification requires that there is sufficient random variation in how treatments are assigned.

Definition 1. A distribution P on A, X satisfies **overlap** if $0 < P(A = 1|x) < 1$ for all x . It satisfies **strict overlap** if $\epsilon < P(A = 1|x) < 1 - \epsilon$ for all x and some $\epsilon > 0$.

Overlap is the requirement that any unit could have either received the treatment or not.

To see the necessity of overlap, consider estimating the effectiveness of a drug in a study where patient sex is a confounder, but the drug was only ever prescribed to male patients. Then, conditional on a patient being female, we would know that patient was assigned to control. Without further assumptions, it's impossible to know the effect of the drug on a population with female patients, because there would be no data to inform the expected outcome for treated female patients, that is, $\mathbb{E}[Y | A = 1, X = \text{female}]$. In this case, the statistical estimand equation 36.19 would not be identifiable. In the same vein, strict overlap ensures that the conditional distributions at each stratum of X can be estimated in finite samples.

Overlap can be particularly limiting in settings where we are adjusting for a large number of covariates (in an effort to satisfy no unobserved confounding). Then, certain combinations of traits may be very highly predictive of treatment assignment, even if individual traits are not. E.g., male patients over age 70 with BMI greater than 25 are very rarely assigned the drug. If such groups represent a significant fraction of the target population, or have significantly different treatment effects, then this issue can be problematic. In this case, the strict overlap assumption puts very strong restrictions on observational studies: for an observational study to satisfy overlap, most dimensions of the confounders X would need to closely mimic the balance we would expect in an RCT [D'A+21].

36.4.2 ATE Estimation with Observed Confounders

We now return to estimating the ATE using observed—i.e., not experimental—data. We've shown that in the case where we observe all common causes of the treatment and outcome, the ATE is

¹ causally identified with a statistical estimand τ . We now consider several strategies for estimating this
² quantity using a finite data sample. Broadly, these techniques are known as backdoor adjustment.²
³ Recall that the defining characteristic of a confounding variable is that it affects both treatment
⁴ and outcome. Thus, an adjustment strategy may aim to account for the influence of confounders on
⁵ the observed outcome, the influence of confounders on treatment, or both. We discuss each of these
⁶ strategies in turn.
⁷

⁸

⁹ 36.4.2.1 Outcome Model Adjustment

¹⁰ We begin with an approach to covariate adjustment that relies on modeling the conditional expectation
¹¹ of the outcome Y given treatment A and confounders X . This strategy is often referred to as g-
¹² computation or outcome adjustment.³ To begin, we define
¹³

¹⁴ **Definition 2.** *The conditional expected outcome is the function Q given by*

$$\underline{16} \quad Q(a, x) = \mathbb{E}[Y|A = a, X = x]. \quad (36.25)$$

¹⁷ Substituting this definition into the definition of our estimand τ , Equation (36.19), we have
¹⁸ $\tau = \mathbb{E}[Q(1, x) - Q(0, x)]$. This suggests a procedure for estimating τ : fit a model \hat{Q} for Q and then
¹⁹ report
²⁰

$$\underline{21} \quad \hat{\tau}^Q \triangleq \frac{1}{n} \sum_i \hat{Q}(1, x_i) - \hat{Q}(0, x_i). \quad (36.26)$$

²⁴ To fit \hat{Q} , recall that $E[Y|a, x] = \operatorname{argmin}_Q \mathbb{E}[(Y - Q(A, X))^2]$. That is, the minimizer (among all
²⁵ functions) of the squared loss risk is the conditional expected outcome.⁴ So, to approximate Q , we
²⁶ simply use mean squared error to fit a predictor that predicts Y from A and X .

²⁷ The estimation procedure takes several steps. We first fit a model \hat{Q} to predict Y . Then, for each
²⁸ unit i , we predict that unit's outcome had they received treatment $\hat{Q}(1, x_i)$ and we predict their
²⁹ outcome had they not received treatment $\hat{Q}(0, x_i)$.⁵ If the unit actually did receive treatment ($a_i = 1$)
³⁰ then $\hat{Q}(0, x_i)$ is our guess about what would have happened in the counterfactual case that they
³¹ did not. The estimated expected gain from treatment for this individual is $\hat{Q}(1, x_i) - \hat{Q}(0, x_i)$ —the
³² difference in expected outcome between being treated and not treated. Finally, we estimate the outer
³³ expectation with respect to $P(X)$ —the true population distribution of the confounders—using the
³⁴ empirical distribution $\hat{P}(X) = 1/n \sum_i \delta_{x_i}$. In effect, this means we substitute the expectation (over
³⁵ an unknown distribution) by an average over the observed data.
³⁶

³⁷ **Linear regression** It's worth saying something more about the special case where Q is modeled
³⁸ as a linear function of both the treatment and all the covariates. That is, the case where we assume
³⁹ the identification conditions of Theorem 2 and we additionally assume that the true, causal law
⁴⁰

⁴¹ 2. As we discuss in Section 36.8, this backdoor adjustment references the estimand returned by the do-calculus to
⁴² eliminate confounding from a backdoor path. This also generalizes the approaches discussed here to some cases where
⁴³ we do not observe all common causes.

⁴⁴ 3. The “g” stands for generalized, for now-inscrutable historical reasons [Rob86].

⁴⁵ 4. To be precise, this definition applies when X and Y are square-integrable, and the minimization taken over measurable
⁴⁵ functions.

⁴⁶ 5. this interpretation is justified by the same conditions as Theorem 2

⁴⁷

(the SCM) governing Y yields: $Q(A, X) = \mathbb{E}[Y|A, X] = \mathbb{E}[f_Y(A, X, \xi)|A, X] = \beta_0 + \beta_A A + \beta_X X$. Plugging in, we see that $Q(1, X) - Q(0, X) = \beta_A$ (and so also $\tau = \beta_A$). Then, the estimator for the average treatment effect reduces to the estimator for the regression coefficient β_A . This “fit linear regression and report the regression coefficient” remains a common way of estimating the association between two variables in practice. The expected-outcome-adjustment procedure here may be viewed as a generalization of this procedure that removes the linear parametric assumption.

36.4.2.2 Propensity Score Adjustment

Outcome model adjustment relies on modeling the relationship between the confounders and the outcome. A popular alternative is to model the relationship between the confounders and the treatment. This strategy adjusts for confounding by directly addressing sampling bias in the treated and control groups. This bias arises from the relationship between the confounders and the treatment. Intuitively, the effect of confounding may be viewed as due to the difference between $P(X|A = 1)$ and $P(X|A = 0)$ —e.g., the population of people who rode horses as children is wealthier than the population of people who did not. When we observe all confounding variables X , this degree of over- or under-representation can be adjusted away by reweighting samples such that the confounders X have the same distribution in the treated and control groups. When the confounders are balanced between the two groups, then any differences between them must be attributable to the treatment.

A key quantity for balancing treatment and control groups is the **propensity score**, which summarises the relationship between confounders and treatment.

Definition 3. *The propensity score is the function g given by $g(x) = P(A = 1|X = x)$.*

To make use of the propensity score in adjustment, we first rewrite the estimand τ in a suggestive form, leveraging the fact that $A \in \{0, 1\}$:

$$\tau = \mathbb{E}\left[\frac{YA}{g(X)} - \frac{Y(1-A)}{1-g(X)}\right]. \quad (36.27)$$

This identity can be verified by noting that $\mathbb{E}[YA|X] = \mathbb{E}[Y|A = 1, X]P(A = 1|X) + 0$, rearranging for $\mathbb{E}[Y|A = 1, X]$, doing the same for $\mathbb{E}[Y|A = 0, X]$, and substituting in to Equation (36.19). Note that the identity is just a mathematical fact about the statistical estimand—it does not rely on any causal assumptions, and holds whether or not τ can be interpreted as a causal effect.

This expression suggests the **inverse probability of treatment weighted estimator**, or IPTW estimator:

$$\hat{\tau}^{\text{IPTW}} \triangleq \frac{1}{n} \sum_i \frac{Y_i A_i}{\hat{g}(X_i)} - \frac{Y_i(1-A_i)}{1-\hat{g}(X_i)}. \quad (36.28)$$

Here, \hat{g} is an estimate of the propensity score function. Recall from Section 14.2.1 that if a model is well-specified and the loss function is a proper scoring rule then risk minimizer $g^* = \operatorname{argmin}_g \mathbb{E}[L(A, g(X))]$ will be $g^*(X) = P(A = 1|X)$. That is, we can estimate the propensity score by fitting a model that predicts A from X . Cross-entropy and squared loss are both proper scoring rules, so we may use standard supervised learning methods.

In summary, the procedure is to estimate the propensity score function (with machine learning), and then to plug the estimated propensity scores $\hat{g}(x_i)$ into Equation (36.28). The IPTW estimator

1 computes a difference of weighted averages between the treated and untreated group. The effect is to
2 upweight the outcomes of units who were unlikely to be treated but who nevertheless actually, by
3 chance, received treatment (and similarly for untreated). Intuitively, such units are typical for the
4 untreated population. So, their outcomes under treatment are informative about what would have
5 happened had a typical untreated unit received treatment.

6 A word of warning is in order. Although the IPTW is asymptotically valid and popular in practice,
7 it can be very unstable in finite samples. If estimated propensity scores are extreme for some values
8 of x (that is, very close to 0 or 1), then the corresponding IPTW weights can be very large, resulting
9 in a high-variance estimator. In some cases, this instability can be mitigated by instead using the
10 Hajek version of the estimator.

$$\hat{\tau}^{\text{h-IPTW}} \triangleq \sum_i Y_i A_i \frac{1/\hat{g}(X_i)}{\sum_i A_i/\hat{g}(X_i)} - \sum_i Y_i (1 - A_i) \frac{1/(1-\hat{g}(X_i))}{\sum_i (1-A_i)/(1-\hat{g}(X_i))}. \quad (36.29)$$

11 However, extreme weights can persist even after self-normalization, either because there are truly
12 strata of X where treatment assignment is highly imbalanced, or because the propensity score
13 estimation method has overfit. In such cases, it is common to apply heuristics such as weight clipping.
14 See Khan and Ugander [KU21] for a longer discussion of inverse-propensity type estimators,
15 including some practical improvements.

16

17 36.4.2.3 Double Machine Learning

18 We have seen how to estimate the average treatment effect using either the relationship between
19 confounders and outcome, or the relationship between confounders and treatment. In each case,
20 we follow a two step estimation procedure. First, we fit models for the expected outcome or the
21 propensity score. Second, we plug these fitted models into a downstream estimator of the effect.

22 Unsurprisingly, the quality of the estimate of τ depends on the quality of the estimates \hat{Q} or \hat{g} . This
23 is problematic because Q and g may be complex functions that require large numbers of samples to
24 estimate. Even though we're only interested in the 1-dimensional parameter τ , the naive estimators
25 described thus far can have very slow rates of convergence. This leads to unreliable inference or very
26 large confidence intervals.

27 Remarkably, there are strategies for combining Q and g in estimators that, in principle, do better
28 than using either Q or g alone. The **Augmented Inverse Probability of Treatment Weighted**
29 **Estimator (AIPTW)** is one such estimator. It is defined as

$$\hat{\tau}^{\text{AIPTW}} \triangleq \frac{1}{n} \sum_i \hat{Q}(1, X_i) - \hat{Q}(0, X_i) + A_i \frac{Y_i - \hat{Q}(1, X_i)}{\hat{g}(x_i)} - (1 - A_i) \frac{Y_i - \hat{Q}(0, X_i)}{1 - \hat{g}(X_i)}. \quad (36.30)$$

30

31 That is, $\hat{\tau}^{\text{AIPTW}}$ is the outcome adjustment estimator plus a stabilization term that depends on
32 the propensity score. This estimator is a particular case of a broader class of estimators that are
33 referred to as **semi-parametrically efficient** or **double machine-learning** estimators [Che+17e;
34 Che+17d]. We'll use the later terminology here.

35 We now turn to understanding the sense in which double machine learning estimators are robust
36 to misestimation of the **nuisance functions** Q and g . To this end, we define the **influence curve**
37

¹ of τ to be the function ϕ defined by⁶

$$\phi(X_i, A_i, Y_i; Q, g, \tau) \triangleq Q(1, X_i) - Q(0, X_i) + A_i \frac{Y_i - Q(1, X_i)}{g(x_i)} - (1 - A_i) \frac{Y_i - Q(0, X_i)}{1 - g(X_i)} - \tau. \quad (36.31)$$

⁷ By design, $\hat{\tau}^{\text{AIPTW}} - \tau = \frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \tau)$. We begin by considering what would happen if we
⁸ simply knew Q and g , and didn't have to estimate them. In this case, the estimator would be
⁹ $\hat{\tau}^{\text{ideal}} = \frac{1}{n} \sum_i \phi(\mathbf{X}_i; Q, g, \tau)$ and, by the central limit theorem, we would have:

$$\sqrt{n}(\hat{\tau}^{\text{ideal}} - \tau) \xrightarrow{d} \text{Normal}(0, \mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]). \quad (36.32)$$

¹⁰ This result characterizes the estimation uncertainty in the best possible case. If we knew Q and g ,
¹¹ we could rely on this result for, e.g., finding confidence intervals for our estimate.

¹² The question is: what happens when Q and g need to be estimated? For general estimators and
¹³ nuisance function models, we don't expect the \sqrt{n} -rate of Equation (36.32) to hold. For instance,
¹⁴ $\sqrt{n}(\hat{\tau}^Q - \tau)$ only converges if $\sqrt{n}\mathbb{E}[(\hat{Q} - Q)^2]^{\frac{1}{2}} \rightarrow 0$. That is, for the naive estimator we only get the
¹⁵ \sqrt{n} rate for estimating τ if we can also estimate Q at the \sqrt{n} rate—a much harder task! This is the
¹⁶ issue that the double machine learning estimator helps with.

¹⁷ To understand how, we decompose the error in estimating τ as follows:

$$\sqrt{n}(\hat{\tau}^{\text{AIPTW}} - \tau) \quad (36.33)$$

$$= \frac{1}{\sqrt{n}} \sum_i \phi(\mathbf{X}_i; Q, g, \tau) \quad (36.34)$$

$$+ \frac{1}{\sqrt{n}} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}_i; Q, g, \tau) - \mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}; Q, g, \tau)] \quad (36.35)$$

$$+ \sqrt{n}\mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}; Q, g, \tau)] \quad (36.36)$$

²⁹ We recognize the first term, Equation (36.34), as $\sqrt{n}(\hat{\tau}^{\text{ideal}} - \tau)$, the estimation error in the optimal
³⁰ case where we know Q and g . Ideally, we'd like the error of $\hat{\tau}^{\text{AIPTW}}$ to be asymptotically equal to
³¹ this ideal case—which will happen if the other two terms go to 0.

³² The second term, Equation (36.35), is a penalty we pay for using the same data to estimate Q, g
³³ and to compute τ . For many model classes, it can be shown that such “empirical process” terms go
³⁴ to 0. This can also be guaranteed in general by using different data for fitting the nuisance functions
³⁵ and for computing the estimator (see the next section).

³⁶ The third term, Equation (36.36), captures the penalty we pay for misestimating the nuisance
³⁷ functions. This is where the particular form of the AIPTW is key. With a little algebra, we can show
³⁸ that

$$\mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}) - \phi(\mathbf{X}; Q, g)] = \mathbb{E}\left[\frac{1}{g(X)}(\hat{g}(X) - g(X))(\hat{Q}(1, X) - Q(1, X))\right] \quad (36.37)$$

$$+ \frac{1}{1 - g(X)}(\hat{g}(X) - g(X))(\hat{Q}(0, X) - Q(0, X)). \quad (36.38)$$

⁴⁴ 6. Influence curves are the foundation of what follows, and the key to generalizing the analysis beyond the ATE.
⁴⁵ Unfortunately, going into the general mathematics would require a major digression, so we omit it. However, see
⁴⁶ references at the end of the chapter for some pointers to the relevant literature.

¹ The important point is that estimation errors of Q and g are multiplied together. Using the Cauchy-Schwarz inequality, we find that $\sqrt{n}\mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}) - \phi(\mathbf{X}; Q, g)] \rightarrow 0$ as long as $\sqrt{n} \max_a \mathbb{E}[(\hat{Q}(a, X) - Q(a, X))^2]^{\frac{1}{2}} \mathbb{E}[(\hat{g}(X) - g(X))^2]^{\frac{1}{2}} \rightarrow 0$. That is, the misestimation penalty will vanish so long as the product of the misestimation errors is $o(\sqrt{n})$. For example, this means that that τ can be estimated at the (optimal) \sqrt{n} rate even when the estimation error of each of Q and g only decreases as $o(n^{-1/4})$.

⁷ The upshot here is that the double machine learning estimator has the special property that the weak condition $\sqrt{n}\mathbb{E}[(\hat{Q}(T, X) - Q(T, X))^2]\mathbb{E}[(\hat{g}(X) - g(X))^2] \rightarrow 0$ suffices to imply that

$$\sqrt{n}(\hat{\tau}^{\text{AIPTW}} - \tau) \xrightarrow{d} \text{Normal}(0, \mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]) \quad (36.39)$$

¹¹(though strictly speaking this requires some additional technical conditions we haven't discussed).
¹²This is *not* true for the earlier estimators we discussed, which require a much faster rate of convergence
¹³for the nuisance function estimation.

¹⁴ The AIPTW estimator has two further nice properties that are worth mentioning. First, it is
¹⁵**non-parametrically efficient**. This means that this estimator has the smallest possible variance
¹⁶of any estimator that does not make parametric assumptions; namely, $\mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]$. This means,
¹⁷for example, that this estimator yields the smallest confidence intervals of any approach that does not
¹⁸rely on parametric assumptions. Second, it is **doubly robust**: the estimator is consistent (converges
¹⁹to the true τ as $n \rightarrow \infty$) as long as at least one of either \hat{Q} or \hat{g} is consistent.
²⁰

²¹36.4.2.4 Cross Fitting

²³The term Equation (36.35) in the error decomposition above is the penalty we pay for reusing the
²⁴same data to both fit Q, g and to compute the estimator. For many choices of model for Q, g , this
²⁵term goes to 0 quickly as n gets large and we achieve the (best case) \sqrt{n} error rate. However, this
²⁶property doesn't always hold.

²⁷ As an alternative, we can always randomly split the available data and use one part for model fitting,
²⁸and the other to compute the estimator. Effectively, this means the nuisance function estimation and
²⁹estimator computation are done using independent samples. It can then be shown that the reuse
³⁰penalty will vanish. However, this comes at the price of reducing the amount of data available for
³¹each of nuisance function estimation and estimator computation.

³² This strategy can be improved upon by a **cross fitting** approach. We divide the data into K
³³folds. For each fold j we use the other $K - 1$ folds to fit the nuisance function models $\hat{Q}^{-j}, \hat{g}^{-j}$.
³⁴Then, for each datapoint i in fold j , we take $\hat{Q}(a_i, x_i) = \hat{Q}^{-j}(a_i, x_i)$ and $\hat{g}(x_i) = \hat{g}^{-j}(x_i)$. That is,
³⁵the estimated conditional outcomes and propensity score for each datapoint are predictions from a
³⁶model that was not trained on that datapoint. Then, we estimate τ by plugging $\{\hat{Q}(a_i, x_i), \hat{g}(x_i)\}_i$
³⁷into Equation (36.30). It can be shown that this cross fitting procedure has the same asymptotic
³⁸guarantee—the central limit theorem at the \sqrt{n} rate—as described above.
³⁹

⁴⁰36.4.3 Uncertainty Quantification

⁴¹In addition to the point estimate $\hat{\tau}$ of the average treatment effect, we'd also like to report a measure
⁴³of the uncertainty in our estimate. For example, in the form of a confidence interval. The asymptotic
⁴⁴normality of $\sqrt{n}\hat{\tau}$ (Equation (36.39)) provides a means for this quantification. Namely, we could
⁴⁵base confidence intervals and similar on the limiting variance $\mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]$. Of course, we don't
⁴⁶actually know any of Q, g , or τ . However, it turns out that it suffices to estimate the asymptotic
⁴⁷

variance with $\frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2$ [Che+17e]. That is, we can estimate the uncertainty by simply plugging in our fitted nuisance models and our point estimate of τ into

$$\hat{\mathbb{V}}[\hat{\tau}] = 1/n \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2. \quad (36.40)$$

This estimated variance can then be used to compute confidence intervals in the usual manner. E.g., we'd report a 95% confidence interval for τ as $\hat{\tau} \pm 1.96\sqrt{\hat{\mathbb{V}}[\hat{\tau}]/n}$.

Alternatively, we could quantify the uncertainty by bootstrapping. Note, however, that this would require refitting the nuisance functions with each bootstrap model. Depending on the model and data, this can be prohibitively computationally expensive.

36.4.4 Matching

One particularly popular approach to adjustment-based causal estimation is **matching**. Intuitively, the idea is to match each treated to unit to an untreated unit that has the same (or at least similar) values of the confounding variables and then compare the observed outcomes of the treated unit and its matched control. If we match on the full set of common causes, then the difference in outcomes is, intuitively, a noisy estimate of the effect the treatment had on that treated unit. We'll now build this up a bit more carefully. In the process we'll see that matching can be understood as, essentially, a particular kind of outcome model adjustment.

For simplicity, consider the case where X is a discrete random variable. Define \mathcal{A}_x to be the set of treated units with covariate value x , and \mathcal{C}_x to be the set of untreated units with covariate value x . In this case, the matching estimator is:

$$\hat{\tau}^{\text{matching}} = \sum_x \hat{P}(x) \left(\frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} Y_i - \frac{1}{|\mathcal{C}_x|} \sum_{j \in \mathcal{C}_x} Y_j \right), \quad (36.41)$$

where $\hat{P}(x)$ is an estimator of $P(X = x)$ —e.g., the fraction of units with $X = x$. Now, we can rewrite $Y_i = Q(A_i, X_i) + \xi_i$ where ξ_i is a unit-specific noise term defined by the equation. In particular, we have that $\mathbb{E}[\xi_i | A_i, X_i] = 0$. Substituting this in, we have:

$$\hat{\tau}^{\text{matching}} = \sum_x \hat{P}(x) (Q(1, x) - Q(0, x)) + \sum_x \frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} \xi_i - \frac{1}{|\mathcal{C}_x|} \sum_{j \in \mathcal{C}_x} \xi_j. \quad (36.42)$$

We can recognize the first term as an estimator of usual target parameter τ (it will be equal to τ if $\hat{P}(x) = P(x)$). The second term is a difference of averages of random variables with expectation 0, and so each term will converge to 0 as long as $|\mathcal{A}_x|$ and $|\mathcal{C}_x|$ each go to infinity as we see more and more data. Thus, we see that the matching estimator is a particular way of estimating the parameter τ . The procedure can be extended to continuous covariates by introducing some notion of values of X being close, and then matching close treatment and control variables.

There are two points we should emphasize here. First, notice that the argument here has nothing to do with causal identification. Matching is a particular technique for estimating the observational parameter τ . Whether or not τ can be interpreted as an average treatment effect is determined by the conditions of Theorem 2—the particular estimation strategy doesn't say anything about this. Second, notice that in essence matching amounts to a particular choice of model for \hat{Q} . Namely,

$\frac{1}{2}$ $\hat{Q}(1, x) = \frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} Y_i$ and similarly for $\hat{Q}(0, x)$. That is, we estimate the conditional expected
 $\frac{3}{4}$ outcome as a sample mean over units with the same covariate value. Whether this is a good idea
 $\frac{5}{6}$ depends on the quality of our model for Q . In situations where better models are possible (e.g., a
 $\frac{7}{8}$ machine-learning model fits the data well), we might expect to get a more accurate estimate by using
 $\frac{9}{10}$ the conditional expected outcome predictor directly.

$\frac{11}{12}$ There is another important case we mention in passing. In general, when using adjustment based
 $\frac{13}{14}$ identification, it suffices to adjust for any function $\phi(X)$ of X such that $A \perp\!\!\!\perp X | \phi(X)$. To see that
 $\frac{15}{16}$ adjusting for only $\phi(X)$ suffices, first notice that $g(X) = P(A = 1 | X) = P(A = 1 | \phi(X))$ only depends
 $\frac{17}{18}$ on $\phi(X)$, and then recall that can write the target parameter as $\tau = \mathbb{E}\left[\frac{YA}{g(X)} - \frac{Y(1-A)}{1-g(X)}\right]$, whence
 $\frac{20}{21}$ τ only depends on X through $g(X)$. That is: replacing X by a reduced version $\phi(X)$ such that
 $\frac{22}{23}$ $g(X) = P(A = 1 | \phi(X))$ can't make any difference to τ . Indeed, the most popular choice of $\phi(X)$ is
 $\frac{25}{26}$ the propensity score itself, $\phi(X) = g(X)$. This leads to **propensity score matching**, a two step
 $\frac{27}{28}$ procedure where we first fit a model for the propensity score, and then run matching based on the
 $\frac{30}{31}$ estimated propensity score values for each unit. Again, this is just a particular estimation procedure
 $\frac{33}{34}$ for the observational parameter τ , and says nothing about whether it's valid to interpret τ as a
 $\frac{36}{37}$ causal effect.

$\frac{38}{39}$

$\frac{40}{41}$ 36.4.5 Practical Considerations and Procedures

$\frac{43}{44}$ when performing causal analysis, many issues can arise in practice, some of which we discuss below.

$\frac{51}{52}$ 36.4.5.1 What to adjust for

$\frac{54}{55}$ Choosing which variables to adjust for is a key detail in estimating causal effects using covariate
 $\frac{57}{58}$ adjustment. The criterion is clear when one has a full causal graph relating A , Y , and all covariates
 $\frac{60}{61}$ X to each other. Namely, adjust for all variables that are actually causal parents of A and Y . In
 $\frac{63}{64}$ fact, with access to the full graph, this criterion can be generalized somewhat—see Section 36.8.

$\frac{66}{67}$ In practice, we often don't actually know the full causal graph relating all of our variables.
 $\frac{69}{70}$ As a result, it is common to apply simple heuristics to determine which variables to adjust for.
 $\frac{72}{73}$ Unfortunately, these heuristics have serious limitations. However, exploring these is instructive.

$\frac{75}{76}$ A key condition in Theorem 2 is that the covariates X that we adjust for must include all the
 $\frac{78}{79}$ common causes. In the absence of a full causal graph, it is tempting to condition on as many observed
 $\frac{81}{82}$ variables as possible to try to ensure this condition holds. However, this can be problematic. For
 $\frac{84}{85}$ instance, suppose that M is a mediator of the effect of A on Y —i.e., M lies on one of the directed
 $\frac{87}{88}$ paths between A and Y . Then, conditioning on M will block this path, removing some of the causal
 $\frac{90}{91}$ effect. Note that this does not always result in an attenuated, or smaller-magnitude, effect estimate.
 $\frac{93}{94}$ The effect through a given mediator may run in the opposite direction of other causal pathways
 $\frac{96}{97}$ from the treatment; thus conditioning on a mediator can inflate or even flip the sign of a treatment
 $\frac{99}{100}$ effect. Alternatively, if C is a collider between A and Y —a variable that is caused by both—then
 $\frac{103}{104}$ conditioning on C will induce an extra statistical dependency between A and Y .

$\frac{106}{107}$ Both pitfalls of the “condition on everything” heuristic discussed above both involve conditioning
 $\frac{109}{110}$ on variables that are downstream of the treatment A . A natural response is to this is to limit
 $\frac{112}{113}$ conditioning to all pre-treatment variables, or those that are causally upstream of the treatment.
 $\frac{115}{116}$ Importantly, if there is a valid adjustment set in the observed covariates X , then there will also be a
 $\frac{118}{119}$ valid adjustment set among the pre-treatment covariates. This is because any open backdoor path
 $\frac{121}{122}$

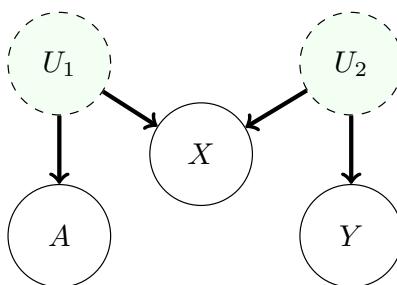


Figure 36.4: The M-bias causal graph. Here, A and Y are not confounded. However, conditioning on the covariate X opens a backdoor path, passing through U_1 and U_2 (because X is a collider). Thus, adjusting for X creates bias. This is true even though X need not be a pre-treatment variable.

between A and Y must include a parent of A , and the set of pre-treatment covariates includes these parents. However, it is still possible that conditioning on the full set of pre-treatment variables can induce new backdoor paths between A and Y through colliders. In particular, if there is a covariate D that is separately confounded with the treatment A and the outcome Y then D is a collider, and conditioning on D opens a new backdoor path. This phenomenon is known as m-bias because of the shape of the graph [Pea09c], see Figure 36.4.

A practical refinement of the pre-treatment variable heuristic is given in VanderWeele TJ [VT11]. Their heuristic suggests conditioning on all pre-treatment variables that are causes of the treatment, outcome, or both. The essential qualifier in this heuristic is that the variable is causally upstream of treatment and/or outcome. This eliminates the possibility of conditioning on covariates that are only confounded with treatment and outcome, avoiding m-bias. Notably, this heuristic requires more causal knowledge than the above heuristics, but does not require detailed knowledge of how different covariates are causally related to each other.

The VanderWeele TJ [VT11] criterion is a useful rule of thumb, but other practical considerations often arise. For example, if one has more knowledge about the causal structure among covariates, it is possible to optimize adjustment sets to minimize the variance of the resulting estimator [RS20]. One important example of reducing variance by pruning adjustment sets is the exclusion of variables that are known to only be a parent of the treatment, and not of the outcome (so called instruments, as discussed in Section 36.5).

Finally, adjustment set selection criteria operate under the assumption that there actually exists a valid adjustment set among observed covariates. When there is no set of observed covariates in X that block all backdoor paths, then any adjusted estimate will be biased. Importantly, in this case, the bias does not necessarily decrease as one conditions on more variables. For example, conditioning on an instrumental variable often results in an estimate that has higher bias, in addition to the higher variance discussed above. This phenomenon is known as bias amplification or z-bias; see Section 36.7.2. A general rule of thumb is that variables that explain away much more variation in the treatment than in the outcome can potentially amplify bias, and should be treated with caution.

1
2 **36.4.5.2 Overlap**

3 Recall that in addition to no-unobserved-confounders, identification of the average treatment effect
4 requires overlap: the condition that $0 < P(A = 1|x) < 1$ for the population distribution P . With
5 infinite data, any amount of overlap will suffice for estimating the causal effect. In realistic settings,
6 even near failures can be problematic. Equation (36.39) gives an expression for the (asymptotic)
7 variance of our estimate: $\mathbb{E}[\phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2]/n$. Notice that $\phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2$ involves terms that are
8 proportional to $1/g(x)$ and $1/(1 - g(x))$. Accordingly, the variance of our estimator will balloon
9 if there are units where $g(x) \approx 0$ or $g(x) \approx 1$ (unless such units are rare enough that they don't
10 contribute much to the expectation).

11 In practice, a simple way to deal with potential overlap violation is to fit a model \hat{g} for the
12 treatment assignment probability—which we need to do anyways—and check that the values $\hat{g}(x)$
13 are not too extreme. In the case that some values are too extreme, the simplest resolution is to cheat.
14 We can simply exclude all the data with extreme values of $\hat{g}(x)$. This is equivalent to considering
15 the average treatment effect over only the subpopulation where overlap is satisfied. This changes
16 the interpretation of the estimand. The restricted subpopulation ATE may or may not provide
17 a satisfactory answer to the real-world problem at hand, and this needs to be justified based on
18 knowledge of the real-world problem.

19

20
21 **36.4.5.3 Choice of Estimand and Average Treatment Effect on the Treated**

22 Usually, our goal in estimating a causal effect is qualitative. We want to know what the sign of the
23 effect is, and whether it's large or small. The utility of the ATE is that it provides a concrete query
24 we can use to get a handle on the qualitative question. However, it is not sacrosanct; sometimes
25 we're better off choosing an alternative causal estimand that still answers the qualitative question but
26 which is easier to estimate statistically. The **average treatment effect on the treated** or ATT,
27

28 $\text{ATT} \triangleq \mathbb{E}_{X|A=1}[\mathbb{E}[Y|X, \text{do}(A = 1)] - \mathbb{E}[Y|X, \text{do}(A = 0)]],$ (36.43)
29

30 is one such an estimand that is frequently useful.

31 The ATT is useful when many members of the population are very unlikely to receive treatment,
32 but the treated units had a reasonably high probability of receiving the control. This can happen if,
33 e.g., we sample control units from the general population, but the treatment units all self-selected
34 into treatment from a smaller subpopulation. In this case, it's not possible to (non-parametrically)
35 determine the treatment effect for the control units where no similar unit took treatment. The ATT
36 solves this obstacle by simply omitting such units from the average.

37 If we have the causal structure Figure 36.3, and the overlap condition $P(A = 1|X = x) < 1$ for all
38 $X = x$ then the ATT is causally identified as

39
40 $\tau^{\text{ATT}} = \mathbb{E}_{X|A=1}[\mathbb{E}[Y|A = 1, X] - \mathbb{E}[Y|A = 0, X]].$ (36.44)
41

42 Note that the required overlap condition here is weaker than for identifying the ATE. (The proof is
43 the same as Theorem 2.)

44 The estimation strategies for the ATE translate readily to estimation strategies for the ATT.
45 Namely, estimate the nuisance functions the same way and then simply replace averages over all
46 data points by averages over the treated datapoints only. In principle, it's possible to do a little
47

better than this by making use of the untreated datapoints as well. A corresponding double machine learning estimator is

$$\hat{\tau}^{\text{ATT-AIPTW}} \triangleq \frac{1}{n} \sum_i \frac{A_i}{\text{P}(A=1)} (Y - \hat{Q}(0, X_i)) - \frac{(1 - A_i)g(X)}{\text{P}(A=1)(1 - g(X))} (Y - \hat{Q}(0, X_i)). \quad (36.45)$$

. The variance of this estimator can be estimated by

$$\begin{aligned} \phi^{\text{ATT}}(\mathbf{X}_i; Q, g, \tau) &\triangleq \frac{1}{n} \sum_i \left[\frac{A_i}{\text{P}(A=1)} (Y - \hat{Q}(0, X_i)) \right. \\ &\quad \left. - \frac{(1 - A_i)g(X)}{\text{P}(A=1)(1 - g(X))} (Y - \hat{Q}(0, X_i) - \frac{A\tau}{\text{P}(A=1)}) \right] \end{aligned} \quad (36.46)$$

$$\hat{\mathbb{V}}[\hat{\tau}^{\text{ATT-AIPTW}}] \triangleq \frac{1}{n} \sum_i \phi^{\text{ATT}}(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau}^{\text{ATT-AIPTW}}). \quad (36.47)$$

Notice that the estimator for the ATT doesn't require estimating $Q(1, X)$. This can be a considerable advantage when the treated units are rare.

See Chernozhukov et al. [Che+17e] for details.

36.4.6 Summary and Practical Advice

We have seen a number of estimators that follow the general procedure:

1. Fit statistical or machine-learning models $\hat{Q}(a, x)$ as a predictor for Y , and/or $\hat{g}(x)$ as a predictor for A
2. Compute the predictions $\hat{Q}(0, x_i), \hat{Q}(1, x_i), \hat{g}(x_i)$ for each data point, and
3. Combine these predictions into an estimate of the average treatment effect.

Importantly, no single estimation approach is a silver bullet. For example, the double machine-learning estimator has appealing theoretical properties, such as asymptotic efficiency guarantees and a recipe for estimating uncertainty without needing to bootstrap the model fitting. However, in terms of the quality of point estimates, the double ML estimators can sometimes underperform their more naive counterparts [KS07]. In fact, there are cases where each of outcome regression, propensity weighting, or doubly robust methods will outperform the others.

One difficulty in choosing an estimator in practice is that there are fewer guardrails in causal inference than there are in standard predictive modeling. In predictive modeling, we construct a train-test split and validate our prediction models using the true labels or outcomes in the held-out dataset. However, for causal problems, the causal estimands are functionals of a different data-generating process from the one that we actually observed. As a result, it is impossible to empirically validate many aspects of causal estimation using standard techniques.

The effectiveness of a given approach is often determined by how much we trust the specification of our propensity score or outcome regression models $\hat{g}(x)$ and $\hat{Q}(a, x)$, and how well the treatment and control groups overlap in the dataset. Using flexible models for the nuisance functions g and Q can alleviate some of the concerns about model misspecification, but our freedom to use such models is

¹ often constrained by dataset size. When we have the luxury of large data, we can use flexible models;
² on the other hand, when the dataset is relatively small, we may need to use a smaller parametric
³ family or stringent regularization to obtain stable estimates of Q and g . Similarly, if overlap is poor
⁴ in some regions of the covariate space, then flexible models for Q may be highly variable, and inverse
⁵ propensity score weights may be large. In these cases, IPTW or AIPTW estimates may fluctuate
⁶ wildly as a function of large weights. Meanwhile, outcome regression estimates will be sensitive to
⁷ the specification of the Q model and its regularization, and can incur bias that is difficult to measure
⁸ if the specification or regularization does not match the true outcome process.

⁹ There are a number of practical steps that we can take to sanity-check causal estimates. The
¹⁰ simplest check is to compute many different ATE estimators (e.g., outcome regression, IPTW, doubly
¹¹ robust) using several comparably complex estimators of Q and g . We can then check whether they
¹² agree, at least qualitatively. If they do agree then this can provide some peace of mind (although it
¹³ is not a guarantee of accuracy). If they disagree, caution is warranted, particularly in choosing the
¹⁴ specification of the Q and g models.

¹⁵ It is also important to check for failures of overlap. Often, issues such as disagreement between
¹⁶ alternative estimators can be traced back to poor overlap. A common way to do this, particularly
¹⁷ with high-dimensional data, is to examine the estimated (ideally cross-fitted) propensity scores $\hat{g}(x_i)$.
¹⁸ This is a useful diagnostic, even if the intention is to use an outcome regression approach that only
¹⁹ incorporates and estimated outcome regression function $\hat{Q}(a, x_i)$. If overlap issues are relevant, it
²⁰ may be better to instead estimate either the average treatment effect on the treated, or the “trimmed”
²¹ estimand given by discarding units with extreme propensities.

²² Uncertainty quantification is also an essential part of most causal analyses. This frequently take
²³ the form of an estimate of the estimator’s variance, or a confidence interval. This may be important
²⁴ for downstream decision-making, and can also be a useful diagnostic. We can calculate variance either
²⁵ by bootstrapping the entire procedure (including refitting the models in each bootstrap replicate),
²⁶ or computing analytical variance estimates from the AIPTW estimator. Generally, large variance
²⁷ estimates may indicate issues with the analysis. For example, poor overlap will often (although
²⁸ not always) manifest as extremely large variances under either of these methods. Small variance
²⁹ estimates should be treated with caution, unless other checks, such as overlap checks, or stability
³⁰ across different Q and g models, also pass.

³¹ The previous advice only addresses the statistical problem of estimating τ from a data sample. It
³² does not speak to whether or not τ can reasonably be interpreted as an average treatment effect.
³³ Considerable care should be devoted to whether or not the assumption that there are no unobserved
³⁴ confounders is reasonable. There are several methods for assessing the sensitivity of the ATE estimate
³⁵ to violations of this assumption. See Section 36.7. Bias due to unobserved confounding can be
³⁶ substantial in practice—often overwhelming bias due to estimation error—so it is wise to conduct
³⁷ such an analysis.

³⁸

³⁹

⁴⁰

⁴¹ 36.5 Instrumental Variable Strategies

⁴²

⁴³ Adjustment-based methods rely on observing all confounders affecting the treatment and outcome.
⁴⁴ In some situations, it is possible to identify interesting causal effects even when there are unobserved
⁴⁵ confounders. We now consider strategies based on **instrumental variables**. The instrumental
⁴⁶ variable graph is shown in Figure 36.5. The key ingredient is the instrumental variable Z , a variable
⁴⁷

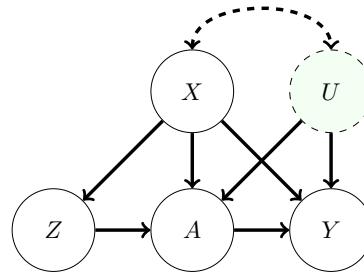


Figure 36.5: Causal graph illustrating the Instrumental Variable setup. The treatment A and outcome Y are both influenced by unobserved confounder U . Nevertheless, identification is sometimes possible due to the presence of the instrument Z . We also allow for observed covariates X that we may need to adjust for. The dashed arrow between U and X indicates a statistical dependency where we remain agnostic to the particular causal relationship.

that has a causal effect on Y only through its causal effect on A . Informally, the identification strategy is to determine the causal effect of Z on Y , the causal effect of Z on A , and then combine these into an estimate of the causal effect of A on Y .

For this identification strategy to work the instrument must satisfy three conditions. There are observed variables (confounders) X such that:

1. **Instrument Relevance** $Z \not\perp\!\!\!\perp A|X$: the instrument must actually affect the treatment assignment.
2. **Instrument Unconfoundedness** Any backdoor path between Z and Y is blocked by X , even conditional on A .
3. **Exclusion Restriction** All directed paths from Z to Y pass through A . That is, the instrument affects the outcome *only* through its effect on A .

(It may help conceptually to first think through the case where X is the empty set—i.e., where the only confounder is the unobserved U). These assumptions are necessary for using instrumental variables for causal identification, but they are not quite sufficient. In practice, they must be supplemented by an additional assumption that depends more closely on the details of the problem at hand. Historically, this additional assumption was usually that both the instrument-treatment and treatment-outcome relationships are linear. We'll examine some less restrictive alternatives below.

Before moving on to how to use instrumental variables for identification, let's consider how we might encounter instruments in practice. The key is that it's often possible to find, and measure, variables that affect treatment and that are assigned (as if) at random. For example, suppose we are interested in measuring the effect of taking a drug A on some health outcome Y . The challenge is that whether a study participant actually takes the drug can be confounded with Y —e.g., sicker people may be more likely to take their medication, but have worse outcomes. However, the assignment of treatments to patients can be randomized and this random assignment can be viewed as an instrument. This **random assignment with non-compliance** scenario is common in practice. The random assignment—the instrument—satisfies relevance (so long as assigning the drug affects the probability of the patient taking the drug). It also satisfies unconfoundedness (because the

¹ instrument is randomized). And, it plausibly satisfies exclusion restriction: telling (or not telling)
² a patient to take a drug has no effect on their health outcome except through influencing whether
³ or not they actually take the drug. As a second example, the **judge fixed effects** research design
⁴ uses the identity of the judge assigned to each criminal case to infer the effect of incarceration on
⁵ some life outcome of interest (e.g., total lifetime earnings). Relevance will be satisfied so long as
⁶ different judges have different propensities to hand out severe sentences. The assignment of trial
⁷ judges to cases is randomized, so unconfoundedness will also be satisfied. And, exclusion restriction
⁸ is also plausible: the particular identity of the judge assigned to your case has no bearing on your
⁹ years-later life outcomes, except through the particular sentence that you're subjected to.

¹⁰ It's important to note that these assumptions require some care, particularly exclusion restriction.
¹¹ Relevance can be checked directly from the data, by fitting a model to predict the treatment from the
¹² instrument (or vice versa). Unconfoundedness is often satisfied by design: the instrument is randomly
¹³ assigned. Even when literal random assignment doesn't hold, we often restrict to instruments
¹⁴ where unconfoundedness is "obviously" satisfied—e.g., using number of rainy days in a month as
¹⁵ an instrument for sun exposure. Exclusion restriction is trickier. For example, it might fail in the
¹⁶ drug assignment case if patients who are not told to take a drug respond by seeking out alternative
¹⁷ treatment. Or, it might fail in the judge fixed effects case if judges hand out additional, unrecorded,
¹⁸ punishments in addition to incarceration. Assessing the plausibility of exclusion restriction requires
¹⁹ careful consideration based on domain expertise.

²⁰ We now return to the question of how to make use of an instrument once we have it in hand. As
²¹ previously mentioned, getting causal identification using instrumental variables requires supplementing
²² the IV assumptions with some additional assumption about the causal process.

²³

²⁴ 36.5.1 Additive Unobserved Confounding

²⁵ We first consider **additive unobserved confounding**. That is, we assume that the structural causal
²⁶ model for the outcome has the form:⁷
²⁷

$$\begin{aligned} \text{36} \quad Y &\leftarrow f(A, X) + f_U(U). \end{aligned} \tag{36.48}$$

²⁸

²⁹ In words, we assume that there are no interaction effects between the treatment and the unobserved
³⁰ confounder—everyone responds to treatment in the same way. With this additional assumption, we
³¹ see that $\mathbb{E}[Y|X, \text{do}(A = a)] - \mathbb{E}[Y|X, \text{do}(A = a')] = f(a, X) - f(a', X)$. In this setting, our goal is to
³² learn this contrast.

³³ **Theorem 3** (Additive Confounding Identification). *If the instrumental variables assumptions hold*
³⁴ *and also additive unobserved confounding holds, then there is a function $\tilde{f}(a, x)$ where*
³⁵

$$\begin{aligned} \text{36} \quad \mathbb{E}[Y|x, \text{do}(A = a)] - \mathbb{E}[Y|x, \text{do}(A = a')] &= \tilde{f}(a, x) - \tilde{f}(a', x), \end{aligned} \tag{36.49}$$

³⁷

³⁸ for all x, a, a' and such that \tilde{f} satisfies

³⁹

$$\begin{aligned} \text{40} \quad \mathbb{E}[Y|z, x] &= \int \tilde{f}(a, x)p(a|z, x)da. \end{aligned} \tag{36.50}$$

⁴¹ Here, $p(a|z, x)$ is the conditional probability density of treatment.

⁴²

⁴³ 7. We roll the unit-specific variables ξ into U to avoid notational overload.

⁴⁴

In particular, if there is a unique function g that satisfies

$$\mathbb{E}[Y|z, x] = \int g(a, x)p(a|z, x)da, \quad (36.51)$$

then $g = \tilde{f}$ and this relation identifies the target causal effect.

Before giving the proof, let's understand the point of this identification result. The key insight is that both the left hand side of Equation (36.51) and $p(a|z, x)$ (appearing in the integrand) are identified by the data, since they involve only observational relationships between observed variables. So, \tilde{f} is identified implicitly as one of the functions that makes Equation (36.51) true. If there is a unique such function, then this fully identifies the causal effect.

Proof. With the additive unobserved confounding assumption, the instrument unconfoundedness implies that $U \perp\!\!\!\perp Z|X$. Then, we have that:

$$\mathbb{E}[Y|Z, X] = \mathbb{E}[f(A, X)|Z, X] + \mathbb{E}[f_U(U)|Z, X] \quad (36.52)$$

$$= \mathbb{E}[f(A, X)|Z, X] + \mathbb{E}[f_U(U)|X] \quad (36.53)$$

$$= \mathbb{E}[\tilde{f}(A, X)|Z, X], \quad (36.54)$$

where $\tilde{f} = f(A, X) + \mathbb{E}[f_U(U)|X]$. Now, identifying just \tilde{f} would suffice for us, because we could then identify contrasts between treatments: $f(a, x) - f(a', x) = \tilde{f}(a, x) - \tilde{f}(a', x)$. (The term $\mathbb{E}[f_U(U)|x]$ cancels out). Accordingly, we rewrite Equation (36.54) as:

$$\mathbb{E}[Y|z, x] = \int \tilde{f}(a, x)p(a|z, x)da. \quad (36.55)$$

□

It's worth dwelling briefly on how the IV assumptions come into play here. The exclusion restriction is implied by the additive unobserved confounding assumption, which we use explicitly. We also use the unconfoundedness assumption to conclude $U \perp\!\!\!\perp Z|X$. However, we do not use relevance. The role of relevance here is in ensuring that few functions solve the relation Equation (36.51). Informally, the solution g is constrained by the requirement that it hold for all values of Z . However, different values of Z only add non-trivial constraints if $p(a|z, x)$ differ depending on the value of z —this is exactly the relevance condition.

Estimation The basic estimation strategy is to fit models for $\mathbb{E}[Y|z, x]$ and $p(a|z, x)$ from the data, and then solve the implicit equation Equation (36.51) to find g consistent with the fitted models. The procedures for doing this can vary considerably depending on the particulars of the data (e.g., if Z is discrete or continuous) and the choice of modeling strategy. We omit a detailed discussion, but [see e.g. NP03; Dar+11; Har+17; SSG19; BKS19; Mua+20; Dik+20] for various concrete approaches.

It's also worth mentioning an additional nuance to the general procedure. Even if relevance holds, there will often be more than one function that satisfies Equation (36.51). So, we have only identified \tilde{f} as a member of this set of functions. In practice, this ambiguity is defeated by making some additional structural assumption about \tilde{f} . For example, we model \tilde{f} with a neural network, and then choose the network satisfying Equation (36.51) that has minimum l_2 -norm on the parameters (i.e., we pick the l_2 -regularized solution).

1 **36.5.2 Instrument Monotonicity and Local Average Treatment Effect**

2 We now consider an alternative assumption to additive unobserved confounding that is applicable
3 when both the instrument and treatment are binary. It will be convenient to conceptualize the
4 instrument as assignment-to-treatment. Then, the population divides into four subpopulations:
5

6 1. Compliers, who take the treatment if assigned to it, and who don't take the treatment otherwise.

7 2. Always takers, who take the treatment no matter their assignment

8 3. Never takers, who refuse the treatment no matter their assignment

9 4. Defiers, who refuse the treatment if assigned to it, and who take the treatment if not assigned.

10 Our goal in this setting will be to identify the average treatment effect among the compliers. The
11 **local average treatment effect (or complier average treatment effect)** is defined to be⁸

$$\underline{16} \quad \text{LATE} = \mathbb{E}[Y|\text{do}(A = 1), \text{complier}] - \mathbb{E}[Y|\text{do}(A = 0), \text{complier}]. \quad (36.56)$$

17 The LATE requires an additional assumption for identification. Namely, **instrument monotonicity**:
18 being assigned (not assigned) the treatment only increases (decreases) the probability that each
19 unit will take the treatment. Equivalently, $P(\text{defier}) = 0$.

20 We can then write down the identification result.

21 **Theorem 4.** *Given the instrumental variable assumptions and instrument monotonicity, the local
22 average treatment is identified as a parameter τ^{LATE} of the observational distributional; that is,
23 $\text{LATE} = \tau^{\text{LATE}}$. Namely,*

$$\underline{26} \quad \tau^{\text{LATE}} = \frac{\mathbb{E}[\mathbb{E}[Y|X, Z = 1] - \mathbb{E}[Y|X, Z = 0]]}{\mathbb{E}[P(A = 1|X, Z = 1) - P(A = 1|X, Z = 0)]}. \quad (36.57)$$

27 *Proof.* We now show that, given the IV assumptions and monotonicity, $\text{LATE} = \tau^{\text{LATE}}$. First, notice
28 that

$$\underline{31} \quad \tau^{\text{LATE}} = \frac{\mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)]}{P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))}. \quad (36.58)$$

32 This follows from backdoor adjustment, Theorem 2, applied to the numerator and denominator
33 separately. Our strategy will be to decompose $\mathbb{E}[Y|\text{do}(Z = z)]$ into the contributions from the
34 compliers, the units that ignore the instrument (the always/never takers), and the defiers. To that
35 end, note that $P(\text{complier}|\text{do}(Z = z)) = P(\text{complier})$ and similarly for always/never takers and
36 defiers—interventions on the instrument don't change the composition of the population. Then,

$$\underline{39} \quad \mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)] \quad (36.59)$$

$$\underline{40} \quad = (\mathbb{E}[Y|\text{complier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{complier}, \text{do}(Z = 0)])P(\text{complier}) \quad (36.60)$$

$$\underline{41} \quad + (\mathbb{E}[Y|\text{always/never}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{always/never}, \text{do}(Z = 0)])P(\text{always/never}) \quad (36.61)$$

$$\underline{42} \quad + (\mathbb{E}[Y|\text{defier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{defier}, \text{do}(Z = 0)])P(\text{defier}). \quad (36.62)$$

43 ⁸ We follow the econometrics literature in using “LATE” because “CATE” is already commonly used for conditional
44 average treatment effect.

The key is the effect on the complier subpopulation, Equation (36.60). First, by definition of the complier population, we have that:

$$\mathbb{E}[Y|\text{complier, do}(Z = z)] = \mathbb{E}[Y|\text{complier, do}(A = z)]. \quad (36.63)$$

That is, the causal effect of the treatment is the same as the causal effect of the instrument in this subpopulation—this is the core reason why access to an instrument allows identification of the local average treatment effect. This means that

$$\text{LATE} = \mathbb{E}[Y|\text{complier, do}(Z = 1)] - \mathbb{E}[Y|\text{complier, do}(Z = 0)]. \quad (36.64)$$

Further, we have that $P(\text{complier}) = P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))$. The reason is simply that, by definition of the subpopulations,

$$P(A = 1|\text{do}(Z = 1)) = P(\text{complier}) + P(\text{always taker}) \quad (36.65)$$

$$P(A = 1|\text{do}(Z = 0)) = P(\text{always taker}). \quad (36.66)$$

Now, plugging the expression for $P(\text{complier})$ and Equation (36.64) into Equation (36.60) we have that:

$$(\mathbb{E}[Y|\text{complier, do}(Z = 1)] - \mathbb{E}[Y|\text{complier, do}(Z = 0)])P(\text{complier}) \quad (36.67)$$

$$= \text{LATE} \times (P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))) \quad (36.68)$$

This gives us an expression for the local average treatment effect in terms of the effect of the instrument on the compliers and the probability that a unit takes the treatment when assigned/not-assigned.

The next step is to show that the remaining instrument effect decomposition terms, Equations (36.61) and (36.62), are both 0. Equation (36.61) is the causal effect of the instrument on the always/never takers. It's equal to 0 because, by definition of this subpopulation, the instrument has no causal effect in the subpopulation—they ignore the instrument! Mathematically, this is just $\mathbb{E}[Y|\text{always/never, do}(Z = 1)] = \mathbb{E}[Y|\text{always/never, do}(Z = 0)]$. Finally, Equation (36.62) is 0 by the instrument monotonicity assumption: we assumed that $P(\text{defier}) = 0$.

In totality, we now have that Equations (36.60) to (36.62) reduces to:

$$\mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)] \quad (36.69)$$

$$= \text{LATE} \times (P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))) + 0 + 0 \quad (36.70)$$

Rearranging for LATE and plugging in to Equation (36.58) gives claimed identification result. \square

36.5.2.1 Estimation

For estimating the local average treatment effect under the monotone instrument assumption, there is a double-machine learning approach that works with generic supervised learning approaches. Here, we want an estimator $\hat{\tau}^{\text{LATE}}$ for the parameter

$$\hat{\tau}^{\text{LATE}} = \frac{\mathbb{E}[\mathbb{E}[Y|X, Z = 1] - \mathbb{E}[Y|X, Z = 0]]}{\mathbb{E}[P(A = 1|X, Z = 1) - P(A = 1|X, Z = 0)]}. \quad (36.71)$$

1 To define the estimator, it's convenient to introduce some additional notation. First, we define the
2 nuisance functions:

3 $\mu(z, x) = \mathbb{E}[Y|z, x]$ (36.72)

4 $m(z, x) = P(A = 1|x, z)$ (36.73)

5 $p(x) = P(Z = 1|x).$ (36.74)

6

7 We also define the score ϕ by:

8 $\phi_{Z \rightarrow Y}(\mathbf{X}; \mu, p) \triangleq \mu(1, X) - \mu(0, X) + \frac{Z(Y - \mu(1, X))}{p(X)} - \frac{(1 - Z)(Y - \mu(0, X))}{1 - p(X)}$ (36.75)

9 $\phi_{Z \rightarrow A}(\mathbf{X}; m, p) \triangleq m(1, X) - m(0, X) + \frac{Z(A - m(1, X))}{p(X)} - \frac{(1 - Z)(A - m(0, X))}{1 - p(X)}$ (36.76)

10 $\phi(\mathbf{X}; \mu, m, p, \tau) \triangleq \phi_{Z \rightarrow Y}(\mathbf{X}; \mu, p) - \phi_{Z \rightarrow A}(\mathbf{X}; m, p) \times \tau$ (36.77)

11 Then, the estimator is defined by a two stage procedure:

12 1. Fit models $\hat{\mu}, \hat{m}, \hat{p}$ for each of μ, m, p (using supervised machine learning).

13 2. Define $\hat{\tau}^{\text{LATE}}$ as the solution to $\frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{\mu}, \hat{m}, \hat{p}, \hat{\tau}^{\text{LATE}}) = 0.$ That is,

14
$$\hat{\tau}^{\text{LATE}} = \frac{\frac{1}{n} \sum_i \phi_{Z \rightarrow Y}(\mathbf{X}_i; \hat{\mu}, \hat{p})}{\frac{1}{n} \sum_i \phi_{Z \rightarrow A}(\mathbf{X}_i; \hat{m}, \hat{p})}$$
 (36.78)

15 It may help intuitions to notice that the double machine learning estimator of the LATE is effectively
16 the double machine learning estimator of the average treatment effect of Z on Y divided by the
17 double machine learning estimator of the average treatment effect of Z on $A.$

18 Similarly to Section 36.4, the nuisance functions can be estimated by:

19

20 1. fit a model $\hat{\mu}$ that predicts Y from Z, X by minimizing mean square error

21 2. fit a model \hat{m} that predicts A from Z, X by minimizing mean cross-entropy

22 3. fit a model \hat{p} that predicts Z from X by minimizing mean cross-entropy.

23

24 As in Section 36.4, reusing the same data for model fitting and computing the estimator can
25 potentially cause problems. This can be avoided with use a cross-fitting procedure as described in
26 Section 36.4.2.4. In this case, we split the data into K folds and, for each fold k , use all the but
27 the k th fold to compute estimates $\hat{\mu}_{-k}, \hat{m}_{-k}, \hat{p}_{-k}$ of the nuisance parameters. Then we compute
28 the nuisance estimates for each datapoint i in fold k by predicting the required quantity using the
29 nuisance model fit on the other folds. That is, if unit i is in fold k , we compute $\hat{\mu}(z_i, x_i) \triangleq \hat{\mu}^{-k}(z_i, x_i)$
30 and so forth.

31 The key result is that if we use the cross-fit version of the estimator and the estimators for the
32 nuisance functions converge to their true values in the sense that

33 1. $\mathbb{E}(\hat{\mu}(Z, X) - \mu(Z, X))^2 \rightarrow 0$, $\mathbb{E}(\hat{m}(Z, X) - m(Z, X))^2 \rightarrow 0$, and $\mathbb{E}(\hat{p}(X) - p(X))^2 \rightarrow 0$

34 2. $\sqrt{\mathbb{E}[(\hat{p}(X) - p(X))^2]} \times (\sqrt{\mathbb{E}[(\hat{\mu}(Z, X) - \mu(Z, X))^2]} + \sqrt{\mathbb{E}[(\hat{m}(Z, X) - m(Z, X))^2]}) = o(\sqrt{n})$

35

then (with some omitted technical conditions) we have asymptotic normality at the \sqrt{n} -rate:

$$\sqrt{n}(\hat{\tau}^{\text{LATE}-\text{cf}} - \tau^{\text{LATE}}) \xrightarrow{d} \text{Normal}(0, \frac{\mathbb{E}[\phi(\mathbf{X}; \mu, m, p, \tau^{\text{LATE}})^2]}{\mathbb{E}[m(1, X) - m(0, X)]^2}). \quad (36.79)$$

As with double machine learning for the confounder adjustment strategy, the key point here is that we can achieve the (optimal) \sqrt{n} rate for estimating the LATE under a relatively weak condition on how well we estimate the nuisance functions—what matters is the *product* of the error in p and the errors in μ, m . So, for example, a very good model for how the instrument is assigned (p) can make up for errors in the estimation of the treatment-assignment (m) and outcome (μ) models.

The double machine learning estimator also gives a recipe for quantifying uncertainty. To that end, define

$$\hat{\tau}_{Z \rightarrow A} \triangleq \frac{1}{n} \sum_i \phi_{Z \rightarrow A}(\mathbf{X}_i; \hat{m}, \hat{p}) \quad (36.80)$$

$$\hat{\mathbb{V}}[\hat{\tau}^{\text{LATE}}] \triangleq \frac{1}{\hat{\tau}_{Z \rightarrow A}^2} \frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{\mu}, \hat{m}, \hat{p}, \hat{\tau}^{\text{LATE}})^2. \quad (36.81)$$

Then, subject to suitable technical conditions, $\hat{\mathbb{V}}[\hat{\tau}^{\text{LATE}-\text{cf}}]$ can be used as an estimate of the variance of the estimator. More precisely,

$$\sqrt{n}(\hat{\tau}^{\text{LATE}} - \tau^{\text{LATE}}) \xrightarrow{d} \text{Normal}(0, \hat{\mathbb{V}}[\hat{\tau}^{\text{LATE}}]). \quad (36.82)$$

Then, confidence intervals or p -values can be computed using this variance in the usual way. The main extra condition required for the variance estimator to be valid is that the nuisance parameters must all converge at rate $O(n^{-1/4})$ (so an excellent estimator for one can't fully compensate for terrible estimators of the others). In fact, even this condition is unnecessary in certain special cases—e.g., when p is known exactly, which occurs when the instrument is randomly assigned. See Chernozhukov et al. [Che+17e] for technical details.

36.5.3 Two Stage Least Squares

Commonly, the IV assumptions are supplemented with the following linear model assumptions:

$$A_i \leftarrow \alpha_0 + \alpha Z_i + \delta_A X_i + \gamma_A X_i + \xi_i^A \quad (36.83)$$

$$Y_i \leftarrow \beta_0 + \beta A_i + \delta_Y X_i + \gamma_Y X_i + \xi_i^Y \quad (36.84)$$

That is, we assume that the real-world process for treatment assignment and the outcome are both linear. In this case, plugging Equation (36.83) into Equation (36.84) yields

$$Y_i \leftarrow \tilde{\beta}_0 + \beta \alpha Z_i + \tilde{\delta} X_i + \tilde{\gamma} X_i + \tilde{\xi}_i. \quad (36.85)$$

The point is that β , the average treatment effect of A on Y , is equal to the coefficient $\beta \alpha$ of the instrument in the outcome-instrument model divided by the coefficient α of the instrument in the treatment-instrument model. So, to estimate the treatment effect, we simply fit both linear models and divide the estimated coefficients. This procedure is called **two stage least squares**.

¹ The simplicity of this procedure is seductive. However, the required linearity assumptions are hard
² to satisfy in practice and frequently lead to severe issues. A particularly pernicious version of this
³ is that linear-model misspecification together with weak relevance can yield standard errors for the
⁴ estimate that are far too small. In practice, this can lead us to find large, significant estimates from
⁵ two stage least squares when the truth is actually a weak or null effect. See [Rei16; You19; ASS19;
⁶ Lal+21] for critical evaluations of two stage least squares in practice.
⁷

⁸

⁹ 36.6 Difference in Differences

¹⁰

¹¹ Unsurprisingly, time plays an important role in causality. Causes precede effects, and we should be
¹² able to incorporate this knowledge into causal identification. We now turn to a particular strategy
¹³ for causal identification that relies on observing each unit at multiple time points. Data of this kind
¹⁴ is sometimes called **panel data**. We'll consider the simplest case. There are two time periods. In
¹⁵ the first period, none of the units are treated, and we observe an outcome Y_{0i} for each unit. Then,
¹⁶ a subset of the units are treated, denoted by $A_i = 1$. In the second time period, we again observe
¹⁷ the outcomes Y_{1i} for each unit, where now the outcomes of the treated units are affected by the
¹⁸ treatment. Our goal is to determine the average effect receiving the treatment had on the treated
¹⁹ units. That is, we want to know the average difference between the outcomes we actually observed
²⁰ for the treated units, and the outcomes we would have observed on those same units if they had not
²¹ been treated. The general strategy we look at is called **difference in differences**.⁹

²² As a concrete motivating example, consider trying to determine the effect raising minimum wage
²³ on employment. The concern here is that, in an efficient labor market, increasing the price of workers
²⁴ will reduce the demand for them, thereby driving down employment. As such, it seems increasing
²⁵ minimum wage may hurt the people the policy is nominally intended to help. The question is: how
²⁶ strong is this effect in practice? Card and Krueger [CK94a] studied this effect using difference in
²⁷ differences. The Philadelphia metropolitan area includes regions in both Pennsylvania and New
²⁸ Jersey (different US states). On April 1st 1992, New Jersey raised its minimum wage from \$4.25 to
²⁹ \$5.05. In Pennsylvania, the wage remained constant at \$4.25. The strategy is to collect employment
³⁰ data from fast food restaurants (which pay many employees minimum wage) in each state before
³¹ and after the change in minimum wage. In this case, for restaurant i , we have Y_{0i} , the number of
³² full time employees in February 1992, and Y_{1i} , the number of full time employees in November 1992.
³³ The treatment is simply $A_i = 1$ if the restaurant was located in New Jersey, and $A_i = 0$ if located in
³⁴ Pennsylvania. Our goal is to estimate the average effect of the minimum wage hike on employment
³⁵ in the restaurants affected by it (i.e., the ones in New Jersey).

³⁶ The assumption in classical difference-in-differences is the following structural equation:

³⁷

$$\text{38 } Y_{ti} \leftarrow W_i + S_t + \tau A_i \mathbb{I}(t = 1) + \xi_{ti}, \quad (36.86)$$

³⁹ with $\mathbb{E}[\xi_{ti}|W_i, S_t, A_i] = 0$. Here, W_i is a unit specific effect that is constant across time (e.g., the
⁴⁰ location of the restuarant or competence of the management) and S_t is a time-specific effect that
⁴¹ applies to all units (e.g., the state of the US economy at each time). Both of these quantities are
⁴² treated as unobserved, and not explicitly accounted for. The parameter τ captures the target causal
⁴³ effect. The (strong) assumption here is that unit, time, and treatment effects are all additive. This
⁴⁴

⁴⁵ 9. See github.com/vveitch/causality-tutorials/blob/main/difference_in_differences.ipynb.

⁴⁶⁴⁷

assumption is called **parallel trends**, because it is equivalent to assuming that, in the absence of treatment, the trend over time would be the same in both groups. It's easy to see that under this assumption, we have:

$$\tau = \mathbb{E}[Y_{1i} - Y_{0i}|A = 1] - \mathbb{E}[Y_{1i} - Y_{0i}|A = 0]. \quad (36.87)$$

That is, the estimand first computes the difference across time for both the treated and untreated group, and then computes the difference between these differences across the groups. The obvious estimator is then

$$\hat{\tau} = \frac{1}{n_A} \sum_{i:A_i=1} Y_{1i} - Y_{0i} - \frac{1}{n - n_A} \sum_{i:A_i=0} Y_{1i} - Y_{0i}, \quad (36.88)$$

where n_A is the number of treated units.

The root identification problem addressed by difference-in-differences is that $\mathbb{E}[W_i|A_i = 1] \neq \mathbb{E}[W_i|A_i = 0]$. That is, restaurants in New Jersey may be systematically different from restaurants in Pennsylvania in unobserved ways that affect employment.¹⁰ This is why we can't simply compare average outcomes for the treated and untreated. The identification assumption is that this unit-specific effect is the only source of statistical association with treatment; in particular we assume the time-specific effect has no such issue: $\mathbb{E}[S_{1i} - S_{0i}|A_i = 1] = \mathbb{E}[S_{1i} - S_{0i}|A_i = 0]$. Unfortunately, this assumption can be too strong. For instance, administrative data shows employment in Pennsylvania falling relative to employment in New Jersey between 1993 and 1996 [AP08, §5.2]. Although this doesn't directly contradict the parallel trends assumption used for identification, which needs to hold only in 1992, it does make it seem less credible.

To weaken the assumption, we'll look at a version that requires parallel trends to hold only after adjusting for covariates. To motivate this, we note that there were several different types of fast food restaurant included in the employment data. These vary, e.g., in the type of food they serve, and in cost per meal. Now, it seems reasonable the trend in employment may depend on the type of restaurant. For example, more expensive chains (such as Kentucky Fried Chicken) might be more affected by recessions than cheaper chains (such as McDonald's). If expensive chains are more common in New Jersey than in Pennsylvania, this effect can create a violation of parallel trends—if there's recession affecting both states, we'd expect employment to go down more in New Jersey than in Pennsylvania. However, we may find it credible that McDonald's restaurants in New Jersey have the same trend as McDonald's in Pennsylvania, and similarly for Kentucky Fried Chicken.

The next step is to give a definition of the target causal effect that doesn't depend on a parametric model, and a non-parametric statement of the identification assumption to go with it. In words, the causal estimand will be the average treatment effect on the units that received the treatment. To make sense of this mathematically, we'll introduce a new piece of notation:

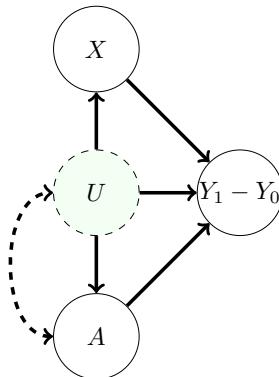
$$\mathbb{P}^{A=1}(Y|\text{do}(A = a)) \triangleq \int \mathbb{P}(Y|A = a, \text{parents of } Y)d\mathbb{P}(\text{parents of } Y|A = 1) \quad (36.89)$$

$$\mathbb{E}^{A=1}[Y|\text{do}(A = a)] \triangleq \mathbb{E}_{\mathbb{P}^{A=1}(Y|\text{do}(A=a))}[Y]. \quad (36.90)$$

In words: recall that the ordinary do operator works by replacing $\mathbb{P}(\text{parents}|A = a)$ by the marginal distribution $\mathbb{P}(\text{parents})$, thereby breaking the backdoor associations. Now, we're replacing the

¹⁰ 10. This is similar to the issue that arises from unobserved confounding, except W_i need not be a cause of the treatment assignment.

1
2
3
4
5
6
7
8
9
10
11
12
13



14 *Figure 36.6: Causal graph assumed for the difference-in-differences setting. Here, the outcome of interest*
 15 *is the difference between the pre- and post-treatment period, $Y_1 - Y_0$. This difference is influenced by the*
 16 *treatment, unobserved factors U , and observed covariates X . The dashed arrow between U and A indicates a*
 17 *statistical dependency between the variables, but where we remain agnostic to the precise causal mechanism.*
 18 *For example, in the minimum wage example, U might be the average income in restaurant's neighbourhood,*
 19 *which is dependent on the state, and hence also the treatment.*

20
21

22 distribution $P(\text{parents}|A = a)$ by $P(\text{parents}|A = 1)$, irrespective of the actual treatment value. This
 23 still breaks all backdoor associations, but is a better match for our target of estimating the treatment
 24 effect only among the treated units.

25 To formalize a causal estimand using the do calculus, we need to assume some partial causal
 26 structure. We'll use the graph in Figure 36.6. With this in hand, our causal estimand is the average
 27 treatment effect on the units that received the treatment, namely:

28

$$29 \quad \text{ATT}^{\text{DiD}} = \mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A = 1)] - \mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A = 0)] \quad (36.91)$$

30 In the minimum wage example, this is the average effect of the minimum wage hike on employment
 31 in the restaurants affected by it (i.e., the ones in New Jersey).

32 Finally, we formalize the identification assumption that, conditional on X , the trends in the treated
 33 and untreated groups are the same. The **conditional parallel trends** assumption is:

34

$$35 \quad \mathbb{E}^{A=1}[Y_1 - Y_0|X, \text{do}(A = 0)] = \mathbb{E}[Y_1 - Y_0|X, A = 0]. \quad (36.92)$$

36 In words, this says that for treated units with covariates X , the trend we would have seen had we not
 37 assigned treatment is the same as the trend we actually saw for the untreated units with covariates
 38 X . That is, if New Jersey had not raised its minimum wage, then McDonald's in New Jersey would
 39 have the same expected change in employment as McDonald's in Pennsylvania.

40 With this in hand, we can give the main identification result:

41 **Theorem 5** (Difference in Differences Identification). *We observe $A, Y_0, Y_1, X \sim P$. Suppose that*

42 *1. (Causal Structure) The data follows the causal graph in Figure 36.6.*

43

44 *2. (Conditional Parallel Trends) $\mathbb{E}^{A=1}[Y_1 - Y_0|X, \text{do}(A = 0)] = \mathbb{E}[Y_1 - Y_0|X, A = 0]$.*

45

46

1 3. (*Overlap*) $P(A = 1) > 0$ and $P(A = 1|X = x) < 1$ for all values of x in the sample space. That is,
2 there are no covariate values that only exist in the treated group.

4 Then, the average treatment effect on the treated is identified as $\text{ATT}^{\text{DiD}} = \tau^{\text{DiD}}$, where
5

$$\tau^{\text{DiD}} = \mathbb{E}[\mathbb{E}[Y_1 - Y_0|A = 1, X] - \mathbb{E}[Y_1 - Y_0|A = 0, X]|A = 1]. \quad (36.93)$$

7 *Proof.* First, by unrolling definitions, we have that

$$\mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A = 1), X] = \mathbb{E}[Y_1 - Y_0|A = 1, X]. \quad (36.94)$$

11 The interpretation is the near-tautology that the average effect among the treated under treatment
12 is equal to the actually observed average effect among the treated. Next,

$$\mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A = 0), X] = \mathbb{E}[Y_1 - Y_0|A = 0, X]. \quad (36.95)$$

15 is just the conditional parallel trends assumption. The result follows immediately.

16 (The overlap assumption is required to make sure all the conditional expectations are well
17 defined). \square

19 36.6.1 Estimation

21 With the identification result in hand, the next task is to estimate the observational estimand
22 Equation (36.93). To that end, we define $\tilde{Y} \triangleq Y_1 - Y_0$. Then, we've assumed that $\tilde{Y}, X, A \stackrel{\text{iid}}{\sim} P$ for
23 some unknown distribution P , and our target estimand is $\mathbb{E}[\mathbb{E}[\tilde{Y}|A = 1, X] - \mathbb{E}[\tilde{Y}|A = 0, X]|A = 1]$.
24 We can immediately recognize this as the observational estimand that occurs in estimating the
25 average treatment effect through adjustment, described in Section 36.4.5.3. That is, even though
26 the causal situation and the identification argument are different between the adjustment setting
27 and the difference in differences setting, the statistical estimation task we end up with is the same.
28 Accordingly, we can use all of the estimation tools we developed for adjustment. That is, all of the
29 techniques there—expected outcome modeling, propensity score methods, double machine learning,
30 and so forth—were purely about the *statistical* task, which is the same between the two scenarios.

31 So, we're left with the same general recipe for estimation we saw in Section 36.4.6. Namely,

- 33 1. Fit statistical or machine-learning models $\hat{Q}(a, x)$ as a predictor for $\tilde{Y} = Y_1 - Y_0$, and/or $\hat{g}(x)$ as
34 a predictor for A
- 35 2. Compute the predictions $\hat{Q}(0, x_i), \hat{Q}(1, x_i), \hat{g}(x_i)$ for each data point, and
- 37 3. Combine these predictions into an estimate of the average treatment effect on the treated.

39 The estimator in the third step can be the expected outcome model estimator, the propensity weighted
40 estimator, the double machine learning estimator, or any other strategy that's valid in the adjustment
41 setting.

43 36.7 Credibility Checks

45 Once we've chosen an identification strategy, fit our models, and produced an estimate, we're faced
46 with a basic question: should we believe it? Whether the reported estimate succeeds in capturing

1 the true causal effect depends on whether the assumptions required for causal identification hold, the
2 quality of the machine learning models, and the variability in the estimate due to only having access
3 to a finite data sample. The latter two problems are already familiar from machine learning and
4 statistical practice. We should, e.g., assess our models by checking performance on held out data,
5 examining feature importance, and so forth. Similarly, we should report measures of the uncertainty
6 due to finite sample (e.g., in the form of confidence intervals). Because these procedures are already
7 familiar practice, we will not dwell on them further. However, model evaluation and uncertainty
8 quantification are key parts of any credible causal analysis.

9 Assessing the validity of identification assumptions is trickier. First, there are assumptions that
10 can in fact be checked from data. For example, overlap should be checked in analysis using backdoor
11 adjustment or difference in differences, and relevance should be checked in the instrumental variable
12 setting. Again, checking these conditions is absolutely necessary for a credible causal analysis. But,
13 again, this involves only familiar data analysis, so we will not discuss it further. Next, there are the
14 causal assumptions that cannot be verified from data; e.g., no unobserved confounding in backdoor
15 adjustment, the exclusion restriction in IV, and conditional parallel trends in DiD. Ultimately, the
16 validity of these assumptions must be assessed using substantive causal knowledge of the particular
17 problem under consideration. However, it is possible to conduct some supplementary analyses that
18 make the required judgement easier. We now discuss two such techniques.

19

20

21

22 36.7.1 Placebo Checks

23

24 In many situations we may be able to find a variable that can be interpreted as a “treatment” that is
25 known to have no effect on the outcome, but which we expect to be confounded with the outcome in
26 a very similar fashion to the true treatment of interest. For example, if we’re trying to estimate the
27 efficacy of a COVID vaccine in preventing symptomatic COVID, we might take our placebo treatment
28 to be vaccination against HPV. We do not expect that there’s any causal effect here. However, it
29 seems plausible that latent factors that cause an individual to seek (or avoid) HPV vaccination and
30 COVID vaccination are similar; e.g., health conscientiousness, fear of needles, and so forth. Then, if
31 our identification strategy is valid for the COVID vaccine, we’d also expect it to be valid for
32 HPV vaccination. Accordingly, our estimation procedure we use for estimating the COVID effect
33 should, when applied to HPV, yield $\hat{\tau} \approx 0$. Or, more precisely, the confidence interval should contain
34 0. If this does not happen, then we may suspect that there are still some confounding factors lurking
35 that are not adequately handled by the identification procedure.

36 A similar procedure works when there is a variable that can be interpreted as an outcome which
37 is known to not be affected by the treatment, but that shares confounders with the outcome we’re
38 actually interested in. For example, in the COVID vaccination case, we might take the null outcome
39 to be symptomatic COVID within 7 days of vaccination [Dag+21]. Our knowledge of both the
40 biological mechanism of vaccination and the amount of time it takes to develop symptoms after
41 COVID infection (at least 2 days) lead us to conclude that it’s unlikely that the treatment has a
42 causal effect on the outcome. However, the properties of the treated people that affect how likely they
43 are to develop symptomatic COVID are largely the same in the 7 day and, e.g., 6 month window.
44 That includes factors such as risk aversion, baseline health, and so forth. Again, we can apply our
45 identification strategy to estimate the causal effect of the treatment on the null outcome. If the
46 confidence interval does not include 0, then we should doubt the credibility of the analysis.

47

36.7.2 Sensitivity Analysis to Unobserved Confounding

We now specialize to the case of estimating the average causal effect of a binary treatment by adjusting for confounding variables, as described in Section 36.4. In this case, causal identification is based on the assumption of ‘no unobserved confounding’; i.e., the assumption that the observed covariates include all common causes of the treatment assignment and outcome. This assumption is fundamentally untestable from observed data, but its violation can induce bias in the estimation of the treatment effect—the unobserved confounding may completely or in part explain the observed association. Our aim in this part is to develop a sensitivity analysis tool to aid in reasoning about potential bias induced by unobserved confounding.

Intuitively, if we estimate a large positive effect then we might expect the real effect is also positive, even in the presence of mild unobserved confounding. For example, consider the association between smoking and lung cancer. One could argue that this association arises from a hormone that both predisposes carriers to both an increased desire to smoke and to a greater risk of lung cancer. However, the association between smoking and lung cancer is large—is it plausible that some unknown hormonal association could have a strong enough influence to explain the association? Cornfield et al. [Cor+59] showed that, for a particular observational dataset, such an unmeasured hormone would need to increase the probability of smoking by at least a factor of nine. This is an unreasonable effect size for a hormone, so they conclude it’s unlikely the causal effect can be explained away.

We would like a general procedure to allow domain experts to make judgments about whether plausible confounding is “mild” relative to the “large” effect. In particular, the domain expert must translate judgments about the strength of the unobserved confounding into judgments about the bias induced in the estimate of the effect. Accordingly, we must formalize what is meant by strength of unobserved confounding, and to show how to translate judgments about confounding strength into judgments about bias.

A prototypical example, due to Imbens [Imb03] (building on [RR83]), illustrates the broad approach. As above, the observed data consists of a treatment A , an outcome Y , and covariates X that may causally affect the treatment and outcome. Imbens [Imb03] then posits an additional unobserved binary confounder U for each patient, and supposes that the observed data and unobserved confounder were generated according to the following assumption, known as **Imbens’ Sensitivity Model**:

$$U_i \stackrel{\text{iid}}{\sim} \text{Bern}(1/2) \quad (36.96)$$

$$A_i | X_i, U_i \stackrel{\text{ind}}{\sim} \text{Bern}(\text{sig}(\gamma X_i + \alpha U_i)) \quad (36.97)$$

$$Y_i | X_i, A_i, U_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\tau A_i + \beta X_i + \delta U_i, \sigma^2). \quad (36.98)$$

where sig is the sigmoid function.

If we had observed U_i , we could estimate $(\hat{\tau}, \hat{\gamma}, \hat{\beta}, \hat{\alpha}, \hat{\delta}, \hat{\sigma}^2)$ from the data and report $\hat{\tau}$ as the estimate of the average treatment effect. Since U_i is not observed, it is not possible to identify the parameters from the data. Instead, we make (subjective) judgments about plausible values of α —how strongly U_i affects the treatment assignment—and δ —how strongly U_i affects the outcome. Contingent on plausible $\alpha = \alpha^*$ and $\delta = \delta^*$, the other parameters can be estimated. This yields an estimate of the treatment effect $\hat{\tau}(\alpha^*, \delta^*)$ under the presumed values of the sensitivity parameters.

The approach just outlined has a major drawback: it relies on a parametric model for the full data generating process. The assumed model is equivalent to assuming that, had U been observed, it

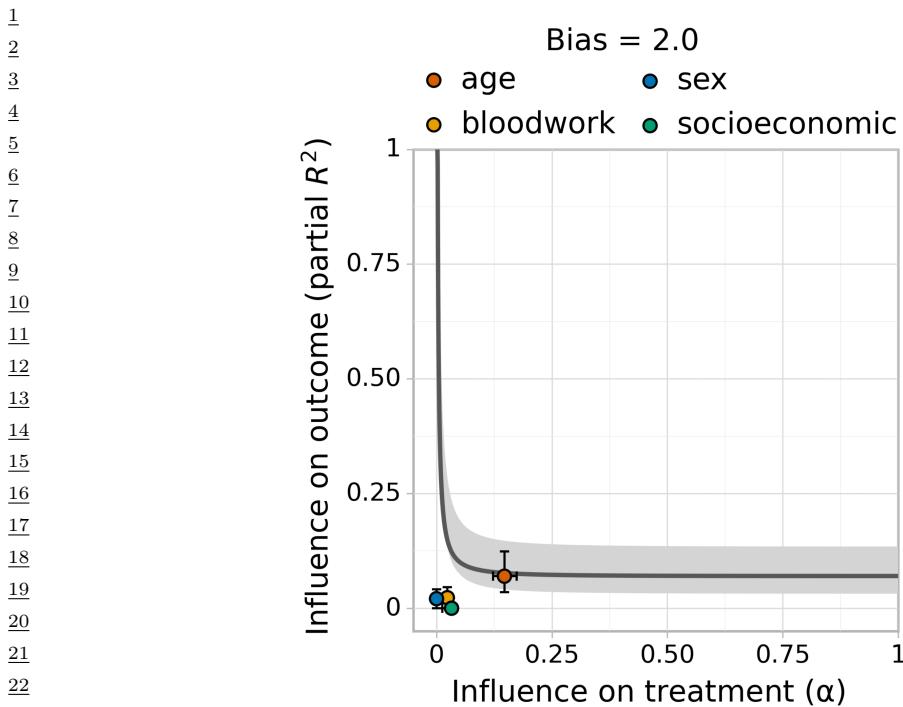


Figure 36.7: Austen plot showing how strong an unobserved confounder would need to be to induce a bias of 2 in an observational study of the effect of combination blood pressure medications on diastolic blood pressure [Dor+16]. We chose this bias to equal the nominal average treatment effect estimated from the data. We model the outcome with Bayesian Additive Regression Trees and the treatment assignment with logistic regression. The curve shows all values treatment and outcome influence that would induce a bias of 2. The colored dots show the influence strength of (groups of) observed covariates, given all other covariates. For example, an unobserved confounder with as much influence as the patient's age might induce a bias of about 2.

would have been appropriate to use logistic regression to model treatment assignment, and linear regression to model the outcome. This assumption also implies a simple, parametric model for the relationships governing the observed data. This restriction is out of step with modern practice, where we use flexible machine-learning methods to model these relationships. For example, the assumption forbids the use of neural networks or random forests, though such methods are often state-of-the-art for causal effect estimation.

Austen plots We now turn to developing an alternative an adaptation of Imbens' approach that fully decouples sensitivity analysis and modeling of the observed data. Namely, the **Austen plots** of [VZ20]. An example Austen plot is shown in Figure 36.7. The high-level idea is to posit a generative model that uses a simple, interpretable parametric form for the influence of the unobserved confounder, but that *puts no constraints on the model for the observed data*. We then use the

parametric part of the model to formalize “confounding strength” and to compute the induced bias as a function of the confounding.

Austen plots further adapt two strategies pioneered by Imbens [Imb03]. First, we find a parameterization of the model so that the sensitivity parameters, measuring strength of confounding, are on a standardized, unitless scale. This allows us to compare the strength of hypothetical unobserved confounding to the strength of observed covariates, measured from data. Second, we plot the curve of all values of the sensitivity parameter that would yield given level of bias. This moves the analyst judgment from “what are plausible values of the sensitivity parameters?” to “are sensitivity parameters this extreme plausible?”

Figure 36.7, an Austen plot for an observational study of the effect of combination medications on diastolic blood pressure, illustrates the idea. A bias of 2 would suffice to undermine the qualitative conclusion that the blood-pressure treatment is effective. Examining the plot, an unobserved confounder as strong as age could induce this amount of confounding, but no other (group of) observed confounders has so much influence. Accordingly, if a domain expert thinks an unobserved confounder as strong as age is unlikely then they may conclude that the treatment is likely effective. Or, if such a confounder is plausible, they may conclude that the study fails to establish efficacy.

Setup The data are generated independently and identically $(Y_i, A_i, X_i, U_i) \stackrel{\text{iid}}{\sim} P$, where U_i is not observed and P is some unknown probability distribution. The approach in Section 36.4 assumes that the observed covariates X contain all common causes of Y and A . If this ‘no unobserved confounding’ assumption holds, then the ATE is equal to parameter, τ , of the observed data distribution, where

$$\tau = \mathbb{E}[\mathbb{E}[Y|X, A = 1] - \mathbb{E}[Y|X, A = 0]]. \quad (36.99)$$

This observational parameter is then estimated from a finite data sample. Recall from Section 36.4 that this involves estimating the conditional expected outcome $Q(A, X) = \mathbb{E}[Y|A, X]$ and the propensity score $g(X) = P(A = 1|X)$, then plugging these into an estimator $\hat{\tau}$.

We are now concerned with the case of possible unobserved confounding. That is, where U causally affects Y and A . If there is unobserved confounding then the parameter τ is not equal to the ATE, so $\hat{\tau}$ is a biased estimate. Inference about the ATE then divides into two tasks. First, the statistical task: estimating τ as accurately as possible from the observed data. And, second, the causal (domain-specific) problem of assessing $\text{bias} = \text{ATE} - \tau$. We emphasize that our focus here is bias due to causal misidentification, not the statistical bias of the estimator. Our aim is to reason about the bias induced by unobserved confounding—the second task—in a way that imposes no constraints on the modeling choices for \hat{Q} , \hat{g} and $\hat{\tau}$ used in the statistical analysis.

Sensitivity Model Our sensitivity analysis should impose no constraints on how the *observed* data is modeled. However, sensitivity analysis demands some assumption on the relationship between the observed data and the *unobserved* confounder. It is convenient to formalize such assumptions by specifying a probabilistic model for how the data is generated. The strength of confounding is then formalized in terms of the parameters of the model (the sensitivity parameters). Then, the bias induced by the confounding can be derived from the assumed model. Our task is to posit a generative model that both yields a useful and easily interpretable sensitivity analysis, and that avoids imposing any assumptions about the observed data.

To begin, consider the functional form of the sensitivity model used by Imbens [Imb03].

$$\text{logit}(\mathbb{P}(A = 1|x, u)) = h(x) + \alpha u \quad (36.100)$$

$$\mathbb{E}[Y|a, x, u] = l(a, x) + \delta u, \quad (36.101)$$

for some functions h and l . That is, the propensity score is logit-linear in the unobserved confounder, and the conditional expected outcome is linear.

By rearranging Equation (36.100) to solve for u and plugging in to Equation (36.101), we see that it's equivalent to assume $\mathbb{E}[Y|t, x, u] = \tilde{l}(t, x) + \delta \text{logit}(\mathbb{P}(A = 1|x, u))$. That is, the unobserved confounder u only influences the outcome through the propensity score. Accordingly, by positing a distribution on $\mathbb{P}(A = 1|x, u)$ directly, we can circumvent the need to explicitly articulate U (and h).

Definition 36.7.1. Let $\tilde{g}(x, u) = \mathbb{P}(A = 1|x, u)$ denote the propensity score given observed covariates x and the unobserved confounder u .

The insight is that we can posit a sensitivity model by defining a distribution on \tilde{g} directly. We choose:

$$\tilde{g}(X, U)|X \sim \text{Beta}(g(X)(1/\alpha - 1), (1 - g(X))(1/\alpha - 1)).$$

That is, the full propensity score $\tilde{g}(X, U)$ for each unit is assumed to be sampled from a Beta distribution centered at the observed propensity score $g(X)$. The sensitivity parameter α plays the same role as in Imbens' model: it controls the influence of the unobserved confounder U on treatment assignment. When α is close to 0 then $\tilde{g}(X, U)|X$ is tightly concentrated around $g(X)$, and the unobserved confounder has little influence. That is, U minimally affects our belief about who is likely to receive treatment. Conversely, when α is close to 1 then \tilde{g} concentrates near 0 and 1; i.e., knowing U would let us accurately predict treatment assignment. Indeed, it can be shown that α is the change in our belief about how likely a unit was to have gotten the treatment, given that they were actually observed to be treated (or not):

$$\alpha = \mathbb{E}[\tilde{g}(X, U)|A = 1] - \mathbb{E}[\tilde{g}(X, U)|A = 0]. \quad (36.102)$$

With the \tilde{g} model in hand, we define the **Austen Sensitivity Model** as follows:

$$\tilde{g}(X, U)|X \sim \text{Beta}(g(X)(1/\alpha - 1), (1 - g(X))(1/\alpha - 1)) \quad (36.103)$$

$$A|X, U \sim \text{Bern}(\tilde{g}(X, U)) \quad (36.104)$$

$$\mathbb{E}[Y|A, X, U] = Q(A, X) + \delta(\text{logit}\tilde{g}(X, U) - \mathbb{E}[\text{logit}\tilde{g}(X, U)|A, X]). \quad (36.105)$$

This model has been constructed to satisfy the requirement that the propensity score and conditional expected outcome are the g and Q actually present in the observed data:

$$\mathbb{P}(A = 1|X) = \mathbb{E}[\mathbb{E}[T|X, U]|X] = \mathbb{E}[\tilde{g}(X, U)|X] = g(X)$$

$$\mathbb{E}[Y|A, X] = \mathbb{E}[\mathbb{E}[Y|A, X, U]|A, X] = Q(A, X).$$

The sensitivity parameters are α , controlling the dependence between the unobserved confounder the treatment assignment, and δ , controlling the relationship with the outcome.

1 **Bias** We now turn to calculating the bias induced by unobserved confounding. By assumption, X
2 and U together suffice to render the average treatment effect identifiable as:

3 $\text{ATE} = \mathbb{E}[\mathbb{E}[Y|A=1, X, U] - \mathbb{E}[Y|A=0, X, U]].$

4 Plugging in our sensitivity model yields,

5 $\text{ATE} = \mathbb{E}[Q(1, X) - Q(0, X)] + \delta(\mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A=1] - \mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A=0]).$

6 The first term is the observed-data estimate τ , so

7 $\text{bias} = \delta(\mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A=1] - \mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A=0]).$

8 Then, by invoking Beta-Bernoulli conjugacy and standard Beta identities,¹¹ we arrive at,

9 **Theorem 6.** *Under the Austen sensitivity model, Equation (36.105), an unobserved confounder with
10 influence α and δ induces bias in the estimated treatment effect equal to*

11
$$\text{bias} = \frac{\delta}{1/\alpha - 1} \mathbb{E}\left[\frac{1}{g(X)} + \frac{1}{1-g(X)}\right].$$

12 That is, the amount of bias is determined by the sensitivity parameters and by the *realized*
13 propensity score. Notice that more extreme propensity scores lead to more extreme bias in response
14 to unobserved confounding. This means, in particular, that conditioning on a covariate that affects
15 the treatment but that does not directly affect the outcome (an instrument) will increase any bias
16 due to unobserved confounding. This general phenomena is known as **z-bias**.

17 **Sensitivity Parameters** The Austen model provides a formalization of confounding strength
18 in terms of the parameters α and δ and tells us how much bias is induced by a given strength of
19 confounding. This lets us translate judgments about confounding strength to judgments about bias.
20 However, it is not immediately obvious how to translate qualitative judgements such as “I think any
21 unobserved confounder would be much less important than Age” to judgements about the possible
22 values of the sensitivity parameters.

23 First, because the scale of δ is not fixed, it may be difficult to compare the influence of potential
24 unobserved confounders to the influence of reference variables. To resolve this, we reexpress the
25 outcome-confounder strength in terms of the (non-parametric) partial coefficient of determination:

26
$$R_{Y,\text{par}}^2(\alpha, \delta) = 1 - \frac{\mathbb{E}(Y - \mathbb{E}[Y|A, X, U])^2}{\mathbb{E}(Y - Q(A, X))^2}.$$

27 The key to computing the reparameterization is the following result

28 **Theorem 7.** *Under the Austen sensitivity model, Equation (36.105), the outcome influence is*

29
$$R_{Y,\text{par}}^2(\alpha, \delta) = \delta^2 \sum_{a=0}^1 \frac{\mathbb{E}[\psi_1(g(X)^a(1-g(X))^{1-a}(1/\alpha - 1) + 1[A=a])]}{\mathbb{E}[(Y - Q(A, X))^2]},$$

30 where ψ_1 is the trigamma function.

31 11. We also use the recurrence relation $\psi(x+1) - \psi(x) = 1/x$, where ψ is the digamma function.

¹ See Veitch and Zaveri [VZ20] for the proof.
²

³ By design, α —the strength of confounding influence on treatment assignment—is already on
⁴ a fixed, unitless scale. However, because the measure is tied to the model it may be difficult to
⁵ interpret, and it is not obvious how to compute reference confounding strength values from the
⁶ observed data. The next result clarifies these issues.

⁷ **Theorem 8.** *Under the Austen sensitivity model, Equation (36.105),*
⁸

$$\frac{9}{10} \quad \alpha = 1 - \frac{\mathbb{E}[\tilde{g}(X, U)(1 - \tilde{g}(X, U))]}{\mathbb{E}[g(X)(1 - g(X))]}.$$
¹¹

¹² See Veitch and Zaveri [VZ20] for the proof. That is, the sensitivity parameter α measures how
¹³ much more extreme the propensity scores become when we condition on U . That is, α is a measure
¹⁴ of the extra predictive power U adds for A , above and beyond the predictive power in X . It may
¹⁵ also be insightful to notice that

$$\frac{16}{17} \quad \alpha = R_{A,\text{par}}^2 = 1 - \frac{\mathbb{E}[(A - \tilde{g}(X, U))^2]}{\mathbb{E}[(A - g(X))^2]}. \quad (36.106)$$
¹⁸

¹⁹ That is, α is just the (non-parametric) partial coefficient of determination of U on A —the same
²⁰ measure used for the outcome influence. (To see this, just expand the expectations conditional on
²¹ $A = 1$ and $A = 0$).
²²

²³ **Estimating bias** In combination, Theorems 6 and 7 yield an expression for the bias in terms of α
²⁴ and $R_{Y,\text{par}}^2$. In practice, we can estimate the bias induced by confounding by fitting models for \hat{Q}
²⁵ and \hat{g} and replacing the expectations by means over the data.
²⁶

²⁷ 36.7.2.1 Calibration using observed data

²⁹ The analyst must make judgments about the influence a hypothetical unobserved confounder might
³⁰ have on treatment assignment and outcome. To calibrate such judgments, we'd like to have a reference
³¹ point for how much the observed covariates influence the treatment assignment and outcome. In the
³² sensitivity model, the degree of influence is measured by partial R_Y^2 and α . We want to measure the
³³ degree of influence of an observed covariate Z given the other observed covariates $X \setminus Z$.

³⁴ For the outcome, this can be measured as:
³⁵

$$\frac{36}{37} \quad R_{Y \cdot Z|T, X \setminus Z}^2 \triangleq 1 - \frac{\mathbb{E}(Y - Q(A, X))^2}{\mathbb{E}(Y - \mathbb{E}[Y|A, X \setminus Z])^2}.$$
³⁸

³⁹ In practice, we can estimate the quantity by fitting a new regression model \hat{Q}_Z that predicts Y from
⁴⁰ A and $X \setminus Z$. Then we compute

$$\frac{41}{42} \quad R_{Y \cdot Z|T, X \setminus Z}^2 = 1 - \frac{\frac{1}{n} \sum_i (y_i - \hat{Q}(t_i, x_i))^2}{\frac{1}{n} \sum_i (y_i - \hat{Q}_Z(t_i, x_i \setminus z_i))^2}.$$
⁴³

⁴⁴ Using Theorem 8, we can measure influence of observed covariate Z on treatment assignment given
⁴⁵ $X \setminus Z$ in an analogous fashion to the outcome. We define $g_{X \setminus Z}(X \setminus Z) = P(A = 1|X \setminus Z)$, then fit a
⁴⁶

model for $g_{X \setminus Z}$ by predicting A from $X \setminus Z$, and estimate

$$\hat{\alpha}_{Z|X \setminus Z} = 1 - \frac{\frac{1}{n} \sum_i \hat{g}(x_i)(1 - \hat{g}(x_i))}{\frac{1}{n} \sum_i \hat{g}_{X \setminus Z}(x_i \setminus z_i)(1 - \hat{g}_{X \setminus Z}(x_i \setminus z_i))}.$$

Grouping covariates The estimated values $\hat{\alpha}_{X \setminus Z}$ and $\hat{R}_{Y,X \setminus Z}^2$ measure the influence of Z conditioned on all the other confounders. In some cases, this can be misleading. For example, if some piece of information is important but there are multiple covariates providing redundant measurements, then the estimated influence of each covariate will be small. To avoid this, group together related or strongly dependent covariates and compute the influence of the entire group in aggregate. For example, grouping income, location, and race as ‘socioeconomic variables’.

36.7.2.2 Practical Use

We now have sufficient results to produce Austen plots such as Figure 36.7. At a high level, the procedure is:

1. Produce an estimate $\hat{\tau}$ using any modeling tools. As a component of this, estimate the propensity score \hat{g} and conditional outcome model \hat{Q}
2. Pick a level of bias that would suffice to change the qualitative interpretation of the estimate (e.g., the lower bound of a 95% confidence interval).
3. Plot the values of α and $R_{Y,\text{par}}^2$ that would suffice to induce that much bias. This is the black curve on the plot. To calculate these values, use Theorems 6 and 7 together with the estimated \hat{g} and \hat{Q} .
4. Finally, compute reference influence level for (groups of) observed covariates. In particular, this requires fitting reduced models for the conditional expected outcome and propensity that do not use the reference covariate as a feature.

In practice, an analyst only needs to do the model fitting parts themselves. The bias calculations, reference value calculations, and plotting can be done automatically with standard libraries.¹²

Austen plots are predicated on Equation (36.105). This assumption replaces the purely parametric Equation (36.98) with a version that eliminates any parametric requirements on the observed data. However, we emphasize that Equation (36.105) does, implicitly, impose some parametric assumption on the structural causal relationship between U and A, Y . Ultimately, any conclusion drawn from the sensitivity analysis depends on this assumption, which is not justified on any substantive grounds. Accordingly, such sensitivity analyses can only be used to informally guide domain experts. They do not circumvent the need to thoroughly adjust for confounding. This reliance on a structural assumption is a generic property of sensitivity analysis.¹³ Indeed, there are now many sensitivity analysis models that allow the use of any machine learning model in the data analysis [e.g., RRS00; FDF19; She+11; HS13; BK19; Ros10; Yad+18; ZSB19; Sch+21a]. However, none of these are yet in

¹² See github.com/vveitch/causality-tutorials/blob/main/Sensitivity_Analysis.ipynb.

¹³ In extreme cases, there can be so little unexplained variation in A or Y that only a very weak confounder could be compatible with the data. In this case, essentially assumption free sensitivity analysis is possible [Man90].

¹ routine use in practice. We have presented Austen plots here not because they make an especially
² virtuous modeling assumption, but because they are (relatively) easy to understand and interpret.
³

⁴ Austen plots are most useful in situations where the conclusion from the plot would be ‘obvious’
⁵ to a domain expert. For instance, in Figure 36.7, we can be confident that an unobserved confounder
⁶ similar to socioeconomic status would not induce enough bias to change the qualitative conclusion.
⁷ By contrast, Austen plots should not be used to draw conclusions such as, “I think a latent confounder
⁸ could only be 90% as strong as ‘age’, so there is evidence of a small non-zero effect”. Such nuanced
⁹ conclusions might depend on issues such as the particular sensitivity model we use, or finite-sample
¹⁰ variation of our bias and influence estimates, or on incautious interpretation of the calibration dots.
¹¹ These issues are subtle, and it would be difficult resolve them to a sufficient degree that a sensitivity
¹² analysis would make an analysis credible.

¹³

¹⁴ **Calibration using observed data** The interpretation of the observed-data calibration requires
¹⁵ some care. The sensitivity analysis requires the analyst to make judgements about the strength
¹⁶ of influence of the unobserved confounder U , *conditional on the observed covariates X* . However,
¹⁷ we report the strength of influence of observed covariate(s) Z , *conditional on the other observed*
¹⁸ *covariates $X \setminus Z$* . The difference in conditioning sets can have subtle effects.

¹⁹ Cinelli and Hazlett [CH20] give an example where Z and U are identical variables in the true
²⁰ model, but where influence of U given A, X is larger than the influence of Z given $A, X \setminus Z$. (The
²¹ influence of Z given $X \setminus Z, U$ would be the same as the influence of U given X). Accordingly, an
²² analyst is *not* justified in a judgement such as, “I know that U and Z are very similar. I see Z has
²³ substantial influence, but the dot is below the line. Thus, U will not undo the study conclusions.” In
²⁴ essence, if the domain expert suspects a strong interaction between U and Z then naively eyeballing
²⁵ the dot-vs-line position may be misleading. A particular subtle case is when U and Z are independent
²⁶ variables that both strongly influence A and Y . The joint influence on A creates an interaction effect
²⁷ between them when A is conditioned on (the treatment is a collider). This affects the interpretation
²⁸ of $R^2_{Y \cdot U | X, A}$. Indeed, we should generally be skeptical of sensitivity analysis interpretation when it is
²⁹ expected that a strong confounder has been omitted. In such cases, our conclusions may depend
³⁰ substantively on the particular form of our sensitivity model, or other unjustifiable assumptions.

³¹ Although the interaction problem is conceptually important, its practical significance is unclear.
³² We often expect the opposite effect: if U and Z are dependent (e.g., race and wealth) then omitting U
³³ should increase the apparent importance of Z —leading to a conservative judgement (a dot artificially
³⁴ towards the top right part of the plot).

³⁵

³⁶

³⁷ 36.8 The Do Calculus

³⁸

³⁹ We have seen several strategies for identifying causal effects as parameters of observational distribu-
⁴⁰ tions. Confounder adjustment (Section 36.4) relied only on the assumed causal graph (and overlap),
⁴¹ which specified that we observe all common causes of A and Y . On the other hand, instrumental
⁴² variable methods and difference-in-differences each relied on both an assumed causal graph and
⁴³ partial functional form assumptions about the underlying structural causal model. Because functional
⁴⁴ form assumptions can be quite difficult to justify on substantive grounds, it’s natural to ask when
⁴⁵ causal identification is possible from the causal graph alone. That is, when can we be agnostic to the
⁴⁶ particular functional form of the structural causal models?

⁴⁷

There is a general “**calculus of intervention**”, known as the **do-calculus**, that gives a general recipe for determining when the causal assumptions expressed in a causal graph can be used to identify causal effects [Pea09c]. The do-calculus is a set of three rewrite rules that allows us to replace statements where we condition on variables being set by intervention, e.g. $P(Y|\text{do}(A = a))$, with statements involving only observational quantities, e.g. $\mathbb{E}_X[P(Y|A = a, X)]$. When causal identification is possible, we can repeatedly apply the three rules to boil down our target causal parameter into an expression involving only the observational distribution.

36.8.1 The three rules

To express the rules, let X , Y , Z , and W be arbitrary disjoint sets of variables in a causal DAG G .

Rule 1 The first rule allows us to insert or delete observations z :

$$p(y|\text{do}(x), z, w) = p(y|\text{do}(x), w) \text{ if } (Y \perp Z|X, W)_{G_{\overline{X}}} \quad (36.107)$$

where $G_{\overline{X}}$ denotes cutting edges going into X , and $(Y \perp Z|X, W)_{G_{\overline{X}}}$ denotes conditional independence in the mutilated graph. The rule follows from d-separation in the mutilated graph. This rule just says that conditioning on irrelevant variables leaves the distribution invariant (as we would expect).

Rule 2 The second rule allows us to replace $\text{do}(z)$ with conditioning on (seeing) z . The simplest case where we can do this is: if Z is a root of the causal graph (i.e., it has no causal parents) then $p(y|\text{do}(z)) = p(y|z)$. The reason is that the do operator is equivalent to conditioning in the mutilated causal graph where all the edges into Z are removed, but, because Z is a root, the mutilated graph is just the original causal graph. The general form of this rule is:

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), z, w) \text{ if } (Y \perp Z|X, W)_{G_{\overline{X}Z}} \quad (36.108)$$

where $G_{\overline{X}Z}$ cuts edges going into X and out of Z . Intuitively, we can replace $\text{do}(z)$ by z as long as there are no backdoor (non-directed) paths between z and y . If there are in fact no such paths, then cutting all the edges going out of Z will mean there are no paths connecting Z and Y , so that $Y \perp Z$. The rule just generalizes this line of reasoning to allow for extra observed and intervened variables.

Rule 3 The third rule allows us to insert or delete actions $\text{do}(z)$:

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), w) \text{ if } (Y \perp Z|X, W)_{G_{\overline{X}Z^*}} \quad (36.109)$$

where $G_{\overline{X}Z^*}$ cuts edges going into X and Z^* , and where Z^* is the set of Z -nodes that are not ancestors of any W -node in $G_{\overline{X}}$. Intuitively, this condition corresponds to intervening on X , and checking whether the distribution of Y is invariant to *any* intervention that we could apply on Z .

36.8.2 Revisiting Backdoor Adjustment

We begin with a more general form of the adjustment formula we used in Section 36.4.

¹ First, suppose we observe all of A 's parents, call them X . For notational simplicity, we'll assume
² for the moment that X is discrete. Then,
³

$$\begin{aligned} \frac{4}{5} p(Y = y|\text{do}(A = a)) &= \sum_x p(Y = y|x, \text{do}(A = a))p(x|\text{do}(A = a)) \\ \frac{6}{7} &= \sum_x p(Y = y|x, A = a)p(x). \end{aligned} \quad (36.110)$$

$$\frac{8}{9} \quad (36.111)$$

⁹ The first line is just a standard probability relation (marginalizing over z). We are using causal
¹⁰ assumptions in two ways in the second line. First, $p(x|\text{do}(A = a)) = p(x)$: the treatment has
¹¹ no causal effect on Z , so interventions on A don't change the distribution of Z . This is rule 3,
¹² Equation (36.109). Second, $p(Y = y|z, \text{do}(A = a)) = p(Y = y|z, A = a)$. This equality holds because
¹³ conditioning on the parents blocks all non-directed paths from A to Y , reducing the causal effect to
¹⁴ be the same as the observational effect. The equality is an application of rule 2, Equation (36.108).

¹⁵ Now, what if we don't observe all the parents of A ? The key issue is **backdoor paths**: paths
¹⁶ between A and Y that contain an arrow into A . These paths are the general form of the problem
¹⁷ that occurs when A and Y share a common cause. Suppose that we can find a set of variables S
¹⁸ such that (1) no node in S is a descendant of A ; and (2) S blocks every backdoor path between A
¹⁹ and Y . Such a set is said to satisfy the **backdoor criterion**. In this case, we can use S instead of
²⁰ the parents of X in the adjustment formula, Equation (36.111). That is,

$$\frac{22}{23} \quad p(Y = y|\text{do}(A = a)) = \mathbb{E}_S[p(Y = y|S, A = a)]. \quad (36.112)$$

²⁴ The proof follows the invocation of rules 3 and 2, in the same way as for the case where S is just the
²⁵ parents of A . Notice that requiring S to not contain any descendants of A means that we don't risk
²⁶ conditioning on any variables that mediate the effect, nor any variables that might be colliders—either
²⁷ would undermine the estimate.

²⁸ The backdoor adjustment formula generalizes the adjust-for-parents approach and adjust-for-all-
²⁹ common-causes approach of Section 36.4. That's because both the parents of A and the common
³⁰ causes satisfy the backdoor criterion.

³¹ In practice, the full distribution $p(Y = y|\text{do}(A = a))$ is rarely used as the causal target. Instead,
³² we try to estimate a low-dimensional parameter of this distribution, such as the average treatment
³³ effect. The adjustment formula immediately translates in the obvious way. If we define

$$\begin{aligned} \frac{34}{35} \quad \tau &= \mathbb{E}_S[\mathbb{E}[Y|A = 1, S] - \mathbb{E}[Y|A = 0, S]], \end{aligned}$$

³⁶ then we have that $\text{ATE} = \tau$ whenever S satisfies the backdoor criteria. The parameter τ can then
³⁷ be estimated from finite data using the methods described in Section 36.4, using S in place of the
³⁸ common causes X .

³⁹

⁴⁰ 36.8.3 Frontdoor Adjustment

⁴¹ Backdoor adjustment is applicable if there's at least one observed variable on every backdoor path
⁴² between A and Y . As we have seen, identification is sometimes still possible even when this condition
⁴³ doesn't hold. Frontdoor adjustment is another strategy of this kind. Figure 36.8 shows the causal
⁴⁴ structure that allows this kind of adjustment strategy. Suppose we're interested in the effect of
⁴⁵ smoking A on developing cancer Y , but we're concerned about some latent genetic confounder U .

⁴⁶

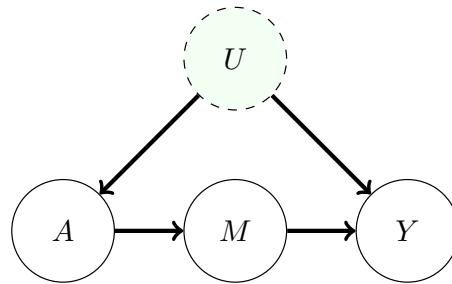


Figure 36.8: Causal graph illustrating the frontdoor criterion setup. The effect of the treatment A on outcome Y is entirely mediated by mediator M . This allows us to infer the causal effect even if the treatment and outcome are confounded by U .

Suppose that all of the directed paths from A to Y pass through some set of variables M . Such variables are called **mediators**. For example, the effect of smoking on lung cancer might be entirely mediated by the amount of tar in the lungs and measured tissue damage. It turns out that if all such mediators are observed, and the mediators do not have an unobserved common cause with A or Y , then causal identification is possible. To understand why this is true, first notice that we can identify the causal effect of A on M and the causal effect of M on A , both by backdoor adjustment. Further, the mechanism of action of A on Y is: A changes M which in turn changes Y . Then, we can combine these as:

$$p(Y|\text{do}(A = a)) = \sum_m p(Y|\text{do}(M = m))p(M = m|\text{do}(A = a)) \quad (36.113)$$

$$= \sum_m \sum_{a'} p(Y|a', m)p(a')p(m|a) \quad (36.114)$$

The second line is just backdoor adjustment applied to identify each of the do expressions (note that A blocks the $M-Y$ backdoor path through U).

Equation (36.114) is called the **front-door formula** [Pea09b, §3.3.2]. To state the result in more general terms, let us introduce a definition. We say a set of variables M satisfies the **front-door criterion** relative to an ordered pair of variables (A, Y) if (1) M intercepts all directed paths from A to Y ; (2) there is no unblocked backdoor path from A to M ; and (3) all backdoor paths from M to Y are blocked by A . If M satisfies this criterion, and if $p(A, M) > 0$ for all values of A and M , then the causal effect of A on Y is identifiable and is given by Equation (36.114).

Let us interpret this theorem in terms of our smoking example. Condition 1 means that smoking A should have no effect on cancer Y except via tar and tissue damage M . Conditions 2 and 3 mean that the genotype U cannot have any effect on M except via smoking A . Finally, the requirement that $p(A, M) > 0$ for all values implies that high levels of tar in the lungs must arise not only due to smoking, but also other factors (e.g., pollutants). In other words, we require $p(A = 0, M = 1) > 0$ so we can assess the impact of the mediator in the untreated setting.

We can now use the do-calculus to derive the frontdoor criterion; following [PM18b, p236]. Assuming

1 the causal graph G shown in Figure 36.8:

$$\begin{aligned}
 \underline{2} \quad p(y|\text{do}(a)) &= \sum_m p(y|\text{do}(a), m)p(m|\text{do}(a)) && (\text{probability axioms}) \\
 \underline{3} \quad &= \sum_m p(y|\text{do}(a), \text{do}(m))p(m|\text{do}(a)) && (\text{rule 2 using } G_{\overline{S}\overline{T}}) \\
 \underline{4} \quad &= \sum_m p(y|\text{do}(a), \text{do}(m))p(m|a) && (\text{rule 2 using } G_{\underline{S}}) \\
 \underline{5} \quad &= \sum_m p(y|\text{do}(m))p(m|a) && (\text{rule 3 using } G_{\overline{S}\overline{T}^*}) \\
 \underline{6} \quad &= \sum_{a'} \sum_m p(y|\text{do}(m), a')p(a'|\text{do}(m))p(m|a) && (\text{probability axioms}) \\
 \underline{7} \quad &= \sum_{a'} \sum_m p(y|m, a')p(a'|\text{do}(m))p(m|a) && (\text{rule 2 using } G_{\underline{T}}) \\
 \underline{8} \quad &= \sum_{a'} \sum_m p(y|m, a')p(a')p(m|a) && (\text{rule 3 using } G_{\overline{T}^*}) \\
 \underline{9} \\
 \underline{10} \\
 \underline{11} \\
 \underline{12} \\
 \underline{13} \\
 \underline{14} \\
 \underline{15} \\
 \underline{16} \\
 \underline{17} \\
 \underline{18} \\
 \underline{19}
 \end{aligned}$$

20 **Estimation** To estimate the causal distribution from data using the frontdoor criterion we need to
21 estimate each of $p(y|m, a)$, $p(a)$, and $p(m|a)$. In practice, we can fit models $\hat{p}(y|m, a)$ by predicting
22 Y from M and A , and $\hat{p}(m|a)$ by predicting M from A . Then, using the empirical distribution to
23 estimate $p(a)$, the final estimate is:

$$\underline{24} \quad \frac{1}{|A|} \sum_{a'} \sum_m \hat{p}(y|m, a') \hat{p}(m|a), \quad (36.115)$$

25 where $|A|$ is the number of treatments.

26 We usually have more modest targets than the full distribution $p(y|\text{do}(a))$. For instance, we may
27 be content with just estimating the average treatment effect. It's straightforward to derive a formula
28 for this using the frontdoor adjustment. Similarly to backdoor adjustment, more advanced estimators
29 of the ATE through frontdoor effect are possible in principle. For example, we might combine fitted
30 models for $\mathbb{E}[Y|m, a]$ and $P(M|a)$. See Fulcher et al. [Ful+20] for an approach to robust estimation
31 via front door adjustment, as well as a generalization of the front door approach to more general
32 settings.

33

34 36.9 Further Reading

35

36 There is an enormous and growing literature on the intersection of causality and machine learning.

37 First, there are many textbooks on theoretical and practical elements of causal inference. These
38 include Pearl [Pea09c], focused on causal graphs, Angrist and Pischke [AP08], focused on econometrics,
39 Hernán and Robins [HR20b], with roots in epidemiology, Imbens and Rubin [IR15], with origin in
40 statistics, and Morgan and Winship [MW15], for a social sciences perspective. The introduction to
41 causality in Shalizi [Sha22, §7] is also recommended, particularly the treatment of matching.

42 Double machine-learning has featured prominently in this chapter. This is a particular instantiation
43 of non-parametric estimation. This topic has substantial theoretical and practical importance in
44

modern causal inference. The double machine learning work includes estimators for many commonly encountered scenarios [Che+17e; Che+17d]. Good references for a lucid explanation of how and why non-parametric estimation works include [Ken16; Ken17; FK21]. Usually, the key guarantees of non-parametric estimator are asymptotic. Generally, there are many estimators that share optimal asymptotic guarantees (e.g. the AIPTW estimator given in Equation (36.30)). Although these are asymptotically equivalent, in finite samples their behavior can be very different. There are estimators that preserve asymptotic guarantees but aim to improve performance in practical finite sample regimes [e.g., vR11].

There is also considerable interest in the estimation of heterogeneous treatment effects. The question here is: what effect would this treatment have when applied to a unit with such-and-such specific characteristics? E.g., what is the effect of this drug on women over the age of 50? The causal identification arguments used here are more-or-less the same as for the estimation of average case effects. However, the estimation problems can be substantially more involved. Some reading includes [Kün+19; NW20; Ken20; Yad+21].

There are several commonly applicable causal identification and estimation strategies beyond the ones we've covered in this chapter. **Regression discontinuity designs** rely on the presence of some sharp, arbitrary non-linearity in treatment assignment. For example, eligibility for some aid programs is determined by whether an individual has income below or above a fixed amount. The effect of the treatment can be studied by comparing units just below and just above this threshold. **Synthetic controls** are a class of methods that try to study the effect of a treatment on a given unit by constructing a synthetic version of that unit that acts as a control. For example, to study the effect of legislation banning smoking indoors in California, we can construct a synthetic California as a weighted average of other states, with weights chosen to balance demographic characteristics. Then, we can compare the observed outcome of California with the outcome of the synthetic control, constructed as the weighted average of the outcomes of the donor states. See Angrist and Pischke [AP08] for a textbook treatment of both strategies. Closely related are methods that use time series modeling to create synthetic outcomes. For example, to study the effect of an advertising campaign beginning at time T on product sales Y_t , we might build a time series model for Y_t using data in the $t < T$ period, and then use this model to predict the values of $(\hat{Y}_t)_{t>T}$ we would have seen had the campaign not been run. We can estimate the causal effect by comparing the factual, realized Y_t to the predicted, counterfactual, \hat{Y}_t . See Brodersen et al. [Bro+15] for an instantiation of this idea.

In this chapter, our focus has been on using machine learning tools to estimate causal effects. There is also a growing interest in using the ideas of causality to improve machine learning tools. This is mainly aimed at building predictors that are robust when deployed in new domains [SS18b; SCS19; Arj+20; Mei18b; PBM16a; RC+18; Zha+13a; Sch+12b; Vei+21] or that do not rely on particular ‘spurious’ correlations in the training data [RPH21; Wu+21; Gar+19; Mit+20; WZ19; KCC20; KHL20; TAH20; Vei+21].

Index

- Q -function, 1118
 f -MAX, 1168
 $\#P$ -hard, 404
Global explanation, 1073
Local explanation, 1072
intrinsic, 1079
“arms”, 1101
“do” operator, 188
1/f, 11
2d lattice model, 141
3-SAT, 404
8-queens problem, 280
80-20 rule, 12
- Sylvester determinant lemma, 664
- A* search, 1016, 1150
A/B test, 1101
A2C, 1145
A3C, 1145
ABC, 537, 846
abduction, 190
absorbing state, 53, 1116
abstain, 702
accept, 468
acquisition function, 273, 275, 1124, 1125
action, 1106
action nodes, 1098
action-value function, 1118
Actionability, 1077
actions, 1095
activation function, 592
active inference, 1166
active learning, 273, 277, 314, 664, 1123
actor critic, 1133
actor-critic, 1144
acyclic directed mixed graph, 173
AD, 231
Adam, 271
adapt, 704
adaptive filtering, 964
adaptive importance sampling, 442
adaptive MCMC, 472
adaptive policies, 1106
adaptive prediction sets, 555
- adaptive rejection sampling, 457
adaptive resampling, 521
adaptive tempering, 533
add one smoothing, 71
add-one smoothing, 51
additive intervention, 188
additive scalar bijection, 795
additive unobserved confounding, 1194
ADF, 373
adjoint sensitivity method, 807
ADMG, 173
admissible, 1150
admixture mixture, 917
admixture model, 27, 915, 917
advantage actor critic, 1145
advantage function, 1118
adversarial stickers, 723
adversarial attack, 723
adversarial autoencoder, 761
adversarial bandit, 1107
adversarial image, 723
adversarial inverse RL, 1168
adversarial risk, 727
adversarial training, 727
ADVI, 326, 434
AEP, 210
affine autoregressive flows, 800
affine flow, 794
affine scalar bijection, 795
affine transformation, 794
affinity propagation, 397
agent, 1095, 1106, 1115
aggregate posterior, 222
aggregated posterior, 228, 756
AI, 221
AIPTW, 1184
AIRL, 1168
Akaike information criterion, 111
aleatoric uncertainty, 67, 550
ALI, 866
alpha divergence, 56
alpha posterior, 622
alphafold, 144
AlphaGo, 1151
AlphaZero, 1151
alternative hypothesis, 106
AMD, 403
- amortization gap, 438
amortized ELBO, 437
amortized inference, 399, 437, 747, 748
analysis by synthesis, 881
analysis-by-synthesis, 1049
ancestor, 341, 519
ancestral graph, 155
ancestral sampling, 128, 286, 835
annealed importance sampling, 460, 533, 741, 811
annotation shift, 697
annulus, 14
anomaly detection, 699, 903, 982
anti-causal prediction, 695
anti-Hebbian term, 815
antiferromagnetic, 142
antithetic sampling, 464
aperiodic, 54
apprenticeship learning, 1166
Approximate Bayesian Computation, 537, 846
approximate dynamic programming, 1120
approximate inference, 322
approximate linear programming, 1120
approximate minimum degree, 403
AR, 731
architectural methods, 721
ARD, 571, 643
ARD kernel, 643
arithmetic circuits, 179
arithmetic coding, 221
ARMA, 976
arrow of time, 984
artificial dynamics, 362
ASOS, 966
ASR, 342, 734
associative Markov model, 144
associative Markov network, 141
associative memory, 145
assumed density filtering, 373
asynchronous advantage actor critic, 1145
asymmetric numeral systems, 221
asymptotic consistency, 1017

- 1
2 asymptotic equipartition property, **210**
3 asynchronous updates, **394**
4 asynchronous value iteration, **1121**
5 ATE, **1170**
6 atoms, **1008**
7 ATT, **1190**
8 attention layer, **597**
9 attention score, **597**
10 attention weight, **597**
11 attribute, **180**
12 audio-visual speech recognition, **956**
13 augmented DAG, **188**
14 Augmented Inverse Probability of Treatment Weighted Estimator, **1184**
15 augmented reality, **528**
16 Austen plots, **1206**
17 Austen Sensitivity Model, **1208**
18 auto-encoding SMC, **538**
19 auto-regressive HMM, **936**
20 auto-regressive model, **783**
21 Auto-Sklearn, **278**
22 Auto-Weka, **278**
23 autoconj, **427**
24 autocorrelation function, **499**
25 autocorrelation time, **499**
26 autodiff, **231**
27 autoencoder, **602**
28 autoencoders, **1050**
29 automatic differentiation, **231**
30 automatic differentiation variational inference, **326**, **434**
31 Automatic Relevance Determination, **643**
32 automatic relevance determination, **903**
33 automatic relevancy determination, **571**, **610**
34 automatic speech recognition, **342**, **734**, **953**
35 autoregressive, **731**
36 autoregressive bijection, **799**
37 autoregressive flows, **798**
38 autoregressive models, **874**
39 auxiliary latent variables, **438**
40 auxiliary variable deep generative model, **439**
41 auxiliary variables, **481**
42 average causal effect, **1179**
43 Average Treatment Effect, **1170**
44 average treatment effect, **1178**
45 average treatment effect on the treated, **1190**
46 axis aligned, **14**
47
48 BA lower bound, **218**
49 back translation, **704**
50 backbone, **601**
51 backcasting, **985**
52 backdoor criterion, **1214**
53 backdoor paths, **1214**
54 backoff smoothing, **103**
55 backpropagation, **240**
56 backup diagram, **1122**
- backwards Kalman gain matrix, **350**
 backwards kernel, **532**
 backwards transfer, **720**
 BACS, **709**
 bagging, **619**
 balance the dataset, **698**
 BALD, **1127**
 Bambi, **587**
 BAMDP, **1135**
 bandit problem, **1106**
 bandwidth, **642**, **736**
 Barlow Twins, **1057**
 base distribution, **791**
 base measure, **30**, **1006**
 base-rate, **1033**
 baseline, **463**
 baseline function, **250**
 basic random variables, **180**
 basin flooding, **281**
 basis functions, **931**
 batch active learning, **1127**
 batch ensemble, **620**
 batch normalization, **595**
 batch optimization, **249**
 batch reinforcement learning, **1157**
 BatchBALD, **1128**
 batched Bayesian optimization, **278**
 Baum-Welch, **941**, **944**
 Baum-Welch algorithm, **135**
 Bayes ball algorithm, **123**
 Bayes by backprop, **613**
 Bayes estimator, **1095**
 Bayes factor, **106**
 Bayes filter, **333**
 Bayes nets, **117**
 Bayes' rule, **64**
 Bayes' rule for Gaussians, **20**
 Bayes-adaptive MDP, **1135**
 Bayes-Newton, **373**
 BayesBiNN, **272**
 Bayesian active learning by disagreement, **1127**
 Bayesian approach, **67**
 Bayesian dark knowledge, **622**
 Bayesian decision theory, **1095**
 Bayesian deep learning, **607**
 Bayesian factor regression, **907**
 Bayesian hypothesis testing, **106**
 Bayesian inference, **1**, **64**
 Bayesian information criterion, **110**
 Bayesian lasso, **569**
 Bayesian learning rule, **265**
 Bayesian model selection, **105**
 Bayesian multi net, **138**
 Bayesian networks, **117**
 Bayesian nonparametric, **1005**
 Bayesian nonparametric models, **507**
 Bayesian Occam's razor, **106**
 Bayesian online changepoint detection, **959**
 Bayesian optimization, **272**, **286**, **641**
- Bayesian Optimization Algorithm, **286**
 Bayesian p-value, **113**
 Bayesian quadrature, **465**
 Bayesian statistics, **63**
 Bayesian structural time series, **977**
 Bayesian transfer learning, **610**
 BayesOpt, **272**
 BBB, **613**
 BBMM, **673**
 BBVI, **428**
 BDL, **607**
 beam search, **1016**
 Bean Machine, **185**
 bearings only tracking problem, **973**
 behavior cloning, **1166**
 behavior policy, **1157**
 behavior-agnostic off-policy, **1157**
 belief networks, **117**
 belief propagation, **383**
 belief state, **333**, **550**, **1109**, **1135**
 belief states, **383**
 belief-state MDP, **1110**
 Bellman backup, **1120**
 Bellman error, **1118**
 Bellman residual, **1118**
 Bellman's optimality equations, **1118**
 Berkson's paradox, **120**, **125**
 Bernoulli bandit, **1110**
 Bernoulli distribution, **5**
 Bernoulli mixture model, **886**
 BERT, **604**, **767**, **1052**
 Bessel function, **27**, **644**
 best arm identification, **272**, **1104**
 best-arm identification, **1108**
 best-first search, **1150**
 beta distribution, **13**
 beta function, **8**
 Beta process, **1023**
 beta-VAE, **758**
 bi-directed graph, **175**
 BIC, **110**, **681**
 BIC loss, **111**
 BIC score, **111**
 big data, **70**
 BiGAN, **866**, **1050**
 bigram model, **49**, **211**
 bigram statistics, **51**
 bijection, **44**, **792**
 bilinear form, **599**
 binary entropy function, **208**
 binary logistic regression, **575**
 binary neural network, **272**
 binomial coefficient, **5**
 binomial distribution, **5**
 binomial regression, **558**
 bit error, **388**
 bits, **198**, **208**
 bits back coding, **224**
 bits per dimension, **740**
 BIVA, **773**
 bivariate Gaussian, **14**
 black box attack, **725**
 black box shift estimation, **708**

- 1
 2 black box variational inference, 428
 3 blackbox, 428
 4 blackbox EP, 446
 5 blackbox matrix-matrix multiplication, 673
 6 Blackwell-MacQueen, 1009
 7 blind inverse problem, 891
 8 blind source separation, 926
 9 block length, 225
 10 block stacking, 1153
 11 blocked Gibbs sampling, 479
 12 BLOG, 185
 13 BN, 595
 14 BNP, 1005
 15 BOA, 286
 16 Bochner's theorem, 650
 17 BOCPD, 959
 18 Boltzmann machine, 145
 19 Boltzmann policy, 1134
 20 bond variables, 483
 21 Bonnet's theorem, 247, 252
 22 boolean satisfiability problems, 280
 23 bootstrap filter, 517
 24 bootstrap sampling, 618
 25 bootstrapping, 1133, 1137
 26 borrow statistical strength, 95
 27 bottom-up inference model, 773
 28 bound optimization, 255, 258
 29 Box-Muller, 455
 30 BP, 383
 31 branching factor, 211
 32 Bregman divergence, 206, 207
 33 Brier score, 546
 34 BRMS, 587
 35 Brownian motion, 490, 491, 645, 1030
 36 Brownian noise, 501
 37 BSS, 926
 38 BSTS, 977
 39 bucket elimination, 399
 40 BUGS, 473
 41 building blocks, 286, 591
 42 burn-in phase, 492
 43 burn-in time, 467
 44 burstiness, 1032
 45 BYOL, 1058
 46 calculus of intervention, 1213
 47 calculus of variations, 39
 48 calibrated, 546
 49 calibration set, 554
 50 canonical correlation analysis, 909, 1043
 51 canonical form, 16, 30
 52 canonical link function, 560
 53 canonical parameters, 16, 30, 33, 151
 54 CAQL, 1119
 55 cart-pole swing-up, 1152
 56 casino HMM, 935
 57 CASP, 144
 58 catastrophic forgetting, 719
 59 categorical, 6
 60 categorical distribution, 72
 61 categorical PCA, 911, 921
 62 CatPCA, 911
 63 Cauchy, 9
 64 Cauchy sequence, 657
 65 causal convolution, 785
 66 causal DAGs, 695
 67 causal discovery, 1003
 68 Causal graphs, 1173
 69 causal hierarchy, 189
 70 causal impact, 191, 985
 71 Causal inference, 1169
 72 causal Markov assumption, 186
 73 causal models, 186
 74 causal prediction, 695
 75 causal representation learning, 1060
 76 causally sufficient, 186
 77 causes of effects, 189
 78 CAVI, 411
 79 cavity distribution, 445
 80 CCA, 909
 81 cdf, 7
 82 CEB, 229
 83 CelebA, 738, 738, 754
 84 Centered kernel alignment (CKA), 1043
 85 centering matrix, 80
 86 central composite design, 678
 87 central limit theorem, 453, 687
 88 central moment, 9
 89 certify, 727
 90 ceteris paribus, 190
 91 chain components, 172
 92 chain compositions, 236
 93 chain graph, 156, 172, 174
 94 chain rule, 235
 95 chance nodes, 1098
 96 change of variables, 44
 97 changepoint detection, 698, 957, 959
 98 channel coding, 193, 225
 99 channel coding theorem, 397
 100 Chapman-Kolmogorov, 48
 101 Chapman-Kolmogorov equation, 334
 102 characteristic length scale, 643
 103 Chernoff-Hoeffding inequality, 1111
 104 chi-squared distance, 57
 105 Chi-squared distribution, 11
 106 children, 117
 107 Chinese restaurant process, 1010
 108 choice theory, 582
 109 Cholesky decomposition, 455
 110 Chomsky normal form, 163, 164
 111 chordal, 169
 112 Chow-Liu algorithm, 286
 113 chromosomes, 283
 114 CI, 117
 115 circuits, 238
 116 circular flow, 806
 117 circular normal, 27
 118 citation matching, 185
 119 CKF, 367
 120 clamped phase, 157, 815
 121 class incremental learning, 719
 122 classical statistics, 63
 123 classifier guidance, 840
 124 classifier-free guidance, 841
 125 Claude Shannon, 221
 126 clausal form, 183
 127 click through rate, 1105, 1108
 128 clinical trials, 1108
 129 CLIP, 789, 1056
 130 clique, 139
 131 clique tree, 405
 132 cliques, 402
 133 closed world assumption, 184, 972
 134 closing the loop, 530
 135 closure, 153
 136 cluster variational method, 395
 137 clustering, 884
 138 clutter problem, 328
 139 CMA-ES, 288, 1154
 140 CNN, 601
 141 co-information, 215
 142 co-parents, 127
 143 coagulation, 1020
 144 coalescence, 518
 145 cocktail party problem, 926
 146 code words, 221
 147 codebook, 778
 148 codebook loss, 779
 149 codewords, 193
 150 coffee, lemon, milk, and tea, 297
 151 cold posterior, 622
 152 cold start problem, 70
 153 collapsed, 135
 154 collapsed Gibbs sampler, 479, 1012
 155 collapsed particles, 524
 156 collective classification, 182
 157 collider, 123, 1173
 158 collocation, 1151
 159 coloring, 473
 160 commitment loss, 780
 161 common corruptions, 693
 162 common random number, 1153
 163 common random numbers, 462
 164 common random numbers trick, 435
 165 compact support, 673
 166 Compactness, Sparsity, 1077
 167 compatible, 1149
 168 complementary log-log, 560
 169 complete, 657
 170 complete data, 156
 171 completely random measures, 1027
 172 Completeness, 1077, 1079
 173 completing the square, 82
 174 complexity penalty, 109
 175 complier average treatment effect, 1196
 176 components, 976
 177 composite likelihood, 158
 178 compositional pattern-producing network, 725
 179 Compression Lemma, 201
 180 computation graph, 591
 181 computation tree, 392
 182 concave, 37
 183 concentration inequality, 1111
 184 concentration of measure, 728
 185 concentration parameter, 1006

- 1
 2 concept drift, 717
 2 concept shift, 636, 697
 3 concrete distribution, 253
 4 condensation, 517
 5 conditional entropy bottleneck, 229
 5 conditional expected outcome, 1182
 6 conditional GAN, 865
 7 conditional generative model, 731, 734
 8 Conditional generative models, 864
 10 conditional independence, 117
 11 conditional KL divergence, 197
 11 conditional moments, 371
 12 conditional parallel trends, 1202
 12 conditional probability distribution, 118
 14 conditional probability table, 119
 15 conditional random field, 160, 1097
 16 conditional random fields, 139, 140
 16 Conditional shift, 697
 17 conditional value at risk, 1131
 18 conditionally conjugate, 80
 19 conditioner, 798, 799
 19 conditioning case, 119
 20 conditioning matrix, 241
 21 conductance, 493
 21 confidence score, 700
 22 conformal prediction, 553, 699, 700
 23 conformal score, 553
 24 conformalized quantile regression, 556
 25 confounder, 175
 26 confounders, 1179
 27 conical combination, 796
 27 conjugate, 70, 321
 28 conjugate gradients, 672
 29 conjugate prior, 20, 29, 70
 29 conjunction of features, 812
 30 conjunctive normal form, 183
 31 consensus sequence, 941
 31 conservative policy iteration, 1147
 32 consistent, 1006
 33 constant symbols, 183
 33 contact map, 144
 34 content addressable memory, 145
 35 content constrained, 726
 35 context free grammar, 163
 36 context variables, 710
 37 Context., 1065
 37 contextual bandit, 1107
 38 continual learning, 630, 703, 716
 39 continuation method, 509
 39 continuing task, 1116
 40 continuous task-agnostic learning, 719
 41 continuous-ranked probability score, 989
 42 continuous-time flows, 806
 43 contraction, 1120
 44 contrastive divergence, 158, 814
 45 contrastive learning, 1054
 45 Contrastive Multiview Coding, 1056
 46
 47
 Contrastiveness, 1079
 control, 1101
 control as inference, 1163
 control theory, 1116
 control variate, 463, 502
 control variates, 250, 429
 controller, 1116
 converge, 493
 conversions, 1105
 convex BP, 395
 convex combination, 71
 ConvNeXt, 601
 convolution, 593
 convolutional layer, 594
 convolutional Markov model, 785
 convolutional neural network, 601, 784
 convolutional neural networks, 165
 cooling schedule, 510
 cooperative cut, 303
 coordinate ascent variational inference, 411
 coresets, 663
 correlated topic model, 921
 correlation coefficient, 14
 correspondence, 970
 cosine distance, 27
 cosine kernel, 646
 count-based exploration, 1134
 Counterfactual queries, 1176
 counterfactual question, 189
 counterfactual reasoning, 984
 coupled HMM, 956
 coupling flows, 797
 coupling layer, 797
 covariance function, 639
 covariance graph, 175, 220
 covariance matrix, 13
 covariate shift, 696
 coverage, 554
 Cox process, 661
 CPD, 118
 CPPN, 725
 CPT, 119
 CQR, 556
 credible interval, 66, 78
 CRF, 139, 160
 CRFs, 140
 critic, 848
 critical temperature, 142, 484
 cross correlation, 593
 cross entropy, 204, 211
 cross entropy method, 285, 1151
 cross fitting, 1186
 cross validation, 107
 cross-entropy method, 287, 537
 crossover, 283
 crowd computing, 1123
 CRP, 1010
 CRPS, 989
 CTR, 1108
 cubature, 451
 cubature Kalman filter, 367
 cubatures, 367
 cumulants, 35
 cumulative distribution function, 7
 cumulative regret, 1114
 cumulative reward, 1107
 curse of dimensionality, 736, 1120
 curse of horizon, 1160
 curved exponential family, 30
 CV, 429
 CVI, 268
 cycle consistency, 873
 cyclical annealing, 770
 d-separated, 123
 D4PG, 1149
 DAGs, 117
 DALL-E, 788
 DALL-E 2, 790, 842
 damped updates, 415
 damping, 393, 447
 dark knowledge, 622
 DARN, 121
 data assimilation, 361
 data association, 970
 data augmentation, 584, 704
 data cleaning, 699
 Data compression, 221
 data compression, 193, 735, 740
 data generating process, 1, 696
 data processing inequality, 202, 213
 data tempering, 513
 data-driven MCMC, 472
 dataset shift, 693
 daydream phase, 444
 DBL, 607
 DBN, 148, 957
 DCGAN, 867
 DDIM, 842
 DDP, 1151
 DDPG, 1149
 de Finetti's theorem, 64
 dead leaves, 890
 decision diagram, 1098
 decision nodes, 1098
 decision tree, 1100
 declarative approach, 185
 decoder, 602, 747, 883
 decomposable, 157, 169, 403
 decompose, 131
 decoupled EKF, 633
 deep autoregressive network, 121
 deep Bayesian learning, 607
 deep belief network, 148
 deep Boltzmann machine, 148
 deep Boltzmann network, 148
 deep CCA, 909
 deep deterministic policy gradient, 1149
 deep ensembles, 618
 Deep Factors, 989
 deep fakes, 734, 871
 deep Gaussian process, 689
 deep generative model, 121
 deep generative models, 731
 deep image prior, 611
 Deep kernel learning, 684
 deep latent Gaussian model, 747
 deep latent variable model, 747
 deep learning, 1, 591
 deep Markov model, 990

- 1
 2 deep neural network, 591
 3 deep PILCO, 1153
 4 deep Q-network, 1140
 5 deep state-space model, 989
 6 deep submodular function, 303
 7 deep unfolding, 399
 8 DeepAR, 989
 9 DeepGLO, 989
 10 DeepSSM, 989
 11 default prior, 88
 12 deformable parts model, 166
 13 degenerate kernel, 649
 14 degree of normality, 8
 15 degrees of freedom, 8, 75, 110
 16 deleted interpolation, 103
 17 delta function, 67
 18 delta method, 253
 19 delta VAE, 770
 20 demand forecasting, 982
 21 denoising autoencoder, 1052
 22 denoising diffusion GAN, 842
 23 denoising diffusion models, 831
 24 denoising diffusion probabilistic
 model, 838
 25 denoising diffusion probabilistic
 models, 822
 26 Denoising Score Matching, 819
 27 DenseFlow, 733
 28 density estimation, 735
 29 density model, 701
 30 density ratio estimation, 62
 31 Derivative free optimization, 279
 32 derivative function, 233
 33 derivative operator, 233
 34 detailed balance, 469
 35 detailed balance equations, 55
 36 determinantal point process, 343
 37 Determinantal point processes,
 1035
 38 determinantal projection point pro-
 cesses, 1037
 39 deterministic ADVI, 435
 40 deterministic annealing, 947
 41 deterministic inducing conditional,
 666
 42 deterministic policy gradient theo-
 rem, 1148
 43 deterministic training conditional,
 667
 44 deviance, 110
 45 DFO, 279
 46 DGP, 1, 689
 47 diagonal covariance matrix, 14
 48 diameter, 389
 49 DIC, 666
 50 dictionary learning, 931
 51 diffeomorphism, 792
 52 difference in differences, 1200
 53 differentiable CEM, 287
 54 differentiable simulators, 846
 55 differential dynamic programming,
 1151
 56 differential entropy, 208
 57 diffuse prior, 88
 58 diffusion matrix, 491
 59 diffusion models, 731, 831
 60 diffusion process, 831
 61 diffusion term, 490, 490
 62 DiffWave, 842
 63 digamma function, 422
 64 dilated convolution, 785
 65 diminishing returns, 300
 66 DINO, 1058
 67 direct coupling analysis, 144
 68 direct method, 1157
 69 directed acyclic graphs, 117
 70 directed Gaussian graphical model,
 122
 71 Dirichlet, 25
 72 Dirichlet distribution, 72
 73 Dirichlet process, 997, 999, 1005
 74 Dirichlet process mixture models,
 424
 75 discount factor, 1107, 1117, 1131
 76 discount parameter, 1018
 77 discrete task-agnostic learning,
 719
 78 discrete with probability one, 1008
 79 discriminative model, 543
 80 discriminative reranking, 343
 81 discriminator, 848
 82 disease mapping, 662
 83 disease transmission, 995
 84 disentangled, 758, 930, 1059
 85 disentangled representation learn-
 ing, 1042
 86 dispersion parameter, 38
 87 distill, 842
 88 distillation, 622
 89 distortion, 222
 90 distributed representation, 147,
 955
 91 distribution free, 553
 92 distribution shift, 693, 727, 1156
 93 distributional particles, 524
 94 distributional RL, 1142, 1149
 95 distributionally robust optimiza-
 tion, 704
 96 diverged, 487
 97 divergence metric, 55
 98 DLGM, 747
 99 DLM, 977
 100 DLVM, 747
 101 DM, 831
 102 DNN, 591
 103 do calculus, 1172
 104 do-calculus, 1213
 105 do-notation, 1176
 106 domain adversarial learning, 707
 107 domain drift, 716
 108 domain generalization, 636, 710
 109 domain randomization, 1155
 110 domain shift, 696
 111 domains, 710
 112 donor, 986
 113 Donsker Varadhan lower bound,
 219
 114 Donsker-Varadhan, 202
 115 dot product attention, 597
 116 double DQN, 1142
 117 double loop algorithms, 393
 118 double machine-learning, 1184
 119 double Q-learning, 1140
 120 double sided exponential, 9
 121 doubly intractable, 156
 122 doubly reparameterized gradient
 estimator, 441
 123 doubly robust, 1159, 1186
 124 doubly stochastic, 428
 125 downstream, 739
 126 DPGM, 117
 127 DPPs, 1035
 128 DQN, 1140
 129 DRE, 62
 130 Dreamer, 1155
 131 dropout, 595
 132 DTC, 667
 133 dual EKF, 362
 134 DualDICE, 1161
 135 dueling DQN, 1142
 136 Dutch book, 65
 137 dyna, 1151
 138 dynamic Bayesian network, 957
 139 dynamic embedded topic model,
 922
 140 dynamic linear model, 963, 977
 141 dynamic programming, 322, 383,
 1119
 142 dynamic programming., 405
 143 dynamic topic model, 921
 144 dynamic VAE, 766
 145 dynamical variational autoen-
 coders, 989
 146 E step, 258, 259
 147 EA, 282
 148 earning while learning, 1107
 149 EB, 101
 150 EBM, 731, 811
 151 ECE, 547
 152 ECM, 265
 153 ECME, 265
 154 EDA, 285
 155 edge potentials, 151
 156 edit distance, 942
 157 effective dimensionality, 625
 158 effective sample size, 497, 499,
 521
 159 effects of causes, 189
 160 EI, 276
 161 eigenfunction, 649
 162 eigengap, 493
 163 eight schools, 98
 164 Einstein summation, 406
 165 einsum, 403, 406
 166 einsum networks, 179
 167 EKF, 357
 168 elastic weight consolidation, 631,
 721
 169 ELBO, 258, 325, 410, 749
 170 elementwise flow, 794
 171 eligibility traces, 1138
 172 elimination order, 401
 173 elite set, 283
 174 ELPD, 108
 175 EM, 258
 176 EMNA, 286
 177 empirical Bayes, 101, 610
 178 empirical distribution, 204
 179 empirical Fisher, 245
 180 empirical MBR, 1098

- 1
- 2 empirical risk, 544
empirical risk minimization, 544
- 3 empirical risk minimization, 265
- 4 emulate, 537
encoder, 602
- 5 End-task., 1065
- 6 endogeneous, 186
energy, 140
- 7 energy based models, 731
- 8 energy disaggregation, 955
energy function, 282, 324, 811, 1034
- 10 energy score, 700
energy-based model, 140
- 11 Energy-based models, 811
- 12 EnKF, 361
ensemble, 595, 1126
- 13 ensemble Kalman filter, 361
- 14 entity resolution, 181
- 15 entropy, 208, 258
entropy sampling, 1125
- 16 entropy search, 277
environment, 1107, 1115
- 17 environments, 710
- 18 EP, 445
epidemiology, 518
- 19 episodic task, 1116
- 20 epistemic uncertainty, 67, 550
epistemological uncertainty, 997
- 21 EPLL, 889
- 22 epsilon-greedy, 1134
- 23 episode, 1116
- 24 equal odds, 1098
- 25 equal opportunity, 1098
- 26 equilibrium distribution, 52
- 27 ergodic, 54
- 26 ERM, 265, 544
error correcting codes, 225, 396
- 27 error correction, 193
error-correcting codes, 178
- 29 ES-RNN, 988
- 29 ESS, 497, 521
estimated potential scale reduction, 497
- 31 estimation of distribution, 285
- 32 estimation of multivariate normal algorithm, 286
- 33 estimator, 63
- 34 Etsy, 942
- 35 EUBO, 441
Euler approximation, 501
- 36 Euler's method, 486, 807
evidence, 64, 106, 325, 749
- 37 evidence lower bound, 258, 325, 410, 749
- 38 evidence maximization, 610
evidence upper bound, 441
- 40 Evolution strategies, 287
- 41 evolutionary algorithm, 282
evolutionary programming, 285
- 42 evolutionary search, 279
- 43 evolutionary strategies, 272
- 44 EWC, 631
excess kurtosis, 9
- 45 exchangeable, 136, 922
exchangeable process, 1010
- 46 exchangeable with, 95
- 47
- Exclusion Restriction, 1193
exclusive KL, 199
execution traces, 185
exogenous, 186
exp-sine-squared kernel, 646
expanded parameterization, 887
expectation backpropagation, 633
expectation maximization, 258
expectation propagation, 380, 445
expected calibration error, 547
expected complete data log likelihood, 259
expected free energy, 1166
expected improvement, 276
expected LPPD, 108
expected patch log likelihood, 889
expected sufficient statistics, 134, 259
- experience replace, 1140
experience replay, 721
explainability, 189
explaining away, 80, 120, 125, 129, 216
explicit duration HMM, 952
explicit layers, 600
explicit probabilistic models, 845
exploration bonus, 1111, 1135
exploration-exploitation tradeoff, 1102, 1108, 1134
exponential cooling schedule, 510
exponential dispersion family, 38
exponential distribution, 10
exponential family, 29, 30, 39, 86
exponential family factor analysis, 910
exponential family harmonium, 147
exponential family PCA, 910
exponential family state-space model, 974
Exponential linear unit, 593
exponentiated quadratic, 642
extended Kalman filter, 357, 358, 974
extended Kalman smoother, 359
extended particle filter, 522
extended RTS smoother, 359
external field, 143
external validity, 695
extrapolation, 675
extrinsic variables, 419
- f-divergence, 56
f-Divergence Max-Ent IRL, 1168
Facebook, 998
factor, 139
factor analysis, 137
factor graph, 176, 390, 397
factor loading matrix, 137, 893
factor of variation, 758
factor rotations problem, 897
factorial HMM, 438, 955
factorization property, 128
FAIRL, 1168
fairness, 189, 1098
Faithfulness, Fidelity, 1077
family marginal, 134
- fan-in, 238
fan-out, 235, 238
fantasy data, 816
Fast Fourier Transform, 340
fast geometric ensembles, 616
fast gradient sign, 724
fast ICA, 929
fast weights, 620
FastSLAM, 531
FB, 337
feature induction, 150
feature-based, 301
feedback loop, 1109
feedforward neural network, 600
ferromagnetic, 141
few-shot learning, 705
Feynman-Kac, 511
FFBS, 335
FFG, 177
FFJORD, 807
FFNN, 600
FGS, 724
FIC, 668
FID, 742
fill-in edges, 402
FiLM, 600
filter, 593
filter response normalization, 595
filtering distribution, 331
filtering SMC, 538
filtering variational objective, 538
FIM, 39
fine-tuning, 1041
finite horizon, 1107
finite horizon problem, 1117
finite state machine, 1116
finite sum objective, 249
finite-state Markov chain, 48
first-order delta method, 253
first-order logic, 182
Fisher divergence, 818
Fisher information, 90
Fisher information matrix, 36, 39, 90
- FITC, 667
fitness, 283
fitness proportionate selection, 283
fitted value iteration, 1142
FIVO, 538
fixed effects, 587
Fixed lag smoothing, 345
fixed-form VI, 427
fixed-lag smoothing distribution, 331
flat minima, 491, 623
flow cytometry, 737
folded over, 7
folds, 107
FOO-VB, 635
fooling images, 725
force, 504
forest plot, 97
fork, 123
Forney factor graph, 177
Forney factor graphs, 176
forward adversarial inverse RL, 1168

- 1
 2 forward transfer, 720
 3 forward-mode automatic differentiation, 236
 4 forwards algorithm, 404, 941
 5 forwards filtering backwards sampling, 343
 6 forwards filtering backwards smoothing, 335
 7 forwards kernel, 532
 8 forwards mapping, 881
 9 forwards process, 831
 10 forwards-backwards, 337, 405
 11 founder variables, 898
 12 Fourier basis, 931
 13 Fourier transform, 650
 Fréchet Inception Distance, 742
 fragmentation, 1020
 Frechet inception distance, 61
 free bits, 771
 free energy, 410
 free energy principle, 1166
 free-form VI, 413
 freeze-thaw algorithm, 279
 frequentist sampling distribution, 565
 frequentist statistics, 63
 friction, 504
 front-door criterion, 1215
 front-door formula, 1215
 frustrated, 484
 frustrated system, 142
 full conditional, 127, 473
 full conditionals, 327
 full conformal prediction, 554
 full covariance matrix, 14
 fully connected CRF, 165
 fully connected layer, 592
 fully independent conditional, 668
 fully independent training conditional, 667
 function space, 654
 functional, 231
 functional causal model, 186
 fundamental problem of causal inference, 190
 funnel shape, 99
 funnel transformer, 767
 fuzzy clustering, 997
- 34
 35 g-prior, 565
 GA, 283
 GAE, 1146
 GAIL, 1168
 GAM, 680
 Gamma, 952
 gamma distribution, 10
 gamma function, 8
 GAN, 731
 GANs, 845
 GaP, 914
 gap, 1114
 gated recurrent unit, 604
 Gauss-Hermite integration, 365, 367
 Gauss-Hermite Kalman filter, 367
 Gaussian bandit, 1110, 1112
 Gaussian Bayes net, 122
- 47 Gaussian copula, 989
 Gaussian distribution, 7
 Gaussian filtering, 374
 Gaussian graphical model, 151
 Gaussian kernel, 643
 Gaussian mixture model, 884
 Gaussian MRF, 151
 Gaussian process, 274, 639, 1005
 Gaussian processes, 689, 1152
 Gaussian scale mixture, 262, 569, 886
 Gaussian soap bubble, 14
 Gaussian sum filter, 375
 Gaussian VI, 433
 GELU, 593
 GEM, 264, 1009
 Gen, 185
 GenDICE, 1161
 general Gaussian filter, 366
 Generality, 1079
 generalization, 740
 generalized additive model, 680, 988
 generalized advantage estimation, 1146
 generalized Bayesian inference, 267, 545
 generalized belief propagation, 395
 generalized CCA, 909
 generalized EM, 264
 generalized Gauss-Newton, 613
 generalized IPLF, 372
 generalized linear mixed model, 587
 generalized linear model, 557
 generalized low rank models, 910
 generalized policy improvement, 1122
 generalized pseudo Bayes filter, 375
 generalized statistical linear regression, 372
 generalized variational continual learning, 635
 generate and test, 472
 generative adversarial imitation learning, 1168
 Generative Adversarial Networks, 845
 generative adversarial networks, 731
 generative design, 739
 generative model, 731
 generative neural samplers, 846
 Generative representation learning, 1049
 generative weights, 926
 generator networks, 846
 genetic algorithm, 283
 genetic programming, 283
 geometric distribution, 7, 951
 geometric path, 533
 GGF, 366
 GGM, 151
 GGN, 613
 GHKF, 367
 Gibbs distribution, 140, 811
- Gibbs point processes, 1034
 Gibbs posterior, 545
 Gibbs sampling, 142, 327, 473
 GIPLF, 372
 gist, 922
 Gittins index, 1110
 Glauber dynamics, 473
 GLIDE, 790
 GLIE, 1139
 GLM, 557
 GLM predictive distribution, 621
 GLMM, 587
 global balance equations, 53
 global latent variables, 320
 global Markov property, 123, 153
 global variables, 131
 globally normalized, 161, 172
 Glorot initialization, 608
 GMM, 884
 GMM-HMM, 936
 GNNs, 605
 goodness-of-fit, 55
 GP, 274
 GP-LVM, 912
 GP-MPC, 1153
 GP-UCB, 276
 GPT, 787
 GPT-2, 787
 GPT-3, 787
 gradient descent, 241
 gradient EM algorithm, 265
 gradient sign reversal, 708
 gradient-based meta-learning, 715
 Gram matrix, 642
 grammar VAE, 769
 grammars, 163
 graph neural networks, 605
 graph surgery, 188, 1173
 graphical lasso, 151, 1002
 greatest common divisor, 54
 greedy action, 1119
 greedy search, 280
 grid approximation, 322
 grid search, 281
 grid world, 1119
 ground network, 180
 ground set, 1035
 ground states, 142
 ground terms, 180
 group lasso, 570
 group normalization, 595
 GRU, 604
 GSM, 886
 guided cost learning, 1167
 guided particle filter, 522
 Gumbel distribution, 253
 Gumbel-Max trick, 253
 Gumbel-Softmax, 781
 Gumbel-Softmax distribution, 253
- half Cauchy, 9
 half-Cauchy distribution, 571
 half-edge, 177
 half-normal distribution, 7
 Halton sequence, 465
 Hamilton's equations, 485
 Hamiltonian, 485

- 1
 2 Hamiltonian mechanics, 484
 3 Hamiltonian Monte Carlo, 327, 484, 503, 581
 4 Hamiltonian Variational Inference, 439
 5 Hammersley-Clifford theorem, 140
 6 Hamming distance, 69
 7 Hamming loss, 69
 7 hard clustering, 885
 8 hard EM, 264
 9 hard negative mining, 1056
 9 hard tanh, 254
 10 hardcore repulsive process, 1035
 10 harmonic mean estimator, 106
 11 harmonium, 147
 12 Hastings correction, 468
 12 Hawkes processes, 1032
 13 hazard function, 958
 14 hBOA, 286
 14 HDP-HMM, 949, 1020
 15 heads, 601
 16 heat bath, 473
 16 heavy tailed, 11
 17 heavy tails, 11
 18 Hebb's rule, 815
 19 Hebbian learning, 815
 20 Hellinger distance, 57
 21 Helmholtz machine, 748
 21 Hessian free optimization, 246
 22 heteroskedastic, 601
 22 heuristic function, 1150
 23 heuristic search, 1150
 24 hidden Gibbs random field, 159
 24 hidden Markov model, 331, 783, 934
 26 hidden semi-Markov model, 952
 26 hidden state, 783
 27 hidden variables, 129, 258
 28 hierarchical Bayesian model, 94, 586, 635
 29 hierarchical Bayesian models, 500
 30 hierarchical Dirichlet process, 1019
 31 hierarchical GLM, 587
 32 hierarchical HMM, 953
 33 hierarchical kernel learning, 679
 33 hierarchical VAE, 773
 34 hierarchical variational model, 439
 35 Hilbert space, 657
 35 hill climbing, 280
 36 hindsight, 337
 36 hinge-loss MRFs, 184
 37 Hinton diagram, 50
 38 Hinton diagrams, 904
 39 histogram binning, 548
 39 HMC, 97, 327, 484
 40 HMM, 331, 934
 41 HMM filter, 346
 41 homogeneous, 47
 42 homogeneous Poisson process, 1030
 43 Hopfield network, 144
 44 horseshoe distribution, 887
 45 horseshoe prior, 570
 45 HSMM, 952
 46 Huffman coding, 221
 47
- Hungarian algorithm, 971
 Hutchinson trace estimator, 805
 Hutter prize, 221
 hybrid MC, 484
 hybrid system, 968
 hypernetwork, 599
 hyperparameters, 94
 hypothesis test, 1104
- I-map, 123
 I-projection, 200
 IAF, 803
 IBIS, 535
 ICA, 927
 ICM, 145
 ID, 699
 identifiable, 1059, 1174
 identification strategy, 1174
 identified, 1174
 identity uncertainty, 181
 image captioning, 734
 image deblurring, 889
 image denoising, 889
 image imputation, 907
 image inpainting, 889
 image super-resolution, 889
 image to image translation, 872
 Image-text supervision, 1048
 image-to-image, 734
 image-to-text, 734
 Imagen, 790, 842
 imagination-augmented agents, 1151
 Imbens' Sensitivity Model, 1205
 Imitation learning, 1166
 IMM, 377
 implicit generative models, 741, 845
 implicit layers, 600
 implicit models, 537, 731
 implicit probabilistic model, 845
 implicit probabilistic models, 845
 implicit probability distributions, 439
 implicit probability model, 760
 importance ratio, 1158
 importance sampling, 326, 458
 importance weighted autoencoder, 440
 importance weighted autoencoders, 460
 importance weights, 459
 imputation, 736
 in-context learning, 705, 787
 in-distribution, 699
 in-domain, 694
 in-painting, 736
 Inception, 61
 Inception Score, 742
 inclusive KL, 198
 income inequality, 12
 incomplete data, 159
 incremental EM, 265
 incremental importance weights, 515
 incumbent, 275
 independence sampler, 470
- independent components analysis, 927
 induced width, 403
 inducing points, 665
 inference, 64, 129, 319
 inference compilation, 437
 inference network, 437, 591, 747
 infinite hidden relational model, 999
 infinite horizon, 1107
 infinite mixture model, 1012
 infinite relational model, 997, 999
 infinite-state hidden Markov models, 1020
 infinitely divisible distribution, 1029
 infinitely exchangeable, 64
 infinitely wide neural networks, 610
 influence curve, 1184
 influence diagram, 189, 1098
 influence model, 956
 infomax, 930
 InfoNCE, 219
 InfoNCE loss, 1057
 information arc, 1099
 information bottleneck, 1060
 information bottleneck principle, 226
 information criterion, 110
 information diagram, 212, 216
 information diagrams, 213, 1128
 information extraction, 162
 information filter, 350
 information form, 16, 33, 151
 information gain, 195, 277, 1126
 information processing, 193
 information projection, 200, 374
 information state, 1109
 informative vector machine, 664
 InfoVAE, 759, 772
 inhomogeneous Poisson process, 1030
 injective, 44
 inner product, 656
 innovation term, 346
 input nodes, 238
 inside outside, 955
 inside-outside algorithm, 163
 instance normalization, 595
 instantaneous ELBO, 437
 instrument monotonicity, 1196
 Instrument Relevance, 1193
 Instrument Unconfoundedness, 1193
 instrumental variables, 1192
 integer-valued random measure, 1028
 integral probability metric, 57
 integrating out, 66
 inter-causal reasoning, 125
 interaction information, 215
 interactive multiple models, 377
 Interactivity., 1078
 interpolated Kneser-Ney, 105
 interpolator, 651
 intervention, 190

- 1
 2 interventions, 188, 1172
 3 Interventional queries, 1176
 4 intrinsic uncertainty, 67
 5 intrinsic variables, 419
 6 invariant, 90, 469
 7 invariant causal prediction, 712
 8 invariant distribution, 52
 9 invariant risk minimization, 712
 10 inventory, 982
 11 inverse autoregressive, 802
 12 inverse chi-squared distribution, 75
 13 inverse Gamma, 74, 563
 14 inverse Gamma distribution, 11
 15 inverse mass matrix, 485, 488
 16 inverse of a partitioned matrix, 17
 17 inverse optimal control, 1167
 18 inverse probability of treatment weighted estimator, 1183
 19 inverse probability theory, 63
 20 inverse probability transform, 454
 21 inverse reinforcement learning, 1167
 22 inverse temperature, 533
 23 inverse Wishart, 29, 79, 80
 24 IPF, 157
 25 IPLF, 370
 26 IPLS, 370
 27 IPM, 57
 28 iResNet, 804
 29 IRLS, 257, 561
 30 IRM, 712, 999
 31 irreducible, 53
 32 Ising model, 141, 143
 33 Ising models, 474
 34 isotonic regression, 548
 35 isotropic covariance matrix, 14
 36 iterated batch importance sampling, 535
 37 iterated EKF, 359, 370
 38 iterated EKS, 360, 370
 39 iterated posterior linearization filter, 370
 40 iterated posterior linearization smoother, 370
 41 iterative amortized inference, 438
 42 iterative conditional modes, 145
 43 iterative proportional fitting, 157
 44 iteratively reweighted least squares, 561
 45 IWAE, 440
 46 IWAE bound, 440
 47 Jacobi, 394
 48 Jacobian, 45
 49 Jacobian determinant, 792
 50 Jacobian vector product, 821
 51 Jacobian-vector product (JVP), 233
 52 JAGS, 473
 53 JamBayes, 1002
 54 Jeffrey's conditionalization rule, 205
 55 Jeffreys prior, 77, 90
 56 Jensen's inequality, 196, 258, 440
 57 JMLB, 968
- join tree, 405
 JPDA, 971
 JTA, 405
 judge fixed effects, 1194
 jump Markov linear system, 968
 junction tree, 405
 junction tree algorithm, 405
- K-means clustering, 886
 Kalahari, 566
 Kalman Filter, 346
 Kalman filter, 344, 350, 962, 979
 Kalman filter algorithm, 20
 Kalman gain matrix, 20, 346
 Kalman smoother, 344, 962
 Kalman smoothing, 350, 921
 KDE, 736
 kernel, 468, 593
 kernel basis function expansion, 573
 kernel density estimation, 736
 kernel function, 641, 657
 Kernel Inception Distance, 743
 kernel inception distance, 61
 kernel mean embedding, 59
 kernel ridge regression, 656, 658
 kernel trick, 59, 649
 keys, 597
 KF, 346
 KFAC, 245, 613
 kinetic energy, 485
 KISS, 680
 KISS-GP, 674
 KL annealing, 770
 KL divergence, 195, 259
 knots, 796
 knowledge gradient, 278
 known knowns, 699
 known unknowns, 699
 Kolmogorov-Smirnov test, 46
 kriging, 640
 Kronecker FActored Curvature, 613
 Krylov subspace methods, 672
 Kullback Leibler divergence, 56
 Kullback-Leibler divergence, 195, 259
 kurtosis, 9, 929
- L-ensembles, 1036
 L0 norm, 568
 L0 regularization, 568
 L2, 657
 L2 loss, 1096
 Lévy processes, 1029
 Lévy subordinators, 1029
 Lévy-Itô decomposition, 1030
 Lévy-Kintchine, 1030
 label bias, 171
 label shift, 697
 label smoothing, 549
 label switching, 948
 ladder network, 774
 lag, 345
 Lagrange multipliers, 38
 Lagrangian, 38
- lambda-return, 1138
 Langevin diffusion, 490, 501
 Langevin MCMC, 814
 Langevin Monte Carlo, 489
 language models, 49
 LapGAN, 869
 Laplace approximation, 323, 578
 Laplace distribution, 9
 Laplace Gaussian filter, 523
 Laplace propagation, 446
 lasso, 568, 569
 latent Dirichlet allocation, 917
 latent factor models, 883
 latent factors, 881
 latent overshooting, 992
 latent semantic analysis, 918
 latent space arithmetic, 738
 latent space interpolation, 738, 767
- latent variable, 883
 latent variable collapse, 776
 latent variable model, 881, 883
 layer, 592
 layer normalization, 595
 layers, 591
 lazy training, 689
 LBP, 389
 LDA, 917
 LDA-HMM, 922
 LDPC code, 396
 LDS, 343, 960
 Leaky ReLU, 593
 leapfrog integrator, 486
 learned loss function, 867
 learning, 1131
 learning from demonstration, 1166
 least confident sampling, 1125
 least mean squares, 963
 leave-one-out cross validation, 107
 LeCun initialization, 608
 left-to-right, 339
 left-to-right transition matrix, 48
 legal reasoning, 189
 leptokurtic, 9
 level sets, 14
 LG-SSM, 343, 960
 life-long learning, 716
 lifted inference, 184
 likelihood function, 64
 likelihood ratio, 106, 701
 likelihood ratio gradient estimator, 250
 likelihood tempering, 513
 likelihood-free inference, 537, 846
 lily pads, 936
 limiting distribution, 54
 linear discriminant analysis, 917
 linear dynamical system, 343, 960
 Linear evaluation, 1040
 linear flow, 794
 linear Gaussian CPD, 122
 linear Gaussian state space model, 343
 linear Gaussian system, 19
 linear layer, 592
 linear programming, 1119
 Linear regression, 562

- 1
 2 linear regression bandit, 1110
 2 linear-Gaussian CPD, 174
 3 linear-Gaussian state-space model,
 960
 4 linear-quadratic-Gaussian, 1151
 5 linearization point, 233
 6 link function, 557, 560
 7 linkage learning, 286
 7 Lipschitz constant, 58
 8 LKJ distribution, 88
 9 LM, 49
 9 LMC, 489
 10 LMS, 963
 11 local and global latent variables,
 321
 12 local average treatment effect,
 1196
 13 local evidence, 414
 14 local factor, 445
 15 local latent variables, 320
 15 local level model, 977
 16 local linear trend, 978
 16 local Markov property, 128
 17 local variables, 131
 18 local+global, 989
 19 localist representation, 147
 19 locally normalized, 161, 170, 172
 20 locally optimal proposal distribu-
 tion, 522
 21 location-scale family, 93
 22 log derivative trick, 250, 428
 22 log loss, 544, 546
 23 log partition function, 30
 24 log-derivative trick, 158
 24 log-linear, 149
 25 log-linear model, 149
 26 log-odds score, 941
 26 log-pointwise predictive-density,
 108
 28 log-returns, 975
 28 log-sum-exp trick, 892
 29 logistic, 575
 30 logistic distribution, 583, 928
 30 logistic normal, 921
 31 Logistic regression, 575
 32 logistic regression bandit, 1110
 32 logits, 576
 33 long short term memory, 604
 34 long tail, 70
 34 long tails, 11
 35 LOO-CV, 107
 36 lookahead function, 538
 36 loopy belief propagation, 389, 398
 37 Lorentz, 9
 38 loss function, 1095
 38 lossless compression, 221
 39 lossy compression, 221
 40 low discrepancy sequences, 465
 40 low variance resampling, 520
 41 low-density parity check code, 396
 42 low-resource languages, 944
 42 LPPD, 108
 43 LQG, 1151
 44 LSTM, 604
 45 LVM, 883
 46 M step, 258
 47
- M-projection, 199
 m-separation, 174
 M2, 765
 M4 forecasting competition, 988
 machine translation, 735
 MADE, 784
 MAE, 1053
 MAF, 802
 Mahalanobis distance, 13
 majorize-minimize, 255
 MALA, 489
 MAML, 715
 manifestation shift, 697
 MAP estimate, 65, 340, 544, 1096
 MAPIE, 553
 MAR, 764
 margin sampling, 1125
 marginal calibration error, 547
 marginal likelihood, 64, 101, 106,
 201
 marginalizing out, 66
 Mariner 10, 350
 marked, 1034
 Markov, 123
 Markov assumption, 47, 783
 Markov blanket, 127, 153, 473
 Markov chain, 47
 Markov chain Monte Carlo, 326,
 467
 Markov decision process, 1115
 Markov kernel, 47
 Markov logic network, 182
 Markov mesh, 139
 Markov model, 47, 783
 Markov model of order n, 49
 Markov network, 139
 Markov random field, 139
 Markovian SCM, 186
 masked attention, 598
 masked autoencoder, 1053
 masked autoregressive flow, 802
 masked convolution, 785
 masked language modeling, 1052
 matching, 1187
 Matern kernel, 644
 matrix determinant lemma, 805
 matrix inversion lemma, 17, 21,
 346
 matrix normal, 28
 matrix normal inverse Wishart,
 574
 matrix vector multiplication, 672
 max margin Markov networks, 168
 max marginals, 69, 387
 max-product belief propagation,
 387
 maxent prior, 89
 maximal clique, 139
 maximal weight bipartite match-
 ing, 971
 maximization bias, 1139
 maximizer of posterior marginals,
 69
 maximizer of the max marginal,
 387
 maximizer of the posterior
 marginal, 387
- maximum a posteriori, 1096
 maximum entropy, 89
 maximum entropy Markov model,
 161, 171
 maximum entropy model, 38, 149
 maximum entropy RL, 1164
 maximum expected utility prin-
 ciple, 1095
 maximum likelihood estimation,
 544
 maximum mean discrepancy, 58,
 58, 760, 853
 MBIE, 1135
 MBIE-EB, 1135
 MBRL, 1149
 MCAR, 764
 MCEM, 264
 MCMC, 326, 467
 MCTS, 1150
 MDL, 110
 MDP, 1115
 mean canonical correlation, 1044
 mean field, 326, 411
 mean function, 557, 639
 mean squared canonical corre-
 lation, 1044
 mean value imputation, 736
 measure, 657
 measurement step, 346
 mediators, 1215
 membership query synthesis, 1123
 memetic algorithm, 285
 MEMM, 171
 MEMO, 709
 memory methods, 721
 memorylessness, 1032
 Mental Model, 1078
 Mercer kernel, 639, 641
 Mercer's theorem, 649
 merit function, 275
 MERL, 1164
 message passing, 383
 message passing algorithms, 383
 message passing schedule, 383
 messages, 338, 383
 meta-data, 710
 meta-learning, 713
 Method., 1066
 metric learning, 1056
 Metrics., 1066
 Metropolis Adjusted Langevin Al-
 gorithm, 489
 Metropolis Hastings, 282, 468, 508
 Metropolis Hastings algorithm,
 326
 Metropolis within Gibbs, 478
 MH, 468
 midi format, 787
 min-fill heuristic, 403
 min-max, 856
 min-max optimization problem,
 704
 min-weight heuristic, 403
 minibatch, 249
 minimal, 30
 minimal I-map, 123
 minimal representation, 31

- 1
2 minimal sufficient statistic, 214,
 226, 226
3 minimally informative prior, 88
4 minimum Bayes risk, 1097
5 minimum description length, 110
6 minimum mean squared error,
 1096
7 minorize-maximize, 255
8 missing at random, 764
9 missing completely at random, 764
10 missing data, 258, 260, 763, 999
11 missing data mechanism, 764
12 mixed effects model, 587
13 mixed graph, 173
14 mixed membership model, 915,
 917
15 mixed membership stochastic block
model, 997
16 mixing matrix, 926
17 mixing time, 467, 492, 493
18 mixture model, 883
19 mixture of Bernoullis, 886
20 mixture of experts, 620, 812, 997
21 mixture of factor analysers, 900
22 mixture of Gaussians, 884
23 mixture of Kalman filters, 525
24 mixture proposal, 471
25 MLE, 544
26 MLP, 600
27 MM, 255
28 MMD, 58, 58, 760, 853
29 MMD VAE, 760
30 MMI, 215
31 MMM, 387
32 MMSE, 1096
33 MNIST, 738
34 Mobius inversion formula, 217
35 MoCo, 1058
36 mode, 1096
37 mode collapse, 859
38 mode connectivity, 625
39 mode hopping, 860
40 mode-covering, 199
41 mode-seeking, 199
42 model checking, 112
43 model predictive control, 1150
44 model-agnostic meta-learning, 715
45 model-based approach, 1
46 model-based RL, 1131, 1133,
 1149
47 model-free RL, 1131
48 Modified Euler's method, 486
49 Modularity, 1078
50 MoG, 884
51 molecular graph structure, 768
52 moment matching, 38, 103, 157,
 199, 366, 376, 446,
 854
53 moment parameters, 16, 30
54 moment projection, 199, 374
55 monference, 748
56 monks, 998
57 Monte Carlo, 69
58 Monte Carlo approximation, 326
59 Monte Carlo control, 1136
60 Monte Carlo dropout, 596, 611
61 Monte Carlo EM, 264
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
999
- mutual information, 212
MuZero, 1151
MVAE, 761
MVM, 672
MVN, 13
myopic, 1117, 1119
N-BEATS, 989
N-best list, 343
n-gram model, 49
NADE, 784
NAGVAC, 614
NAGVAC-1, 433
naive Bayes classifier, 137
named entity extraction, 162
named variable, 231
nats, 198, 208
natural evolution strategies, 288
natural evolutionary strategies, 246
natural exponential family, 30
natural gradient, 43, 243, 929
natural gradient descent, 39, 242, 267, 1147
Natural Gradient Gaussian variational approximation, 433
natural language processing, 161
natural parameters, 16, 30, 33
Neal's funnel, 500, 590
nearest neighbor data association, 971
NEF, 30
negative binomial, 952
negative binomial distribution, 6
negative ELBO, 749
negative log likelihood, 544, 740
negative phase, 157
negative transfer, 710
Negative-free representation learning, 1057
negentropy, 929
neighbors, 280
nested dissection, 403
nested dissection order, 403
nested plates, 138
nested SMC, 524
neural architecture search, 611
neural auto-regressive density estimator, 784
neural bandit, 1110
neural CRF parser, 163
neural enhanced BP, 399
neural net kernel, 687
neural network Gaussian process, 610
Neural ODE, 807
neural process, 715
neural spike trains, 523, 974
neural tangent kernel, 688
neural-linear, 614
neuroevolution, 285
neuron, 592
NeuTra HMC, 488
neutral process, 1028
Newton's method, 241
NGD, 242

- 1
 2 NICE, 808
 3 NIW, 81
 3 NIX, 76
 4 NLDS, 973
 5 NLG-SSM, 973
 5 NLP, 161
 6 NMAR, 764
 7 NMF, 914
 7 NN-GP, 686
 8 No-U-Turn Sampler, 488
 node potentials, 151
 9 Noise Conditional Score Network, 824
 10 Noise Contrastive Estimation, 825
 11 noisy channel, 225
 12 noisy channel model, 942
 noisy nets, 1142
 13 non-centered parameterization, 99,
 14 500, 590
 15 non-contrastive representation
 learning, 1057
 16 non-descendants, 128
 non-factorial prior, 573
 17 non-Gaussian SSM, 974
 18 non-linear squared flow, 795
 non-Markovian models, 512
 19 non-negative matrix factorization, 914
 20 non-null recurrent, 54
 21 non-parametric Bayesian, 1005
 22 non-parametric Bayesian models, 882
 23 non-parametric BP, 392, 531
 24 non-parametric models, 543
 25 non-parametrically efficient, 1186
 25 non-terminals, 163
 26 nondecreasing, 1029
 noninformative, 88
 27 nonlinear dynamical system, 973
 28 nonlinear factor analysis, 912
 29 nonlinear Gaussian SSM, 973
 29 nonparametric copula, 989
 30 nonparametric models, 639
 nonstationary kernel, 687
 31 normal distribution, 7
 32 normal factor graph, 177
 33 normal inverse chi-squared, 76
 33 normal inverse Gamma, 75, 563
 34 Normal-Inverse-Wishart, 1012
 35 Normal-inverse-Wishart, 81
 normalization layers, 595
 36 normalized completely random measures, 1028
 37 normalized occupancy distribution, 1123
 38 normalized random measures (NRMs), 1028
 40 normalized stable process, 1029
 41 normalized weights, 460, 514
 41 normalizes, 791
 42 normalizing flow, 936
 43 Normalizing flows, 439
 43 normalizing flows, 731, 791
 44 not missing at random, 764
 45 noun phrase chunking, 162
 46 noun phrases, 162
 46 Nouveau VAE, 777
 47
 novelty detection, 699
 NP-hard, 404
 NSSM, 974
 NTK, 689
 nuisance functions, 1184
 nuisance variables, 129, 388
 null hypothesis, 106
 numerical integration, 367, 451
 NUTS, 488
 NUV, 571
 NWJ lower bound, 219
 Nyström approximation, 665
 object detection, 165
 objective, 88
 observation function, 933
 observation model, 933, 935
 observation noise, 933, 961, 962
 observation overshooting, 992
 Occam factor, 111
 Occam's razor, 904
 occasionally dishonest casino, 333
 occlusion, 165
 off-policy, 1139
 off-policy policy-gradient, 1158
 offline reinforcement learning, 1157
 offspring, 283
 OGN, 271
 Olivetti faces dataset, 685
 on-policy, 1139
 one-armed bandit, 1107
 one-max, 285
 one-step-ahead predictive distribution, 373
 one-to-one, 44
 online advertising system, 1108
 online Bayesian inference, 962
 online EM, 265
 online EWC, 635
 Online Gauss-Newton, 271
 online learning, 630, 721
 online structured Laplace, 635
 ontological uncertainty, 997
 ontology, 1000
 OOD, 629, 699
 open class, 162
 open set recognition, 703
 open world, 972
 open world classification, 703
 open world recognition, 697
 open-universe probability models, 185
 OpenGAN, 702
 opportunity cost, 1102
 optimal action-value function, 1118
 optimal partial policy, 1121
 optimal policy, 1095, 1118
 optimal resampling, 526
 optimal state-value function, 1118
 optimal transport, 288
 optimism in the face of uncertainty, 1110
 optimization problems, 231
 optimizer's curse, 1139
 Optimus, 767
 oracle, 272
 ordered Markov property, 117, 128
 ordinal regression, 585
 Ornstein-Uhlenbeck process, 645
 orthodox statistics, 63
 orthogonal additive kernel, 680
 orthogonal Monte Carlo, 466
 orthogonal random features, 651
 OUPM, 185
 out-of-distribution detection, 699
 out-of-domain, 694
 outer product method, 145
 outlier detection, 699, 735
 outlier exposure, 699
 over-complete representation, 31
 overcomplete representation, 931
 overdispersed, 494
 overfitting, 544
 overlap, 1181
 Pólya urn, 1009
 PAC-Bayes, 545, 626
 padding, 594
 PageRank, 52
 paired data, 864
 pairwise Markov property, 153
 pairwise potentials, 1034
 panel data, 983, 1200
 parallel prefix scan, 340, 352, 674,
 689
 parallel tempering, 494, 510
 parallel trends, 1201
 parallel wavenet, 803
 parameter learning, 130
 parameter tying, 47, 95, 179
 parametric Bayesian model, 1005
 parametric model, 1005
 parametric models, 543
 parametric prior, 1006
 parent, 283
 parents, 117
 Pareto distribution, 11
 pareto index, 12
 Pareto smoothed importance sampling, 109
 parity check bits, 225
 part of speech, 162, 925
 part of speech tagging, 171
 parti, 790
 partial least squares, 908
 partially directed acyclic graph, 172
 partially observable Markov decision process, 1116
 partially observed data, 260
 partially observed Markov model, 933
 partially pooled model, 586
 particle BP, 392, 531
 particle filter, 513
 particle filtering, 327, 381, 511,
 1016
 particle impoverishment, 515
 particle smoothing, 531
 partition function, 30, 140, 811
 partition of the integers, 1010
 parts, 914

- 1
- patchGAN, **872**
path consistency learning, **1162**
path degeneracy, **518**
path diagram, **187**
path sampling, **441**
pathwise derivative, **251**
patience, **429**
pattern completion, **145**
PBIL, **285**
PBP, **614, 633**
PCFG, **163**
PCL, **1162**
PDAG, **172**
peaks function, **508**
peeling algorithm, **399**
PEGASUS, **1153**
per-decision importance sampling, **1158**
per-sample ELBO, **437**
per-step importance ratio, **1158**
per-step regret, **1114**
perceptual aliasing, **881**
perceptual distance metrics, **742**
perfect elimination ordering, **403**
perfect information, **1104**
perfect intervention, **188**
perfect map, **168**
period, **54**
periodic kernel, **646, 659**
permuted MNIST, **719**
perplexity, **211, 740**
persistent contrastive divergence, **817**
persistent variational inference, **156**
personalized recommendations, **70**
perturb-and-MAP, **407**
perturbation, **233**
PETS, **1153**
PGD, **724**
PGMs, **117**
phase space, **485**
phase transition, **142**
phi-exponential family, **34**
phone, **954**
phosphorylation state, **737**
Picard–Lindelöf theorem, **807**
pictorial structure, **166**
PILCO, **1152**
Pilot Studies, **1087**
pinball loss, **555**
pipe, **123**
Pitman-Koopman-Darmois theorem, **37, 214**
pix2pix, **872**
pixelCNN, **785**
pixelCNN++, **786**
pixelRNN, **786**
PixelSNAIL, **780**
placebo, **1108**
planar flow, **806**
PlaNet, **1155**
planning, **1119, 1131**
planning horizon, **1150**
plant, **1116**
plates, **136**
Platt scaling, **548**
- 2
- platykurtic, **9**
PLS, **908**
plug-in estimator, **1157**
plugin approximation, **67**
plutocratic, **12**
PoE, **812**
point estimate, **63, 66**
point process, **1028**
Poisson, **6**
Poisson process, **1028**
Poisson regression, **559**
policy, **1106, 1115**
policy evaluation, **1119, 1121**
policy gradient, **1133**
policy gradient theorem, **1143**
policy improvement, **1121**
policy iteration, **1121**
policy optimization, **1119**
policy search, **1133, 1142**
Polyak-Ruppert averaging, **616**
polymatroid function, **300**
polynomial kernel, **647**
polynomial regression, **568**
polysemy, **918**
polytrees, **387**
POMDP, **1116**
pool-based-sampling, **1123**
pooled, **95**
pooling layer, **594**
population, **282**
population-based incremental learning, **285**
position-specific scoring matrix, **941**
positive definite, **28**
positive definite kernel, **641**
positive phase, **157**
possible worlds, **180, 183**
post-order, **385**
posterior collapse, **424, 769**
posterior distribution, **64**
posterior expected loss, **1095**
posterior inference, **64, 331**
posterior linearization filter, **370**
posterior marginal, **129**
posterior mean, **1096**
posterior predictive check, **112**
posterior predictive distribution, **66**
posterior-predictive p-value, **113**
potential energy, **485**
potential function, **21, 139**
potential outcome, **190**
potential outcomes, **1176**
Potts model, **144**
Potts models, **474**
power EP, **448**
power law, **11**
power posterior, **622**
PPCA, **898**
PPL, **185**
PPO, **1148**
pre-order, **385**
pre-train and fine-tune, **705**
precision, **7, 73**
precision matrix, **16, 33**
precision-weighted mean, **16**
- 3
- preconditioned SGLD, **502**
predict-update, **333**
prediction, **191**
prediction step, **334**
predictive coding, **1166**
predictive distribution, **331**
predictive entropy search, **277**
predictive model, **543**
predictive sparse decomposition, **931**
predictive state representation, **949**
predictive uncertainty, **66**
prequential analysis, **108**
prescribed probabilistic models, **845**
pretext tasks, **1052**
prevalence shift, **697**
Price's theorem, **247, 252**
primitive nodes, **238**
primitive operations, **236**
principle of insufficient reason, **89**
prior distribution, **64**
prior linearization filter, **369**
prior network, **618**
prior predictive distribution, **567**
prior shift, **697**
prioritized experience replay, **1142**
PrLF, **369**
Probabilistic backpropagation, **633**
probabilistic backpropagation, **614**
probabilistic circuit, **179**
probabilistic ensembles with trajectory sampling, **1153**
Probabilistic graphical models, **117**
probabilistic graphical models, **731**
probabilistic LSA, **918**
probabilistic principal components analysis, **898**
probabilistic programming language, **185**
probabilistic soft logic, **184**
probability integral transform, **45**
probability matching, **1112**
probability of improvement, **275, 1125**
probability simplex, **25**
probit approximation, **579**
probit function, **582**
probit link function, **560**
probit regression, **582**
procedural approach, **185**
process noise, **631, 933, 961**
product density function, **1036**
product of experts, **147, 762, 812**
product partition model, **958**
production rules, **163**
profile HMM, **941**
projected gradient descent, **724**
projecting, **374**
projection, **173**
projection pursuit, **929**
Projection-weighted CCA (PWCCA), **1044**
prompt, **705**
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47

- 1
2 prompt tuning, 705
3 propensity score, 1183
3 propensity score matching, 1188
4 proper scoring rule, 546, 849
Properties., 1066
5 Prophet, 988
6 proposal distribution, 327, 456, 458, 468, 509
7 propose, 468
8 protein sequence alignment, 940
protein structure prediction, 144
9 protein-protein interaction net-
works, 995
10 prototypical networks, 716
11 proximal policy optimization, 1148
12 pseudo counts, 66, 70
13 pseudo inputs, 665
14 pseudo likelihood, 158, 158
15 pseudo random number generator, 454
16 PSIS, 109
pure exploration, 1114
17 pushforward, 792
18 pushing sums inside products, 399
Pyro, 185
19
20 Q-learning, 1133, 1139
21 QKF, 367
QT-Opt, 1119
22 quadratic approximation, 323
23 quadratic kernel, 649
quadratic loss, 1096
24 quadrature, 451
25 quadrature Kalman filter, 367
quadratures, 367
26 Qualitative Studies, 1086
27 quantile loss, 555
28 quantile regression, 555, 989
quantization, 209
29 Quasi Monte Carlo, 465
quasi-Newton EM algorithm, 265
30 queries, 272
31 query, 597
query by committee (QBC), 1126
query nodes, 129
33
34 R-hat, 497
radar, 970
35 radial basis function, 642
36 radon, 588
rainbow, 1142, 1149
37 random accelerations model, 961
38 random assignment with non-
compliance, 1193
39 random effects, 587
40 random finite sets, 973
41 random Fourier features, 651
random measure, 1023
42 random prior deep ensemble, 618
43 random restart, 280
random restart hill climbing, 280
44 random search, 281
random walk kernel, 648
45 random walk Metropolis, 327, 470
46 random walk on the integers, 54
47
- random walk proposal, 509
randomized control trials, 1178
Randomized QMC, 465
Rao-Blackwellisation, 462
Rao-Blackwellised particle filtering, 525
rare event, 536
raster scan, 785
rate, 222
rate distortion curve, 222, 758
rational quadratic, 646, 690
rats, 97
RBF, 642
RBM, 147
RBPF, 525
Real NVP, 808
real-time dynamic programming, 1121
receding horizon control, 1150
recognition network, 437, 747
recognition weights, 927
recombination, 283
recommender system, 180
reconstruction error, 222, 702
record linkage, 181
Rectified linear unit, 593
recurrent, 54
recurrent layer, 598
recurrent neural network, 604, 783
recurrent neural networks, 598
recurrent SSM, 991
recursive, 186
recursive least squares, 345, 631, 962
redundancy, 216, 225
reference prior, 94
refractoriness, 1032
regime switching Markov model, 937, 968
Regression discontinuity designs, 1217
regression estimator, 1157
regression model, 543
regret, 722, 1104, 1113, 1114
regular, 36, 54
regularization methods, 721
regularized evolution, 283
rehearsal, 721
REINFORCE, 428, 1133, 1144
Reinforcement learning, 1131
reject action, 1096
rejection sampling, 455
relational probability models, 179
relational UGMs, 181
relational uncertainty, 181
relative entropy, 195
relative risk, 662
relevance network, 220
relevance vector machine, 573
reliability diagram, 547
Renewal processes, 1032
reparameterization gradient, 251
reparameterization trick, 430, 613, 751
reparameterized VI, 430
reparametrization trick, 859
repeated trials, 63
- representation, 226
Representation learning, 739, 1039
representation learning, 675
Representational similarity analysis (RSA), 1042
representer theorem, 658
Reproducing Kernel Hilbert Space, 657
reproducing property, 657
resample, 517
resample-move, 531
residual belief propagation, 390, 394
residual block, 804
residual connections, 594, 804
residual error, 346
residual flows, 804
residues, 144
ResNet, 601
resource allocation, 1108
response surface model, 272
responsibility, 885
restless bandit, 1108
restricted Boltzmann machine, 147
return, 1117
reverse process, 831
reverse-mode automatic differentiation, 237
reversible jump MCMC, 504, 904
reward, 1101, 1106
reward function, 1116
reward model, 1115
reward-to-go, 1117
reweighted wake-sleep, 442, 443
RFF, 651
rich get richer, 424, 1010
ridge regression, 562, 565, 656
Riemann Manifold HMC, 489
Riemannian manifold, 243
risk, 1095
RJMCMC, 504
RKHS, 657
RL, 1131
RLS, 962
RM-HMC, 489
RMSprop, 271, 271
RNA-Sed, 982
RNADE, 784
RNN, 604
RNN-HMM, 952
robust IEKF, 359, 370
robust IEKS, 370
robust optimization, 727
robust priors, 87
robust regression, 662
robustness, 703
robustness analysis, 87
roll call, 915
roulette wheel selection, 283
row stochastic matrix, 119, 935
RPMs, 179
RStanARM, 587
RTS Smoother, 350
RTSS, 350
Rubin Causal Model, 191
run length, 957

- 1
 2 Russian roulette estimator, 805
 3 RVI, 430
 4 SAC, 1164
 5 safe policy iteration, 1147
 6 SAGA-LD, 503
 7 SAGAN, 867
 8 sample diversity, 740
 9 sample inefficient, 1133
 10 sample quality, 740
 11 sample standard deviation, 77
 12 sampling distribution, 63, 565
 13 sampling with replacement, 6
 14 SARSA, 1133, 1139
 15 satisfying assignment, 404
 16 SBEED, 1162
 17 scale-invariant prior, 93
 18 scaled inverse chi-squared, 11
 19 scaling-binning calibrator, 548
 20 scatter matrix, 80
 21 SCFGs, 955
 22 Schrödinger bridge, 842
 23 Schur complement, 17, 18
 24 SCM, 186
 25 score, 818
 26 score function, 39, 41, 814
 27 score function estimator, 250, 428,
 875
 28 score matching, 818
 29 score-based generative models,
 822, 831
 30 seasonality, 978
 31 second order EKF, 359
 32 second-order delta method, 253
 33 segment, 333
 34 segmental HMM, 952
 35 selection bias, 126, 698
 36 selection function, 283
 37 selective prediction, 702, 1096
 38 self attention, 786
 39 Self Attention GAN, 867
 40 self-attention, 598
 41 self-normalized importance sam-
 pling, 459
 42 self-train, 709
 43 semantic network, 1000
 44 semantic segmentation, 161, 164
 45 semi-amortized VI, 438
 46 semi-Markov model, 952
 47 semi-Markovian SCM, 186
 48 semi-parametric model, 654
 49 semi-parametrically efficient, 1184
 50 semi-supervised learning, 709, 765
 51 semilocal linear trend, 978
 52 sensible PCA, 898
 53 sensitive attribute, 1098
 54 sensitivity analysis, 87
 55 sensor fusion, 21
 56 sequence memoizer, 1020
 57 sequence-to-sequence, 735, 735
 58 sequential Bayesian inference, 327,
 630
 59 sequential Bayesian updating, 333,
 962
 60 sequential decision problem, 1106
 61 sequential importance sampling,
 515
 62 sequential importance sampling
 with resampling, 516
 63 sequential model-based optimiza-
 tion, 273
 64 sequential Monte Carlo, 327, 511
 65 sequential VAE, 766
 66 sFA, 915
 67 SFE, 250
 68 SG-HMC, 503
 69 SGD, 241
 70 SGLD, 490, 502
 71 SGLD-Adam, 502
 72 SGLD-CV, 502
 73 SGPR, 671
 74 SGRLD, 502
 75 shaded nodes, 136
 76 sharp minima, 623
 77 Sherman-Morrison-Woodbury for-
 mula, 17
 78 shift equivariance, 594
 79 shift invariance, 594
 80 shift-invariant kernel, 650
 81 shooting, 1151
 82 shortcut features, 694
 83 shortest path problems, 1121
 84 shrinkage, 74, 97
 85 sigma point BP, 392
 86 sigma point filter, 363
 87 sigma points, 362
 88 sigmoid, 575
 89 sigmoid belief net, 121
 90 signal to noise ratio, 832
 91 signed measure, 212
 92 silent state, 944
 93 SimCLR, 1055
 94 simple regret, 1114
 95 simplex factor analysis, 915
 96 Simplicity, 1079
 97 Simulability, 1078
 98 Simulated annealing, 282
 99 simulated annealing, 272, 508
 100 simulation-based inference, 537,
 846
 101 simultaneous localization and map-
 ping, 528
 102 single site updating, 478
 103 single world intervention graph,
 191
 104 singular statistical model, 111
 105 Singular vector CCA (SVCCA),
 1044
 106 SIS, 515
 107 SISR, 516
 108 site potential, 445
 109 SIXO, 538
 110 sketch-rnn, 768
 111 SKI, 673, 674
 112 SKIP, 674
 113 skip connections, 594, 771, 773
 114 skip-chain CRF, 163
 115 skip-VAE, 771
 116 SLAM, 528
 117 SLDS, 968
 118 sleep phase, 443
 119 slice sampling, 481
 120 sliced Fisher divergence, 820
 121 Sliced Score Matching, 820
 122 sliding window detector, 165
 123 slippage, 530
 124 slot machines, 1107
 125 slow weights, 620
 126 SLR, 368
 127 SLS, 280
 128 SMBO, 273
 129 SMC, 327, 511
 130 SMC sampler, 511
 131 SMC samplers, 531
 132 SMC², 537
 133 SMC-ABC, 537
 134 SMILES, 768
 135 smoothed Bellman error embed-
 ding, 1162
 136 smoothing, 335
 137 smoothing distribution, 331
 138 snapshot ensembles, 616
 139 SNGP, 615
 140 Sobol sequence, 465
 141 social networks, 995
 142 soft actor-critic, 1164
 143 soft clustering, 885, 997
 144 soft constraint, 812
 145 soft Q-learning, 1165
 146 soft-thresholding, 424
 147 softmax, 32
 148 softmax function, 576
 149 Softplus, 593
 150 SOLA, 635
 151 SOR, 666
 152 Soundness, 1079
 153 source coding, 193, 221
 154 source coding theorem, 221
 155 source distribution, 693
 156 source distributions, 709
 157 source domain, 872
 158 space filling, 465
 159 SPADA, 971
 160 sparse, 25, 568
 161 sparse Bayesian learning, 571
 162 sparse coding, 931
 163 sparse factor analysis, 898
 164 sparse GP, 666
 165 sparse GP regression, 671
 166 sparse variational GP, 668
 167 sparsity promoting priors, 610
 168 spectral density, 650, 682
 169 spectral estimation, 949
 170 spectral estimation method, 967
 171 spectral mixture kernel, 682
 172 spectral mixture kernels, 650
 173 speech-to-text, 734
 174 spelling correction, 942
 175 sphere the data, 928
 176 spherical covariance matrix, 14
 177 spherical cubature integration, 367
 178 spherical K-means algorithm, 27
 179 spherical topic model, 27
 180 sphering, 928
 181 spike and slab, 887
 182 spike-and-slab, 568
 183 spin, 141
 184 splines, 796
 185 split conformal prediction, 554
 186 split MNIST, 719
 187 split-Rhat, 497

- 1
2 spurious correlations, 694
3 SQF-RNN, 989
3 square root filter, 350
4 square root information filter, 350
5 square-integrable functions, 657
5 squared error, 1096
6 squared exponential, 642
6 SS, 568
7 SSID, 966
8 SSM, 933
9 Stability, 1077
9 stacking, 620
10 standard Brownian motion, 1030
11 standard error, 66, 453
11 standard error of the mean, 78
12 standard Normal, 7
12 state estimation, 331
13 state of nature, 1095
14 state transition diagram, 48
15 state-space model, 933
15 state-space models, 331
16 state-value function, 1118
17 stateful, 598
17 static calibration error, 548
18 stationary, 47
19 stationary distribution, 52
19 stationary kernels, 642
20 statistical estimand, 1175
21 statistical linear regression, 368
21 statistical parity, 1098
22 statistically linearized filter, 369
23 Statistics, 63
23 steepest ascent, 280
24 steepest descent, 241
25 stepping out, 482
25 stepwise EM, 265
26 stick-breaking construction, 1008
27 stick-breaking process, 1009
27 sticking the landing, 431
28 sticky, 471
29 stochastic approximation, 264, 272
29 stochastic approximation EM, 264
30 stochastic automaton, 48
31 stochastic averaged gradient acceleration, 503
32 stochastic bandit, 1107
33 stochastic block model, 996
33 stochastic computation graph, 254
34 stochastic context free grammars, 955
35 stochastic differential equations, 831
36 stochastic EP, 448
37 stochastic gradient descent, 241
38 Stochastic Gradient Langevin Descent, 490
39 stochastic gradient Langevin dynamics, 502
40 Stochastic Gradient Riemannian Langevin Dynamics, 502
42 stochastic hill climbing, 280
43 stochastic Lanczos quadrature, 672
44 stochastic local search, 279, 280, 282
45 stochastic matrix, 48
47
- stochastic meta descent, 167
Stochastic MuZero, 1151
stochastic process, 1005
stochastic relaxation, 272
stochastic RNN, 990
stochastic variance reduced gradient, 502
stochastic variational inference, 428, 672
stochastic video generation, 994
stochastic volatility model, 975
stochastic weight averaging, 616
stochastically ordered, 1009
stop gradient, 779
stop words, 920
straight-through estimator, 254, 778
stratified resampling, 520
streaks, 332
stream-based selective sampling, 1123
streaming variational Bayes, 634
strict overlap, 1181
strictly monotonic scalar function, 795
string kernel, 647
structural causal models, 186, 188, 1171
structural equation model, 174, 187
structural support vector machine, 168
structural time series, 976
structural zeros, 151
structured kernel interpolation, 673
structured mean field, 438
structured prediction, 160
Structured Prediction Energy Networks, 168
structured prediction model, 1097
STS, 976
Student distribution, 8
Student network, 119
student network, 120, 170, 399
Student t distribution, 8
style transfer, 873
sub-Gaussian, 9
subjective probability, 117
Submodular, 299
submodular, 1128
subphones, 954
Subscale Pixel Network, 786
subset of regressors, 666
subspace identification, 966
subspace neural bandit, 1110
sufficient, 226
sufficient statistic, 47, 214
sufficient statistics, 29, 30
sum of squares, 81
sum-product algorithm, 385
sum-product networks, 179
SupCon, 1055
super-Gaussian, 9
supervised PCA, 907
surjective, 44
surrogate assisted EA, 285
- surrogate function, 255, 272
survival of the fittest, 517
susbet-of-data, 663
SUTVA, 190
SVG, 994
SVGP, 668
SVI, 428
SVRG-LD, 502
SWA, 616
SWAG, 617
Swendsen Wang, 483
SWIG, 191
Swish, 593
swiss roll, 822
switching linear dynamical system, 525, 968
Sylvester flow, 806
symamd, 403
symmetric, 468
synchronous updates, 394
synergy, 216
syntactic sugar, 136
synthetic control, 986
Synthetic controls, 1217
systems biology, 1002
systems identification, 964
systolic array, 389
- T5, 767
tabu search, 280
tabular representation, 1116
tacotron, 785
TAN, 138
target aware Bayesian inference, 460
target distribution, 455, 458, 511, 693, 709
target domain, 872
target function, 458
target network, 1141
target policy, 1157
targeted attack, 724
TASA corpus, 918
task, 719
task incremental learning, 719
task interference, 710
task-aware learning, 719
Taylor series, 324
Taylor series expansion, 357
TD, 1137
TD error, 1133, 1137
TD(λ), 1138
TD3, 1149
telescoping sum, 836
temperature, 491
temperature scaling, 548
tempered posterior, 622
tempering, 635
template, 180
templates, 890
temporal difference, 1133, 1137
tensor decomposition, 949
tensor train decomposition, 674
TENT, 709
terminal state, 1116
terminals, 163
test and roll, 1101

- 1
 2 test statistics, 112
 3 test time adaptation, 708
 4 text generation, 735
 5 text to speech, 785
 6 text-to-image, 734
 7 text-to-speech, 874
 8 the deadly triad, 1162
 9 thermodynamic integration, 441
 10 thermodynamic variational objective, 441
 11 thin shell, 14
 12 thinning theorem, 1031
 13 Thompson sampling, 277, 1112
 14 threat model, 726
 15 tilted distribution, 446
 16 time reversal, 532
 17 time reversible, 55
 18 time series forecasting, 976
 19 time update step, 346
 20 time-invariant, 47
 21 time-series forecasting, 174
 22 Toeplitz, 674
 23 top-down inference model, 773
 24 topic, 1020
 25 topic model, 917
 26 topic vector, 917
 27 topic-RNN, 925
 28 topological order, 128
 29 topological ordering, 117
 30 total correlation, 214, 758
 31 total derivative, 252
 32 total regret, 1114
 33 total variation distance, 61
 34 tournament selection, 283
 35 trace plot, 494
 36 trace rank plot, 495
 37 traceback, 341, 387
 38 track, 969
 39 tracking, 344
 40 tractable substructure, 438
 41 trajectory, 1131
 42 trunkplot, 495
 43 trans-dimensional MCMC, 504
 44 transductive active learning, 552
 45 transductive learning, 699
 46 transfer learning, 705, 1040
 47 transformer, 604, 605
 48 transformer VAE, 767
 49 transformers, 767
 50 transient, 54
 51 transition, 1115
 52 transition function, 47, 933
 53 transition kernel, 47
 54 transition matrix, 48, 48
 55 transition model, 933, 935, 1115
 56 translation invariant, 931
 57 translation-invariant prior, 93
 58 Translucence, 1078
 59 Transparency, 1075
 60 transportable, 695
 61 treatment, 1101
 62 treatment effect, 985
 63 tree-augmented naive Bayes classifier, 138
 64 treewidth, 403
 65 trellis diagram, 340
 66 triggering kernel, 1033
 67
- trigram model, 49
 triplet loss, 1056
 TRPO, 1147
 truncated Gaussian, 584
 truncation selection, 283, 285
 trust region policy optimization, 1147
 TT-GP, 674
 TTA, 708
 TTT, 709
 turbocodes, 396
 Turing, 185
 turning the Bayesian crank, 319
 TVO, 441
 Tweedie's formula, 819
 twin network, 191, 985
 twisted particle filters, 538
 twisting function, 538
 two part code, 224
 two sample tests, 744
 two stage least squares, 1199
 two-filter smoothing, 337, 352
 two-moons, 636
 two-sample test, 55
 two-sample testing, 699
 two-slice marginal, 336
 type II maximum likelihood, 101
 type signature, 179
 typical set, 14, 210
- UCB, 276, 1111
 UCBVI, 1135
 UCRL2, 1135
 UDA, 707
 UKF, 362
 ULA, 490
 ULD, 504
 UMDA, 285
 UME, 60
 Unadjusted Langevin Algorithm, 490
 unary terms, 143, 151
 unbiased, 1158
 uncertainty metric, 700
 uncertainty quantification, 553
 uncertainty sampling, 1125
 unclamped phase, 157, 815
 unconstrained monotonic neural networks, 796
 underdamped Langevin dynamics, 504
 underlying predictive model, 958
 underspecification, 607
 undirected local Markov property, 153
 unfaithful, 154
 unidentifiable, 66
 Unified Medical Language System, 1000
 uniform dequantization, 741
 unigram model, 49
 unigram statistics, 51
 uninformative, 88
 uninformative prior, 564
 units, 190
 univariate marginal distribution algorithm, 285
 unnormalized mean embedding, 60
 unnormalized target distribution, 459
 unnormalized weights, 460, 514
 unpaired data, 872
 unroll, 180, 182
 unrolled, 136
 unscented Kalman filter, 362, 974
 unscented Kalman smoother, 365
 unscented particle filter, 522
 unscented RTS smoother, 365
 unscented transform, 362, 363
 unsupervised domain adaptation, 707
 unsupervised domain translation, 872
 untargeted attack, 724
 update step, 334
 UPGM, 139
 UPM, 958
 upper confidence bound, 276, 1111
 user rating profile model, 915
 User studies, 1080
 utility function, 1095
 utility nodes, 1098
- v-structure, 123, 125
 VAE, 149, 731, 747
 VAE-RNN, 766
 VAFC, 433
 vague prior, 567
 validation set, 107
 value function, 1118
 value iteration, 1120
 value nodes, 1098
 value of perfect information, 1100
 ValueDICE, 1168
 values, 597
 VAR, 174
 variable binding problem, 789
 variable duration HMM, 952
 variable elimination, 399, 1101
 variance exploding, 832
 variance preserving, 832
 Variational Approximation with Factor Covariance, 433
 variational autoencoder, 747
 variational autoencoders, 731, 912
 variational Bayes, 324, 415
 variational Bayes EM, 420
 variational continual learning, 635, 721
 variational diffusion model, 831
 variational diffusion models, 831
 variational EM, 259, 264, 420
 variational free energy, 410, 669, 749, 1166
 variational gap, 749
 variational GP, 439
 variational IB, 227
 variational inference, 29, 263, 324, 380, 409, 545, 741
 variational message passing, 268, 426
 variational method, 409

- 1
- 2 Variational Online Gauss-Newton, **270**, **271**
- 3 variational online Gauss-Newton, **614**
- 4 Variational Online Generalized
- 5 Gauss-Newton, **271**
- 6 variational optimization, **246**, **272**
- 7 variational overpruning, **635**, **769**
- 8 variational parameters, **324**, **409**
- 9 variational pruning, **776**
- 10 variational pruning effect, **424**
- 11 variational RNN, **993**
- 12 variational SMC, **538**
- 13 varimax, **898**
- 14 VIB, **324**
- 15 VCL, **635**
- 16 VD-VAE, **774**
- 17 VDM, **831**
- 18 vector auto-regressive, **174**
- 19 vector quantization, **221**
- 20 vector-Jacobian product (VJP), **233**
- 21 verb phrases, **162**
- 22 very deep VAE, **774**
- 23 VFE, **669**
- 24 VI, **324**
- 25 VIB, **227**
- 26 VIM, **782**
- 27 VIREL, **1165**
- 28 virtual samples, **84**
- 29 visible nodes, **129**
- 30 vision as inverse graphics, **881**
- 31 visual SLAM, **528**
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- Visualization, **1075**
- Viterbi algorithm, **340**, **941**
- Viterbi training, **947**
- VMP, **426**
- VOGNN, **271**
- VOGN, **270**, **271**, **614**
- von Mises, **28**
- von Mises-Fisher, **27**
- VQ-GAN, **782**
- VQ-VAE, **777**
- VRNN, **993**
- wake phase, **442**
- wake-phase q update, **444**
- wake-sleep, **748**
- wake-sleep algorithm, **442**
- warmup, **488**
- Wasserstein-1 distance, **58**
- Watanabe-Akaike information criterion, **111**
- wavegrad, **842**
- wavenet, **785**
- weak marginalization, **376**
- weak prior, **567**
- weak supervision, **184**
- weakly informative, **87**
- weakly-supervised representation learning, **1048**
- wealth, **12**
- website testing, **1105**
- weight degeneracy, **515**
- weight of evidence, **201**
- weight perturbation, **271**
- weight space, **654**
- weighted ERM, **706**
- weighted least squares, **263**
- Weiner process, **1030**
- Weinstein–Aronszajn identity, **805**
- white noise kernel, **659**
- whitebox attack, **724**
- whitened coordinate system, **243**
- whitening, **928**
- widely applicable information criterion, **111**
- Wiener noise, **501**, **504**
- wildcatter, **1098**
- Wishart, **28**
- witness function, **58**
- word error, **388**
- word2vec, **738**
- world model, **739**
- Xavier initialization, **608**
- XMC-GAN, **790**
- z-bias, **1209**
- zero-avoiding, **199**
- zero-forcing, **199**, **445**
- zero-inflated Poisson, **559**, **982**
- zero-one loss, **1096**
- zero-shot transfer, **1048**
- zero-sum losses, **856**
- ZIP, **559**, **982**
- Zipf's law, **12**

Bibliography

- [AA18] D. Amir and O. Amir. “Highlights: Summarizing agent behavior to people”. In: *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. 2018, pp. 1168–1176.
- [AB08] C. Archambeau and F. Bach. “Sparse probabilistic projections”. In: *NIPS*. 2008.
- [AB17] M. Arjovsky and L. Bottou. “Towards principled methods for training generative adversarial networks”. In: (2017).
- [AB21] A. N. Angelopoulos and S. Bates. “A Gentle Introduction to Conformal Prediction and Distribution-Free Uncertainty Quantification”. In: (2021). arXiv: [2107.07511 \[cs.LG\]](https://arxiv.org/abs/2107.07511).
- [Aba] A. Abadie. “Using Synthetic Controls: Feasibility, Data Requirements, and Methodological Aspects”. In: *J. of Economic Literature* () .
- [Abd+18] A. Abdolmaleki, J. T. Springenberg, Y. Tassa, R. Munos, N. Heess, and M. A. Riedmiller. “Maximum a Posteriori Policy Optimisation”. In: *ICLR*. 2018.
- [ABM10] J.-Y. Audibert, S. Bubeck, and R. Munos. “Best Arm Identification in Multi-Armed Bandits”. In: *COLT*. 2010, pp. 41–53.
- [Aba+21] S. Abnar, M. Dehghani, B. Neyshabur, and H. Sedghi. “Exploring the limits of large scale pre-training”. In: *ICLR*. 2021.
- [ABV21] S. Akbayrak, I. Bocharov, and B. de Vries. “Extended Variational Message Passing for Automated Approximate Bayesian Inference”. In: *Entropy* 23.7 (2021).
- [AC17] S. Aminikhanghahi and D. J. Cook. “A Survey of Methods for Time Series Change Point Detection”. en. In: *Knowl. Inf. Syst.* 51.2 (2017), pp. 339–367.
- [AC93] J. Albert and S. Chib. “Bayesian analysis of binary and polychotomous response data”. In: *JASA* 88.422 (1993), pp. 669–679.
- [ACB17] M. Arjovsky, S. Chintala, and L. Bottou. “Wasserstein generative adversarial networks”. In: *ICML*. 2017, pp. 214–223.
- [ACL16] L. Aitchison, N. Corradi, and P. E. Latham. “Zipf’s Law Arises Naturally When There Are Underlying, Unobserved Variables”. en. In: *PLoS Comput. Biol.* 12.12 (2016), e1005110.
- [ACL89] L. Atlas, D. Cohn, and R. Ladner. “Training connectionist networks with queries and selective sampliing”. In: *Advances in neural information processing systems* 2 (1989).
- [ACP87] S. Arnborg, D. G. Corneil, and A. Proskurowski. “Complexity of finding embeddings in a k-tree”. In: *SIAM J. on Algebraic and Discrete Methods* 8 (1987), pp. 277–284.
- [Ada00] L. Adamic. *Zipf, Power-laws, and Pareto - a ranking tutorial*. Tech. rep. 2000.
- [Ada+20] V. Adam, S. Eleftheriadis, N. Durrande, A. Artemev, and J. Hensman. “Doubly Sparse Variational Gaussian Processes”. In: *AISTATS*. 2020.
- [Ade+20a] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. *Sanity Checks for Saliency Maps*. 2020. arXiv: [1810.03292 \[cs.CV\]](https://arxiv.org/abs/1810.03292).
- [Ade+20b] J. Adebayo, M. Muelly, I. Liccardi, and B. Kim. “Debugging tests for model explanations”. In: *arXiv preprint arXiv:2011.05429* (2020).
- [Adl+18] P. Adler, C. Falk, S. A. Friedler, T. Nix, G. Rybeck, C. Scheidegger, B. Smith, and S. Venkatasubramanian. “Auditing black-box models for indirect influence”. In: *Knowledge and Information Systems* 54.1 (2018), pp. 95–122.
- [Ado] *Taking It to the MAX: Adobe Photoshop Gets New NVIDIA AI-Powered Neural Filters*. <https://blogs.nvidia.com/blog/2020/10/20/adobe-max-ai/>. Accessed: 2021-08-12.
- [AE+20] D. Agudelo-España, S. Gomez-Gonzalez, S. Bauer, B. Schölkopf, and J. Peters. “Bayesian Online Prediction of Change Points”. In: *UAI*. Vol. 124. Proceedings of Machine Learning Research. PMLR, 2020, pp. 320–329.
- [AFD01] C. Andrieu, N. de Freitas, and A. Doucet. “Robust Full Bayesian Learning for Radial Basis Networks”. In: *Neural Computation* 13.10 (2001), pp. 2359–2407.
- [AFG19] M. Akten, R. Fiebrink, and M. Grierson. “Learning to See: You Are What You See”. In: *ACM SIGGRAPH 2019 Art Gallery. SIGGRAPH ’19*. Association for Computing Machinery, 2019.
- [AGI11] A. Allahverdyan and A. Galstyan. “Comparative Analysis of Viterbi Training and Maximum Likelihood Estimation for HMMs”. In: *NIPS*. 2011, pp. 1674–1682.
- [AG13] S. Agrawal and N. Goyal. “Further Optimal Regret Bounds for Thompson Sampling”. In: *AISTATS*. 2013.
- [Aga+14] D. Agarwal, B. Long, J. Traupman, D. Xin, and L. Zhang. “LASER: a scalable response prediction platform for online advertising”. In: *WSDM*. 2014.
- [Aga+21a] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun. *Reinforcement Learning: Theory and Algorithms*. 2021.
- [Aga+21b] R. Agarwal, L. Melnick, N. Frost, X. Zhang, B. Lengerich, R. Caruana, and G. Hinton. *Neural Additive Models: Interpretable Machine Learning with Neural Nets*. 2021. arXiv: [2004.13912 \[cs.LG\]](https://arxiv.org/abs/2004.13912).
- [AGM14] P. Agrawal, R. Girshick, and J. Malik. “Analyzing the performance of multilayer neural networks for object recognition”. In: *European conference on computer vision*. Springer, 2014, pp. 329–344.
- [AH09] I. Arasaratnam and S. Haykin. “Cubature Kalman Filters”. In: *IEEE Trans. Automat. Contr.* 54.6 (2009), pp. 1254–1269.
- [AHE07] I. Arasaratnam, S. Haykin, and R. J. Elliott. “Discrete-Time Nonlinear Filtering Algorithms

- 1 Using Gauss–Hermite Quadrature". In: *Proc. IEEE*
 2 95.5 (2007), pp. 953–977.
- 3 [AHG20] L. Ambrogioni, M. Hinne, and M. van Gerven, "Automatic structured variational inference". In: (2020). arXiv: [2002.00643 \[stat.ML\]](#).
- 4 [AHK01] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, "On the Surprising Behavior of Distance Metrics in High Dimensional Space". In: *Database Theory – ICDT 2001*. Springer Berlin Heidelberg, 2001, pp. 420–434.
- 5 [AHK12] A. Anandkumar, D. Hsu, and S. M. Kakade, "A Method of Moments for Mixture Models and Hidden Markov Models". In: *COLT*, Vol. 23. Proceedings of Machine Learning Research. PMLR, 2012, pp. 33.1–33.34.
- 6 [AHK65] K. Abend, T. J. Harley, and L. N. Kanal, "Classification of Binary Random Patterns". In: *IEEE Transactions on Information Theory* 11(4) (1965), pp. 538–544.
- 7 [Ahm+17] S. Ahmad, A. Lavin, S. Purdy, and Z. Agha, "Unsupervised real-time anomaly detection for streaming data". In: *Neurocomputing* 262 (2017), pp. 134–147.
- 8 [AHP+05] P. K. Agarwal, S. Har-Peled, et al. "Geometric approximation via coresets". In: *Combinatorial and computational geometry* 52.1–30 (2005), p. 3.
- 9 [AHS85] D. Ackley, G. Hinton, and T. Sejnowski, "A Learning Algorithm for Boltzmann Machines". In: *Cognitive Science* 9 (1985), pp. 147–169.
- 10 [AHT07] Y. Altun, T. Hofmann, and I. Tschantaridis, "Support Vector Machine Learning for Interdependent and Structured Output Spaces". In: *Predicting Structured Data*. Ed. by G. Bakir, T. Hofmann, B. Scholkopf, A. Smola, B. Taskar, and S. Vishwanathan. MIT Press, 2007.
- 11 [Ahu+21] K. Ahuja, J. Wang, A. Dhurandhar, K. Shannugam, and K. R. Varshney, "Empirical or Invariant Risk Minimization? A Sample Complexity Perspective". In: *ICLR*. 2021.
- 12 [AI 19] AI Artists. *Creative Tools to Generate AI Art*. 2019.
- 13 [Air+08] E. Airoldi, D. Blei, S. Fienberg, and E. Xing, "Mixed-membership stochastic blockmodels". In: *JMLR* 9 (2008), pp. 1981–2014.
- 14 [Ait18] L. Aitchison, "A unified theory of adaptive stochastic gradient descent as Bayesian filtering". In: (2018). arXiv: [1807.07540 \[stat.ML\]](#).
- 15 [Ait21] L. Aitchison, "A statistical theory of cold posteriors in deep neural networks". In: *ICLR*. 2021.
- 16 [Aka74] H. Akaike, "A new look at the statistical model identification". In: *IEEE Trans. on Automatic Control* 19.6 (1974).
- 17 [AKO18] S.-I. Amari, R. Karakida, and M. Oizumi, "Fisher Information and Natural Gradient Learning of Random Deep Networks". In: (2018). arXiv: [1808.07172 \[cs.LG\]](#).
- 18 [AKZK19] B. Amos, V. Koltun, and J. Zico Kolter, "The Limited Multi-Label Projection Layer". In: (2019). arXiv: [1906.08070 \[cs.LG\]](#).
- 19 [AL+16] J. Ala-Luhtala, N. Whiteley, K. Heine, and R. Piche, "An Introduction to Twisted Particle Filters and Parameter Estimation in Non-linear State-space Models". In: *IEEE Trans. Signal Process* 64.18 (2016), pp. 4875–4890.
- 20 [al21] M. A. et al. *Understanding Dataset Shift and Potential Remedies*. Tech. rep. Vector Institute, 2021.
- 21 [Ale+16] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep Variational Information Bottleneck". In: *ICLR*. 2016.
- 22 [Ale+17] A. A. Alemi, I. S. Fischer, J. V. Dillon, and K. P. Murphy, "Deep Variational Information Bottleneck". In: *ArXiv* abs/1612.00410 (2017).
- 23 [Ale+18] A. A. Alemi, B. Poole, I. Fischer, J. V. Dillon, R. A. Saurous, and K. Murphy, "Fixing a broken ELBO". In: *ICML*. 2018.
- 24 [Alq21] P. Alquier, "User-friendly introduction to PAC-Bayes bounds". In: (2021). arXiv: [2110.11216 \[stat.ML\]](#).
- 25 [Als+19] J. Alsing, T. Charnock, S. Feeney, and B. Wandelt, "Fast likelihood-free cosmology with neural density estimators and active learning". In: *Monthly Notices of the Royal Astronomical Society* 488.3 (2019), pp. 4440–4458.
- 26 [ALS20] B. Axelrod, Y. P. Liu, and A. Sidford, "Near-optimal approximate discrete and continuous submodular function minimization". In: *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2020, pp. 837–853.
- 27 [AM07] R. P. Adams and D. J. C. MacKay, "Bayesian Online Changepoint Detection". In: (2007). arXiv: [0710.3742 \[stat.ML\]](#).
- 28 [AM+16] M. Auger-Méthé, C. Field, C. M. Albertsen, A. E. Derocher, M. A. Lewis, I. D. Jonson, and J. Mills Flemming, "State-space models' dirty little secrets: even simple linear Gaussian models can have estimation problems". In: *Sci. Rep.* 6 (2016), p. 26677.
- 29 [AM74] D. Andrews and C. Mallows, "Scale mixtures of Normal distributions". In: *J. of Royal Stat. Soc. Series B* 36 (1974), pp. 99–102.
- 30 [AM89] B. D. Anderson and J. B. Moore, *Optimal Control: Linear Quadratic Methods*. Prentice-Hall International, Inc., 1989.
- 31 [Ama09] S.-I. Amari, "α-Divergence Is Unique, Belonging to Both f -Divergence and Bregman Divergence Classes". In: *IEEE Trans. Inf. Theory* 55.11 (2009), pp. 4925–4931.
- 32 [Ama98] S. Amari, "Natural Gradient Works Efficiently in Learning". In: *Neural Comput.* 10.2 (1998), pp. 251–276.
- 33 [Ame+19] S. Amershi, D. Weld, M. Vorvoreanu, A. Fournier, B. Nushi, P. Collisson, J. Suh, S. Iqbal, P. N. Bennett, K. Inkpen, et al. "Guidelines for human-AI interaction". In: *Proceedings of the 2019 chi conference on human factors in computing systems*. 2019, pp. 1–13.
- 34 [Ami01] E. Amir, "Efficient Approximation for Triangulation of Minimum Treewidth". In: *UAI*. 2001.
- 35 [AMJ18a] D. Alvarez-Melis and T. S. Jaakkola, "On the robustness of interpretability methods". In: *arXiv preprint arXiv:1806.08049* (2018).
- 36 [AMJ18b] D. Alvarez-Melis and T. S. Jaakkola, "Towards robust interpretability with self-explaining neural networks". In: *arXiv preprint arXiv:1806.07538* (2018).
- 37 [Amo+16] D. Amodei, C. Olah, J. Steinhardt, P. F. Christiano, J. Schulman, and D. Mané, "Concrete Problems in AI Safety". In: *CoRR* abs/1606.06565 (2016). arXiv: [1606.06565](#).

- 1 [Amo+18] B. Amos, L. Dinh, S. Cabi, T. Rothörl, S. G.
2 Colmenarejo, A. Muldal, T. Erez, Y. Tassa, N. de Freitas, and M. Denil. “Learning Awareness Models”. In:
3 *ICLR*. 2018.
- 4 [Amo22] B. Amos. “Tutorial on amortized optimization
5 for learning to optimize over continuous domains”. In:
6 (2022). arXiv: [2202.00665 \[cs.LG\]](#).
- 7 [AMO88] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin.
8 “Network flows”. In: (1988).
- 9 [Ana+14] A. Anandkumar, R. Ge, D. Hsu, S. M.
10 Kakade, and M. Telgarsky. “Tensor Decompositions
11 for Learning Latent Variable Models”. In: *JMLR* 15
12 (2014), pp. 2773–2832.
- 13 [And+03] C. Andrieu, N. de Freitas, A. Doucet, and
14 M. Jordan. “An introduction to MCMC for machine
15 learning”. In: *Machine Learning* 50 (2003), pp. 5–43.
- 16 [And+20] O. M. Andrychowicz et al. “Learning dexterous
17 in-hand manipulation”. In: *Int. J. Rob. Res.* 39.1
18 (2020), pp. 3–20.
- 19 [Ang+18] E. Angelino, N. Larus-Stone, D. Alabi, M.
20 Seltzer, and C. Rudin. *Learning Certifiably Optimal
21 Rule Lists for Categorical Data*. 2018. arXiv: [1704.01701 \[stat.ML\]](#).
- 22 [Ang+20] C. Angermueller, D. Dohan, D. Belanger, R.
23 Deshpande, K. Murphy, and L. Colwell. “Model-based
24 reinforcement learning for biological sequence design”.
25 In: *ICLR*. 2020.
- 26 [Ang+21] A. N. Angelopoulos, S. Bates, E. J. Candès,
27 M. I. Jordan, and L. Lei. “Learn then Test: Calibrating
28 Predictive Algorithms to Achieve Risk Control”.
29 In: (2021). arXiv: [2110.01052 \[cs.LG\]](#).
- 30 [Ang88] D. Angluin. “Queries and concept learning”. In:
31 *Machine learning* 2.4 (1988), pp. 319–342.
- 32 [Ani+18] R. Anirudh, J. J. Thiagarajan, B. Kailkhura,
33 and T. Bremer. “An Unsupervised Approach to Solving
34 Inverse Problems using Generative Adversarial
35 Networks”. In: (2018). arXiv: [1805.07281 \[cs.CV\]](#).
- 36 [Ano19] Anonymous. “Neural Tangents: Fast and Easy
37 Infinite Neural Networks in Python”. In: (2019).
- 38 [Ant+22] I. Antonoglou, J. Schrittwieser, S. Ozair,
39 T. K. Hubert, and D. Silver. “Planning in Stochastic
40 Environments with a Learned Model”. In: *ICLR*. 2022.
- 41 [AO03] J.-H. Ahn and J.-H. Oh. “A Constrained EM
42 Algorithm for Principal Component Analysis”. In:
43 *Neural Computation* 15 (2003), pp. 57–65.
- 44 [AOM17] M. G. Azar, I. Osband, and R. Munos. “Min-
45 imax Regret Bounds for Reinforcement Learning”. In:
46 *ICML*. 2017, pp. 263–272.
- 47 [AP08] J. D. Angrist and J.-S. Pischke. *Mostly
48 harmless econometrics: An empiricist’s companion*.
49 Princeton university press, 2008.
- 50 [AP09] J. Angrist and J.-S. Pischke. *Mostly Harmless
51 Econometrics*. 2009.
- 52 [AP19] M. Abadi and G. D. Plotkin. “A simple differen-
53 tiable programming language”. In: *Proceedings of the
54 ACM on Programming Languages* 4.POPL (2019),
55 pp. 1–28.
- 56 [Ara+09] A. Aravkin, B. Bell, J. Burke, and G. Pil-
57 lonetto. *An L1-Laplace Robust Kalman Smoother*.
58 Tech. rep. U. Washington, 2009.
- 59 [Ara10] A. Aravkin. *Student’s t Kalman Smoother*.
60 Tech. rep. U. Washington, 2010.
- 61 [Ara+17] A. Aravkin, J. V. Burke, L. Ljung, A. Lozano,
62 and G. Pillonetto. “Generalized Kalman smoothing:
63 Modeling and algorithms”. In: *Automatica* 86 (2017),
64 pp. 63–86.
- 65 [Arb+18] M. Arbel, D. Sutherland, M. Bińkowski,
66 and A. Gretton. “On gradient regularizers for MMD
67 GANs”. In: *Advances in neural information process-
68 ing systems*. 2018, pp. 6700–6710.
- 69 [Arj+19] M. Arjovsky, L. Bottou, I. Gulrajani, and D.
70 Lopez-Paz. “Invariant Risk Minimization”. In: (2019).
71 arXiv: [1907.02893 \[stat.ML\]](#).
- 72 [Arj+20] M. Arjovsky, L. Bottou, I. Gulrajani, and D.
73 Lopez-Paz. *Invariant Risk Minimization*. 2020. arXiv:
74 [1907.02893 \[stat.ML\]](#).
- 75 [Arn+10] C. W. Arnold, S. M. El-Saden, A. A. Bui, and
76 R. Taira. “Clinical case-based retrieval using latent
77 topic analysis”. In: *AMIA annual symposium proceed-
78 ings*. Vol. 2010. American Medical Informatics Associa-
79 tion. 2010, p. 26.
- 80 [Aro+13] S. Arora, R. Ge, Y. Halpern, D. Mimno, A.
81 Moitra, D. Sontag, Y. Wu, and M. Zhu. “A Practical
82 Algorithm for Topic Modeling with Provable Guarantees”. In: *ICML*. 2013.
- 83 [Aro+19] S. Arora, S. S. Du, W. Hu, Z. Li, R. Salakhut-
84 dinov, and R. Wang. “On Exact Computation with an
85 Infinitely Wide Neural Net”. In: (2019). arXiv: [1904.11955 \[cs.LG\]](#).
- 86 [Aro+21] R. Arora et al. *Theory of deep learning*.
87 2021.
- [ARS13] N. S. Arora, S. Russell, and E. Sudderth.
88 “NET-VISA: Network Processing Vertically Integrated
89 Seismic AnalysisNET-VISA: Network Processing Ver-
90 tically Integrated Seismic Analysis”. In: *Bull. Seismol.
91 Soc. Am.* 103.2A (2013), pp. 709–729.
- [Aru+02] M. Arulampalam, S. Maskell, N. Gordon,
92 and T. Clapp. “A Tutorial on Particle Filters for
93 Online Nonlinear/Non-Gaussian Bayesian Tracking”.
94 In: *IEEE Trans. on Signal Processing* 50.2 (2002),
95 pp. 174–189.
- [Aru+17] K. Arulkumaran, M. P. Deisenroth, M.
96 Brundage, and A. A. Bharath. “A Brief Survey of Deep
97 Reinforcement Learning”. In: *IEEE Signal Processing
98 Magazine, Special Issue on Deep Learning for Image
99 Understanding* (2017).
- [AS17] A. Achille and S. Soatto. “On the Emergence
100 of Invariance and Disentangling in Deep Representa-
101 tions”. In: (2017). arXiv: [1706.01350 \[cs.LG\]](#).
- [AS18] A. Achille and S. Soatto. “On the Emergence
102 of Invariance and Disentangling in Deep Representa-
103 tions”. In: *JMLR* 18 (2018), pp. 1–34.
- [AS66] S. M. Ali and S. D. Silvey. “A General Class
104 of Coefficients of Divergence of One Distribution
105 from Another”. In: *J. R. Stat. Soc. Series B Stat.
106 Methodol.* 28.1 (1966), pp. 131–142.
- [Asa00] C. Asavathiratham. “The Influence Model: A
107 Tractable Representation for the Dynamics of Net-
108 worked Markov Chains”. PhD thesis. MIT, Dept.
109 EECS, 2000.
- [ASD20] A. Agrawal, D. Sheldon, and J. Domke. “Ad-
110 vances in Black-Box VI: Normalizing Flows, Im-
111 portance Weighting, and Optimization”. In: *NIPS*. 2020.
- [ASM17] A. Azuma, M. Shimbo, and Y. Matsumoto.
112 “An Algebraic Formalization of Forward and Forward-
113 backward Algorithms”. In: (2017). arXiv: [1702.06941
114 \[cs.LG\]](#).

- 1
- 2 [ASN20] R. Agarwal, D. Schuurmans, and M. Norouzi.
“An Optimistic Perspective on Offline Reinforcement
3 Learning”. In: *ICML*. 2020.
- 4 [ASS19] I. Andrews, J. H. Stock, and L. Sun. “Weak In-
5 struments in Instrumental Variables Regression: The-
6 ory and Practice”. In: *Annual Review of Economics*
7 11.1 (2019).
- 8 [AT08] C. Andrieu and J. Thoms. “A tutorial on adap-
9 tive MCMC”. In: *Statistical Computing* 18 (2008),
pp. 343–373.
- 10 [AT20] E. Agustsson and L. Theis. “Universally Quan-
11 tized Neural Compression”. 2020.
- 12 [Att00] H. Attias. “A Variational Bayesian Framework
13 for Graphical Models”. In: *NIPS-12*. 2000.
- 14 [Att03] H. Attias. “Planning by Probabilistic Infer-
15 ence”. In: *AI-Stats*. 2003.
- 16 [Aue12] J. E. Auerbach. “Automated evolution of inter-
17 esting images”. In: *Artificial Life* 13. 2012.
- 18 [AWR17] J. Altschuler, J. Weed, and P. Rigollet. “Near-
19 linear time approximation algorithms for optimal
20 transport via Sinkhorn iteration”. In: *arXiv preprint
arXiv:1705.09634* (2017).
- 21 [AXK17] B. Amos, L. Xu, and J. Z. Kolter. “Input con-
22 vex neural networks”. In: *International Conference on
Machine Learning*. PMLR. 2017, pp. 146–155.
- 23 [AY19] B. Amos and D. Yarats. “The Differentiable
24 Cross-Entropy Method”. In: (2019). arXiv: [1909.12830](#)
[cs.LG].
- 25 [AZ17] S. Arora and Y. Zhang. “Do gans actually
26 learn the distribution? an empirical study”. In: *arXiv
preprint arXiv:1706.08224* (2017).
- 27 [Aze+20] E. M. Azevedo, A. Deng, J. L. Montiel Olea,
28 J. Rao, and E. G. Weyl. “A/B Testing with Fat Tails”.
29 In: *J. Polit. Econ.* (2020), pp. 000–000.
- 30 [Azi+15] H. Azizpour, A. Sharif Razavian, J. Sullivan,
31 A. Maki, and S. Carlsson. “From generic to specific
32 deep representations for visual recognition”. In: *Pro-
ceedings of the IEEE conference on computer vision
and pattern recognition workshops*. 2015, pp. 36–45.
- 33 [BA03] D. Barber and F. Agakov. “The IM Algorithm:
A Variational Approach to Information Maximization”.
34 In: *NIPS*. NIPS’03. MIT Press, 2003, pp. 201–208.
- 35 [BA05] W. Bechtel and A. Abrahamsen. “Explanation:
36 A mechanist alternative”. In: *Studies in History and
Philosophy of Science Part C: Studies in History
and Philosophy of Biological and Biomedical Sciences* 36.2 (2005), pp. 421–441.
- 37 [Bac09] F. Bach. “High-Dimensional Non-Linear Vari-
able Selection through Hierarchical Kernel Learning”.
38 In: (2009). arXiv: [0909.0844](#) [cs.LG].
- 39 [Bac+13] F. Bach et al. “Learning with Submodular
Functions: A Convex Optimization Perspective”. In:
Foundations and Trends® in Machine Learning 6.2-
40 3 (2013), pp. 145–373.
- 41 [Bac+15a] S. Bach, A. Binder, G. Montavon, F.
42 Klauschen, K.-R. Müller, and W. Samek. “On pixel-
43 wise explanations for non-linear classifier decisions by
layer-wise relevance propagation”. In: *PloS one* 10.7
44 (2015), e0130140.
- 45 [Bac+15b] S. H. Bach, M. Broecheler, B. Huang, and
L. Getoor. “Hinge-Loss Markov Random Fields and
46
- 47 Probabilistic Soft Logic”. In: (2015). arXiv: [1505.04406](#)
[cs.LG].
- [Bac+18] E. Bach, J. Dusart, L. Hellerstein, and D. D.
Kletenik. “Submodular goal value of Boolean func-
tions”. In: *Discrete Applied Mathematics* 238 (2018),
pp. 1–13.
- [Bad+18] M. A. Badgeley, J. R. Zech, L. Oakden-
Rayner, B. S. Glicksberg, M. Liu, W. Gale, M. V.
McConnell, B. Percha, T. M. Snyder, and J. T. Dud-
ley. “Deep Learning Predicts Hip Fracture using Con-
founding Patient and Healthcare Variables”. In: *CoRR
abs/1811.03695* (2018). arXiv: [1811.03695](#).
- [Bah+20] Y. Bahri, J. Kadmon, J. Pennington, S.
Schoenholz, J. Sohl-Dickstein, and S. Ganguli. “Sta-
tistical Mechanics of Deep Learning”. In: *Annu. Rev.
Condens. Matter Phys.* (2020).
- [Bai+15] R. Bairi, R. Iyer, G. Ramakrishnan, and J.
Bilmes. “Summarization of multi-document topic hier-
archies using submodular mixtures”. In: *Proceedings
of the 53rd Annual Meeting of the Association for
Computational Linguistics and the 7th International
Joint Conference on Natural Language Processing
(Volume 1: Long Papers)*. 2015, pp. 553–563.
- [Bai95] L. C. Baird. “Residual Algorithms: Reinforce-
ment Learning with Function Approximation”. In:
ICML. 1995, pp. 30–37.
- [Bak+17] J. Baker, P. Fearnhead, E. B. Fox, and C.
Nemeth. “Control Variates for Stochastic Gradient
MCMC”. In: (2017). arXiv: [1706.05439](#) [stat.CO].
- [Bal17] S. Baluja. “Learning deep models of optimiza-
tion landscapes”. In: *IEEE Symposium Series on
Computational Intelligence (SSCI)* (2017).
- [Bal+18] D. Balduzzi, S. Racaniere, J. Martens, J. Fo-
erster, K. Tuyls, and T. Graepel. “The mechanics of
n-player differentiable games”. In: *International Con-
ference on Machine Learning*. PMLR. 2018, pp. 354–
363.
- [Ban+05] A. Banerjee, I. S. Dhillon, J. Ghosh, and S.
Sra. “Clustering on the unit hypersphere using von
Mises-Fisher distributions”. In: *JMLR*. 2005, pp. 1345–
1382.
- [Ban06] A. Banerjee. “On bayesian bounds”. In: *ICML*.
2006, pp. 81–88.
- [Ban+18] A. Bansal, S. Ma, D. Ramanan, and Y.
Sheikh. “Recycle-gan: Unsupervised video retarget-
ing”. In: *Proceedings of the European conference on
computer vision (ECCV)*. 2018, pp. 119–135.
- [Bao+22] H. Bao, L. Dong, S. Piao, and F. Wei. “BEiT:
BERT Pre-Training of Image Transformers”. In: *Inter-
national Conference on Learning Representations*.
2022.
- [Baq+16] P. Baqué, T. Bagautdinov, F. Fleuret, and P.
Fua. “Principled parallel mean-field inference for dis-
crete random fields”. In: *CVPR*. 2016, pp. 5848–5857.
- [Bar17] D. Barber. *Evolutionary Optimization as a
Variational Method*. 2017.
- [Bas+01] S. Basu, T. Choudhury, B. Clarkson, and A.
Pentland. *Learning Human Interactions with the In-
fluence Model*. Tech. rep. 539. MIT Media Lab, 2001.
- [Bat+12] D. Batra, P. Yadollahpour, A. Guzman-
Rivera, and G. Shakhnarovich. “Diverse M-Best Solu-
tions in Markov Random Fields”. In: *ECCV*. Springer
Berlin Heidelberg, 2012, pp. 1–16.
- [Bau+17] D. Bau, B. Zhou, A. Khosla, A. Oliva, and
A. Torralba. “Network Dissection: Quantifying Inter-

- 1 preability of Deep Visual Representations". In: *Computer Vision and Pattern Recognition*. 2017.
- 2 [Bau+18] D. Bau, J.-Y. Zhu, H. Strobelt, B. Zhou,
3 J. B. Tenenbaum, W. T. Freeman, and A. Torralba. "Gan dissection: Visualizing and understanding
4 generative adversarial networks". In: *arXiv preprint arXiv:1811.10597* (2018).
- 5 [Bau+20] D. Bau, J.-Y. Zhu, H. Strobelt, A. Lapedriza,
6 B. Zhou, and A. Torralba. "Understanding the role of
7 individual units in a deep neural network". In: *Proceedings of the National Academy of Sciences* (2020).
- 8 [Bau+70] L. E. Baum, T. Petrie, G. Soules, and N.
9 Weiss. "A maximization technique occurring in the statistical analysis of probabilistic functions in Markov
10 chains". In: *The Annals of Mathematical Statistics* 41 (1970), pp. 164–171.
- 11 [Bau74] B. G. Baumgart. "Geometric modeling for computer vision." In: 1974.
- 12 [Bax00] J. Baxter. "A Model of Inductive Bias Learning". In: *JAIR* (2000).
- 13 [Bay+15] A. G. Baydin, B. A. Pearlmutter, A. A.
14 Radul, and J. M. Siskind. "Automatic differentiation in machine learning: a survey". In: (2015). arXiv: [1502.05767 \[cs.SC\]](https://arxiv.org/abs/1502.05767).
- 15 [BB12] J. Bergstra and Y. Bengio. "Random Search for Hyper-Parameter Optimization". In: *JMLR* 13 (2012), pp. 281–305.
- 16 [BB15a] A. Bendale and T. Boult. "Towards Open
17 World Recognition". In: *CVPR*. 2015.
- 18 [BB15b] J. Bornschein and Y. Bengio. "Reweighted
19 Wake-Sleep". In: *ICLR*. 2015.
- 20 [BB17] J. Bilmes and W. Bai. "Deep Submodular Func-
21 tions". In: *Arxiv abs/1701.08939* (2017).
- 22 [BB18] W. Bai and J. Bilmes. "Greed is Still Good:
23 Maximizing Monotone Submodular+Supermodular
24 (BP) Functions". In: *International Conference on Machine Learning (ICML)*. <http://proceedings.mlr.press/v80/bai18a.html>. Stockholm, Sweden, 2018.
- 25 [BBH12] C. Blundell, J. Beck, and K. A. Heller. "Mod-
26 elling reciprocating relationships with Hawkes processes". In: *Advances in Neural Information Processing Systems*. 2012, pp. 2600–2608.
- 27 [BBM07] A. Banerjee, S. Basu, and S. Merugu. "Multi-
28 way Clustering on Relation Graphs". In: *Proc. SIAM Intl. Conf. on Data Mining (SDM)*. 2007.
- 29 [BBM10] N. Bhatnagar, C. Bogdanov, and E. Mossel.
30 *The Computational Complexity of Estimating Con-
31 vergence Time*. Tech. rep. arxiv, 2010.
- 32 [BBS09] J. O. Berger, J. M. Bernardo, and D. Sun.
33 "The Formal Definition of Reference Priors". In: *Ann. Stat.* 37.2 (2009), pp. 905–938.
- 34 [BBS95] A. G. Barto, S. J. Bradtke, and S. P. Singh.
35 "Learning to act using real-time dynamic program-
36 ming". In: *AIJ* 72.1 (1995), pp. 81–138.
- 37 [BBV11a] R. Benassi, J. Bect, and E. Vazquez.
38 "Bayesian optimization using sequential Monte Carlo".
39 In: (2011). arXiv: [1111.4802 \[math.OC\]](https://arxiv.org/abs/1111.4802).
- 40 [BBV11b] R. Benassi, J. Bect, and E. Vazquez. "Rob-
41 ust Gaussian Process-Based Global Optimization Us-
42 ing a Fully Bayesian Expected Improvement Criterion". In: *Intl. Conf. on Learning and Intelligent
43 Optimization (LION)*. 2011, pp. 176–190.
- 44 [BC07] D. Barber and S. Chiappa. "Unified inference
45 for variational Bayesian linear Gaussian state space
46 models". In: *NIPS*. 2007.
- 47 [BC08] M. Bădoiu and K. L. Clarkson. "Optimal coresets
48 for balls". In: *Computational Geometry* 40.1 (2008), pp. 14–22.
- 49 [BC14] J. Ba and R. Caruana. "Do Deep Nets Really
50 Need to be Deep?" In: *Advances in Neural Informa-
51 tion Processing Systems* 27 (2014).
- 52 [BC17] D. Beck and T. Cohn. "Learning Kernels over
53 Strings using Gaussian Processes". In: *Proceedings of the Eighth International Joint Conference on Na-
54 tural Language Processing (Volume 2: Short Papers)*. Vol. 2. 2017, pp. 67–73.
- 55 [BC89] D. P. Bertsekas and D. A. Castanon. "The auc-
56 tion algorithm for the transportation problem". In: *Annals of Operations Research* 20.1 (1989), pp. 67–
57 96.
- 58 [BC93] B. M. Bell and F. W. Cathey. "The iterated
59 Kalman filter update as a Gauss-Newton method". In: *IEEE Trans. Automat. Contr.* 38.2 (Feb. 1993),
60 pp. 294–297.
- 61 [BC94] P. Baldi and Y. Chauvin. "Smooth online learn-
62 ing algorithms for Hidden Markov Models". In: *Neural Computation* 6 (1994), pp. 305–316.
- 63 [BC95] S. Baluja and R. Caruana. "Removing the
64 Genetics from the Standard Genetic Algorithm". In:
65 *ICML*. 1995, pp. 38–46.
- 66 [BCF10] E. Brochu, V. M. Cora, and N. de Freitas. "A
67 Tutorial on Bayesian Optimization of Expensive Cost
68 Functions, with Application to Active User Modeling
69 and Hierarchical Reinforcement Learning". In: (2010). arXiv: [1012.2599 \[cs.LG\]](https://arxiv.org/abs/1012.2599).
- 70 [BCH20] M. Briers, M. Charalambides, and C. Holmes.
71 "Risk scoring calculation for the current NHSx contact
72 tracing app". In: (2020). arXiv: [2005.11057 \[cs.CY\]](https://arxiv.org/abs/2005.11057).
- 73 [BCJ20] A. Buchholz, N. Chopin, and P. E. Ja-
74 cob. "Adaptive Tuning Of Hamiltonian Monte Carlo
75 Within Sequential Monte Carlo". In: *Bayesian Anal-*
76 (2020).
- 77 [BCN18] L. Bottou, F. E. Curtis, and J. Nocedal. "Opt-
78 imization Methods for Large-Scale Machine Learn-
79 ing". In: *SIAM Rev.* 60.2 (2018), pp. 223–311.
- 80 [BCNM06] C. Buciluă, R. Caruana, and A. Niculescu-
81 Mizil. "Model compression". In: *Proceedings of the 12th ACM SIGKDD international conference on
82 Knowledge discovery and data mining*. 2006, pp. 535–
83 541.
- 84 [BCV13] Y. Bengio, A. Courville, and P. Vincent. "Rep-
85 resentation learning: A review and new perspectives".
86 In: *IEEE transactions on pattern analysis and ma-
87 chine intelligence* 35.8 (2013), pp. 1798–1828.
- 88 [BD+10] S. Ben-David, J. Blitzer, K. Crammer, A.
89 Kulesza, F. Pereira, and J. W. Vaughan. "A theory
90 of learning from different domains". In: *Mach. Learn.*
91 79.1 (May 2010), pp. 151–175.
- 92 [BD11] A. Bhattacharya and D. B. Dunson. "Simplex
93 factor models for multivariate unordered categorical
94 data". In: *JASA* (2011).
- 95 [BD87] G. Box and N. Draper. *Empirical Model-
96 Building and Response Surfaces*. Wiley, 1987.

- 1
- 2 [BD92] D. Bayer and P. Diaconis. “Trailing the dovetail
3 shuffle to its lair”. In: *The Annals of Applied Probability* 2.2 (1992), pp. 294–313.
- 4 [BD97] S. Baluja and S. Davies. “Using Optimal
5 Dependency-Trees for Combinatorial Optimization:
6 Learning the Structure of the Search Space”. In: *ICML*.
7 1997.
- 8 [BDM09] R. Burkard, M. Dell’Amico, and S. Martello.
9 *Assignment Problems*. SIAM, 2009.
- 10 [BDM10] M. Briers, A. Doucet, and S. Maskel.
11 “Smoothing algorithms for state-space models”. In: *Annals of the Institute of Statistical Mathematics*
12 62.1 (2010), pp. 61–89.
- 13 [BDM17] M. G. Bellemare, W. Dabney, and R. Munos.
14 “A Distributional Perspective on Reinforcement Learning”. In: *ICML*. 2017.
- 15 [BDM18] N. Brosse, A. Durmus, and E. Moulines. “The
16 promises and pitfalls of Stochastic Gradient Langevin
17 Dynamics”. In: *NIPS*. 2018.
- 18 [BDS18] A. Brock, J. Donahue, and K. Simonyan.
19 “Large Scale GAN Training for High Fidelity Natural
20 Image Synthesis”. In: (2018). arXiv: [1809.11096 \[cs.LG\]](#).
- 21 [Bea03] M. Beal. “Variational Algorithms for Approximate Bayesian Inference”. PhD thesis. Gatsby Unit,
22 2003.
- 23 [Bea19] M. A. Beaumont. “Approximate Bayesian
24 Computation”. In: *Annual Review of Statistics and Its Application* 6.1 (2019), pp. 379–403.
- 25 [Béd08] M. Bédard. “Optimal acceptance rates for
26 Metropolis algorithms: Moving beyond 0.234”. In:
27 *Stochastic Process. Appl.* 118.12 (2008), pp. 2198–
2222.
- 28 [Beh+19] J. Behrmann, W. Grathwohl, R. T. Q. Chen,
29 D. Duvenaud, and J.-H. Jacobsen. “Invertible Residual
30 Networks”. In: *ICML*. 2019.
- 31 [Bel03] A. J. Bell. “The co-information lattice”. In: *ICA conference*. 2003.
- 32 [Bel+16] M. G. Bellemare, S. Srinivasan, G. Ostrovski,
33 T. Schaul, D. Saxton, and R. Munos. “Unifying
34 Count-Based Exploration and Intrinsic Motivation”. In: *NIPS*. 2016.
- 35 [Bel+18] M. I. Belghazi, A. Baratin, S. Rajeshwar, S.
36 Ozair, Y. Bengio, A. Courville, and D. Hjelm. “Mutual
37 Information Neural Estimation”. In: *ICML*. Ed. by J.
38 Dy and A. Krause. Vol. 80. Proceedings of Machine
39 Learning Research. PMLR, 2018, pp. 531–540.
- 40 [Bel+19] D. Belanger, S. Vora, Z. Mariet, R. Deshpande,
41 D. Dohan, C. Angermueller, K. Murphy, O. Chapelle,
42 and L. Colwell. “Biological Sequence Design using
43 Batched Bayesian Optimization”. In: *NIPS workshop on ML for the sciences*. 2019.
- 44 [Bel94] B. M. Bell. “The Iterated Kalman Smoother as
45 a Gauss–Newton Method”. In: *SIAM J. Optim.* 4.3
46 (Aug. 1994), pp. 626–636.
- 47 [Ben13] Y. Bengio. “Estimating or Propagating Gradients Through Stochastic Neurons”. In: (2013). arXiv:
48 [1305.2982 \[cs.LG\]](#).
- 49 [Ben+20] K. Benidis et al. “Neural forecasting: Intro-
50 duction and literature overview”. In: (2020). arXiv:
51 [2004.10240 \[cs.LG\]](#).
- 52 [Bén+21] C. Bénard, G. Biau, S. Veiga, and E. Scornet.
53 “Interpretable random forests via rule extraction”. In:
54 *International Conference on Artificial Intelligence
and Statistics*. PMLR, 2021, pp. 937–945.
- 55 [Ben+21a] Y. Bengio, T. Deleu, E. J. Hu, S. Lahouli,
56 M. Tiwari, and E. Bengio. “GFlowNet Foundations”.
57 In: (Nov. 2021). arXiv: [2111.09266 \[cs.LG\]](#).
- 58 [Ben+21b] G. W. Benton, W. J. Maddox, S. Lotfi, and
59 A. G. Wilson. “Loss Surface Simplexes for Mode Con-
60 nnecting Volumes and Fast Ensembling”. In: *ICML*.
61 2021.
- 62 [Ber05] J. M. Bernardo. “Reference Analysis”. In:
63 *Handbook of Statistics*. Ed. by D. K. Dey and C. R.
64 Rao. Vol. 25. Elsevier, 2005, pp. 17–90.
- 65 [Ber15] D. Bertsekas. *Convex Optimization Algo-
66 rithms*. Athena Scientific, 2015.
- 67 [Ber16] D. Bertsekas. *Nonlinear Programming*. Third.
68 Athena Scientific, 2016.
- 69 [Ber+18] R. van den Berg, L. Hasenclever, J. M. Tom-
70 czak, and M. Welling. “Sylvester normalizing flows for
71 variational inference”. In: *AISTATS*. 2018.
- 72 [Ber+19] H. Berard, G. Gidel, A. Almahairi, P. Vin-
73 cent, and S. Lacoste-Julien. “A Closer Look at the
74 Optimization Landscapes of Generative Adversarial
75 Networks”. In: *International Conference on Learning
76 Representations*. 2019.
- 77 [Ber19] D. Bertsekas. *Reinforcement learning and op-
78 timal control*. Athena Scientific, 2019.
- 79 [Ber+21] J. Berner, P. Grohs, G. Kutyniok, and P. Pe-
80 tersen. “The Modern Mathematics of Deep Learning”.
81 In: (2021). arXiv: [2105.04026 \[cs.LG\]](#).
- 82 [Ber85] J. Berger. “Bayesian Salesmanship”. In:
83 *Bayesian Inference and Decision Techniques with
84 Applications: Essays in Honor of Bruno deFinetti*. Ed.
85 by P. K. Goel and A. Zellner. North-Holland,
86 1985.
- 87 [Ber96] J. Bertoin. *Lévy processes*. Vol. 121. Cam-
88 bridge university press Cambridge, 1996.
- 89 [Ber97] D. Bertsekas. *Parallel and Distribution Com-
90 putation: Numerical Methods*. Athena Scientific,
91 1997.
- 92 [Ber99] A. Berchtold. “The double chain Markov
93 model”. In: *Comm. Stat. Theor. Methods* 28 (1999),
94 pp. 2569–2589.
- 95 [Bes75] J. Besag. “Statistical analysis of non-lattice
96 data”. In: *The Statistician* 24 (1975), pp. 179–196.
- 97 [Bet13] M. Betancourt. “A General Metric for Rieman-
98 nian Manifold Hamiltonian Monte Carlo”. In: *Geomet-
99 ric Science of Information*. Springer Berlin Heidelberg,
100 2013, pp. 327–334.
- 101 [Bet17] M. Betancourt. “A Conceptual Introduction to
102 Hamiltonian Monte Carlo”. In: (2017). arXiv: [1701.02434 \[stat.ME\]](#).
- 103 [BFH75] Y. Bishop, S. Fienberg, and P. Holland. *Dis-
crete Multivariate Analysis: Theory and Practice*. MIT Press, 1975.
- 104 [BFY20] T. D. Barfoot, J. R. Forbes, and D. Yoon.
105 “Exactly Sparse Gaussian Variational Inference with
106 Application to Derivative-Free Batch Nonlinear State
107 Estimation”. In: *Intl. J. of Robotics Research* (2020).
- 108 [BG06] M. Beal and Z. Ghahramani. “Variational
109 Bayesian Learning of Directed Graphical Models with
110 Hidden Variables”. In: *Bayesian Analysis* 1.4 (2006).
- 111 [BG13] M. J. Betancourt and M. Girolami. “Hamilto-
112 nian Monte Carlo for Hierarchical Models”. In: (2013).
113 arXiv: [1312.0906 \[stat.ME\]](#).

- [BG73] A. Björck and G. H. Golub. “Numerical methods for computing angles between linear subspaces”. In: *Mathematics of computation* 27.123 (1973), pp. 579–594.
- [BG96] A. Becker and D. Geiger. “A sufficiently fast algorithm for finding close to optimal junction trees”. In: *UAI*. 1996.
- [BGHM17] J. Boyd-Graber, Y. Hu, and D. Mimno. “Applications of Topic Models”. In: *Foundations and Trends® in Information Retrieval* 11.2-3 (2017), pp. 143–296.
- [BGM17] J. Ba, R. Grosse, and J. Martens. “Distributed Second-Order Optimization using Kronecker-Factored Approximations”. In: *ICLR*. openreview.net, 2017.
- [BGS16] Y. Burda, R. Grosse, and R. Salakhutdinov. “Importance Weighted Autoencoders”. In: *ICLR*. 2016.
- [BGT93] C. Berrou, A. Glavieux, and P. Thitimajshima. “Near Shannon limit error-correcting coding and decoding: Turbo codes”. In: *Proc. IEEE Intl. Comm. Conf.* (1993).
- [BH11] M.-F. Balcan and N. J. Harvey. “Learning submodular functions”. In: *Proceedings of the forty-third annual ACM symposium on Theory of computing*. 2011, pp. 793–802.
- [BH20] M. T. Bahadori and D. Heckerman. “Debiasing Concept-based Explanations with Causal Analysis”. In: *International Conference on Learning Representations*. 2020.
- [BH92] D. Barry and J. A. Hartigan. “Product partition models for change point problems”. In: *Annals of statistics* 20 (1992), pp. 260–279.
- [Bha+19] A. Bhadra, J. Datta, N. G. Polson, and B. T. Willard. “Lasso Meets Horseshoe: a survey”. In: *Bayesian Anal.* 34.3 (2019), pp. 405–427.
- [Bha+21] K. Bhatia, N. Kuang, Y. Ma, and Y. Wang. *Statistical and computational tradeoffs in variational Bayes: a case study of inferential model selection*. Tech. rep. 2021.
- [Bha+22] K. Bhatia, N. L. Kuang, Y.-A. Ma, and Y. Wang. “Statistical and Computational Trade-offs in Variational Inference: A Case Study in Inferential Model Selection”. In: (July 2022). arXiv: [2207.11208 \[stat.ML\]](#).
- [BHB19] P. Bachman, R. D. Hjelm, and W. Buchwalter. *Learning Representations by Maximizing Mutual Information Across Views*. 2019. arXiv: [1906.00910 \[cs.LG\]](#).
- [BHC19] M. Binkowski, D. Hjelm, and A. Courville. “Batch weight for domain adaptation with mass shift”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 1844–1853.
- [BHP02] M. Bădoiu, S. Har-Peled, and P. Indyk. “Approximate clustering via core-sets”. In: *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*. 2002, pp. 250–257.
- [BHW16] P. G. Bissiri, C. Holmes, and S. Walker. “A General Framework for Updating Belief Distributions”. In: *JRSSB* 78.5 (2016), 1103–1130.
- [Bic09] D. Bickson. “Gaussian Belief Propagation: Theory and Application”. PhD thesis. Hebrew University of Jerusalem, 2009.
- [Bie06] G. J. Bierman. *Factorization Methods for Discrete Sequential Estimation (Dover Books on Mathematics)*. en. Illustrated edition. Dover Publications, 2006.
- [Big+11] B. Biggio, G. Fumera, I. Pillai, and F. Roli. “A survey and experimental evaluation of image spam filtering techniques”. In: *Pattern recognition letters* 32.10 (2011), pp. 1436–1446.
- [Bil01] J. A. Bilmes. *Graphical Models and Automatic Speech Recognition*. Tech. rep. UWEETR-2001-0005. Univ. Washington, Dept. of Elec. Eng., 2001.
- [Bil22] J. Bilmes. “Submodularity In Machine Learning and Artificial Intelligence”. In: (2022). arXiv: [2202.00132 \[cs.LG\]](#).
- [Biń+18] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton. “Demystifying MMD GANs”. In: *ICLR*. 2018.
- [Bin+18] M. Binkowski, D. J. Sutherland, M. Arbel, and A. Gretton. “Demystifying MMD GANs”. In: *International Conference on Learning Representations*. 2018.
- [Bin+19] E. Bingham, J. P. Chen, M. Jankowiak, F. Obermeyer, N. Pradhan, T. Karaletsos, R. Singh, P. Szerlip, P. Horschall, and N. D. Goodman. “Pyro: Deep Universal Probabilistic Programming”. In: *JMLR* 20.28 (2019), pp. 1–6.
- [Biń+19] M. Bińkowski, J. Donahue, S. Dieleman, A. Clark, E. Elsen, N. Casagrande, L. C. Cobo, and K. Simonyan. “High Fidelity Speech Synthesis with Adversarial Networks”. In: *International Conference on Learning Representations*. 2019.
- [Bin+97] J. Binder, D. Koller, S. J. Russell, and K. Kanazawa. “Adaptive Probabilistic Networks with Hidden Variables”. In: *Machine Learning* 29 (1997), pp. 213–244.
- [Bis06] C. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [Bis99] C. Bishop. “Bayesian PCA”. In: *NIPS*. 1999.
- [Bit16] S. Bitzer. *The UKF exposed: How it works, when it works and when it's better to sample*. 2016.
- [Bit+21] J. Bitterwolf, A. Meinke, M. Augustin, and M. Hein. “Revisiting out-of-distribution detection: A simple baseline is surprisingly effective”. In: *ICML Workshop on Uncertainty in Deep Learning (UDL)*. 2021.
- [BJ05] F. Bach and M. Jordan. *A probabilistic interpretation of canonical correlation analysis*. Tech. rep. 688. U. C. Berkeley, 2005.
- [BJ+06] D. M. Blei, M. I. Jordan, et al. “Variational inference for Dirichlet process mixtures”. In: *Bayesian analysis* 1.1 (2006), pp. 121–143.
- [BJ06] W. Buntine and A. Jakulin. “Discrete Component Analysis”. In: *Subspace, Latent Structure and Feature Selection: Statistical and Optimization Perspectives*. Workshop. 2006.
- [BJV97] J. S. D. Bonet, C. L. I. Jr., and P. A. Viola. “MIMIC: Finding Optima by Estimating Probability Densities”. In: *NIPS*. MIT Press, 1997, pp. 424–430.
- [BK01] Y. Boykov and V. Kolmogorov. “An Experimental Comparison of Min-Cut/Max-Flow Algorithms for Energy Minimization in Computer Vision”. In: *Third International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition*. 2001.
- [BK10] R. Bardeh et and B. Kegl. “Surrogating the surrogate: accelerating Gaussian-process-based global optimization with a mixture cross-entropy algorithm”. In: *ICML*. 2010.

- 1
- 2 [BK15] D. Belanger and S. Kakade. “A Linear Dynamical System Model for Text”. en. In: *ICML*. 2015, pp. 833–842.
- 3
- 4 [BK19] M. Bonvini and E. H. Kennedy. “Sensitivity Analysis via the Proportion of Unmeasured Confounding”. In: *arXiv e-prints*. arXiv:1912.02793 (Dec. 2019), arXiv:1912.02793. arXiv: 1912.02793 [stat.ME].
- 5
- 6 [BKB17] O. Bastani, C. Kim, and H. Bastani. “Interpreting blackbox models via model extraction”. In: *arXiv preprint arXiv:1705.08504* (2017).
- 7
- 8 [BKH16] J. L. Ba, J. R. Kiros, and G. E. Hinton. “Layer Normalization”. In: (2016). arXiv: 1607.06450 [stat.ML].
- 9
- 10 [Bkj13] T. Broderick, B. Kulis, and M. Jordan. “MAD-Bayes: MAP-based asymptotic derivations from Bayes”. In: *International Conference on Machine Learning*. PMLR. 2013, pp. 226–234.
- 11
- 12 [BKM16] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. “Variational Inference: A Review for Statisticians”. In: *JASA* (2016).
- 13
- 14 [BKS19] A. Bennett, N. Kallus, and T. Schnabel. “Deep Generalized Method of Moments for Instrumental Variable Analysis”. In: *Advances in Neural Information Processing Systems*. 2019, pp. 3564–3574.
- 15
- 16 [BL06] D. Blei and J. Lafferty. “Dynamic topic models”. In: *ICML*. 2006, pp. 113–120.
- 17
- 18 [BL07a] C. M. Bishop and J. Lasserre. “Generative or discriminative? Getting the best of both worlds”. In: *Bayesian Statistics 8*. 2007.
- 19
- 20 [BL07b] D. Blei and J. Lafferty. “A Correlated Topic Model of “Science””. In: *Annals of Applied Stat.* 1.1 (2007), pp. 17–35.
- 21
- 22 [BLC18] A. J. Bose, H. Ling, and Y. Cao. “Adversarial Contrastive Estimation”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2018, pp. 1021–1032.
- 23
- 24 [Ble12] D. M. Blei. “Probabilistic topic models”. In: *Commun. ACM* 55.4 (2012), pp. 77–84.
- 25
- 26 [Ble17] D. Blei. *Variational inference: foundations and innovations (Lecture)*. Simons Institute Lecture. 2017.
- 27
- 28 [BLM16] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2016.
- 29
- 30 [Blo+16] A. Bloniarz, H. Liu, C.-H. Zhang, J. S. Sekhon, and B. Yu. “Lasso adjustments of treatment effect estimates in randomized experiments”. In: *Proceedings of the National Academy of Sciences* 113.27 (2016), pp. 7383–7390.
- 31
- 32 [BLS11] G. Blanchard, G. Lee, and C. Scott. “Generalizing from several related classification tasks to a new unlabeled sample”. In: *NIPS*. 2011.
- 33
- 34 [BLS17] J. Ballé, V. Laparra, and E. P. Simoncelli. “End-to-end Optimized Image Compression”. In: *ICLR*. 2017.
- 35
- 36 [Blu+15] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. “Weight Uncertainty in Neural Networks”. In: *ICML*. 2015.
- 37
- 38 [BM+18] G. Barth-Maron, M. W. Hoffman, D. Buden, W. Dabney, D. Horgan, T. B. Dhruva, A. Muldal, N. Heess, and T. Lillicrap. “Distributed Distributional Deterministic Policy Gradients”. In: *ICLR*. 2018.
- 39
- 40 [BM19] Y. Blau and T. Michaeli. “Rethinking Lossy Compression: The Rate-Distortion-Perception Trade-off”. In: *ICML*. 2019.
- 41
- 42 [BM+73] D. Blackwell, J. B. MacQueen, et al. “Ferguson distributions via Pólya urn schemes”. In: *The annals of statistics* 1.2 (1973), pp. 353–355.
- 43
- 44 [BMK20] Z. Boros, M. Mutny, and A. Krause. “Coresets via Bilevel Optimization for Continual Learning and Streaming”. In: *Advances in Neural Information Processing Systems* 33 (2020).
- 45
- 46 [BMP19] A. Brunel, D. Mazza, and M. Pagani. “Backpropagation in the Simply Typed Lambda-Calculus with Linear Negation”. In: *Proc. ACM Program. Lang. 4.POPL* (2019).
- 47
- 48 [BMR97] J. Binder, K. Murphy, and S. Russell. “Space-efficient inference in dynamic probabilistic networks”. In: *IJCAI*. 1997.
- 49
- 50 [BMS11] S. Bubeck, R. Munos, and G. Stoltz. “Pure Exploration in Finitely-armed and Continuous-armed Bandits”. In: *Theoretical Computer Science* 412.19 (2011), pp. 1832–1852.
- 51
- 52 [BNJ03a] D. Blei, A. Ng, and M. Jordan. “Latent Dirichlet allocation”. In: *JMLR* 3 (2003), pp. 993–1022.
- 53
- 54 [BNJ03b] D. M. Blei, A. Y. Ng, and M. I. Jordan. “Latent dirichlet allocation”. In: *JMLR* 3 (2003), pp. 993–1022.
- 55
- 56 [BO14] J. Bayer and C. Osendorfer. “Learning Stochastic Recurrent Networks”. In: *Workshop on Advances in Variational Inference*. 2014.
- 57
- 58 [Boe+05] P.-T. de Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein. “A Tutorial on the Cross-Entropy Method”. en. In: *Ann. Oper. Res.* 134.1 (2005), pp. 19–67.
- 59
- 60 [Boh92] D. Bohnig. “Multinomial logistic regression algorithm”. In: *Annals of the Inst. of Statistical Math.* 44 (1992), pp. 197–200.
- 61
- 62 [Bol89] K. Bollen. *Structural Equation Models with Latent Variables*. John Wiley & Sons, 1989.
- 63
- 64 [Bon64] G. Bonnet. “Transformations des signaux aléatoires à travers les systèmes non linéaires sans mémoire”. In: *Annales des Télécommunications* 19 (1964).
- 65
- 66 [Bor] A. Borodin. “Determinantal point processes”. In: *The Oxford Handbook of Random Matrix Theory*.
- 67
- 68 [Bös+17] J.-H. Böse, V. Flunkert, J. Gasthaus, T. Januschowski, D. Lange, D. Salinas, S. Schelter, M. Seeger, and Y. Wang. “Probabilistic Demand Forecasting at Scale”. In: *Proceedings VLDB Endowment* 10.12 (2017), pp. 1694–1705.
- 69
- 70 [Bot+13] L. Bottou, J. Peters, J. Quiñonero-Candela, D. X. Charles, D. M. Chickering, E. Portugaly, D. Ray, P. Simard, and E. Snelson. “Counterfactual Reasoning and Learning Systems: The Example of Computational Advertising”. In: *JMLR* 14 (2013), pp. 3207–3260.
- 71
- 72 [Bow+16a] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, and S. Bengio. “Generating Sentences from a Continuous Space”. In: *CONLL*. 2016.
- 73
- 74 [Bow+16b] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, and S. Bengio. “Generating Sentences from a Continuous Space”. In: *CONLL*. 2016.

- [Box80] G. E. P. Box. "Sampling and Bayes' Inference in Scientific Modelling and Robustness". In: *J. of Royal Stat. Soc. Series A* 143.4 (1980), pp. 383–430.
- [BP13] K. A. Bollen and J. Pearl. "Eight Myths About Causality and Structural Equation Models". In: *Handbook of Causal Analysis for Social Research*. Ed. by S. L. Morgan. Springer Netherlands, 2013, pp. 301–328.
- [BP16] E. Bareinboim and J. Pearl. "Causal inference and the data-fusion problem". en. In: *Proc. Natl. Acad. Sci. U. S. A.* 113.27 (2016), pp. 7345–7352.
- [BP21] T. Bricken and C. Pehlevan. "Attention Approximates Sparse Distributed Memory". In: *NIPS*. 2021.
- [BPK16] S. Bulo, L. Porzi, and P. Kortschieder. "Distillation dropout". In: *ICML*. 2016.
- [BPL21a] R. Balestrieri, J. Pesenti, and Y. LeCun. "Learning in High Dimension Always Amounts to Extrapolation". In: (Oct. 2021). arXiv: [2110.09485](https://arxiv.org/abs/2110.09485) [cs.LG].
- [BPL21b] A. Bardes, J. Ponce, and Y. LeCun. "Vicreg: Variance-invariance-covariance regularization for self-supervised learning". In: *arXiv preprint arXiv:2105.04906* (2021).
- [BPS16] A. G. Baydin, B. A. Pearlmutter, and J. M. Siskind. "DiffSharp: An AD library for .NET languages". In: *arXiv preprint arXiv:1611.03423* (2016).
- [BR05] H. Bang and J. M. Robins. "Doubly Robust Estimation in Missing Data and Causal Inference Models". In: *Biometrics* 61.4 (2005), pp. 962–973.
- [BR18] B. Biggio and F. Roli. "Wild patterns: Ten years after the rise of adversarial machine learning". In: *Pattern Recognition* 84 (2018), pp. 317–331.
- [BR98] S. Brooks and G. Roberts. "Assessing convergence of Markov Chain Monte Carlo algorithms". In: *Statistics and Computing* 8 (1998), pp. 319–335.
- [Bra+18] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, and S. Wanderman-Milne. *JAX: composable transformations of Python+NumPy programs*. Version 0.1.55. 2018.
- [Bra96] M. Brand. *Coupled hidden Markov models for modeling interacting processes*. Tech. rep. 405. MIT Lab for Perceptual Computing, 1996.
- [Bre01] L. Breiman. "Statistical modeling: The two cultures (with comments and a rejoinder by the author)". In: *Statistical science* 16.3 (2001), pp. 199–231.
- [Bre+17] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and regression trees*. Routledge, 2017.
- [Bre+20a] J. Brehmer, G. Louppe, J. Pavéz, and K. Cranmer. "Mining gold from implicit models to improve likelihood-free inference". en. In: *Proc. Natl. Acad. Sci. U. S. A.* 117.10 (2020), pp. 5242–5249.
- [Bre+20b] R. Brekelmans, V. Masrani, F. Wood, G. Ver Steeg, and A. Galstyan. "All in the Exponential Family: Bregman Duality in Thermodynamic Variational Inference". In: *ICML*. 2020.
- [Bre67] L. M. Bregman. "The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming". In: *USSR Computational Mathematics and Mathematical Physics* 7.3 (1967), pp. 200–217.
- [Bre91] Y. Brenier. "Polar factorization and monotone rearrangement of vector-valued functions". In: *Communications on pure and applied mathematics* 44.4 (1991), pp. 375–417.
- [Bre92] J. Breese. "Construction of belief and decision networks". In: *Computational Intelligence* 8 (1992), 624–647.
- [Bre96] L. Breiman. "Stacked regressions". In: *Mach. Learn.* 24.1 (1996), pp. 49–64.
- [BRG20] R. Bai, V. Rockova, and E. I. George. "Spike-and-slab meets LASSO: A review of the spike-and-slab LASSO". In: (2020). arXiv: [2010.06451](https://arxiv.org/abs/2010.06451) [stat.ME].
- [Bri50] G. W. Brier. "Verification of forecasts expressed in terms of probability". In: *Monthly Weather Review* 78.1 (1950), pp. 1–3.
- [Bro09] G. Brown. "A new perspective on information theoretic feature selection". In: *AISTATS*. 2009.
- [Bro+13] T. Broderick, N. Boyd, A. Wibisono, A. C. Wilson, and M. I. Jordan. "Streaming Variational Bayes". In: *NIPS*. 2013.
- [Bro+15] K. H. Brodersen, F. Gallusser, J. Koehler, N. Remy, and S. L. Scott. "Inferring causal impact using Bayesian structural time-series models". In: *Ann. Appl. Stat.* 9.1 (2015), pp. 247–274.
- [Bro18] T. Broderick. *Tutorial: Variational Bayes and Beyond*. 2018.
- [Bro19] J. Brownlee. *Generative Adversarial Networks with Python*. Accessed: 2019-8-27. Machine Learning Mastery, 2019.
- [Bro+20a] D. Brookes, A. Busia, C. Fannjiang, K. Murphy, and J. Listgarten. "A view of estimation of distribution algorithms through the lens of expectation-maximization". In: *GECCO '20*. Association for Computing Machinery, 2020, pp. 189–190.
- [Bro+20b] P. Brouillard, S. Lachapelle, A. Lacoste, S. Lacoste-Julien, and A. Drouin. "Differentiable Causal Discovery from Interventional Data". In: *NIPS*. July 2020.
- [Bro+20c] D. Brown, R. Coleman, R. Srinivasan, and S. Nieckum. "Safe imitation learning via fast bayesian reward inference from preferences". In: *International Conference on Machine Learning*. PMLR. 2020, pp. 1165–1177.
- [Bro+20d] T. B. Brown et al. "Language Models are Few-Shot Learners". In: (2020). arXiv: [2005 . 14165](https://arxiv.org/abs/2005.14165) [cs.CL].
- [BRS17] E. Balkanski, A. Rubinstein, and Y. Singer. "The Limitations of Optimization from Samples". In: *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*. STOC 2017. Montreal, Canada: Association for Computing Machinery, 2017, 1016–1027.
- [BRSS18] N. Bou-Rabee and J. M. Sanz-Serna. "Geometric integrators and the Hamiltonian Monte Carlo method". In: *Acta Numer.* (2018).
- [Bru+18] M. Brundage et al. "The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation". In: (2018). arXiv: [1802.07228](https://arxiv.org/abs/1802.07228) [cs.AI].
- [BS17] E. Balkanski and Y. Singer. "Minimizing a Submodular Function from Samples". In: *NIPS*. 2017, pp. 814–822.
- [BS18] S. Barratt and R. Sharma. "A note on the inception score". In: *arXiv preprint arXiv:1801.01973* (2018).

- 1
- 2 [BS20] E. Balkanski and Y. Singer. “A lower bound for
3 parallel submodular minimization”. In: *Proceedings of
the 52nd Annual ACM SIGACT Symposium on Theory
of Computing*. 2020, pp. 130–139.
- 4
- 5 [BS+20] S. Bobadilla-Suarez, C. Ahlheim, A. Mehrotra,
6 A. Panos, and B. C. Love. “Measures of neural
similarity”. In: *Computational Brain & Behavior* 3.4
(2020), pp. 369–383.
- 7
- 8 [BS95a] A. J. Bell and T. J. Sejnowski. “An infor-
9 mation maximisation approach to blind separation
and blind deconvolution”. In: *Neural Computation* 7.6
(1995), pp. 1129–1159.
- 10 [BS95b] A. J. Bell and T. J. Sejnowski. “An infor-
11 mation-maximization approach to blind separa-
12 tion and blind deconvolution”. In: *Neural computation*
7.6 (1995), pp. 1129–1159.
- 13 [BSA83] A. G. Barto, R. S. Sutton, and C. W. Ander-
14 son. “Neuronlike adaptive elements that can solve dif-
ficult learning control problems”. In: *SMC* 13.5 (1983),
pp. 834–846.
- 15
- 16 [BSF88] Y. Bar-Shalom and T. Fortmann. *Tracking
and data association*. Academic Press, 1988.
- 17
- 18 [BSL93] Y. Bar-Shalom and X. Li. *Estimation and
Tracking: Principles, Techniques and Software*.
Artech House, 1993.
- 19
- 20 [BSWT11] Y. Bar-Shalom, P. K. Willett, and X. Tian.
[BSWT11] *Tracking and Data Fusion: A Handbook of Algo-
rithms*. en. Yaakov Bar-Shalom, 2011.
- 21
- 22 [BT00] C. Bishop and M. Tipping. “Variational rele-
vance vector machines”. In: *UAI*. 2000.
- 23
- 24 [BT04] G. Bouchard and B. Triggs. “The tradeoff be-
tween generative and discriminative classifiers”. In:
*IASC International Symposium on Computational
Statistics (COMPSTAT '04)*. 2004.
- 25
- 26 [BT12] M. Botvinick and M. Toussaint. “Planning as
inference”. en. In: *Trends Cogn. Sci.* 16.10 (2012),
pp. 485–488.
- 27
- 28 [BT73] G. Box and G. Tiao. *Bayesian inference in sta-
tistical analysis*. Addison-Wesley, 1973.
- 29
- 30 [BTEGN09] A. Ben-Tal, L. El Ghaoui, and A. Ne-
mirovski. *Robust optimization*. Vol. 28. Princeton Uni-
versity Press, 2009.
- 31
- 32 [Buc+12] N. Buchbinder, M. Feldman, J. Naor, and R.
Schwartz. “A tight (1/2) linear-time approximation
to unconstrained submodular maximization”. In: *In
FOCS* (2012).
- 33
- 34 [Buc+17] C. L. Buckley, C. S. Kim, S. McGregor, and
A. K. Seth. “The free energy principle for action and
perception: A mathematical review”. In: *J. Math. Psy-
chol.* 81 (2017), pp. 55–79.
- 35
- 36 [Bud+21] K. Budhathoki, D. Janzing, P. Bloebaum,
and H. Ng. “Why did the distribution change?” In:
AISTATS. Ed. by A. Banerjee and K. Fukumizu.
Vol. 130. Proceedings of Machine Learning Research.
PMLR, 2021, pp. 1666–1674.
- 37
- 38 [Bul11] A. D. Bull. “Convergence rates of efficient
global optimization algorithms”. In: *JMLR* 12 (2011),
pp. 2879–2904.
- 39
- 40 [Bul+20] S. Bulusu, B. Kailkhura, B. Li, P. K. Varshney,
and D. Song. “Anomalous Example Detection in
Deep Learning: A Survey”. In: *IEEE Access* 8 (2020),
pp. 132330–132347.
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- [BV04] S. Boyd and L. Vandenberghe. *Convex optimi-
zation*. Cambridge, 2004.
- [BVHP18] S. Beery, G. Van Horn, and P. Perona.
“Recognition in terra incognita”. In: *Proceedings of the
European Conference on Computer Vision (ECCV)*.
2018, pp. 456–473.
- [BVW02] H. Bui, S. Venkatesh, and G. West. “Policy
Recognition in the Abstract Hidden Markov Model”.
In: *JAIR* 17 (2002), pp. 451–499.
- [BVZ01] Y. Boykov, O. Veksler, and R. Zabih. “Fast
Approximate Energy Minimization via Graph Cuts”.
In: *IEEE PAMI* 23.11 (2001).
- [BVZ99] Y. Boykov, O. Veksler, and R. Zabih. “Fast
Approximate Energy Minimization via Graph Cuts”.
In: *ICCV* (1). 1999, pp. 377–384.
- [BW20] G. J. J. van den Burg and C. K. I. Williams.
“An Evaluation of Change Point Detection Algo-
rithms”. In: (2020). arXiv: 2003.06222 [stat.ML].
- [BW21] G. J. van den Burg and C. K. Williams. “On
Memorization in Probabilistic Deep Generative Mod-
els”. In: *NIPS*. 2021.
- [BWM18] A. Buchholz, F. Wenzel, and S. Mandt.
“Quasi-Monte Carlo Variational Inference”. In: *ICML*.
2018.
- [BWR16] M. Bauer, M. van der Wilk, and C. E. Ras-
mussen. “Understanding Probabilistic Sparse Gaus-
sian Process Approximations”. In: *NIPS*. 2016,
pp. 1533–1541.
- [BYH20] O. Bohdal, Y. Yang, and T. Hospedales.
“Flexible Dataset Distillation: Learn Labels Instead of
Images”. In: *arXiv preprint arXiv:2006.08572* (2020).
- [BYM17] D. Belanger, B. Yang, and A. McCallum.
“End-to-End Learning for Structured Prediction En-
ergy Networks”. In: *ICML*. Ed. by D. Precup and
Y. W. Teh. Vol. 70. Proceedings of Machine Learning
Research. PMLR, 2017, pp. 429–439.
- [Byr+16] R Byrd, S Hansen, J Nocedal, and Y Singer.
“A Stochastic Quasi-Newton Method for Large-Scale
Optimization”. In: *SIAM J. Optim.* 26.2 (2016),
pp. 1008–1031.
- [BZ20] A. Barbu and S.-C. Zhu. *Monte Carlo Methods*.
en. Springer, 2020.
- [CA13] E. F. Camacho and C. B. Alba. *Model predic-
tive control*. Springer, 2013.
- [Cac+18] M. Caccia, L. Caccia, W. Fedus, H.
Larochelle, J. Pineau, and L. Charlin. “Language
GANs Falling Short”. In: *CoRR* abs/1811.02549
(2018). arXiv: 1811.02549.
- [CAII20] V. Coscrito, M. H. de Almeida Inácio, and
R. Izbicki. “The NN-Stacking: Feature weighted linear
stacking through neural networks”. In: *Neurocomput-
ing* (2020).
- [Cal+07] G. Calinescu, C. Chekuri, M. Pál, and J. Von-
drák. “Maximizing a submodular set function sub-
ject to a matroid constraint”. In: *Proceedings of the
12th International Conference on Integer Program-
ming and Combinatorial Optimization (IPCO)*. 2007,
pp. 182–196.
- [Cal20] O. Calin. *Deep Learning Architectures: A
Mathematical Approach*. en. 1st ed. Springer, 2020.
- [Cam+21] A. Campbell, Y. Shi, T. Rainforth, and A.
Doucet. “Online Variational Filtering and Parameter
Learning”. In: *NIPS*. 2021.

- [Can04] J. Canny. “GaP: a factor model for discrete data”. In: *SIGIR*. 2004, pp. 122–129.
- [Cao+15] Y. Cao, M. A Brubaker, D. J Fleet, and A. Hertzmann. “Efficient Optimization for Sparse Gaussian Process Regression”. In: *IEEE PAMI* 37.12 (2015), pp. 2415–2427.
- [Car03] P. Carbonetto. “Unsupervised Statistical Models for General Object Recognition”. MA thesis. University of British Columbia, 2003.
- [Car+15] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. “Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission”. In: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015, pp. 1721–1730.
- [Car+19] N. Carlini, A. Athalye, N. Papernot, W. Brendel, J. Rauber, D. Tsipras, I. Goodfellow, A. Madry, and A. Kurakin. “On evaluating adversarial robustness”. In: *arXiv preprint arXiv:1902.06705* (2019).
- [Car+21] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. *Emerging Properties in Self-Supervised Vision Transformers*. 2021. arXiv: [2104.14294 \[cs.CV\]](https://arxiv.org/abs/2104.14294).
- [Car97] R. Caruana. “Multitask Learning”. In: *Machine Learning* 28.1 (1997), pp. 41–75.
- [Cat08] A. Caticha. *Lectures on Probability, Entropy, and Statistical Physics*. 2008. arXiv: [0808 . 0012 \[physics.data-an\]](https://arxiv.org/abs/0808.0012).
- [Cat+11] A. Caticha, A. Mohammad-Djafari, J.-F. Bercher, and P. Bessière. “Entropic Inference”. In: *AIP Conference Proceedings* 1305.1 (2011), pp. 20–29. eprint: <https://doi.org/10.1063/1.3573619>.
- [CB20] Y. Chen and P. Bühlmann. “Domain adaptation under structural causal models”. In: *JMLR* (2020).
- [CBL20] K. Cranmer, J. Brehmer, and G. Louppe. “The frontier of simulation-based inference”. In: *Proceedings of the National Academy of Sciences* 117.48 (2020), pp. 30055–30062.
- [CBR20] Z. Chen, Y. Bei, and C. Rudin. “Concept whitening for interpretable image recognition”. In: *Nature Machine Intelligence* 2.12 (2020), pp. 772–782.
- [CC84] M. Conforti and G. Cornuejols. “Submodular set functions, matroids and the greedy algorithm: tight worst-case bounds and some generalizations of the Rado-Edmonds theorem”. In: *Discrete Applied Mathematics* 7.3 (1984), pp. 251–274.
- [CC96] M. Cowles and B. Carlin. “Markov Chain Monte Carlo Convergence Diagnostics: A Comparative Review”. In: *JASA* 91 (1996), pp. 883–904.
- [CCS22] A. Corenflos, N. Chopin, and S. Särkkä. “De-Sequentialized Monte Carlo: a parallel-in-time particle smoother”. In: (2022). arXiv: [2202.02264 \[stat.CO\]](https://arxiv.org/abs/2202.02264).
- [CDC15] C. Chen, N. Ding, and L. Carin. “On the Convergence of Stochastic Gradient MCMC Algorithms with High-Order Integrators”. In: *NIPS*. 2015.
- [CDS02] M. Collins, S. Dasgupta, and R. E. Schapire. “A Generalization of Principal Components Analysis to the Exponential Family”. In: *NIPS-14*. 2002.
- [CDS19] A. Clark, J. Donahue, and K. Simonyan. “Adversarial video generation on complex datasets”. In: *arXiv preprint arXiv:1907.06571* (2019).
- [Cér+12] F. Cérou, P. Del Moral, T. Furon, and A. Guyader. “Sequential Monte Carlo for rare event estimation”. In: *Stat. Comput.* 22.3 (2012), pp. 795–808.
- [CFG14] T. Chen, E. B. Fox, and C. Guestrin. “Stochastic Gradient Hamiltonian Monte Carlo”. In: *ICML*. 2014.
- [CG00] S. S. Chen and R. A. Gopinath. “Gaussianization”. In: *NIPS*. 2000, pp. 423–429.
- [CG15] X. Chen and A. Gupta. “Webly supervised learning of convolutional networks”. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 1431–1439.
- [CG18] C. Ceylan and M. U. Gutmann. “Conditional Noise-Contrastive Estimation of Unnormalised Models”. In: *International Conference on Machine Learning*. 2018, pp. 726–734.
- [CG96] S. Chen and J. Goodman. “An empirical study of smoothing techniques for language modeling”. In: *Proc. 34th ACL*. 1996, pp. 310–318.
- [CG98] S. Chen and J. Goodman. *An empirical study of smoothing techniques for language modeling*. Tech. rep. TR-10-98. Dept. Comp. Sci., Harvard, 1998.
- [CGR06] S. R. Cook, A. Gelman, and D. B. Rubin. “Validation of Software for Bayesian Models Using Posterior Quantiles”. In: *J. Comput. Graph. Stat.* 15.3 (2006), pp. 675–692.
- [CGS15] T. Chen, I. Goodfellow, and J. Shlens. “Net2net: Accelerating learning via knowledge transfer”. In: *International Conference on Learning Representations*. 2015.
- [CH20] C. Cinelli and C. Hazlett. “Making sense of sensitivity: extending omitted variable bias”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 82.1 (2020), pp. 39–67. eprint: <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/rssb.12348>.
- [Cha12] K. M. A. Chai. “Variational Multinomial Logit Gaussian Process”. In: *JMLR* 13.Jun (2012), pp. 1745–1808.
- [Cha14] N. Chapados. “Effective Bayesian Modeling of Groups of Related Count Time Series”. In: *ICML*. 2014.
- [Cha+14] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. “Return of the devil in the details: Delving deep into convolutional nets”. In: *British Machine Vision Conference*. 2014.
- [Cha+17] D. Chakrabarty, Y. T. Lee, A. Sidford, and S. C.-w. Wong. “Subquadratic submodular function minimization”. In: *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*. 2017, pp. 1220–1231.
- [Cha+18] N. S. Chatterji, N. Flammarion, Y.-A. Ma, P. L. Bartlett, and M. I. Jordan. “On the theory of variance reduction for stochastic gradient Monte Carlo”. In: *ICML*. 2018.
- [Cha+19a] C. Chan, S. Ginosar, T. Zhou, and A. A. Efros. “Everybody dance now”. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2019, pp. 5933–5942.
- [Cha+19b] T. Chavdarova, G. Gidel, F. Fleuret, and S. Lacoste-Julien. “Reducing noise in GAN training with variance reduced extragradient”. In: *Advances in Neural Information Processing Systems*. 2019, pp. 393–403.
- [Cha21] S. H. Chan. *Introduction to Probability for Data Science*. Michigan Publishing, 2021.

- 1
2 [Che+05] G. Chechik, A. Globerson, N. Tishby, and
Y. Weiss. “Information Bottleneck for Gaussian Vari-
ables”. In: *JMLR* 6.Jan (2005), pp. 165–188.
- 3
4 [Che+15] L.-C. Chen, G. Papandreou, I. Kokkinos, K.
Murphy, and A. L. Yuille. “Semantic Image Segmen-
tation with Deep Convolutional Nets and Fully Con-
nected CRFs”. In: *ICLR*. 2015.
- 5
6 [Che+16] X. Chen, Y. Duan, R. Houthooft, J. Schul-
man, I. Sutskever, and P. Abbeel. “InfoGAN: In-
terpretable Representation Learning by Information
Maximizing Generative Adversarial Nets”. In: *NIPS*.
2016.
- 7
8 [Che+17] T. Che, Y. Li, R. Zhang, R. D. Hjelm, W.
Li, Y. Song, and Y. Bengio. “Maximum-likelihood aug-
mented discrete generative adversarial networks”. In:
arXiv preprint arXiv:1702.07983 (2017).
- 9
10 [Che17] C. Chelba. *Language Modeling in the Era of
Abundant Data*. AI with the best. 2017.
- 11
12 [Che+17a] L.-C. Chen, G. Papandreou, I. Kokkinos, K.
Murphy, and A. L. Yuille. “DeepLab: Semantic Image
Segmentation with Deep Convolutional Nets, Atrous
Convolution, and Fully Connected CRFs”. In: *IEEE
PAMI* (2017).
- 13
14 [Che+17b] X. Chen, D. P. Kingma, T. Salimans, Y.
Duan, P. Dhariwal, J. Schulman, I. Sutskever, and P.
Abbeel. “Variational Lossy Autoencoder”. In: *ICLR*.
2017.
- 15
16 [Che+17c] X. Chen, N. Mishra, M. Rohaninejad, and
P. Abbeel. “PixelSNAIL: An Improved Autoregres-
sive Generative Model”. In: (2017). arXiv: 1712.09763
[cs.LG].
- 17
18 [Che+17d] V. Chernozhukov, D. Chetverikov, M.
Demirer, E. Duflo, C. Hansen, and W. Newey. “Double/
Debiased/Neyman Machine Learning of Treat-
ment Effects”. In: *American Economic Review* 107.5
(2017), pp. 261–65.
- 19
20 [Che+17e] V. Chernozhukov, D. Chetverikov, M.
Demirer, E. Duflo, C. Hansen, W. Newey, and
J. Robins. “Double/Debiased Machine Learning for
Treatment and Structural parameters”. In: *The
Econometrics Journal* (2017).
- 21
22 [Che+18a] C. Chen, W. Wang, Y. Zhang, Q. Su, and
L. Carin. “A convergence analysis for a class of practi-
cal variance-reduction stochastic gradient MCMC”. In:
Sci. China Inf. Sci. 62.1 (2018), p. 12101.
- 23
24 [Che+18b] C. Chen, O. Li, C. Tao, A. J. Barnett, J. Su,
and C. Rudin. “This looks like that: deep learning for
interpretable image recognition”. In: arXiv preprint
arXiv:1806.10574 (2018).
- 25
26 [Che+18c] R. T. Q. Chen, Y. Rubanova, J. Bettens-
court, and D. Duvenaud. “Neural Ordinary Differen-
tial Equations”. In: *NIPS*. 2018.
- 27
28 [Che+18d] X. Cheng, N. S. Chatterji, P. L. Bartlett,
and M. I. Jordan. “Underdamped Langevin MCMC:
A non-asymptotic analysis”. In: *COLT*. 2018.
- 29
30 [Che+19] R. T. Q. Chen, J. Behrmann, D. Duvenaud,
and J.-H. Jacobsen. “Residual Flows for Invertible
Generative Modeling”. In: *NIPS*. 2019.
- 31
32 [Che+20a] M. Chen, A. Radford, R. Child, J. Wu, H.
Jun, D. Luan, and I. Sutskever. “Generative pretrain-
ing from pixels”. In: *International Conference on Ma-
chine Learning*. PMLR. 2020, pp. 1691–1703.
- 33
34 [Che+20b] M. Chen, Y. Wang, T. Liu, Z. Yang, X. Li,
Z. Wang, and T. Zhao. “On computation and gener-
alization of generative adversarial imitation learning”.
In: arXiv preprint arXiv:2001.02792 (2020).
- 35
36 [Che+20c] T. Chen, S. Kornblith, M. Norouzi, and G.
Hinton. “A simple framework for contrastive learning
of visual representations”. In: *ICML*. 2020.
- 37
38 [Che+21] N. Chen, Y. Zhang, H. Zen, R. J. Weiss, M.
Norouzi, N. Dehak, and W. Chan. “WaveGrad 2: Iter-
ative refinement for text-to-speech synthesis”. In: *In-
terspeech 2021*. ISCA: ISCA, Aug. 2021.
- 39
40 [Che+22] D. Chen, D. Wang, T. Darrell, and S. E.
Ebrahimi. “Contrastive Test-Time Adaptation”. In:
CVPR. Apr. 2022.
- 41
42 [Che95] Y. Cheng. “Mean shift, mode seeking, and clus-
tering”. In: *IEEE PAMI* 17.8 (1995).
- 43
44 [Chi14] S. Chiappa. “Explicit-Duration Markov Switch-
ing Models”. In: *Foundations and Trends in Machine
Learning* 7.6 (2014), pp. 803–886.
- 45
46 [Chi21a] R. Child. “Very Deep VAEs Generalize Au-
toregressive Models and Can Outperform Them on
Images”. In: *ICLR*. 2021.
- 47
48 [Chi21b] R. Child. “Very Deep VAEs Generalize Au-
toregressive Models and Can Outperform Them on
Images”. In: ArXiv abs/2011.10650 (2021).
- 49
50 [CHL+05] S. Chopra, R. Hadsell, Y. LeCun, et al.
“Learning a similarity metric discriminatively, with
application to face verification”. In: *CVPR*. 2005,
pp. 539–546.
- 51
52 [CHM97] D. M. Chickering, D. Heckerman, and C.
Meek. “A Bayesian Approach to Learning Bayesian
Networks with Local Structure”. In: *UAI*. UAI’97.
1997, pp. 80–89.
- 53
54 [Cho02] N. Chopin. “A Sequential Particle Filter
Method for Static Models”. In: *Biometrika* 89.3
(2002), pp. 539–551.
- 55
56 [Cho11] M. J. Choi. “Trees and Beyond: Exploiting and
Improving Tree-Structured Graphical Models”. PhD
thesis. MIT, 2011.
- 57
58 [Cho15] Y. Chow, A. Tamar, S. Mannor, and M.
Pavone. “Risk-Sensitive and Robust Decision-Making:
A CVaR Optimization Approach”. In: *NIPS*. 2015,
pp. 1522–1530.
- 59
60 [Cho21] F. Chollet. *Deep learning with Python (sec-
ond edition)*. Manning, 2021.
- 61
62 [Cho57] N. Chomsky. *Syntactic Structures*. Mouton,
1957.
- 63
64 [Chr+21] R. Christiansen, N. Pfister, M. E. Jakobsen,
N. Gnecco, and J. Peters. “A causal framework for dis-
tribution generalization”. In: *IEEE PAMI* (2021).
- 65
66 [Chu+15] J. Chung, K. Kastner, L. Dinh, K. Goel, A.
Courville, and Y. Bengio. “A Recurrent Latent Vari-
able Model for Sequential Data”. In: *NIPS*. 2015.
- 67
68 [Chu+18] K. Chua, R. Calandra, R. McAllister, and
S. Levine. “Deep Reinforcement Learning in a Hand-
ful of Trials using Probabilistic Dynamics Models”. In:
NIPS. 2018.
- 69
70 [Chu+19] G. Chuang, G. DeSalvo, L. Karydas, J.-F.
Kagy, A. Rostamizadeh, and A. Theeraphol. “Active
Learning Empirical Study”. In: *NeurIPS LIRE Work-
shop*. 2019.
- 71
72 [Chw+15] K. Chwialkowski, A. Ramdas, D. Sejdinovic,
and A. Gretton. “Fast Two-Sample Testing with An-
alytic Representations of Probability Measures”. In:
NIPS. 2015.

- 1 [CI80] D. R. Cox and V. Isham. *Point processes*. Vol. 12. CRC Press, 1980.
- 2 [CJ21] A. D. Cobb and B. Jalaian. “Scaling Hamiltonian Monte Carlo inference for Bayesian neural networks with symmetric splitting”. In: *UAI*. Vol. 161. Proceedings of Machine Learning Research. PMLR, 2021, pp. 675–685.
- 3 [CK05] M. Collins and T. Koo. “Discriminative Reranking for Natural Language Parsing”. In: *Proc. ACL*. 2005.
- 4 [CK07] J. J. F. Commandeur and S. J. Koopman. *An Introduction to State Space Time Series Analysis (Practical Econometrics)*. en. 1st ed. Oxford University Press, 2007.
- 5 [CK21] D. Chakrabarty and S. Khanna. “Better and simpler error analysis of the Sinkhorn–Knopp algorithm for matrix scaling”. In: *Mathematical Programming* 188.1 (2021), pp. 395–407.
- 6 [CK94a] D. Card and A. B. Krueger. “Minimum Wages and Employment: A Case Study of the Fast-Food Industry in New Jersey and Pennsylvania”. In: *American Economic Review* 84.4 (1994), pp. 772–793.
- 7 [CK94b] C. Carter and R. Kohn. “On Gibbs sampling for state space models”. In: *Biometrika* 81.3 (1994), pp. 541–553.
- 8 [CK96] C. Carter and R. Kohn. “Markov Chain Monte Carlo in Conditionally Gaussian State Space Models”. In: *Biometrika* 83 (1996), pp. 589–601.
- 9 [CKK17] L. Chen, A. Krause, and A. Karbasi. “Interactive Submodular Bandit”. In: *NIPS*. 2017, pp. 141–152.
- 10 [CL00] R. Chen and S. Liu. “Mixture Kalman Filters”. In: *J. Royal Stat. Soc. B* (2000).
- 11 [CL07] L. Carvahlo and C. Lawrence. “Centroid estimation in discrete high-dimensional spaces with applications in biology”. In: *PNAS* 105.4 (2007).
- 12 [CL11] O. Chapelle and L. Li. “An empirical evaluation of Thompson sampling”. In: *NIPS*. 2011.
- 13 [CL18] Z. Chen and B. Liu. *Lifelong Machine Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan Claypool, 2018.
- 14 [CL96] B. P. Carlin and T. A. Louis. *Bayes and Empirical Bayes Methods for Data Analysis*. Chapman and Hall, 1996.
- 15 [Cla20] P. Clavier. “Sum-Product Network in the context of missing data”. en. MA thesis. KTH, 2020.
- 16 [Cla21] A. Clayton. *Bernoulli's Fallacy: Statistical Illogic and the Crisis of Modern Science*. en. Columbia University Press, 2021.
- 17 [CLD18] C. Cremer, X. Li, and D. Duvenaud. “Inference Suboptimality in Variational Autoencoders”. In: *ICML*. 2018.
- 18 [Clo+19] J. R. Clough, I. Oksuz, E. Puyol-Antón, B. Ruijsink, A. P. King, and J. A. Schnabel. “Global and local interpretability for cardiac MRI classification”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 656–664.
- 19 [Clo20] Cloudera. *Causality for ML*. 2020.
- 20 [CLV19] M. Cox, T. van de Laar, and B. de Vries. “A factor graph approach to automated design of Bayesian signal processing algorithms”. In: *Int. J. Approx. Reason.* 104 (2019), pp. 185–204.
- 21 [CLW18] Y. Chen, L. Li, and M. Wang. “Scalable Bilinear π -Learning Using State and Action Features”. In: *ICML*. 2018, pp. 833–842.
- 22 [CM09] O. Cappé and E. Moulines. “Online EM Algorithm for Latent Data Models”. In: *J. of Royal Stat. Soc. Series B* 71.3 (2009), pp. 593–613.
- 23 [CMD17] C. Cremer, Q. Morris, and D. Duvenaud. “Reinterpreting Importance-Weighted Autoencoders”. In: *ICLR Workshop*. 2017.
- 24 [CMR05] O. Cappe, E. Moulines, and T. Ryden. *Inference in Hidden Markov Models*. Springer, 2005.
- 25 [CMR12] C. Cortes, M. Mohri, and A. Rostamizadeh. “Algorithms for learning kernels based on centered alignment”. In: *The Journal of Machine Learning Research* 13.1 (2012), pp. 795–828.
- 26 [CMS12] D. Ciregan, U. Meier, and J. Schmidhuber. “Multi-column deep neural networks for image classification”. In: *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 2012, pp. 3642–3649.
- 27 [CN01] H. Choset and K. Nagatani. “Topological simultaneous localization and mapping (SLAM): toward exact localization without explicit localization”. In: *IEEE Trans. Robotics and Automation* 17.2 (2001).
- 28 [CNW20] M. Collier, A. Nazabal, and C. K. I. Williams. “VAEs in the Presence of Missing Data”. In: *ICML Workshop on the Art of Learning with Missing Values*. 2020.
- 29 [CO06] N. Chater and M. Oaksford. “Mental mechanisms”. In: *Information sampling and adaptive cognition* (2006), pp. 210–236.
- 30 [COB18] L. Chizat, E. Oyallon, and F. Bach. “On Lazy Training in Differentiable Programming”. In: (2018). arXiv: [1812.07956 \[math.OC\]](https://arxiv.org/abs/1812.07956).
- 31 [Cor+12] G. Cormode, M. Garofalakis, P. J. Haas, and C. Jermaine. “Synopses for massive data: Samples, histograms, wavelets, sketches”. In: *Foundations and Trends in Databases* 4.1–3 (2012), pp. 1–294.
- 32 [Cor17] G. Cormode. “Data sketching”. In: *Communications of the ACM* 60.9 (2017), pp. 48–55.
- 33 [Cor+59] J. Cornfield, W. Haenszel, E. C. Hammond, A. M. Lilienfeld, M. B. Shimkin, and E. L. Wynder. “Smoking and Lung Cancer: Recent Evidence and a Discussion of Some Questions”. In: *JNCI: Journal of the National Cancer Institute* 22.1 (Jan. 1959), pp. 173–203. eprint: <https://academic.oup.com/jnci/article-pdf/22/1/173/2704718/22-1-173.pdf>.
- 34 [Cor+87] A. Corana, M. Marchesi, C. Martini, and S. Ridella. “Minimizing Multimodal Functions of Continuous Variables with the ‘Simulated Annealing’ Algorithm”. In: *ACM Trans. Math. Softw.* 13.3 (1987), pp. 262–280.
- 35 [Cou16] Council of European Union. *General Data Protection Regulation*. 2016.
- 36 [Cou+20] J. Courts, A. Wills, T. Schön, and B. Ninness. “Variational System Identification for Nonlinear State-Space Models”. In: (2020). arXiv: [2012.05072 \[stat.ML\]](https://arxiv.org/abs/2012.05072).
- 37 [Cou+21] J. Courts, J. Hendriks, A. Wills, T. Schön, and B. Ninness. “Variational State and Parameter Estimation”. In: *19th IFAC Symposium on System Identification SYSID 2021*. 2021.
- 38 [Cov99] T. M. Cover. *Elements of information theory*. John Wiley & Sons, 1999.

- 1
- 2 [CP20a] R. Chen and I. C. Paschalidis. *Distributionally Robust Learning*. NOW Foundations and Trends
3 in Optimization, 2020.
- 4 [CP20b] N. Chopin and O. Papaspiliopoulos. *An Introduction to Sequential Monte Carlo*. en. 1st ed.
5 Springer, 2020.
- 6 [CPD17] P. Constantinou and A Philip Dawid. “Extended conditional independence and applications in
7 causal inference”. en. In: *Ann. Stat.* 45.6 (2017),
8 pp. 2618–2653.
- 9 [CPS10] C. Carvalho, N. Polson, and J. Scott. “The
10 horseshoe estimator for sparse signals”. In: *Biometrika*
10 97.2 (2010), p. 465.
- 11 [Cri+02] N. Cristianini, J. Shawe-Taylor, A. Elisseeff,
12 and J. S. Kandola. “On Kernel-Target Alignment”. In:
13 *Advances in Neural Information Processing Systems*.
2002, pp. 367–373.
- 14 [CRK19] J. M. Cohen, E. Rosenfeld, and J. Z.
15 Kolter. “Certified adversarial robustness via randomiz-
16 ed smoothing”. In: *arXiv preprint arXiv:1902.02918*
16 (2019).
- 17 [Cro+11] D. F. Crouse, P. Willett, K. Pattipati, and
17 L. Svensson. “A look at Gaussian mixture reduction
18 algorithms”. In: *14th International Conference on In-
formation Fusion*. 2011, pp. 1–8.
- 19 [CS04] I. Csiszár and P. C. Shields. “Information the-
20 ory and statistics: A tutorial”. In: (2004).
- 21 [CS09] Y. Cho and L. K. Saul. “Kernel Methods for
22 Deep Learning”. In: *NIPS*. 2009, pp. 342–350.
- 22 [CS18] P. Chaudhari and S. Soatto. “Stochastic gra-
23 dient descent performs variational inference, converges
24 to limit cycles for deep networks”. In: *ICLR*. 2018.
- 25 [CSF16] A. W. Churchill, S. Sigtia, and C. Fernando.
26 “Learning to Generate Genotypes with Neural Net-
26 works”. In: (2016). arXiv: [1604.04153 \[cs.NE\]](https://arxiv.org/abs/1604.04153).
- 27 [Csi67] I. Csiszar. “Information-Type Measures of Dif-
28 ference of Probability Distributions and Indirect Ob-
29 servations”. In: *Acta Scientiarum Mathematicarum
Hungarica* 2 (1967), pp. 299–318.
- 30 [CSN21] J. Couillon, L. South, and C. Nemeth.
31 “Stochastic Gradient MCMC with Multi-Armed Ban-
31 dit Tuning”. In: (2021). arXiv: [2105.13059 \[stat.CO\]](https://arxiv.org/abs/2105.13059).
- 32 [CT06] T. M. Cover and J. A. Thomas. *Elements of
Information Theory*. 2nd edition. John Wiley, 2006.
- 33 [CT+19] M. F. Cusumano-Towner, F. A. Saad, A. K.
34 Lew, and V. K. Mansinghka. “Gen: a general-purpose
35 probabilistic programming system with programmable
36 inference”. In: *Proceedings of the 40th ACM SIG-
PLAN Conference on Programming Language De-
sign and Implementation*. PLDI 2019. Association for
Computing Machinery, 2019, pp. 221–236.
- 37 [CT91] T. M. Cover and J. A. Thomas. *Elements of
Information Theory*. John Wiley, 1991.
- 38 [CTM17] M. F. Cusumano-Towner and V. K. Mans-
39 inghka. “AIDE: An algorithm for measuring the accu-
40 racy of probabilistic inference algorithms”. In: *NIPS*.
41 2017.
- 42 [CTN17] Y. Chali, M. Tanvee, and M. T. Nayeem. “To-
43 wards abstractive multi-document summarization us-
44 ing submodular function-based framework, sentence
44 compression and merging”. In: *Proceedings of the
45 Eighth International Joint Conference on Natural
Language Processing (Volume 2: Short Papers)*. 2017,
45 pp. 418–424.
- 46
- 47
- [CTS78] C. Cannings, E. A. Thompson, and M. H. Skol-
nick. “Probability functions in complex pedigrees”. In:
Advances in Applied Probability 10 (1978), pp. 26–61.
- [CUH16] D.-A. Clevert, T. Unterthiner, and S. Hochreiter. “Fast and Accurate Deep Network Learning by
Exponential Linear Units (ELUs)”. In: *ICLR*. 2016.
- [Cui+18] Y. Cui, Y. Song, C. Sun, A. Howard, and S.
Belongie. “Large scale fine-grained categorization and
domain-specific transfer learning”. In: *Proceedings of
the IEEE conference on computer vision and pattern
recognition*. 2018, pp. 4109–4118.
- [Cun22] Y. L. Cun. *A path towards autonomous AI*.
2022.
- [Cun83] W. H. Cunningham. “Decomposition of sub-
modular functions”. In: *Combinatorica* 3.1 (1983),
pp. 53–68.
- [Cut13] M. Cuturi. “Sinkhorn Distances: Lightspeed
Computation of Optimal Transportation Distances”.
In: *NIPS*. 2013.
- [CW07] C. M. Carvalho and M. West. “Dynamic
Matrix-Variate Graphical Models”. In: *Bayesian Anal-
ysis* 2.1 (2007), pp. 69–98.
- [CW16] T. Cohen and M. Welling. “Group Equiva-
lent Convolutional Networks”. en. In: *ICML*. 2016,
pp. 2990–2999.
- [CWG20] D. C. Castro, I. Walker, and B. Glocker.
“Causality matters in medical imaging”. en. In: *Nat.
Commun.* 11.1 (2020), p. 3673.
- [CWS21] J. Courts, A. G. Wills, and T. B. Schön.
“Gaussian Variational State Estimation for Nonlinear
State-Space Models”. In: *IEEE Trans. Signal Process.*
69 (2021), pp. 5979–5993.
- [CXH21] X. Chen, S. Xie, and K. He. “An empirical
study of training self-supervised vision transformers”.
In: *arXiv preprint arXiv:2104.02057* (2021).
- [CY20] G. Cormode and K. Yi. *Small Summaries for
Big Data*. Cambridge University Press, 2020.
- [Cza+20] J. Czarnowski, T. Laidlow, R. Clark, and
A. J. Davison. “DeepFactors: Real-Time Probabilistic
Dense Monocular SLAM”. In: *ICRA*. 2020.
- [CZG20] B. Charpentier, D. Zügner, and S. Günne-
mann. “Posterior network: Uncertainty estimation
without ood samples via density-based pseudo-counts”.
In: *arXiv preprint arXiv:2006.09239* (2020).
- [D'A+20] A. D'Amour et al. “Underspecification
Presents Challenges for Credibility in Modern Ma-
chine Learning”. In: (2020). arXiv: [2011.03395 \[cs.LG\]](https://arxiv.org/abs/2011.03395).
- [D'A+21] A. D'Amour, P. Ding, A. Feller, L. Lei, and
J. Sekhon. “Overlap in observational studies with high-
dimensional covariates”. In: *Journal of Econometrics*
221.2 (2021), pp. 644–654.
- [Dag+21] N. Dagan, N. Barda, E. Kepten, O. Miron,
S. Perchik, M. A. Katz, M. A. Hernán, M. Lipsitch,
B. Reis, and R. D. Balicer. “BNT162b2 mRNA Covid-
19 Vaccine in a Nationwide Mass Vaccination Setting”.
In: *New England Journal of Medicine* 384.15 (2021),
pp. 1412–1423. eprint: <https://doi.org/10.1056/NEJMoa2101765>.
- [Dai+17] H. Dai, B. Dai, Y.-M. Zhang, S. Li, and
L. Song. “Recurrent Hidden Semi-Markov Model”. In:
ICLR. 2017.
- [Dai+18] B. Dai, A. Shaw, L. Li, L. Xiao, N. He, Z.
Liu, J. Chen, and L. Song. “SBEED: Convergent Rein-

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
- forcement Learning with Nonlinear Function Approximation". In: *ICML*. 2018, pp. 1133–1142.
- [Dai+19a] B. Dai, H. Dai, A. Gretton, L. Song, D. Schuurmans, and N. He. "Kernel exponential family estimation via doubly dual embedding". In: *AISTATS*. PMLR. 2019, pp. 2321–2330.
- [Dai+19b] B. Dai, Z. Liu, H. Dai, N. He, A. Gretton, L. Song, and D. Schuurmans. "Exponential family estimation via adversarial dynamics embedding". In: *Advances in Neural Information Processing Systems*. 2019, pp. 10979–10990.
- [Dai+20a] C. Dai, J. Heng, P. E. Jacob, and N. Whiteley. "An invitation to sequential Monte Carlo samplers". In: (2020). arXiv: [2007.11936 \[stat.CO\]](https://arxiv.org/abs/2007.11936).
- [Dai+20b] Z. Dai, G. Lai, Y. Yang, and Q. V. Le. "Funnel-Transformer: Filtering out Sequential Redundancy for Efficient Language Processing". In: *NIPS*. 2020.
- [Dal+04] N. Dalvi, P. Domingos, S. Sanghavi, and D. Verma. "Adversarial classification". In: *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*. 2004, pp. 99–108.
- [Dar03] A. Darwiche. "A Differential Approach to Inference in Bayesian Networks". In: *J. ACM* 50.3 (2003), pp. 280–305.
- [Dar09] A. Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge, 2009.
- [Dar+11] S. Darolles, Y. Fan, J.-P. Florens, and E. Renault. "Nonparametric instrumental regression". In: *Econometrica* 79.5 (2011), pp. 1541–1565.
- [Dar80] R. A. Darton. "Rotation in Factor Analysis". In: *Journal of the Royal Statistical Society. Series D (The Statistician)* 29.3 (1980), pp. 167–194.
- [Dau07] H. Daume. "Fast search for Dirichlet process mixture models". In: *AISTATS*. 2007.
- [Dav+04] T. A. Davis, J. R. Gilbert, S. I. Larimore, and E. G. Ng. "A Column Approximate Minimum Degree Ordering Algorithm". In: *ACM Trans. Math. Softw.* 30.3 (2004), pp. 353–376.
- [Dav+18] T. R. Davidson, L. Falorsi, N. De Cao, T. Kipf, and J. M. Tomczak. "Hyperspherical Variational Auto-Encoders". In: *UAI*. 2018.
- [Daw00] A. P. Dawid. "Causal Inference Without Counterfactuals". In: *JASA* 95.450 (2000), pp. 407–424.
- [Daw02] A. P. Dawid. "Influence diagrams for causal modelling and inference". In: *Intl. Stat. Review* 70 (2002). Corrections p437, pp. 161–189.
- [Daw15] A. P. Dawid. "Statistical Causality from a Decision-Theoretic Perspective". In: *Annu. Rev. Stat. Appl.* 2.1 (2015), pp. 273–303.
- [Daw82] A. P. Dawid. "The Well-Calibrated Bayesian". In: *JASA* 77.379 (1982), pp. 605–610.
- [Dax+21] E. Daxberger, A. Kristiadi, A. Immer, R. Eschenhagen, M. Bauer, and P. Hennig. "Laplace Redux—Effortless Bayesian Deep Learning". In: *NIPS*. 2021.
- [Day+95] P. Dayan, G. Hinton, R. Neal, and R. Zemel. "The Helmholtz machine". In: *Neural Networks* 9.8 (1995).
- [DB+13] P. De Blasi, S. Favaro, A. Lijoi, R. H. Mena, I. Prünster, and M. Ruggiero. "Are Gibbs-type priors the most natural generalization of the Dirichlet process?". In: *IEEE PAMI* 37.2 (2013), pp. 212–229.
- [DB+21] V. De Bortoli, J. Thornton, J. Heng, and A. Doucet. "Diffusion Schrödinger Bridge with Applications to Score-Based Generative Modeling". In: *NIPS*. June 2021.
- [DBB20] S. Daulton, M. Balandat, and E. Bakshy. "Differentiable Expected Hypervolume Improvement for Parallel Multi-Objective Bayesian Optimization". In: *NIPS*. 2020.
- [DPB19] C. Durkan, A. Bekasov, and I. M. G. Papamakarios. "Neural Spline Flows". In: *NIPS*. 2019.
- [DBW20] I. A. Delbridge, D. S. Bindel, and A. G. Wilson. "Randomly Projected Additive Gaussian Processes for Regression". In: *International Conference on Machine Learning*. 2020.
- [DCF+15] E. Denton, S. Chintala, R. Fergus, et al. "Deep generative image models using a Laplacian pyramid of adversarial networks". In: *NIPS*. 2015.
- [DD22] S. Doyen and N. B. Dadario. "12 Plagues of AI in Healthcare: A Practical Guide to Current Issues With Using Machine Learning in a Medical Context". en. In: *Front Digit Health* 4 (May 2022), p. 765406.
- [DDL97] S. DellaPietra, V. DellaPietra, and J. Lafferty. "Inducing features of random fields". In: *IEEE PAMI* 19.4 (1997).
- [DE00] R. Dahlhaus and M. Eichler. "Causality and graphical models for time series". In: *Highly structured stochastic systems*. Ed. by P. Green, N. Hjort, and S. Richardson. Oxford University Press, 2000.
- [DE04] J. Dow and J. Endersby. "Multinomial probit and multinomial logit: a comparison of choice models for voting research". In: *Electoral Studies* 23.1 (2004), pp. 107–122.
- [Dec96] R. Dechter. "Bucket elimination: a unifying framework for probabilistic inference". In: *UAI*. 1996.
- [DeG70] M. DeGroot. *Optimal Statistical Decisions*. McGraw-Hill, 1970.
- [DEL20] H. M. Dolatabadi, S. Erfani, and C. Leckie. "Invertible Generative Modeling using Linear Rational Splines". In: *AISTATS*. 2020, pp. 4236–4246.
- [Del+21] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardi, G. Slabaugh, and T. Tuytelaars. "A continual learning survey: Defying forgetting in classification tasks". en. In: *IEEE PAMI* (2021).
- [Den+02] D. Denison, C. Holmes, B. Mallick, and A. Smith. *Bayesian methods for nonlinear classification and regression*. Wiley, 2002.
- [Den+20] Y. Deng, A. Bakhtin, M. Ott, A. Szlam, and M. Ranzato. "Residual Energy-Based Models for Text Generation". In: *International Conference on Learning Representations*. 2020.
- [Dev+18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. "Bert: Pre-training of deep bidirectional transformers for language understanding". In: *arXiv preprint arXiv:1810.04805* (2018).
- [Dev+19] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding". In: *NAACL*. 2019.
- [Dev+21] L. Devlin, P. Horridge, P. L. Green, and S. Maskell. "The No-U-Turn Sampler as a Proposal Distribution in a Sequential Monte Carlo Sampler with a Near-Optimal L-Kernel". In: (2021). arXiv: [2108.02498 \[stat.CO\]](https://arxiv.org/abs/2108.02498).

- 1
- 2 [Dev85] P. A. Devijver. “Baum’s forward-backward al-
gorithm revisited”. In: *Pattern Recognition Letters*
3 3.6 (1985), pp. 369–373.
- 4 [Dex] Dex: Research language for array processing in
the Haskell/ML family. [https://github.com/google-
research/dex-lang](https://github.com/google-research/dex-lang). 2019.
- 5 [DF18] E. Denton and R. Fergus. “Stochastic Video
Generation with a Learned Prior”. In: *ICML*. 2018.
- 6 [DF19] X. Ding and D. J. Freedman. “Learning Deep
Generative Models with Annealed Importance Sam-
pling”. In: (2019). arXiv: 1906.04904 [stat.ML].
- 7 [DF21] F. D’Angelo and V. Fortuin. “Repulsive Deep
Ensembles are Bayesian”. In: *NIPS*. May 2021.
- 8 [dff21] S. Dozinski, U. Feige, and M. Feldman. “Are
Gross Substitutes a Substitute for Submodular Valua-
tions?”. In: *arXiv preprint arXiv:2102.13343* (2021).
- 9 [DFO20] M. Deisenroth, A. Faisal, and C. S. Ong.
Mathematics for machine learning. Cambridge, 2020.
- 10 [DFR15] M. P. Deisenroth, D. Fox, and C. E. Ras-
mussen. “Gaussian Processes for Data-Efficient Learn-
ing in Robotics and Control”. In: *IEEE PAMI* 37.2
(2015), pp. 408–423.
- 11 [DFS16] A. Daniely, R. Frostig, and Y. Singer. “Toward
Deeper Understanding of Neural Networks: The Power
of Initialization and a Dual View on Expressivity”. In:
NIPS. 2016, pp. 2253–2261.
- 12 [DG17] P. Dabkowski and Y. Gal. “Real time image
saliency for black box classifiers”. In: *NeurIPS* (2017).
- 13 [DG84] P. J. Diggle and R. J. Gratton. “Monte Carlo
methods of inference for implicit statistical models”.
In: *Journal of the Royal Statistical Society. Series B
(Methodological)* (1984), pp. 193–227.
- 14 [DGA00] A. Doucet, S. Godsill, and C. Andrieu.
“On sequential Monte Carlo Sampling Methods for
Bayesian Filtering”. In: *Statistics and Computing*
10.3 (2000), pp. 197–208.
- 15 [DGE15] C. Doersch, A. Gupta, and A. A. Efros. “Un-
supervised visual representation learning by context
prediction”. In: *Proceedings of the IEEE interna-
tional conference on computer vision*. 2015, pp. 1422–
1430.
- 16 [DGF20] Y. Dubois, J. Gordon, and A. Y. Foong. *Neu-
ral Process Family*. [http://yanndubs.github.io/
Neural-Process-Family/](http://yanndubs.github.io/
Neural-Process-Family/). 2020.
- 17 [DGK01] A. Doucet, N. Gordon, and V. Krishna-
murthy. “Particle Filters for State Estimation of Jump
Markov Linear Systems”. In: *IEEE Trans. on Signal
Processing* 49.3 (2001), pp. 613–624.
- 18 [DH22] F. Dellaert and S. Hutchinson. *Introducion to
Robotics and Perception*. 2022.
- 19 [DHK14] A. Deshpande, L. Hellerstein, and D.
Kletenik. “Approximation algorithms for stochastic
boolean function evaluation and stochastic submodu-
lar set cover”. In: *Proceedings of the twenty-fifth
annual ACM-SIAM Symposium on Discrete Algo-
rithms*. SIAM. 2014, pp. 1453–1466.
- 20 [Dia88] P. Diaconis. “Sufficiency as statistical symme-
try”. In: *Proceedings of the AMS Centennial Sympo-
sium*. 1988, pp. 15–26.
- 21 [Die00] T. G. Dietterich. “Ensemble Methods in Ma-
chine Learning”. In: *Multiple Classifier Systems*.
Springer Berlin Heidelberg, 2000, pp. 1–15.
- 22 [Die+07] M. Diehl, H. G. Bock, H. Diedam, and P.-B.
Wieber. “Fast Direct Multiple Shooting Algorithms for
Optimal Robot Control”. In: *Lecture Notes in Control
and Inform. Sci.* 340 (2007).
- 23 [Die10] L. Dietz. *Directed Factor Graph Notation for
Generative Models*. Tech. rep. MPI, 2010.
- 24 [Die+17] A. B. Dieng, C. Wang, J. Gao, and J. Paisley.
“TopicRNN: A Recurrent Neural Network with Long-
Range Semantic Dependency”. In: *ICLR*. 2017.
- 25 [Die+19a] A. B. Dieng, Y. Kim, A. M. Rush, and D. M.
Blei. “Avoiding Latent Variable Collapse With Genera-
tive Skip Models”. In: *AISTATS*. 2019.
- 26 [Die+19b] A. B. Dieng, F. J. Ruiz, D. M. Blei, and
M. K. Titsias. “Prescribed generative adversarial net-
works”. In: *arXiv preprint arXiv:1910.04302* (2019).
- 27 [Die22] S. Dieleman. *Guidance: a cheat code for diffu-
sion models*. 2022.
- 28 [Dik+20] N. Dikkala, G. Lewis, L. Mackey, and V.
Syrgkanis. “Minimax Estimation of Conditional Mo-
ment Models”. In: *Advances in Neural Information
Processing Systems*. 2020.
- 29 [Din+17] L. Dinh, R. Pascanu, S. Bengio, and Y. Ben-
gio. “Sharp Minima Can Generalize For Deep Nets”.
In: (2017). arXiv: 1703.04933 [cs.LG].
- 30 [DJ11] A. Doucet and A. M. Johansen. “A Tutorial on
Particle Filtering and Smoothing: Fifteen years later”.
In: *Handbook of Nonlinear Filtering*. Ed. by D Crisan
and B Rozovsk. 2011.
- 31 [DJ21] K. Desai and J. Johnson. *VirTex: Learning
Visual Representations from Textual Annotations*.
2021. arXiv: 2006.06666 [cs.CV].
- 32 [DJK18] J. Djolonga, S. Jegelka, and A. Krause.
“Provable Variational Inference for Constrained Log-
Submodular Models”. In: *NeurIPS*. 2018.
- 33 [DK12] J. Durbin and S. J. Koopman. *Time Series
Analysis by State Space Methods: Second Edition*. en.
Revised ed. edition. Oxford University Press, 2012.
- 34 [DK14] J. Djolonga and A. Krause. “From map to
marginals: Variational inference in bayesian submodu-
lar models”. In: *Advances in Neural Information Pro-
cessing Systems*. 2014, pp. 244–252.
- 35 [DK15a] J. Djolonga and A. Krause. “Scalable varia-
tional inference in log-supermodular models”. In: *In-
ternational Conference on Machine Learning*. PMLR.
2015, pp. 1804–1813.
- 36 [DK15b] G. Durrett and D. Klein. “Neural CRF Pars-
ing”. In: *Proc. ACL*. 2015.
- 37 [DKB15] L. Dinh, D. Krueger, and Y. Bengio. “NICE:
Non-linear Independent Components Estimation”. In:
ICLR. 2015.
- 38 [DKD16] J. Donahue, P. Krähenbühl, and T. Dar-
rell. “Adversarial feature learning”. In: *arXiv preprint
arXiv:1605.09782* (2016).
- 39 [DKD17] J. Donahue, P. Krähenbühl, and T. Dar-
rell. *Adversarial Feature Learning*. 2017. arXiv: 1605.
09782 [cs.LG].
- 40 [DKS13] J. Dick, F. Y. Kuo, and I. H. Sloan. “High-
dimensional integration: the quasi-Monte Carlo way”.
In: *Acta Numerica* 22 (2013), 133–288.
- 41 [DL09] P. Domingos and D. Lowd. *Markov Logic: An
Interface Layer for AI*. Morgan & Claypool, 2009.
- 42
- 43
- 44
- 45
- 46
- 47

- 1
- [DL10] J. V. Dillon and G. Lebanon. “Stochastic Composite Likelihood”. In: *J. Mach. Learn. Res.* 11 (2010), pp. 2597–2633.
- 2
- [DL13] A. Damianou and N. Lawrence. “Deep Gaussian Processes”. In: *AISTATS*. 2013, pp. 207–215.
- 3
- [DL93] P. Dagum and M. Luby. “Approximating probabilistic inference in Bayesian belief networks is NP-hard”. In: *Artificial Intelligence* 60 (1993), pp. 141–153.
- 4
- [DLB17] C. Dann, T. Lattimore, and E. Brunskill. “Unifying PAC and Regret: Uniform PAC Bounds for Episodic Reinforcement Learning”. In: *NIPS*. 2017, pp. 5717–5727.
- 5
- [DLM09] H. Daumé III, J. Langford, and D. Marcu. “Search-based Structured Prediction”. In: *MLJ* 75.3 (2009), pp. 297–325.
- 6
- [DLM99] B. Delyon, M. Lavielle, and E. Moulines. “Convergence of a stochastic approximation version of the EM algorithm”. In: *Annals of Statistics* 27.1 (1999), pp. 94–128.
- 7
- [DLR77] A. P. Dempster, N. M. Laird, and D. B. Rubin. “Maximum likelihood from incomplete data via the EM algorithm”. In: *J. of the Royal Statistical Society, Series B* 34 (1977), pp. 1–38.
- 8
- [DLT21] M. De Lange and T. Tuytelaars. “Continual Prototype Evolution: Learning Online from Non-Stationary Data Streams”. In: *ICCV*. 2021.
- 9
- [DM01] D. van Dyk and X.-L. Meng. “The Art of Data Augmentation”. In: *J. Computational and Graphical Statistics* 10.1 (2001), pp. 1–50.
- 10
- [DM19a] Y. Du and I. Mordatch. “Implicit Generation and Generalization in Energy-Based Models”. In: (2019). arXiv: [1903.08689 \[cs.LG\]](https://arxiv.org/abs/1903.08689).
- 11
- [DM19b] Y. Du and I. Mordatch. “Implicit Generation and Modeling with Energy Based Models”. In: *NIPS*. 2019, pp. 3608–3618.
- 12
- [DM22] A. P. Dawid and M. Musio. “Effects of Causes and Causes of Effects”. In: *Annu. Rev. Stat. Appl.* 9.1 (Mar. 2022), pp. 261–287.
- 13
- [DMDJ12] P. Del Moral, A. Doucet, and A. Jasra. “An adaptive sequential Monte Carlo method for approximate Bayesian computation”. In: *Stat. Comput.* 22.5 (2012), pp. 1009–1020.
- 14
- [DMKM22] G. Duran-Martin, A. Kara, and K. Murphy. “Efficient Online Bayesian Inference for Neural Bandits”. In: *AISTATS*. 2022.
- 15
- [DMM17] A. P. Dawid, M. Musio, and R. Murtas. “The probability of causation”. In: *Law, Probability and Risk* 16.4 (2017), pp. 163–179.
- 16
- [DMP18] C. Donahue, J. McAuley, and M. Puckette. “Adversarial Audio Synthesis”. In: *International Conference on Learning Representations*. 2018.
- 17
- [DMV15] S. Dash, D. M. Malioutov, and K. R. Varshney. “Learning interpretable classification rules using sequential rowsampling”. In: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE. 2015, pp. 3337–3341.
- 18
- [DN21] P. Dhariwal and A. Q. Nichol. “Diffusion Models Beat GANs on Image Synthesis”. In: *NIPS*. May 2021.
- 19
- [Dnp] .
- 20
- [DNR11] D. Duvenaud, H. Nickisch, and C. E. Rasmussen. “Additive Gaussian Processes”. In: *NIPS*. 2011.
- 21
- [Dom+06] P. Domingos, S. Kok, H. Poon, M. Richardson, and P. Singla. “Unifying Logical and Statistical AI”. In: *IJCAI*. 2006.
- 22
- [Dom+19] A.-K. Dombrowski, M. Alber, C. J. Anders, M. Ackermann, K.-R. Müller, and P. Kessel. “Explanations can be manipulated and geometry is to blame”. In: *arXiv preprint arXiv:1906.07983* (2019).
- 23
- [Don+14] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. “Decaf: A deep convolutional activation feature for generic visual recognition”. In: *International conference on machine learning*. 2014, pp. 647–655.
- 24
- [Don+17a] K. Dong, D. Eriksson, H. Nickisch, D. Bindel, and A. G. Wilson. “Scalable Log Determinants for Gaussian Process Kernel Learning”. In: *NIPS*. 2017, pp. 6327–6337.
- 25
- [Don+17b] Y. Dong, H. Su, J. Zhu, and F. Bao. “Towards interpretable deep neural networks by leveraging adversarial examples”. In: *arXiv preprint arXiv:1708.05493* (2017).
- 26
- [Don+21] X. Dong, J. Bao, T. Zhang, D. Chen, W. Zhang, L. Yuan, D. Chen, F. Wen, and N. Yu. “Peco: Perceptual codebook for bert pre-training of vision transformers”. In: *arXiv preprint arXiv:2111.12710* (2021).
- 27
- [Dor+16] V. Dorie, M. Harada, N. B. Carnegie, and J. Hill. “A flexible, interpretable framework for assessing sensitivity to unmeasured confounding”. In: *Statistics in Medicine* 35.20 (2016), pp. 3453–3470. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/sim.6973>.
- 28
- [Dos+21] A. Dosovitskiy et al. “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale”. In: *International Conference on Learning Representations*. 2021.
- 29
- [Doy+07] K. Doya, S. Ishii, A. Pouget, and R. P. N. Rao, eds. *Bayesian Brain: Probabilistic Approaches to Neural Coding*. MIT Press, 2007.
- 30
- [DR11] M. P. Deisenroth and C. E. Rasmussen. “PILCO: A Model-Based and Data-Efficient Approach to Policy Search”. In: *ICML*. 2011.
- 31
- [DR17] G. K. Dziugaite and D. M. Roy. “Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data”. In: *UAI*. 2017.
- 32
- [DRB19] A. B. Dieng, F. J. R. Ruiz, and D. M. Blei. “The Dynamic Embedded Topic Model”. In: (2019). arXiv: [1907.05545 \[cs.CL\]](https://arxiv.org/abs/1907.05545).
- 33
- [DRG15] G. K. Dziugaite, D. M. Roy, and Z. Ghahramani. “Training generative neural networks via Maximum Mean Discrepancy optimization”. In: *ICML*. 2015.
- 34
- [Dru08] J. Drugowitsch. *Bayesian linear regression*. Tech. rep. U. Rochester, 2008.
- 35
- [DS18] J. Domke and D. R. Sheldon. “Importance Weighting and Variational Inference”. In: *NIPS*. Curran Associates, Inc., 2018, pp. 4474–4483.
- 36
- [DS19] J. Donahue and K. Simonyan. “Large scale adversarial representation learning”. In: *arXiv preprint arXiv:1907.02544* (2019).
- 37
- [DSDB17] L. Dinh, J. Sohl-Dickstein, and S. Bengio. “Density estimation using Real NVP”. In: *ICLR*. 2017.

- 1
- 2 [DSK16] V. Dumoulin, J. Shlens, and M. Kudlur. “A
3 Learned Representation For Artistic Style”. In: (2016).
- 4 [DSZ16] A. Datta, S. Sen, and Y. Zick. “Algorithmic
5 transparency via quantitative input influence:
6 Theory and experiments with learning systems”. In:
7 *2016 IEEE symposium on security and privacy (SP)*.
8 IEEE. 2016, pp. 598–617.
- 9 [DTK16] J. Djolonga, S. Tschiatschek, and A. Krause.
10 “Variational inference in mixed probabilistic submodular
11 models”. In: *Advances in Neural Information Processing Systems*. 2016, pp. 1759–1767.
- 12 [Du+16] N. Du, H. Dai, R. Trivedi, U. Upadhyay, M.
13 Gomez-Rodriguez, and L. Song. “Recurrent marked
14 temporal point processes: Embedding event history to
15 vector”. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016, pp. 1555–1564.
- 16 [Du+18] J. Du, S. Ma, Y.-C. Wu, S. Kar, and J. M. F.
17 Moura. “Convergence Analysis of Distributed Inference
18 with Vector-Valued Gaussian Belief Propagation”. In:
19 *JMLR* 18.172 (2018), pp. 1–38.
- 20 [Du+19] S. S. Du, K. Hou, B. Póczos, R. Salakhutdinov,
21 R. Wang, and K. Xu. “Graph Neural Tangent Kernel: Fusing Graph Neural Networks with Graph Kernels”. In: (2019). arXiv: [1905.13192 \[cs.LG\]](https://arxiv.org/abs/1905.13192).
- 22 [Du+20] Y. Du, S. Li, J. Tenenbaum, and I. Mordatch.
23 “Improved Contrastive Divergence Training of Energy Based Models”. In: *arXiv preprint arXiv:2012.01316* (2020).
- 24 [Du+21] C. Du, Z. Gao, S. Yuan, L. Gao, Z. Li, Y. Zeng,
25 X. Zhu, J. Xu, K. Gai, and K.-C. Lee. “Exploration in Online Advertising Systems with Deep Uncertainty-Aware Learning”. In: *KDD*. KDD ’21. Association for Computing Machinery, 2021, pp. 2792–2801.
- 26 [Dua+20] S. Duan, N. Watters, L. Matthey, C. P.
27 Burgess, A. Lerchner, and I. Higgins. “A Heuristic for Unsupervised Model Selection for Variational Disentangled Representation Learning”. In: *ArXiv abs/1905.12614* (2020).
- 28 [Dua+87] S. Duane, A. Kennedy, B. Pendleton, and D.
29 Roweth. “Hybrid Monte Carlo”. In: *Physics Letters B* 195.2 (1987), pp. 216–222.
- 30 [Dub+16] A. Dubey, S. J. Reddi, B. Póczos, A. J.
31 Smola, E. P. Xing, and S. A. Williamson. “Variance Reduction in Stochastic Gradient Langevin Dynamics”. In: *NIPS*. Vol. 29. 2016, pp. 1154–1162.
- 32 [Dud13] J. Duda. “Asymmetric numeral systems: entropy coding combining speed of Huffman coding with compression rate of arithmetic coding”. In: (2013). arXiv: [1311.2540 \[cs.IT\]](https://arxiv.org/abs/1311.2540).
- 33 [Duf02] M. Duff. “Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes”. PhD thesis. U. Mass. Dept. Comp. Sci., 2002.
- 34 [Dum+16] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville.
35 “Adversarially learned inference”. In: *arXiv preprint arXiv:1606.00704* (2016).
- 36 [Dum+17] V. Dumoulin, I. Belghazi, B. Poole, O. Mastropietro, A. Lamb, M. Arjovsky, and A. Courville.
37 “Adversarially Learned Inference”. 2017. arXiv: [1606.00704 \[stat.ML\]](https://arxiv.org/abs/1606.00704).
- 38 [Dur+19] C. Durkan, A. Bekasov, I. Murray, and G.
39 Papamakarios. “Cubic-Spline Flows”. In: *ICML Workshop on Invertible Neural Networks and Normalizing Flows*. 2019.
- 40 [Dur+98] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, 1998.
- 41 [Duv+13] D. Duvenaud, J. Lloyd, R. Grosse, J. Tenenbaum, and G. Zoubin. “Structure Discovery in Nonparametric Regression through Compositional Kernel Search”. In: *ICML*. 2013, pp. 1166–1174.
- 42 [Duv14] D. Duvenaud. “Automatic Model Construction with Gaussian Processes”. PhD thesis. Computational and Biological Learning Laboratory, University of Cambridge, 2014.
- 43 [dV04] D. de Farias and B. Van Roy. “On Constraint Sampling in the Linear Programming Approach to Approximate Dynamic Programming”. In: *Mathematics of Operations Research* 29.3 (2004), pp. 462–478.
- 44 [DV+17] H. De Vries, F. Strub, J. Mary, H. Larochelle, O. Pietquin, and A. C. Courville. “Modulating early visual processing by language”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 6594–6604.
- 45 [DV75] M. D. Donsker and S. S. Varadhan. “Asymptotic evaluation of certain Markov process expectations for large time, I”. In: *Communications on Pure and Applied Mathematics* 28.1 (1975), pp. 1–47.
- 46 [DV99] A. P. Dawid and V. Vovk. “Prequential probability: Principles and properties”. In: *Bernoulli* 5 (1999), pp. 125–162.
- 47 [DKV17] F. Doshi-Velez and B. Kim. *Towards A Rigorous Science of Interpretable Machine Learning*. 2017. eprint: [1702.08608](https://arxiv.org/abs/1702.08608) (stat.ML).
- 48 [DW19] B. Dai and D. Wipf. “Diagnosing and Enhancing VAE Models”. In: *ICLR*. 2019.
- 49 [DWS12] T. Degris, M. White, and R. S. Sutton. “Off-Policy Actor-Critic”. In: *ICML*. 2012.
- 50 [DWW19] B. Dai, Z. Wang, and D. Wipf. “The Usual Suspects? Reassessing Blame for VAE Posterior Collapse”. In: (2019). arXiv: [1912.10702 \[cs.LG\]](https://arxiv.org/abs/1912.10702).
- 51 [DWW20] B. Dai, Z. Wang, and D. Wipf. “The Usual Suspects? Reassessing Blame for VAE Posterior Collapse”. In: *ICML*. Ed. by H. D. Iii and A. Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 2313–2322.
- 52 [DY79] P. Diaconis and D. Ylvisaker. “Conjugate priors for exponential families”. In: vol. 7. 1979, pp. 269–281.
- 53 [ECM18] P. Etoori, M. Chinnakotla, and R. Mamidi. “Automatic Spelling Correction for Resource-Scarce Languages using Deep Learning”. In: *Proceedings of ACL 2018, Student Research Workshop*. Association for Computational Linguistics, 2018, pp. 146–152.
- 54 [ED05] D. Earl and M. Deem. “Parallel tempering: Theory, applications, and new perspectives”. In: *Phys. Chem. Chem. Phys.* 7 (2005), p. 3910.
- 55 [Edm69] H. P. Edmundson. “New methods in automatic extracting”. In: *Journal of the ACM (JACM)* 16.2 (1969), pp. 264–285.
- 56 [Edm70] J. Edmonds. “Matroids, submodular functions, and certain polyhedra”. In: *Combinatorial Structures and Their Applications* (1970), pp. 69–87.
- 57 [EFL04] E. Erosheva, S. Fienberg, and J. Lafferty. “Mixed-membership models of scientific publications”. In: *PNAS* 101 (2004), pp. 5220–2227.

- [Efr86] B. Efron. "Why Isn't Everyone a Bayesian?" In: *The American Statistician* 40.1 (1986).
- [EGW05] D. Ernst, P. Geurts, and L. Wehenkel. "Tree-Based Batch Mode Reinforcement Learning". In: *JMLR* 6 (2005), pp. 503–556.
- [Eis16] J. Eisner. "Inside-Outside and Forward-Backward Algorithms Are Just Backprop (Tutorial Paper)". In: *EMNLP Workshop on Structured Prediction for NLP*. 2016.
- [Eke+13] M. Ekeberg, C. Lökvist, Y. Lan, M. Weigt, and E. Aurell. "Improved contact prediction in proteins: using pseudolikelihoods to infer Potts models". en. In: *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.* 87.1 (2013), p. 012707.
- [EM07] D. Eaton and K. Murphy. "Exact Bayesian structure learning from uncertain interventions". In: *AI/Statistics*. 2007.
- [EMH19] T. Elsken, J. H. Metzen, and F. Hutter. "Neural Architecture Search: A Survey". In: *JMLR* 20 (2019), pp. 1–21.
- [EMK06] G. Elidan, I. McGraw, and D. Koller. "Residual belief propagation: Informed scheduling for asynchronous message passing". In: *UAI*. 2006.
- [Eng+18] J. Engel, K. K. Agrawal, S. Chen, I. Gulrajani, C. Donahue, and A. Roberts. "GANSynth: Adversarial Neural Audio Synthesis". In: *International Conference on Learning Representations*. 2018.
- [Erh+09] D. Erhan, Y. Bengio, A. Courville, and P. Vincent. "Visualizing higher-layer features of a deep network". In: *University of Montreal* 1341.3 (2009), p. 1.
- [Erm+13] S. Ermon, C. P. Gomes, A. Sabharwal, and B. Selman. "Taming the Curse of Dimensionality: Discrete Integration by Hashing and Optimization". In: *ICML*. Feb. 2013.
- [ERO21] P. Esser, R. Rombach, and B. Ommer. "Taming Transformers for High-Resolution Image Synthesis". In: *CVPR*. 2021.
- [Eva+18] R. Evans et al. "De novo structure prediction with deep-learning based scoring". In: (2018).
- [Eve09] G. Evensen. *Data Assimilation: The Ensemble Kalman Filter*. en. 2nd ed. 2009 edition. Springer, 2009.
- [EW10] El-Yaniv and Wiener. "On the Foundations of Noise-free Selective Classification". In: *JMLR* (2010).
- [EW95] M. D. Escobar and M. West. "Bayesian density estimation and inference using mixtures". In: *JASA* 90.430 (1995), pp. 577–588.
- [Ewe72] W. J. Ewens. "The sampling theory of selectively neutral alleles". In: *Theoretical population biology* 3.1 (1972), pp. 87–112.
- [Ewe90] W. Ewens. "Population genetics theory - the past and the future". In: *Mathematical and Statistical Developments of Evolutionary Theory*. Ed. by S. Lessard. Reidel, 1990, pp. 177–227.
- [EY09] M. Elad and I. Yavneh. "A plurality of sparse representations is better than the sparsest one alone". In: *IEEE Trans. on Info. Theory* 55.10 (2009), pp. 4701–4714.
- [Eyk+18] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song. "Robust Physical-World Attacks on Deep Learning Models". In: *CVPR*. 2018.
- [Eys+21] B. Eysenbach, A. Khazatsky, S. Levine, and R. Salakhutdinov. "Mismatched No More: Joint Model-Policy Optimization for Model-Based RL". In: (2021). arXiv: [2110.02758 \[cs.LG\]](https://arxiv.org/abs/2110.02758).
- [Fas20] I. Fischer and A. A. Alemi. "CEB Improves Model Robustness". In: *Entropy* 22.10 (2020), p. 1081.
- [Fad+20] S. G. Fadel, S. Mair, R. da S. Torres, and U. Brefeld. "Principled Interpolation in Normalizing Flows". In: (2020). arXiv: [2010.12059 \[stat.ML\]](https://arxiv.org/abs/2010.12059).
- [Fag+18] F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler. "VSE++: Improving Visual-Semantic Embeddings with Hard Negatives". In: *BMVC*. 2018.
- [FAL17] C. Finn, P. Abbeel, and S. Levine. "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks". In: *ICML*. 2017.
- [Fan+17] H. Fang, N. Tian, Y. Wang, M. Zhou, and M. A. Haile. "Nonlinear Bayesian Estimation: From Kalman Filtering to a Broader Horizon". In: (2017). arXiv: [1712.01406 \[cs.SY\]](https://arxiv.org/abs/1712.01406).
- [Fau+18] L. Faury, F. Vasile, C. Calauzènes, and O. Fercoq. "Neural Generative Models for Global Optimization with Gradients". In: (2018). arXiv: [1805.08594 \[cs.NE\]](https://arxiv.org/abs/1805.08594).
- [FB19] E. M. Feit and R. Berman. "Test & Roll: Profit-Maximizing A/B Tests". In: *Marketing Science* 38.6 (2019), pp. 1038–1058.
- [FBW21] M. Finzi, G. Benton, and A. G. Wilson. "Residual Pathway Priors for Soft Equivariance Constraints". In: *NIPS*. 2021.
- [FC03] P. Fearnhead and P. Clifford. "On-line inference for hidden Markov models via particle filters". en. In: *J. of Royal Stat. Soc. Series B* 65.4 (2003), pp. 887–899.
- [FD07a] B. Frey and D. Dueck. "Clustering by Passing Messages Between Data Points". In: *Science* 315 (2007), 972–976.
- [FD07b] B. J. Frey and D. Dueck. "Clustering by passing messages between data points". In: *science* 315.5814 (2007), pp. 972–976.
- [FDF19] A. M. Franks, A. D'Amour, and A. Feller. "Flexible Sensitivity Analysis for Observational Studies Without Observable Implications". In: *Journal of the American Statistical Association* 0.0 (2019), pp. 1–33. eprint: <https://doi.org/10.1080/01621459.2019.1604369>.
- [FDZ19] A. Fasano, D. Durante, and G. Zanella. "Scalable and Accurate Variational Bayes for High-Dimensional Binary Regression Models". In: (2019). arXiv: [1911.06743 \[stat.ME\]](https://arxiv.org/abs/1911.06743).
- [FE73] M. Fischler and R. Elschlager. "The representation and matching of pictorial structures". In: *IEEE Trans. on Computer* 22.1 (1973).
- [Fed+18] W. Fedus, M. Rosca, B. Lakshminarayanan, A. M. Dai, S. Mohamed, and I. Goodfellow. "Many Paths to Equilibrium: GANs Do Not Need to Decrease a Divergence At Every Step". In: *International Conference on Learning Representations*. 2018.
- [Fei98] U. Feige. "A threshold of $\ln n$ for approximating set cover". In: *Journal of the ACM (JACM)* (1998).
- [Fel+10] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. "Object Detection with Discriminatively Trained Part Based Models". In: *IEEE PAMI* 32.9 (2010).

- 1
- 2 [Fel+19] M. Fellows, A. Mahajan, T. G. J. Rudner, and
S. Whiteson. “VIREL: A Variational Inference Frame-
3 work for Reinforcement Learning”. In: *NeurIPS*. 2019,
pp. 7120–7134.
- 4
- 5 [Fen+21] S. Y. Feng, V. Gangal, J. Wei, S. Chandar,
S. Vosoughi, T. Mitamura, and E. Hovy. “A Survey of
6 Data Augmentation Approaches for NLP”. In: (2021).
arXiv: [2105.03075 \[cs.CL\]](https://arxiv.org/abs/2105.03075).
- 7
- 8 [Fer73] T. S. Ferguson. “A Bayesian analysis of some
nonparametric problems”. In: *The Annals of Statistics* (1973), pp. 209–230.
- 9
- 10 [Feu+15] M. Feurer, A. Klein, K. Eggensperger, J.
Springenberg, M. Blum, and F. Hutter. “Efficient
and Robust Automated Machine Learning”. In: *NIPS*.
2015, pp. 2962–2970.
- 11
- 12 [FG15] N. Fournier and A. Guillin. “On the rate of con-
vergence in Wasserstein distance of the empirical mea-
sure”. In: *Probability Theory and Related Fields* 162.3
(2015), pp. 707–738.
- 13
- 14 [FG18] S. Farquhar and Y. Gal. “Towards Robust Evalu-
ations of Continual Learning”. In: (2018). arXiv: [1805.09733 \[stat.ML\]](https://arxiv.org/abs/1805.09733).
- 15
- 16 [FGG97] N. Friedman, D. Geiger, and M. Goldszmidt.
“Bayesian network classifiers”. In: *MLJ* 29 (1997),
pp. 131–163.
- 17
- 18 [FH12] P. F. Felzenszwalb and D. P. Huttenlocher.
“Distance Transforms of Sampled Functions”. In: *The-
ory of Computing* 8.19 (2012), pp. 415–428.
- 19
- 20 [FH17] N. Frosst and G. Hinton. *Distilling a Neural
Network Into a Soft Decision Tree*. 2017. arXiv: [1711.09784 \[cs.LG\]](https://arxiv.org/abs/1711.09784).
- 21
- 22 [FH20] E. Fong and C. Holmes. “On the marginal
likelihood and cross-validation”. In: *Biometrika* 107.2
(2020).
- 23
- 24 [FH75] K. Fukunaga and L. Hostetler. “The estimation
of the gradient of a density function, with applications
in pattern recognition”. In: *IEEE Trans. Inf. Theory*
21.1 (1975), pp. 32–40.
- 25
- 26 [FH97] B. J. Frey and G. Hinton. “Efficient stochastic
source coding and an application to a Bayesian net-
work source model”. In: *Computer Journal* (1997).
- 27
- 28 [FHDV20] J. Futoma, M. C. Hughes, and F. Doshi-
Velez. “Popcorn: Partially observed prediction con-
strained reinforcement learning”. In: *AISTATS*
(2020).
- 29
- 30 [FKH03] P. Felzenszwalb, D. Huttenlocher, and J.
Kleinberg. “Fast Algorithms for Large State Space
HMMs with Applications to Web Usage Analysis”. In:
NIPS. 2003.
- 31
- 32 [FHL19] S. Fort, H. Hu, and B. Lakshminarayanan.
“Deep Ensembles: A Loss Landscape Perspective”. In:
(2019). arXiv: [1912.02757 \[stat.ML\]](https://arxiv.org/abs/1912.02757).
- 33
- 34 [FHM18] S. Fujimoto, H. van Hoof, and D. Meger.
“Addressing Function Approximation Error in Actor-
Critic Methods”. In: *ICLR*. 2018.
- 35
- 36 [FHT08] J. Friedman, T. Hastie, and R. Tibshirani.
“Sparse inverse covariance estimation the graphical
lasso”. In: *Biostatistics* 9.3 (2008), pp. 432–441.
- 37
- 38 [FI10] A. Fischer and C. Igel. “Empirical analysis of
the divergence of Gibbs sampling based learning al-
gorithms for restricted Boltzmann machines”. In: *Inter-
national conference on artificial neural networks*.
Springer. 2010, pp. 208–217.
- 39
- 40 [Fie70] S. Fienberg. “An Iterative Procedure for Esti-
mation in Contingency Tables”. In: *Annals of Mathe-
matical Statistics* 41.3 (1970), pp. 907–917.
- 41
- 42 [Fin+16] C. Finn, P. Christiano, P. Abbeel, and S.
Levine. “A connection between generative adversarial
networks, inverse reinforcement learning, and energy-
based models”. In: *arXiv preprint arXiv:1611.03832*
(2016).
- 43
- 44 [Fis20] I. Fischer. “The Conditional Entropy Bot-
tle-neck”. In: *Entropy* 22.9 (2020).
- 45
- 46 [Fis25] R. Fisher. *Statistical Methods for Research
Workers*. Biological monographs and manuals. Oliver
and Boyd, 1925.
- 47
- 48 [FJ02] M. A. T. Figueiredo and A. K. Jain. “Unsuper-
vised Learning of Finite Mixture Models”. In: *IEEE
PAMI* 24.3 (2002), pp. 381–396.
- 49
- 50 [FJL18] R. Frostig, M. J. Johnson, and C. Leary. “Com-
piling machine learning programs via high-level trac-
ing”. In: *Machine Learning and Systems (MLSys)*
(2018).
- 51
- 52 [FK13a] V. Feldman and P. Kothari. “Learning Cov-
erage Functions”. In: *CoRR* abs/1304.2079 (2013).
arXiv: [1304.2079](https://arxiv.org/abs/1304.2079).
- 53
- 54 [FK13b] M. Frei and H. R. Künsch. “Bridging the en-
semble Kalman and particle filters”. In: *Biometrika*
100.4 (2013), pp. 781–800.
- 55
- 56 [FK14] V. Feldman and P. Kothari. “Learning Cov-
erage Functions and Private Release of Marginals”.
In: *Proceedings of The 27th Conference on Learn-
ing Theory*. Ed. by M. F. Balcan, V. Feldman, and
C. Szepesvári. Vol. 35. Proceedings of Machine Learn-
ing Research. Barcelona, Spain: PMLR, 2014, pp. 679–
702.
- 57
- 58 [FK21] A. Fisher and E. H. Kennedy. “Visually Com-
municating and Teaching Intuition for Influence Func-
tions”. In: *The American Statistician* 75.2 (2021),
pp. 162–172. eprint: <https://doi.org/10.1080/00031305.2020.1717620>.
- 59
- 60 [FKH17] S. Falkner, A. Klein, and F. Hutter. “Combin-
ing Hyperband and Bayesian Optimization”. In: *NIPS
2017 Bayesian Optimization Workshop*. 2017.
- 61
- 62 [FKV13] V. Feldman, P. Kothari, and J. Vondrák.
“Representation, Approximation and Learning of Sub-
modular Functions Using Low-rank Decision Trees”.
In: *Proceedings of the 26th Annual Conference on
Learning Theory*. Ed. by S. Shalev-Shwartz and I.
Steinwart. Vol. 30. Proceedings of Machine Learning
Research. Princeton, NJ, USA: PMLR, 2013, pp. 711–
740.
- 63
- 64 [FKV14] V. Feldman, P. Kothari, and J. Vondrák.
“Nearly tight bounds on ℓ_1 approximation of self-
bounding functions”. In: *CoRR*, abs/1404.4702 1
(2014).
- 65
- 66 [FKV17] V. Feldman, P. Kothari, and J. Vondrák.
“Tight Bounds on ℓ_1 Approximation and Learning of
Self-Bounding Functions”. In: *International Confer-
ence on Algorithmic Learning Theory*. PMLR. 2017,
pp. 540–559.
- 67
- 68 [FKV20] V. Feldman, P. Kothari, and J. Vondrák.
“Tight bounds on ℓ_1 approximation and learning of
self-bounding functions”. In: *Theoretical Computer
Science* 808 (2020), pp. 86–98.
- 69
- 70 [FL07] P. Fearnhead and Z. Liu. “Online Inference for
Multiple Changepoint Problems”. In: *J. of Royal Stat.
Soc. Series B* 69 (2007), pp. 589–605.

- 1
- 2 [FL11] P. Fearnhead and Z. Liu. “Efficient Bayesian
3 analysis of multiple changepoint models with depen-
4 dence across segments”. In: *Statistics and Computing*
5 21.2 (2011), pp. 217–229.
- 6 [FL+18] V. François-Lavet, P. Henderson, R. Islam,
7 M. G. Bellemare, and J. Pineau. “An Introduction to
8 Deep Reinforcement Learning”. In: *Foundations and
9 Trends in Machine Learning* 11.3 (2018).
- 10 [FLA16] C. Finn, S. Levine, and P. Abbeel. “Guided
11 Cost Learning: Deep Inverse Optimal Control via Policy
12 Optimization”. In: *ICML*. 2016, pp. 49–58.
- 13 [Fla+16] S. Flaxman, D. Sejdinovic, J. P. Cunningham,
14 and S. Filippi. “Bayesian Learning of Kernel Embed-
15 dings”. In: *UAI*. 2016.
- 16 [FLL18] J. Fu, K. Luo, and S. Levine. “Learning Rob-
17 bust Rewards with Adversarial Inverse Reinforce-
18 ment Learning”. In: *ICLR*. 2018.
- 19 [FLL19] Y. Feng, L. Li, and Q. Liu. “A Kernel Loss
20 for Solving the Bellman Equation”. In: *NeurIPS*. 2019,
pp. 15430–15441.
- 21 [FLMM21] D. T. Frazier, R. Loaiza-Maya, and G. M.
22 Martin. “Variational Bayes in State Space Models: Inferential and Predictive Accuracy”. In: (June 2021).
arXiv: [2106.12262 \[stat.ME\]](#).
- 23 [FMM18] M. Figurnov, S. Mohamed, and A. Mnih. “Im-
24 plicit Reparameterization Gradients”. In: *NIPS*. 2018.
- 25 [FMP19] S. Fujimoto, D. Meger, and D. Precup. “Off-
26 Policy Deep Reinforcement Learning without Exploration”. In: *ICML*. 2019, pp. 2052–2062.
- 27 [FNG00] N. de Freitas, M. Niranjan, and A. Gee. “Hier-
28 archical Bayesian models for regularisation in sequen-
29 tial learning”. In: *Neural Computation* 12.4 (2000),
pp. 955–993.
- 30 [FNW78] M. Fisher, G. Nemhauser, and L. Wolsey.
“An analysis of approximations for maximizing sub-
31 modular set functions—II”. In: *Polyhedral combinatorics* (1978), pp. 73–87.
- 32 [FO20] F. Farnia and A. Ozdaglar. “Do GANs always
33 have Nash equilibria?” In: *Proceedings of the 37th International Conference on Machine Learning*. Ed. by
H. D. III and A. Singh. Vol. 119. Proceedings of Ma-
34 chine Learning Research. PMLR, 2020, pp. 3029–3039.
- 35 [Fon+21] D. Fontanel, F. Cermelli, M. Mancini, and B.
Caputo. “On the Challenges of Open World Recognition
36 under Shifting Visual Domains”. In: *ICRA*. 2021.
- 37 [For01] G. D. Forney. “Codes on graphs: normal real-
38 izations”. In: *IEEE Trans. Inf. Theory* 47.2 (2001),
pp. 520–548.
- 39 [For+18a] V. Fortuin, G. Dresdner, H. Strathmann,
40 and G. Rätsch. “Scalable Gaussian Processes on
41 Discrete Domains”. In: (2018). arXiv: [1810 . 10368 \[stat.ML\]](#).
- 42 [For+18b] M. Fortunato et al. “Noisy Networks for Ex-
ploration”. In: *ICLR*. 2018.
- 43 [For+19] N. Ford, J. Gilmer, N. Carlini, and D. Cubuk.
“Adversarial Examples Are a Natural Consequence of
44 Test Error in Noise”. In: (2019). arXiv: [1901 . 10513 \[cs.LG\]](#).
- 45 [For21] V. Fortuin. “Priors in Bayesian Deep Learning:
A Review”. In: (2021). arXiv: [2105 . 06868 \[stat.ML\]](#).
- 46 [For+95] J. Forbes, T. Huang, K. Kanazawa, and S.
Russell. “The BATmobile: Towards a Bayesian Auto-
47 mated Taxi”. In: *IJCAI*. 1995.
- [Fot+14] N. Foti, J. Xu, D. Laird, and E. Fox. “Stochas-
tic variational inference for hidden Markov models”. In:
NIPS. 2014, pp. 3599–3607.
- [Fox+08] E. Fox, E. Sudderth, M. Jordan, and A. Wil-
sky. “An HDP-HMM for Systems with State Persis-
tence”. In: *ICML*. 2008.
- [Fox09] E. B. Fox. “Bayesian nonparametric learning
of complex dynamical phenomena”. PhD thesis. Mas-
sachusetts Institute of Technology, 2009.
- [FP08] N. Friel and A. N. Pettitt. “Marginal Likelihood
Estimation via Power Posteriors”. In: *J. of Royal Stat.
Soc. Series B* 70.3 (2008), pp. 589–607.
- [FP69] D. C. Fraser and J. E. Potter. “The optimum
linear smoother as a combination of two optimum lin-
ear filters”. In: *IEEE Trans. on Automatical Control*
(1969), pp. 387–390.
- [FPD09] P. Frazier, W. Powell, and S. Dayanik. “The
knowledge-gradient policy for correlated normal be-
liefs”. In: *INFORMS J. on Computing* 21.4 (2009),
pp. 599–613.
- [Fra08] A. Fraser. *Hidden Markov Models and Dynam-
ical Systems*. SIAM Press, 2008.
- [Fra+16] M. Fraccaro, S. K. Sønderby, U. Paquet, and
O. Winther. “Sequential Neural Models with Stochas-
tic Layers”. In: *NIPS*. 2016.
- [Fra18] P. I. Frazier. “Bayesian Optimization”. In: *Re-
cent Advances in Optimization and Modeling of Con-
temporary Problems*. INFORMS TutORials in Opera-
tions Research. INFORMS, 2018, pp. 255–278.
- [Fre14] A. A. Freitas. “Comprehensible classification
models: a position paper”. In: *ACM SIGKDD explo-
rations newsletter* 15.1 (2014), pp. 1–10.
- [Fre+22] M. Freitag, D. Grangier, Q. Tan, and B.
Liang. “High Quality Rather than High Model Prob-
ability: Minimum Bayes Risk Decoding with Neural
Metrics”. In: *TACL*. 2022.
- [Fre98] B. Frey. *Graphical Models for Machine Learn-
ing and Digital Communication*. MIT Press, 1998.
- [Fre99] R. M. French. “Catastrophic forgetting in con-
nectionist networks”. In: *Trends in Cognitive Science*
(1999).
- [Fri09] K. Friston. “The free-energy principle: a rough
guide to the brain?” en. In: *Trends Cogn. Sci.* 13.7
(2009), pp. 293–301.
- [Fro+13] A. Frome, G. Corrado, J. Shlens, S. Bengio,
J. Dean, M. Ranzato, and T. Mikolov. “Devise: A deep
visual-semantic embedding model”. In: (2013).
- [Fro+21] R. Frostig, M. Johnson, D. Maclaurin, A.
Paszke, and A. Radul. “Decomposing reverse-mode au-
tomatic differentiation”. In: *LAFI workshop at POPL*
2021. 2021.
- [FS07] S. Frühwirth-Schnatter. *Finite Mixture and
Markov Switching Models*. Springer, 2007.
- [FSF10] S. Frühwirth-Schnatter and R. Frühwirth.
“Data Augmentation and MCMC for Binary and
Multinomial Logit Models”. In: *Statistical Modelling
and Regression Structures*. Ed. by T. Kneib and G.
Tutz. Springer, 2010, pp. 111–132.
- [FST98] S. Fine, Y. Singer, and N. Tishby. “The Hier-
archical Hidden Markov Model: Analysis and Applica-
tions”. In: *Machine Learning* 32 (1998), p. 41.

- 1
- 2 [FT05] M. Fashing and C. Tomasi. “Mean shift is a
3 bound optimization”. en. In: *IEEE Trans. Pattern
Anal. Mach. Intell.* 27.3 (2005), pp. 471–474.
- 4 [FT19] A. Finke and A. H. Thiery. “On the relationship
5 between variational inference and adaptive impor-
tance sampling”. In: (2019). arXiv: 1907 . 10477
6 [[stat.ML](#)].
- 7 [FT74] J. H. Friedman and J. W. Tukey. “A Projection
8 Pursuit Algorithm for Exploratory Data Analysis”. In:
9 *IEEE Trans. Comput.* C-23.9 (1974), pp. 881–890.
- 10 [Fu15] M. Fu, ed. *Handbook of Simulation Optimiza-
tion*. 1st ed. Springer-Verlag New York, 2015.
- 11 [Fu+17] Z. Fu, X. Tan, N. Peng, D. Zhao, and R. Yan.
“Style transfer in text: Exploration and evaluation”. In:
12 *arXiv preprint arXiv:1711.06861* (2017).
- 13 [Fu+19] H. Fu, C. Li, X. Liu, J. Gao, A. Celikyilmaz,
14 and L. Carin. “Cyclical Annealing Schedule: A Simple
15 Approach to Mitigating KL Vanishing”. In: *NAACL*.
2019.
- 16 [Fu+20] J. Fu, A. Kumar, O. Nachum, G. Tucker, and
17 S. Levine. *D4RL: Datasets for Deep Data-Driven Re-
inforcement Learning*. arXiv:2004.07219. 2020.
- 18 [Fuj05] S. Fujishige. *Submodular functions and opti-
mization*. Vol. 58. Elsevier Science, 2005.
- 19 [Ful+20] I. R. Fulcher, I. Shpitser, S. Marealle, and
20 E. J. Tchetgen Tchetgen. “Robust inference on popula-
21 tion indirect causal effects: the generalized front door
criterion”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 82.1 (2020),
22 pp. 199–214. eprint: <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/rssb.12345>.
- 23 [FV15] V. Feldman and J. Vondrák. “Tight Bounds
24 on Low-Degree Spectral Concentration of Submodu-
lar and XOS Functions”. In: *2015 IEEE 56th An-
nual Symposium on Foundations of Computer Sci-
ence*. 2015, pp. 923–942.
- 25 [FV16] V. Feldman and J. Vondrák. “Optimal bounds
26 on approximation of submodular and XOS functions
by junta”. In: *SIAM Journal on Computing* 45.3
27 (2016), pp. 1129–1170.
- 28 [FV17] R. C. Fong and A. Vedaldi. “Interpretable ex-
planations of black boxes by meaningful perturbation”.
In: *Proceedings of the IEEE international conference
on computer vision*. 2017, pp. 3429–3437.
- 29 [FW12] N. Friel and J. Wyse. “Estimating the evidence
30 – a review”. In: *Stat. Neerl.* 66.3 (2012), pp. 288–308.
- 31 [FW21] R. Friedman and Y. Weiss. “Posterior Sam-
pling for Image Restoration using Explicit Patch Priors”. In: (2021). arXiv: 2104.09895 [[cs.CV](#)].
- 32 [FWW21] M. Finzi, M. Welling, and A. G. Wilson. “A
33 Practical Method for Constructing Equivariant Multi-
layer Perceptrons for Arbitrary Matrix Groups”. In:
34 *ICML*. 2021.
- 35 [Gaj+19] A. Gajewski, J. Clune, K. O. Stanley, and J.
Lehman. “Evolvability ES: Scalable and Direct Opti-
mization of Evolvability”. In: *Proc. of the Conf. on
Genetic and Evolutionary Computation*. 2019.
- 36 [Gan07] M. Ganapathiraju. “Application of language
37 technologies in biology: Feature extraction and model-
ing for transmembrane helix prediction”. en. PhD the-
sis. 2007.
- 38 [Gan+16a] Y. Ganin, E. Ustinova, H. Ajakan, P. Ger-
main, and others. “Domain-adversarial training of neu-
ral networks”. In: *JMLR* (2016).
- 39 [Gan+16b] Y. Ganin, E. Ustinova, H. Ajakan, P. Ger-
main, H. Larochelle, F. Laviolette, M. Marchand, and
V. Lempitsky. “Domain-adversarial training of neural
networks”. In: *The journal of machine learning re-
search* 17.1 (2016), pp. 2096–2030.
- 40 [Gao+18] R. Gao, J. Xie, S.-C. Zhu, and Y. N. Wu.
“Learning Grid-like Units with Vector Representa-
tion of Self-Position and Matrix Representation of
Self-Motion”. In: *arXiv preprint arXiv:1810.05597*
2018.
- 41 [Gao+20] R. Gao, E. Nijkamp, D. P. Kingma, Z. Xu,
A. M. Dai, and Y. N. Wu. “Flow contrastive esti-
mation of energy-based models”. In: *Proceedings of
the IEEE/CVF Conference on Computer Vision and
Pattern Recognition*. 2020, pp. 7518–7528.
- 42 [Gär03] T. Gärtner. “A Survey of Kernels for Struc-
tured Data”. In: *SIGKDD Explor. Newsl.* 5.1 (2003),
pp. 49–58.
- 43 [GAR16] R. B. Grosse, S. Ancha, and D. M. Roy. “Mea-
suring the reliability of MCMC inference with bidirec-
tional Monte Carlo”. In: *NIPS*. 2016.
- 44 [Gar+18a] J. Gardner, G. Pleiss, K. Q. Weinberger,
D. Bindel, and A. G. Wilson. “GPyTorch: Blackbox
Matrix-Matrix Gaussian Process Inference with GPU
Acceleration”. In: *NIPS*. Ed. by S Bengio, H Wallach,
H Larochelle, K Grauman, N Cesa-Bianchi, and R Gar-
nett. Curran Associates, Inc., 2018, pp. 7576–7586.
- 45 [Gar+18b] J. R. Gardner, G. Pleiss, R. Wu, K. Q. Wein-
berger, and A. G. Wilson. “Product Kernel Interpolation
for Scalable Gaussian Processes”. In: *AISTATS*.
2018.
- 46 [Gar+18c] T. Garipov, P. Izmailov, D. Podoprikhin, D.
Vetrov, and A. G. Wilson. “Loss Surfaces, Mode Con-
nectivity, and Fast Ensembling of DNNs”. In: *NIPS*.
2018.
- 47 [Gar+18d] M. Garnelo, D. Rosenbaum, C. Maddison,
T. Ramalho, D. Saxton, M. Shanahan, Y. W. Teh, D.
Rezende, and S. M. A. Eslami. “Conditional Neural
Processes”. In: *ICML*. Ed. by J. Dy and A. Krause.
Vol. 80. Proceedings of Machine Learning Research.
PMLR, 2018, pp. 1704–1713.
- 48 [Gar+18e] M. Garnelo, J. Schwarz, D. Rosenbaum, F.
Viola, D. J. Rezende, S. M. Ali Eslami, and Y. W. Teh.
“Neural Processes”. In: *ICML workshop on Theoret-
ical Foundations and Applications of Deep Genera-
tive Models*. 2018.
- 49 [Gar+19] S. Garg, V. Perot, N. Limtiaco, A. Taly, E. H.
Chi, and A. Beutel. “Counterfactual Fairness in Text
Classification through Robustness”. In: *Proceedings of
the 2019 AAAI/ACM Conference on AI, Ethics, and
Society*. AIES ’19. Association for Computing Machinery,
2019, pp. 219–226.
- 50 [Gar22] R. Garnett. *Bayesian Optimization*. in prepara-
tion. Cambridge University Press, 2022.
- 51 [Gas+19] J. Gasthaus, K. Benidis, Y. Wang, S. S.
Rangapuram, D. Salinas, V. Flunkert, and T.
Januschowski. “Probabilistic Forecasting with Spline
Quantile Function RNNs”. In: *ICML*. Vol. 89. Pro-
ceedings of Machine Learning Research. PMLR, 2019,
pp. 1901–1910.
- 52 [GAZ19] A. Ghorbani, A. Abid, and J. Zou. “Interpre-
tation of neural networks is fragile”. In: *Proceedings
of the AAAI Conference on Artificial Intelligence*.
Vol. 33. 01. 2019, pp. 3681–3688.
- 53 [GB00] Z. Ghahramani and M. Beal. “Variational in-
ference for Bayesian mixtures of factor analysers”. In:
54 *NIPS-12*. 2000.

- [GB09] A. Guillory and J. Bilmes. "Label Selection on Graphs". In: *NIPS*. Vancouver, Canada, 2009.
- [GB10] X. Glorot and Y. Bengio. "Understanding the difficulty of training deep feedforward neural networks". In: *AISTATS*. 2010, pp. 249–256.
- [GB11] A. Guillory and J. Bilmes. "Active Semi-Supervised Learning using Submodular Functions". In: *UAI*. Barcelona, Spain: AUAI, 2011.
- [GB+18] R. Gómez-Bombarelli, J. N. Wei, D. Duvenaud, J. M. Hernández-Lobato, B. Sánchez-Lengeling, D. Sheberla, J. Aguilera-Iparragirre, T. D. Hirzel, R. P. Adams, and A. Aspuru-Guzik. "Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules". en. In: *American Chemical Society Central Science* 4.2 (2018), pp. 268–276.
- [GBB11] X. Glorot, A. Bordes, and Y. Bengio. "Deep Sparse Rectifier Neural Networks". In: *AISTATS*. 2011.
- [GBC16] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
- [GBJ18] R. Giordano, T. Broderick, and M. I. Jordan. "Covariances, Robustness, and Variational Bayes". In: *JMLR* 19.51 (2018), pp. 1–49.
- [GBP18] C. Gurau, A. Bewley, and I. Posner. "Dropout Distillation for Efficiently Estimating Model Confidence". 2018.
- [GC11] M. Girolami and B. Calderhead. "Riemann manifold Langevin and Hamiltonian Monte Carlo methods". In: *J. of Royal Stat. Soc. Series B* 73.2 (2011), pp. 123–214.
- [GC90] R. P. Goldman and E. Charniak. "Dynamic Construction of Belief Networks". In: *UAI*. 1990.
- [GCW19] M. Gerber, N. Chopin, and N. Whiteley. "Negative association, ordering and convergence of resampling methods". In: *Ann. Stat.* 47.4 (2019), pp. 2236–2260.
- [GD20] T. Geffner and J. Domke. "A Rule for Gradient Estimator Selection, with an Application to Variational Inference". In: *AISTATS*. Ed. by S. Chiappa and R. Calandra. Vol. 108. Proceedings of Machine Learning Research. PMLR, 2020, pp. 1803–1812.
- [GDFY16] S. Ghosh, F. M. Delle Fave, and J. Yedidia. "Assumed Density Filtering Methods for Learning Bayesian Neural Networks". In: *AAAI*. 2016.
- [Geb+21] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford. "Datasheets for datasets". In: *Communications of the ACM* 64.12 (2021), pp. 86–92.
- [Ged+20] D. Gedon, N. Wahlström, T. B. Schön, and L. Ljung. "Deep State Space Models for Nonlinear System Identification". In: (Mar. 2020). arXiv: [2003.14162 \[eess.SY\]](https://arxiv.org/abs/2003.14162).
- [Gei+20a] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann. "Shortcut Learning in Deep Neural Networks". In: *arXiv preprint arXiv:2004.07780* (2020).
- [Gei+20b] R. Geirhos, J. Jacobsen, C. Michaelis, R. S. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann. "Shortcut Learning in Deep Neural Networks". In: *CoRR* abs/2004.07780 (2020). arXiv: [2004.07780](https://arxiv.org/abs/2004.07780).
- [Gel+04] A. Gelman, J. Carlin, H. Stern, and D. Rubin. *Bayesian data analysis*. 2nd edition. Chapman and Hall, 2004.
- [Gel06] A. Gelman. "Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper)". en. In: *Bayesian Anal.* 1.3 (2006), pp. 515–534.
- [Gel+08] A. Gelman, A. Jakulin, M. G. Pittau, and Y.-S. Su. "A weakly informative default prior distribution for logistic and other regression models". en. In: *The Annals of Applied Statistics* 2.4 (2008), pp. 1360–1383.
- [Gel+14a] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian Data Analysis, Third Edition*. Third edition. Chapman and Hall/CRC, 2014.
- [Gel+14b] A. Gelman, A. Vehtari, P. Jylänki, C. Robert, N. Chopin, and J. P. Cunningham. "Expectation propagation as a way of life". In: (2014). arXiv: [1412.4869 \[stat.CO\]](https://arxiv.org/abs/1412.4869).
- [Gel+20] A. Gelman, A. Vehtari, D. Simpson, C. C. Margossian, B. Carpenter, Y. Yao, L. Kennedy, J. Gabry, P.-C. Bürkner, and M. Modrák. "Bayesian Workflow". In: (2020). arXiv: [2011.01808 \[stat.ME\]](https://arxiv.org/abs/2011.01808).
- [Gel+22] A. Gelman, J. Hill, B. Goodrich, J. Gabry, D. Simpson, and A. Vehtari. *Applied Regression and Multilevel Models*. To appear. 2022.
- [Gel74] A. Gelb. *Applied Optimal Estimation*. en. The MIT Press, May 1974.
- [Gel90] M. Gelbrich. "On a formula for the L2 Wasserstein metric between measures on Euclidean and Hilbert spaces". In: *Mathematische Nachrichten* 147.1 (1990), pp. 185–203.
- [Gen+19] A. Genevay, L. Chizat, F. Bach, M. Cuturi, and G. Peyré. "Sample complexity of sinkhorn divergences". In: *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019, pp. 1574–1583.
- [Geo+17] D. George et al. "A generative vision model that trains with high data efficiency and breaks text-based CAPTCHAs". In: *Science* 358.6368 (2017), eaag2612, eprint: <https://www.science.org/doi/pdf/10.1126/science.aag2612>.
- [Geo+18] T. George, C. Laurent, X. Bouthillier, N. Balas, and P. Vincent. "Fast Approximate Natural Gradient Descent in a Kronecker Factored Eigenbasis". In: *NIPS*. Curran Associates, Inc., 2018, pp. 9550–9560.
- [Geo88] H.-O. Georgii. *Gibbs Measures and Phase Transitions*. en. Walter De Gruyter Inc, 1988.
- [Ger+15] M. Germain, K. Gregor, I. Murray, and H. Larochelle. "MADE: Masked Autoencoder for Distribution Estimation". In: *ICML*. Ed. by F. Bach and D. Blei. Vol. 37. Proceedings of Machine Learning Research. PMLR, 2015, pp. 881–889.
- [Gér19] A. Géron. *Hands-On Machine Learning with Scikit-Learn and TensorFlow: Concepts, Tools, and Techniques for Building Intelligent Systems (2nd edition)*. en. O'Reilly Media, Incorporated, 2019.
- [Ger19] S. J. Gershman. "What does the free energy principle tell us about the brain?" In: *Neurons, Behavior, Data Analysis, and Theory* (2019).
- [GEY19] Y. Geifman and R. El-Yaniv. "SelectiveNet: A Deep Neural Network with an Integrated Reject Option". In: *ICML*. 2019.
- [Gey92] C. Geyer. "Practical Markov chain Monte Carlo". In: *Statistical Science* 7 (1992), pp. 473–483.

- 1
- 2 [GF00] E. George and D. Foster. “Calibration and em-
pirical Bayes variable selection”. In: *Biometrika* 87.4
(2000), pp. 731–747.
- 3
- 4 [GF09] I. E. Givoni and B. J. Frey. “A Binary Variable
Model for Affinity Propagation”. In: *Neural Computation* 21.6 (2009), pp. 1589–1600.
- 5
- 6 [GF+15] Á. F. García-Fernández, L. Svensson, M. R.
Morelande, and S. Särkkä. “Posterior Linearization
Filter: Principles and Implementation Using Sigma
Points”. In: *IEEE Trans. Signal Process.* 63.20 (Oct.
2015), pp. 5561–5573.
- 7
- 8 [GF17] B. Goodman and S. Flaxman. “European
Union regulations on algorithmic decision-making and
a “right to explanation””. In: *AI magazine* 38.3 (2017),
pp. 50–57.
- 9
- 10 [GFSS17] Á. F. García-Fernández, L. Svensson, and S.
Särkkä. “Iterated Posterior Linearization Smoother”.
In: *IEEE Trans. Automat. Contr.* 62.4 (2017),
pp. 2056–2063.
- 11
- 12 [GFTS19] Á. F. García-Fernández, F. Tronarp, and S.
Särkkä. “Gaussian Process Classification Using Pos-
terior Linearization”. In: *IEEE Signal Process. Lett.*
26.5 (2019), pp. 735–739.
- 13
- 14 [GFWO20] T. Galy-Fajou, F. Wenzel, and M. Opper.
“Automated Augmented Conjugate Inference for Non-
conjugate Gaussian Process Models”. In: *AISTATS*.
2020.
- 15
- 16 [GG11] T. L. Griffiths and Z. Ghahramani. “The In-
dian Buffet Process: An Introduction and Review.” In:
JMLR 12.4 (2011).
- 17
- 18 [GG14] S. J. Gershman and N. D. Goodman. “Amor-
tized Inference in Probabilistic Reasoning”. In: *36th
Annual Conference of the Cognitive Science Society*.
2014.
- 19
- 20 [GG16] Y. Gal and Z. Ghahramani. “Dropout as a
Bayesian Approximation: Representing Model Uncer-
tainty in Deep Learning”. In: *ICML*. 2016.
- 21
- 22 [GG84] S. Geman and D. Geman. “Stochastic Relax-
ation, Gibbs Distributions, and the Bayesian Restora-
tion of Images”. In: *IEEE PAMI* 6.6 (1984).
- 23
- 24 [GGG15] M. Gygli, H. Grabner, and L. Gool. “Video
summarization by learning submodular mixtures of ob-
jectives”. In: *2015 IEEE Conference on Computer
Vision and Pattern Recognition (CVPR)* (2015),
pp. 3090–3098.
- 25
- 26 [GGS21] M. Grcic, I. Grubisic, and S. Segvic. “Dense-
Flow: Official implementation of Densely connected
normalizing flows”. en. In: *NIPS*. 2021.
- 27
- 28 [GH07] A. Gelman and J. Hill. *Data analysis using
regression and multilevel/ hierarchical models*. Cam-
bridge, 2007.
- 29
- 30 [GH10] M. Gutmann and A. Hyvärinen. “Noise-
contrastive estimation: A new estimation principle for
unnormalized statistical models”. In: *Proceedings of the
Thirteenth International Conference on Artificial
Intelligence and Statistics*. 2010, pp. 297–304.
- 31
- 32 [GH12] M. Gutmann and J.-i. Hirayama. “Bregman di-
vergence as general framework to estimate unnormal-
ized statistical models”. In: *arXiv:1202.3727* (2012).
- 33
- 34 [GH96a] Z. Ghahramani and G. Hinton. *Parameter
estimation for linear dynamical systems*. Tech. rep.
CRG-TR-96-2. Dept. Comp. Sci., Univ. Toronto, 1996.
- 35
- 36 [GH96b] Z. Ghahramani and G. Hinton. *The EM Al-
gorithm for Mixtures of Factor Analyzers*. Tech. rep.
Dept. of Comp. Sci., Uni. Toronto, 1996.
- 37
- 38 [GH98] Z. Ghahramani and G. Hinton. “Variational
learning for switching state-space models”. In: *Neural
Computation* 12.4 (1998), pp. 963–996.
- 39
- 40 [Gha+15] M. Ghavamzadeh, S. Mannor, J. Pineau, and
A. Tamar. “Bayesian Reinforcement Learning: A Sur-
vey”. In: *Foundations and Trends in ML* (2015).
- 41
- 42 [Gha+21] A. Ghandeharioun, B. Kim, C.-L. Li, B. Jou,
B. Eoff, and R. W. Picard. *DISSECT: Disentangled
Simultaneous Explanations via Concept Traversals*.
2021. arXiv: [2105.15164 \[cs.LG\]](https://arxiv.org/abs/2105.15164).
- 43
- 44 [GHC20] C. Geng, S.-J. Huang, and S. Chen. “Recent
Advances in Open Set Recognition: A Survey”. In:
IEEE PAMI (2020).
- 45
- 46 [GHC21] S. Gould, R. Hartley, and D. J. Campbell.
“Deep Declarative Networks”. en. In: *IEEE PAMI* PP
(2021).
- 47
- 48 [GHK17] Y. Gal, J. Hron, and A. Kendall. “Concrete
Dropout”. In: (2017). arXiv: [1705.07832 \[stat.ML\]](https://arxiv.org/abs/1705.07832).
- 49
- 50 [Gho+19] A. Ghorbani, J. Wexler, J. Zou, and B.
Kim. *Towards Automatic Concept-based Explanations*.
2019. arXiv: [1902.03129 \[stat.ML\]](https://arxiv.org/abs/1902.03129).
- 51
- 52 [Gho+21] A. Ghosh, A. Honoré, D. Liu, G. E. Henter,
and S. Chatterjee. “Normalizing Flow based Hidden
Markov Models for Classification of Speech Phones
with Explainability”. In: (July 2021). arXiv: [2107.00730 \[cs.LG\]](https://arxiv.org/abs/2107.00730).
- 53
- 54 [GHV14] A. Gelman, J. Hwang, and A. Vehtari.
“Understanding predictive information criteria for
Bayesian models”. In: *Statistics and Computing* 24.6
(2014), pp. 997–1016.
- 55
- 56 [GHV20] A. Gelman, J. Hill, and A. Vehtari. *Regres-
sion and Other Stories*. Cambridge, 2020.
- 57
- 58 [GIG17] Y. Gal, R. Islam, and Z. Ghahramani. “Deep
bayesian active learning with image data”. In: *Inter-
national Conference on Machine Learning*. PMLR.
2017, pp. 1183–1192.
- 59
- 60 [Gil+18a] J. Gilmer, R. P. Adams, I. Goodfellow, D.
Andersen, and G. E. Dahl. “Motivating the rules of
the game for adversarial example research”. In: *arXiv
preprint arXiv:1807.06732* (2018).
- 61
- 62 [Gil+18b] J. Gilmer, L. Metz, F. Faghri, S. S. Schoen-
holz, M. Raghu, M. Wattenberg, and I. Good-
fellow. “Adversarial spheres”. In: *arXiv preprint
arXiv:1801.02774* (2018).
- 63
- 64 [Gil88] J. R. Gilbert. “Some nested dissection order is
nearly optimal”. In: *Inf. Process. Lett.* 26.6 (1988),
pp. 325–328.
- 65
- 66 [Gir+14] R. Girshick, J. Donahue, T. Darrell, and J.
Malik. “Rich feature hierarchies for accurate object
detection and semantic segmentation”. In: *Proceedings
of the IEEE conference on computer vision and pat-
tern recognition*. 2014, pp. 580–587.
- 67
- 68 [Gir+15] R. Girshick, F. Iandola, T. Darrell, and J. Ma-
lik. “Deformable Part Models are Convolutional Neu-
ral Networks”. In: *CVPR*. 2015.
- 69
- 70 [Gir+21] L. Girin, S. Leglaive, X. Bie, J. Diard, T. Hue-
ber, and X. Alameda-Pineda. “Dynamical Variational
Autoencoders: A Comprehensive Review”. In: *Foun-
dations and Trends® in Machine Learning* 15.1–2
(2021), pp. 1–175.

- [Git89] J. Gittins. *Multi-armed Bandit Allocation Indices*. Wiley, 1989.
- [GJ97] Z. Ghahramani and M. Jordan. “Factorial Hidden Markov Models”. In: *Machine Learning* 29 (1997), pp. 245–273.
- [GK19] L. Graesser and W. L. Keng. *Foundations of Deep Reinforcement Learning: Theory and Practice in Python*. en. 1 edition. Addison-Wesley Professional, 2019.
- [GKS05] C. Guestrin, A. Krause, and A. P. Singh. “Near-optimal sensor placements in gaussian processes”. In: *Proceedings of the 22nd international conference on Machine learning*. 2005, pp. 265–272.
- [GL10] K. Gregor and Y. LeCun. “Learning fast approximations of sparse coding”. In: *ICML*. 2010, pp. 399–406.
- [GL97] F. Glover and M. Laguna. Kluwer Academic Publishers, 1997.
- [Gla03] P. Glasserman. *Monte Carlo Methods in Financial Engineering*. 1st ed. Stochastic Modelling and Applied Probability. Springer-Verlag New York, 2003.
- [Gle02] S. Glennan. “Rethinking mechanistic explanation”. In: *Philosophy of science* 69.S3 (2002), S342–S353.
- [GLM15] J. Ghosh, Y. Li, and R. Mitra. “On the Use of Cauchy Prior Distributions for Bayesian Logistic Regression”. In: (2015). arXiv: [1507.07170 \[stat.ME\]](https://arxiv.org/abs/1507.07170).
- [GLP21] I. Gulrajani and D. Lopez-Paz. “In Search of Lost Domain Generalization”. In: *ICLR*. 2021.
- [GLS81] M. Grötschel, L. Lovász, and A. Schrijver. “The ellipsoid method and its consequences in combinatorial optimization”. In: *Combinatorica* 1.2 (1981), pp. 169–197.
- [GM12] G. Gordon and J. McNulty. *Matroids: a geometric introduction*. Cambridge University Press, 2012.
- [GM15] J. Gorham and L. Mackey. “Measuring sample quality with Stein’s method”. In: *Advances in Neural Information Processing Systems*. 2015, pp. 226–234.
- [GM16] R. Grosse and J. Martens. “A Kronecker-factored approximate Fisher matrix for convolution layers”. In: *ICML*. 2016.
- [GM17] P. L. Green and S Maskell. “Estimating the parameters of dynamical systems from Big Data using Sequential Monte Carlo samplers”. In: *Mech. Syst. Signal Process.* 93 (2017), pp. 379–396.
- [GM98] A. Gelman and X.-L. Meng. “Simulating normalizing constants: from importance sampling to bridge sampling to path sampling”. In: *Statistical Science* 13 (1998), pp. 163–185.
- [GMAR16] Y. Gal, R. T. Mc Allister, and C. E. Rasmussen. “Improving PILCO with Bayesian Neural Network Dynamics Models”. In: *ICML workshop on Data-efficient machine learning*. 2016.
- [GMH20] M. Gorinova, D. Moore, and M. Hoffman. “Automatic Reparameterisation of Probabilistic Programs”. In: *ICML*. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 3648–3657.
- [GNM19] D. Greenberg, M. Nonnenmacher, and J. Macke. “Automatic Posterior Transformation for Likelihood-Free Inference”. In: *ICML*. 2019.
- [Goe+09] M. X. Goemans, N. J. Harvey, S. Iwata, and V. Mirrokni. “Approximating submodular functions everywhere”. In: *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms*. SIAM, 2009, pp. 535–544.
- [Gol+04] J. Goldberger, S. T. Roweis, G. E. Hinton, and R. Salakhutdinov. “Neighbourhood Components Analysis”. In: *NIPS*. 2004.
- [Gol+17] N. Gold, M. G. Frasch, C. Herry, B. S. Richardson, and X. Wang. “A Doubly Stochastic Change Point Detection Algorithm for Noisy Biological Signals”. In: *Front. Physiol.* (2017), p. 106088.
- [Gol17] Y. Goldberg. *Neural Network Methods for Natural Language Processing*. Morgan & Claypool Publishers, 2017.
- [Gol+19] Z. Goldfeld, E. van den Berg, K. H. Greenwald, I. Melnyk, N. Nguyen, B. Kingsbury, and Y. Polyanskiy. “Estimating Information Flow in Deep Neural Networks”. In: *ICML*. 2019.
- [Gol89] D. E. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*. 1st. Addison-Wesley Longman Publishing Co., Inc., 1989.
- [Gom+17] A. N. Gomez, M. Ren, R. Urtasun, and R. B. Grosse. “The Reversible Residual Network: Backpropagation Without Storing Activations”. In: *NIPS*. 2017.
- [Gon+11] J. Gonzalez, Y. Low, A. Gretton, and C. Guestrin. “Parallel gibbs sampling: From colored fields to thin junction trees”. In: *AISTATS*. 2011, pp. 324–332.
- [Gon+14] B. Gong, W.-L. Chao, K. Grauman, and F. Sha. “Diverse sequential subset selection for supervised video summarization”. In: *Advances in neural information processing systems* 27 (2014), pp. 2069–2077.
- [Gon+20] P. J. Goncalves et al. “Training deep neural density estimators to identify mechanistic models of neural dynamics”. In: *Elife* 9 (2020).
- [Goo+14] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. “Generative Adversarial Networks”. In: *NIPS*. 2014.
- [Goo16] I. Goodfellow. “NIPS 2016 Tutorial: Generative Adversarial Networks”. In: *NIPS Tutorial*. 2016.
- [Goo85] I. Good. “Weight of evidence: A brief survey”. In: *Bayesian statistics* 2 (1985), pp. 249–270.
- [Gor+14] A. D. Gordon, T. A. Henzinger, A. V. Nori, and S. K. Rajamani. “Probabilistic programming”. In: *Intl. Conf. on Software Engineering*. 2014.
- [Gor+19] J. Gordon, J. Bronskill, M. Bauer, S. Nowozin, and R. E. Turner. “Meta-Learning Probabilistic Inference For Prediction”. In: *ICLR*. 2019.
- [Gor93] N. Gordon. “Novel Approach to Nonlinear/Non-Gaussian Bayesian State Estimation”. In: *IEE Proceedings (F)* 140.2 (1993), pp. 107–113.
- [Gor95] G. J. Gordon. “Stable Function Approximation in Dynamic Programming”. In: *ICML*. 1995, pp. 261–268.
- [Gou+96] C. Gourieroux, M. Gourieroux, A. Monfort, and D. A. Monfort. *Simulation-based econometric methods*. Oxford university press, 1996.
- [Goy+19] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. *Counterfactual Visual Explanations*. 2019. arXiv: [1904.07451 \[cs.LG\]](https://arxiv.org/abs/1904.07451).
- [Goy+22] S. Goyal, M. Sun, A. Raghunathan, and Z. Kolter. “Test-Time Adaptation via Conjugate Pseudo-labels”. In: (July 2022). arXiv: [2207.09640 \[cs.LG\]](https://arxiv.org/abs/2207.09640).

- 1
- 2 [GPS89] D. Greig, B. Porteous, and A. Seheult. “Exact maximum a posteriori estimation for binary images”. In: *J. of Royal Stat. Soc. Series B* 51.2 (1989), pp. 271–279.
- 3
- 4 [GR06a] M. Girolami and S. Rogers. “Variational Bayesian Multinomial Probit Regression with Gaussian Process Priors”. In: *Neural Comput.* 18.8 (2006), pp. 1790–1817.
- 5
- 6 [GR06b] M. Girolami and S. Rogers. “Variational Bayesian multinomial probit regression with Gaussian process priors”. In: *Neural Computation* 18.8 (2006), pp. 1790–1817.
- 7
- 8 [GR07a] T. Gneiting and A. E. Raftery. “Strictly Proper Scoring Rules, Prediction, and Estimation”. In: *JASA* 102.477 (2007), pp. 359–378.
- 9
- 10 [GR07b] T. Gneiting and A. E. Raftery. “Strictly proper scoring rules, prediction, and estimation”. In: *Journal of the American statistical Association* 102.477 (2007), pp. 359–378.
- 11
- 12 [Gra+10] T. Graepel, J. Quinonero-Candela, T. Borchert, and R. Herbrich. “Web-Scale Bayesian Click-Through Rate Prediction for Sponsored Search Advertising in Microsoft’s Bing Search Engine”. In: *ICML*. 2010.
- 13
- 14 [Gra11] A. Graves. “Practical Variational Inference for Neural Networks”. In: *NIPS*. 2011.
- 15
- 16 [Gra+18] W. Grathwohl, R. T. Q. Chen, J. Bettencourt, I. Sutskever, and D. Duvenaud. “FFJORD: Free-form Continuous Dynamics for Scalable Reversible Generative Models”. In: (2018). arXiv: 1810 . 01367 [cs, LG].
- 17
- 18 [Gra+20a] W. Grathwohl, J. Kelly, M. Hashemi, M. Norouzi, K. Swersky, and D. Duvenaud. “No MCMC for me: Amortized sampling for fast and stable training of energy-based models”. In: *arXiv preprint arXiv:2010.04230* (2020).
- 19
- 20 [Gra+20b] W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky. “Your classifier is secretly an energy based model and you should treat it like one”. In: *ICLR*. 2020.
- 21
- 22 [Gra+20c] W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, and R. Zemel. “Cutting out the Middle-Man: Training and Evaluating Energy-Based Models without Sampling”. In: *arXiv preprint arXiv:2002.05616* (2020).
- 23
- 24 [Gre03] P. Green. “Tutorial on trans-dimensional MCMC”. In: *Highly Structured Stochastic Systems*. Ed. by P. Green, N. Hjort, and S. Richardson. OUP, 2003.
- 25
- 26 [Gre+12] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. “A Kernel Two-Sample Test”. In: *JMLR* 13.Mar (2012), pp. 723–773.
- 27
- 28 [Gre+14] K. Gregor, I. Danihelka, A. Mnih, C. Blundell, and D. Wierstra. “Deep AutoRegressive Networks”. In: *ICML*. 2014.
- 29
- 30 [Gre20] F. Greenlee. *Transformer VAE*. 2020.
- 31
- 32 [Gre+22] P. L. Green, R. E. Moore, R. J. Jackson, J. Li, and S. Maskell. “Increasing the efficiency of Sequential Monte Carlo samplers through the use of approximately optimal L-kernels”. In: *Mech. Syst. Signal Process.* 162 (2022).
- 33
- 34 [Gre98] P. Green. “Reversible Jump Markov Chain Monte Carlo computation and Bayesian model determination”. In: *Biometrika* 82 (1998), pp. 711–732.
- 35
- 36 [Gri+04] T. Griffiths, M. Steyvers, D. Blei, and J. Tenenbaum. “Integrating Topics and Syntax”. In: *NIPS*. 2004.
- 37
- 38 [Gri20] T. L. Griffiths. “Understanding Human Intelligence through Human Limitations”. en. In: *Trends Cogn. Sci.* 24.11 (2020), pp. 873–883.
- 39
- 40 [Gri+20] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, et al. “Bootstrap your own latent: A new approach to self-supervised learning”. In: *arXiv preprint arXiv:2006.07733* (2020).
- 41
- 42 [GRS96] W. Gilks, S. Richardson, and D. Spiegelhalter. *Markov Chain Monte Carlo in Practice*. Chapman and Hall, 1996.
- 43
- 44 [GS04] T. Griffiths and M. Steyvers. “Finding scientific topics”. In: *PNAS* 101 (2004), pp. 5228–5235.
- 45
- 46 [GS08] Y. Guo and D. Schuurmans. “Efficient global optimization for exponential family PCA and low-rank matrix factorization”. In: *2008 46th Annual Allerton Conference on Communication, Control, and Computing*. 2008, pp. 1100–1107.
- 47
- 48 [GS13] A. Gelman and C. R. Shalizi. “Philosophy and the practice of Bayesian statistics”. en. In: *Br. J. Math. Stat. Psychol.* 66.1 (Feb. 2013), pp. 8–38.
- 49
- 50 [GS15] R. B. Grosse and R. Salakhutdinov. “Scaling Up Natural Gradient by Sparsely Factorizing the Inverse Fisher Matrix”. In: *ICML*. 2015.
- 51
- 52 [GS90] A. Gelfand and A. Smith. “Sampling-based approaches to calculating marginal densities”. In: *JASA* 85 (1990), pp. 385–409.
- 53
- 54 [GS92] G. Grimmett and D. Stirzaker. *Probability and Random Processes*. Oxford, 1992.
- 55
- 56 [GSA14] M. A. Gelbart, J. Snoek, and R. P. Adams. “Bayesian Optimization with Unknown Constraints”. In: *UAI*. 2014.
- 57
- 58 [GSD13] A. Guez, D. Silver, and P. Dayan. “Scalable and Efficient Bayes-Adaptive Reinforcement Learning Based on Monte-Carlo Tree Search”. In: *JAIR* 48 (2013), pp. 841–883.
- 59
- 60 [GSJ19] A. Gretton, D. Sutherland, and W. Jitkrittum. *NIPS tutorial on interpretable comparison of distributions and models*. 2019.
- 61
- 62 [GSK18] S. Gidaris, P. Singh, and N. Komodakis. “Unsupervised Representation Learning by Predicting Image Rotations”. In: *International Conference on Learning Representations*. 2018.
- 63
- 64 [GSM18] U. Garciarena, R. Santana, and A. Mendiburu. “Expanding Variational Autoencoders for Learning and Exploiting Latent Representations in Search Distributions”. In: *Proc. of the Conf. on Genetic and Evolutionary Computation*. 2018, pp. 849–856.
- 65
- 66 [GSS15] I. J. Goodfellow, J. Shlens, and C. Szegedy. “Explaining and Harnessing Adversarial Examples”. In: *ICLR*. 2015.
- 67
- 68 [GSZ21] P. Grünwald, T. Steinke, and L. Zakynthinou. “PAC-Bayes, MAC-Bayes and Conditional Mutual Information: Fast rate bounds that handle general VC classes”. In: *COLT*. 2021.
- 69
- 70 [GT86] J. R. Gilbert and R. E. Tarjan. “The analysis of a nested dissection algorithm”. In: *Numer. Math.* 50.4 (1986), pp. 377–404.
- 71
- 72 [Gue19] B. Guedj. “A primer on PAC-Bayesian learning”. In: *arXiv preprint arXiv:1901.05353* (2019).

- 1 [Gul+17] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Du-
2 moulin, and A. C. Courville. “Improved training of
3 wasserstein gans”. In: *NIPS*. 2017, pp. 5767–5777.
- 4 [Gul+20] C. Gulcehre et al. *RL Unplugged: Benchmarks for Offline Reinforcement Learning*. arXiv:2006.13888. 2020.
- 5 [Gum54] E. J. Gumbel. *Statistical theory of extreme values and some practical applications;; A series of lectures (United States. National Bureau of Standards. Applied mathematics series)*. en. 1st edition. U.S. Govt. Print. Office, 1954.
- 6 [Guo09] Y. Guo. “Supervised exponential family principal component analysis via convex optimization”. In: *NIPS*. 2009.
- 7 [Guo+17] C. Guo, G. Pleiss, Y. Sun, and K. Q. Wein-
8 berger. “On Calibration of Modern Neural Networks”. In: *ICML*. 2017.
- 9 [Gup+16] M. Gupta, A. Cotter, J. Pfeifer, K. Voevod-
10 ski, K. Canini, A. Mangylov, W. Moczydłowski, and
11 A. van Esbroeck. “Monotonic Calibrated Interpolated
12 Look-Up Tables”. In: *Journal of Machine Learning
13 Research* 17.109 (2016), pp. 1–47.
- 14 [Gur+18] S. Gururangan, S. Swayamdipta, O. Levy, R.
15 Schwartz, S. R. Bowman, and N. A. Smith. “Annotation
16 Artifacts in Natural Language Inference Data”. In: *CoRR* abs/1803.02324 (2018). arXiv: 1803.02324.
- 17 [Gus01] M. Gustafsson. “A probabilistic derivation
18 of the partial least-squares algorithm”. In: *Journal
19 of Chemical Information and Modeling* 41 (2001),
20 pp. 288–294.
- 21 [Gut+14] M. U. Gutmann, R. Dutta, S. Kaski, and
22 J. Corander. “Statistical inference of intractable gener-
23 ative models via classification”. In: *arXiv preprint
24 arXiv:1407.4981* (2014).
- 25 [Gut22] M. U. Gutmann. “Pen and Paper Exercises in
26 Machine Learning”. In: (June 2022). arXiv: 2206.13446
27 [cs.LG].
- 28 [GV17] S. Ghosal and A. van der Vaart. *Fundamen-
29 tals of Nonparametric Bayesian Inference*. en. 1st ed.
30 Cambridge University Press, 2017.
- 31 [GW08] A. Griewank and A. Walther. *Evaluating
32 Derivatives: Principles and Techniques of Algorithmic
33 Differentiation*. Second. Society for Industrial
34 and Applied Mathematics, 2008.
- 35 [GW92] W. Gilks and P. Wild. “Adaptive rejection
36 sampling for Gibbs sampling”. In: *Applied Statistics*
37 41 (1992), pp. 337–348.
- 38 [GXG18] H. Ge, K. Xu, and Z. Ghahramani. “Turing: a
39 language for flexible probabilistic inference”. In: *AIS-
40 TATS*. 2018, pp. 1682–1690.
- 41 [GZG19] S. K. S. Ghasemipour, R. S. Zemel, and S. Gu.
42 “A Divergence Minimization Perspective on Imitation
43 Learning Methods”. In: *CORL*. 2019, pp. 1259–1277.
- 44 [HA21] R. J. Hyndman and G. Athanasopoulos. *Fore-
45 casting: Principles and Practice*. en. 3rd ed. Otexts,
46 2021.
- 47 [Haa+17] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. “Reinforcement learning with deep energy-
48 based policies”. In: *Proceedings of the 34th Interna-
49 tional Conference on Machine Learning-Volume 70*. 2017, pp. 1352–1361.
- 50 [Haa+18a] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochas-
51 tic Actor”. In: *International Conference on Machine
52 Learning*. 2018, pp. 1861–1870.
- 53 [Haa+18b] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor”. In: *ICML*. 2018.
- 54 [Haa+18c] T. Haarnoja et al. “Soft Actor-Critic Al-
55 gorithms and Applications”. In: (2018). arXiv: 1812.
56 05905 [cs.LG].
- 57 [HAB17] J. H. Huggins, R. P. Adams, and T. Broderick.
58 “PASS-GLM: polynomial approximate sufficient statis-
59 tics for scalable Bayesian GLM inference”. In: *NIPS*.
60 2017.
- 61 [Had+20] R. Hadsell, D. Rao, A. A. Rusu, and R. Pascanu. “Embracing Change: Continual Learning in
62 Deep Neural Networks”. en. In: *Trends Cogn. Sci.*
63 24.12 (2020), pp. 1028–1040.
- 64 [HAE16] M. Huh, P. Agrawal, and A. A. Efros. “What
65 makes ImageNet good for transfer learning?” In: *arXiv
66 preprint arXiv:1608.08614* (2016).
- 67 [Haf18] D. Hafner. *Building Variational Auto-
68 Encoders in TensorFlow*. Blog post. 2018.
- 69 [Haf+19] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas,
70 D. Ha, H. Lee, and J. Davidson. “Learning Latent
71 Dynamics for Planning from Pixels”. In: *ICML*. 2019.
- 72 [Haf+20] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. “Dream to Control: Learning Behaviors by
73 Latent Imagination”. In: *ICLR*. 2020.
- 74 [Hag+17] M. Hagen, M. Potthast, M. Gohsen, A.
75 Rathgeber, and B. Stein. “A Large-Scale Query
76 Spelling Correction Corpus”. In: *Proceedings of the
77 40th International ACM SIGIR Conference on Re-
78 search and Development in Information Retrieval*.
79 SIGIR ’17. ACM, 2017, pp. 1261–1264.
- 80 [Háj08] A. Hájek. “Dutch Book Arguments”. In: *The
81 Oxford Handbook of Rational and Social Choice*. Ed.
82 by P. Anand, P. Pattanaik, and C. Puppe. Oxford Uni-
83 versity Press, 2008.
- 84 [Haj88] B. Hajek. “Cooling Schedules for Optimal An-
85 nealing”. In: *Math. Oper. Res.* 13.2 (1988), pp. 311–
86 329.
- 87 [Häl+21] H. Hälvä, S. L. Corff, L. Leh'ericq, J. So, Y.
88 Zhu, E. Gassiat, and A. Hyvärinen. “Disentangling
89 Identifiable Features from Noisy Data with Structured
90 Nonlinear ICA”. In: *NeurIPS*. 2021.
- 91 [Ham90] J. Hamilton. “Analysis of time series subject
92 to changes in regime”. In: *J. Econometrics* 45 (1990),
93 pp. 39–70.
- 94 [Han16] N. Hansen. “The CMA Evolution Strategy: A
95 Tutorial”. In: (2016). arXiv: 1604.00772 [cs.LG].
- 96 [Han+20] K. Han et al. “A Survey on Vision Trans-
97 former”. In: (2020). arXiv: 2012.12556 [cs.CV].
- 98 [Han80] T. S. Han. “Multiple mutual informations and
99 multiple interactions in frequency data”. In: *Informa-
100 tion and Control* 46.1 (1980), pp. 26–45.
- 101 [Har+17] J. Hartford, G. Lewis, K. Leyton-Brown, and
102 M. Taddy. “Deep IV: A flexible approach for counter-
103 factual prediction”. In: *Proceedings of the 34th Interna-
104 tional Conference on Machine Learning-Volume 70*. JMLR.org, 2017, pp. 1414–1423.
- 105 [Har18] K. Hartnett. “To Build Truly Intelligent Ma-
106 chines, Teach Them Cause and Effect”. In: *Quanta
107 Magazine* (2018).

- 1
- 2 [Har+22] W. Harvey, S. Naderiparizi, V. Masrani, C.
3 Weilbach, and F. Wood. “Flexible Diffusion Modeling
4 of Long Videos”. In: (May 2022). arXiv: 2205.11495
[cs.CV].
- 5 [Har90] A. C. Harvey. *Forecasting, Structural Time
Series Models, and the Kalman Filter*. Cambridge
University Press, 1990.
- 6 [Has10] H. van Hasselt. “Double Q-learning”. In: *NIPS*.
Ed. by J. D. Lafferty, C. K. I. Williams, J. Shawe-
Taylor, R. S. Zemel, and A. Culotta. Curran Associates,
Inc., 2010, pp. 2613–2621.
- 7 [Has70] W. Hastings. “Monte Carlo Sampling Meth-
10 ods Using Markov Chains and Their Applications”. In:
Biometrika 57.1 (1970), pp. 97–109.
- 8 [Hau+10] J. R. Hauser, O. Toubia, T. Evgeniou, R. Be-
furt, and D. Dzyabura. “Disjunctions of conjunctions,
cognitive simplicity, and consideration sets”. In: *Jour-
nal of Marketing Research* 47.3 (2010), pp. 485–496.
- 9 [Hau+11] M. Hauschild, M. Pelikan, M. Hauschild, and
M. Pelikan. “An introduction and survey of estimation
15 of distribution algorithms”. In: *Swarm and Evolution-
ary Computation*. 2011.
- 10 [Haw71] A. G. Hawkes. “Point spectra of some mutu-
ally exciting point processes”. In: *Journal of the Royal
Statistical Society: Series B (Methodological)* 33.3
19 (1971), pp. 438–443.
- 11 [HBB10] M. Hoffman, D. Blei, and F. Bach. “Online
learning for latent Dirichlet allocation”. In: *NIPS*.
2010.
- 12 [HBW19] E. Hoogeboom, R. van den Berg, and M.
Welling. “Emerging Convolutions for Generative Nor-
malizing Flows”. In: *ICML*. 2019.
- 13 [HC93] G. Hinton and D. V. Camp. “Keeping Neu-
ral Networks Simple by Minimizing the Description
Length of the Weights”. In: *in Proc. of the 6th
Ann. ACM Conf. on Computational Learning Theory*. ACM Press, 1993, pp. 5–13.
- 14 [HCG20] S. Huang, Y. Cao, and R. Grosse. “Evaluating
Lossy Compression Rates of Deep Generative Models”.
In: *ICML*. 2020.
- 15 [HCL06] R. Hadsell, S. Chopra, and Y. LeCun. “Di-
mensionality Reduction by Learning an Invariant Map-
ping”. In: *2006 IEEE Computer Society Conference
on Computer Vision and Pattern Recognition (CVPR’06)* 2 (2006), pp. 1735–1742.
- 16 [HD19] D. Hendrycks and T. Dietterich. “Benchmark-
ing Neural Network Robustness to Common Corrup-
tions and Perturbations”. In: *ICLR*. 2019.
- 17 [HDL17] D. Ha, A. M. Dai, and Q. V. Le. “HyperNet-
works”. In: *ICLR*. 2017.
- 18 [He+16a] K. He, X. Zhang, S. Ren, and J. Sun. “Deep
Residual Learning for Image Recognition”. In: *CVPR*.
2016.
- 19 [HE16a] J. Ho and S. Ermon. “Generative adversarial
imitation learning”. In: *Proceedings of the 30th Interna-
tional Conference on Neural Information Process-
ing Systems*. 2016, pp. 4572–4580.
- 20 [He+16b] K. He, X. Zhang, S. Ren, and J. Sun. “Iden-
tity Mappings in Deep Residual Networks”. In: *ECCV*.
2016.
- 21 [HE16b] J. Ho and S. Ermon. “Generative Adversarial
Imitation Learning”. In: *NIPS*. 2016, pp. 4565–4573.
- 22 [HE18] D. Ha and D. Eck. “A Neural Representation
of Sketch Drawings”. In: *ICLR*. 2018.
- 23 [He+19] J. He, D. Spokoyny, G. Neubig, and T. Berg-
Kirkpatrick. “Lagging Inference Networks and Poste-
rior Collapse in Variational Autoencoders”. In: *ICLR*.
2019.
- 24 [He+20] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick.
“Momentum contrast for unsupervised visual represen-
tation learning”. In: *CVPR*. 2020, pp. 9729–9738.
- 25 [He+21] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R.
Girshick. *Masked Autoencoders Are Scalable Vision
Learners*. 2021. arXiv: 2111.06377 [cs.CV].
- 26 [He+22] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and
R. Girshick. “Masked autoencoders are scalable vision
learners”. In: *Proceedings of the IEEE/CVF Confer-
ence on Computer Vision and Pattern Recognition*.
2022, pp. 16000–16009.
- 27 [Heg06] P. Heggernes. “Minimal triangulations of
graphs: A survey”. In: *Discrete Math.* 306.3 (2006),
pp. 297–317.
- 28 [Hel17] J. Helske. “KFAS: Exponential Family State
Space Models in R”. In: *J. Stat. Softw.* (2017).
- 29 [Hen+15] J. Hensman, A. Matthews, M. Filippone,
and Z. Ghahramani. “MCMC for Variationally Sparse
Gaussian Processes”. In: *NIPS*. 2015, pp. 1648–1656.
- 30 [Hen+16] L. A. Hendricks, Z. Akata, M. Rohrbach, J.
Donahue, B. Schiele, and T. Darrell. “Generating vis-
ual explanations”. In: *European conference on com-
puter vision*. Springer, 2016, pp. 3–19.
- 31 [Hen+18] G. E. Henter, J. Lorenzo-Trueba, X. Wang,
and J. Yamagishi. “Deep Encoder-Decoder Models for
Unsupervised Learning of Controllable Speech Synthe-
sis”. In: (2018). arXiv: 1807.11470 [eess.AS].
- 32 [Hen+19a] O. J. Henaff, A. Razavi, C. Doersch, S. M.
Ali Esfandi, and A. van den Oord. “Data-Efficient Im-
age Recognition with Contrastive Predictive Coding”.
In: *arXiv [cs.CV]* (2019).
- 33 [Hen+19b] D. Hendrycks, S. Basart, M. Mazeika, A.
Zou, J. Kwon, M. Mostajabi, J. Steinhhardt, and D.
Song. “Scaling Out-of-Distribution Detection for Real-
World Settings”. In: (2019). arXiv: 1911.11132 [cs.CV].
- 34 [Hen+20] D. Hendrycks*, N. Mu*, E. D. Cubuk, B.
Zoph, J. Gilmer, and B. Lakshminarayanan. “AugMix:
A Simple Data Processing Method to Improve Robust-
ness and Uncertainty”. In: *ICLR*. 2020.
- 35 [Hen+21] C. Henning, M. R. Cervera, F. D’Angelo, J.
von Oswald, R. Traber, B. Ehret, S. Kobayashi, B. F.
Gewe, and J. Sacramento. “Posterior Meta-Replay for
Continual Learning”. In: *NIPS*. 2021.
- 36 [Hes00] T. Heskes. “On ‘Natural’ Learning and Prun-
ing in Multilayered Perceptrons”. In: *Neural Comput.*
12.4 (2000), pp. 881–901.
- 37 [Hes+18] M. Hessel, J. Modayil, H. van Hasselt, T.
Schaul, G. Ostrovski, W. Dabney, D. Horgan, B.
Piot, M. Azar, and D. Silver. “Rainbow: Combining
Improvements in Deep Reinforcement Learning”. In:
AAAI. 2018.
- 38 [Heu+17a] M. Heusel, H. Ramsauer, T. Unterthiner, B.
Nessler, and S. Hochreiter. “GANs Trained by a Two
Time-Scale Update Rule Converge to a Local Nash
Equilibrium”. In: *NIPS*. 2017.
- 39 [Heu+17b] M. Heusel, H. Ramsauer, T. Unterthiner, B.
Nessler, and S. Hochreiter. “Gans trained by a two
time-scale update rule converge to a local nash equilib-
rium”. In: *Advances in neural information processing
systems*. 2017, pp. 6626–6637.

- [HFL13] J. Hensman, N. Fusi, and N. D. Lawrence. “Gaussian Processes for Big Data”. In: *UAI*. 2013.
- [HFM17] D. W. Hogg and D. Foreman-Mackey. “Data analysis recipes: Using Markov Chain Monte Carlo”. In: (2017). arXiv: [1710.06068 \[astro-ph.IM\]](#).
- [HG12] D. I. Hastie and P. J. Green. “Model Choice using Reversible Jump Markov Chain Monte Carlo”. In: *Statistica Neerlandica* 66 (2012), pp. 309–338.
- [HG14] M. D. Hoffman and A. Gelman. “The No-U-Turn Sampler: Adaptively Setting Path Lengths in Hamiltonian Monte Carlo”. In: *JMLR* 15 (2014), pp. 1593–1623.
- [HG16] D. Hendrycks and K. Gimpel. “Gaussian Error Linear Units (GELUs)”. In: *arXiv /cs.LG/* (2016).
- [HGMG18] J. Hron, A. G. de G. Matthews, and Z. Ghahramani. “Variational Bayesian dropout: pitfalls and fixes”. In: *ICML*. 2018.
- [HGS16] H. van Hasselt, A. Guez, and D. Silver. “Deep Reinforcement Learning with Double Q-Learning”. In: *AAAI. AAAI'16*. AAAI Press, 2016, pp. 2094–2100.
- [HH06] C. Holmes and L. Held. “Bayesian auxiliary variable models for binary and multinomial regression”. In: *Bayesian Analysis* 1.1 (2006), pp. 145–168.
- [HHC12] J. Hu, P. Hu, and H. S. Chang. “A Stochastic Approximation Framework for a Class of Randomized Optimization Algorithms”. In: *IEEE Trans. Automatic Control* 57.1 (2012).
- [HH09] A. Hyvärinen, J. Hurri, and P. Hoyer. *Natural Image Statistics: a probabilistic approach to early computational vision*. Springer, 2009.
- [HHK19] M. Haukemann, F. A. Hamprecht, and M. Kandemir. “Sampling-Free Variational Inference of Bayesian Neural Networks by Variance Backpropagation”. In: *UAI*. 2019.
- [HHLB11] F. Hutter, H. H. Hoos, and K. Leyton-Brown. “Sequential Model-Based Optimization for General Algorithm Configuration”. In: *Intl. Conf. on Learning and Intelligent Optimization (LION)*. 2011, pp. 507–523.
- [HHLMF18] M. Havasi, J. M. Hernández-Lobato, and J. J. Murillo-Fuentes. “Inference in Deep Gaussian Processes using Stochastic Gradient Hamiltonian Monte Carlo”. In: (2018). arXiv: [1806.05490 \[stat.ML\]](#).
- [Hig+17a] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. “beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In: *ICLR*. 2017.
- [Hig+17b] I. Higgins, L. Matthey, A. Pal, C. P. Burgess, X. Glorot, M. M. Botvinick, S. Mohamed, and A. Lerchner. “beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework”. In: *ICLR*. 2017.
- [Hin02] G. E. Hinton. “Training products of experts by minimizing contrastive divergence”. en. In: *Neural Computation* 14.8 (2002), pp. 1771–1800.
- [Hin10] G. Hinton. *A Practical Guide to Training Restricted Boltzmann Machines*. Tech. rep. U. Toronto, 2010.
- [Hin14] G. Hinton. *Lecture 6e on neural networks (RMSPROP: Divide the gradient by a running average of its recent magnitude)*. 2014.
- [Hin+95] G. E. Hinton, P. Dayan, B. J. Frey, and R. M. Neal. “The “wake-sleep” algorithm for unsupervised neural networks”. en. In: *Science* 268.5214 (1995), pp. 1158–1161.
- [HIY19] K. Hayashi, M. Imaizumi, and Y. Yoshida. “On Random Subsampling of Gaussian Process Regression: A Graphon-Based Analysis”. In: (2019). arXiv: [1901.09541 \[stat.ML\]](#).
- [HJ12] T. Hazan and T. Jaakkola. “On the Partition Function and Random Maximum A-Posteriori Perturbations”. In: *ICML*. June 2012.
- [HJ20] J. Huang and N. Jiang. “From Importance Sampling to Doubly Robust Policy Gradient”. In: *ICML*. 2020.
- [HJA20] J. Ho, A. Jain, and P. Abbeel. “Denoising Diffusion Probabilistic Models”. In: *NIPS*. 2020.
- [Hje+18] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. “Learning deep representations by mutual information estimation and maximization”. In: *International Conference on Learning Representations*. 2018.
- [Hjo+10] N. Hjort, C. Holmes, P. Müller, and S. Walker, eds. *Bayesian Nonparametrics*. Cambridge, 2010.
- [HJT18] M. D. Hoffman, M. J. Johnson, and D. Tran. “Autoconj: Recognizing and Exploiting Conjugacy Without a Domain-Specific Language”. In: *NIPS*. 2018.
- [HKO22] P. Hennig, H. Kersting, and M. Osborne. *Probabilistic Numerics: Computation as Machine Learning*. 2022.
- [HKP91] J. Hertz, A. Krogh, and R. G. Palmer. *An Introduction to the Theory of Neural Computation*. Addison-Wesley, 1991.
- [HKZ12] D. Hsu, S. Kakade, and T. Zhang. “A spectral algorithm for learning hidden Markov models”. In: *J. of Computer and System Sciences* 78.5 (2012), pp. 1460–1480.
- [HL04] D. R. Hunter and K. Lange. “A Tutorial on MM Algorithms”. In: *The American Statistician* 58 (2004), pp. 30–37.
- [HL+16a] J. Hernandez-Lobato, Y. Li, M. Rowland, T. Bui, D. Hernandez-Lobato, and R. Turner. “Black-Box Alpha Divergence Minimization”. en. In: *ICML*. 2016, pp. 1511–1520.
- [HL+16b] M. Hernandez-Lobato, M. A. Gelbart, R. P. Adams, M. W. Hoffman, and Z. Ghahramani. “A General Framework for Constrained Bayesian Optimization using Information-based Search”. In: *JMLR* (2016).
- [HL20] X. Hu and J. Lei. “A Distribution-Free Test of Covariate Shift Using Conformal Prediction”. In: (2020). arXiv: [2010.07147 \[stat.ME\]](#).
- [HLA15a] J. M. Hernández-Lobato and R. P. Adams. “Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks”. In: *ICML*. 2015.
- [HLA15b] J. M. Hernández-Lobato and R. P. Adams. “Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks”. In: *ICML*. 2015.
- [HLA19] J. Ho, E. Lohn, and P. Abbeel. “Compression with Flows via Local Bits-Back Coding”. In: *NeurIPS*. 2019.
- [HLC19] Y.-P. Hsieh, C. Liu, and V. Cevher. “Finding mixed nash equilibria of generative adversarial networks”. In: *International Conference on Machine Learning*. 2019, pp. 2810–2819.

- 1
- 2 [HLHG14] J. Hernandez-Lobato, M. W. Hoffman, and
Z. Ghahramani. "Predictive entropy search for efficient
global optimization of black-box functions". In: *NIPS*.
2014.
- 3
- 4 [HLR16] K. Hofmann, L. Li, and F. Radlinski. "On-
line Evaluation for Information Retrieval". In: *Founda-
tions and Trends in Information Retrieval* 10.1
(2016), pp. 1–117.
- 5
- 6 [HLRW14] J. R. Hershey, J. Le Roux, and F. Weninger.
"Deep Unfolding: Model-Based Inspiration of Novel
Deep Architectures". In: (2014). arXiv: [1409 . 2574 \[cs.LG\]](https://arxiv.org/abs/1409.2574).
- 7
- 8 [HLS03] R. Herbrich, N. D. Lawrence, and M. Seeger.
"Fast Sparse Gaussian Process Methods: The Infor-
mative Vector Machine". In: *NIPS*. MIT Press, 2003,
pp. 625–632.
- 9
- 10 [HM81] R. Howard and J. Matheson. "Influence dia-
grams". In: *Readings on the Principles and Appli-
cations of Decision Analysis, volume II*. Ed. by R.
Howard and J. Matheson. Strategic Decisions Group,
1981.
- 11
- 12 [HMD18] J. C. Higuera, D Meger, and G Dudek. "Syn-
thesizing Neural Network Controllers with Probabilistic
Model-Based Reinforcement Learning". In: *IROS*.
2018, pp. 2538–2544.
- 13
- 14 [HMD19] D. Hendrycks, M. Mazeika, and T. Dietterich.
"Deep Anomaly Detection with Outlier Exposure". In:
ICLR. 2019.
- 15
- 16 [HMK04] D. Heckerman, C. Meek, and D. Koller. *Prob-
abilistic Models for Relational Data*. Tech. rep. MSR-
TR-2004-30. Microsoft Research, 2004.
- 17
- 18 [HNBK18] J. He, G. Neubig, and T. Berg-Kirkpatrick.
"Unsupervised Learning of Syntactic Structure with
Invertible Neural Projections". In: *EMNLP*. 2018.
- 19
- 20 [HNP09] A. Halevy, P. Norvig, and F. Pereira. "The un-
reasonable effectiveness of data". In: *IEEE Intelligent
Systems* 24.2 (2009), pp. 8–12.
- 21
- 22 [HO00] A. Hyvärinen and E. Oja. "Independent compo-
nent analysis: algorithms and applications". In: *Neural
Networks* 13 (2000), pp. 411–430.
- 23
- 24 [Ho+21] J. Ho, C. Saharia, W. Chan, D. J. Fleet, M.
Norouzi, and T. Salimans. "Cascaded Diffusion Models
for High Fidelity Image Generation". In: (May 2021).
arXiv: [2106 . 15282 \[cs.CV\]](https://arxiv.org/abs/2106.15282).
- 25
- 26 [Ho+22] J. Ho, T. Salimans, A. Gritsenko, W. Chan, M.
Norouzi, and D. J. Fleet. "Video Diffusion Models". In:
(Apr. 2022). arXiv: [2204 . 03458 \[cs.CV\]](https://arxiv.org/abs/2204.03458).
- 27
- 28 [HO48] C. G. Hempel and P. Oppenheim. "Studies in
the Logic of Explanation". In: *Philosophy of science*
15.2 (1948), pp. 135–175.
- 29
- 30 [Hob69] A. Hobson. "A new theorem of information
theory". In: *Journal of Statistical Physics* 1.3 (1969),
pp. 383–391.
- 31
- 32 [Hoe12] J. M. ver Hoef. "Who invented the delta
method?" In: *The American Statistician* 66.2 (2012),
pp. 124–127.
- 33
- 34 [Hoe+21] T. Hoefer, D. Alistarh, T. Ben-Nun, N. Dry-
den, and A. Peste. "Sparsity in Deep Learning: Prun-
ing and growth for efficient inference and training
in neural networks". In: (2021). arXiv: [2102 . 00554 \[cs.LG\]](https://arxiv.org/abs/2102.00554).
- 35
- 36 [Hof09] P. D. Hoff. *A First Course in Bayesian Sta-
tistical Methods*. Springer, 2009.
- 37
- 38 [Hof+13] M. D. Hoffman, D. M. Blei, C. Wang, and J.
Paisley. "Stochastic Variational Inference". In: *JMLR*
14 (2013), pp. 1303–1347.
- 39
- 40 [Hof17] M. D. Hoffman. "Learning Deep Latent Gaus-
sian Models with Markov Chain Monte Carlo". In:
ICML. 2017, pp. 1510–1519.
- 41
- 42 [Hof+18] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu,
P. Isola, K. Saenko, A. Efros, and T. Darrell. "Cy-
cada: Cycle-consistent adversarial domain adaptation".
In: *International conference on machine learning*.
PMLR. 2018, pp. 1989–1998.
- 43
- 44 [Hof+19] M. Hoffman, P. Sountsov, J. V. Dillon, I.
Langmore, D. Tran, and S. Vasudevan. "NeuTra-lizing
Bad Geometry in Hamiltonian Monte Carlo Using
Neural Transport". In: (2019). arXiv: [1903 . 03704 \[stat.CO\]](https://arxiv.org/abs/1903.03704).
- 45
- 46 [Hof99] T. Hofmann. "Probabilistic latent semantic in-
dexing". In: *Research and Development in Informa-
tion Retrieval* (1999), pp. 50–57.
- 47
- 48 [Hoh+20] F. Hohman, M. Conlen, J. Heer, and D. H. P.
Chau. "Communicating with interactive articles". In:
Distill 5.9 (2020), e28.
- 49
- 50 [Hol86] P. W. Holland. "Statistics and Causal Infer-
ence". In: *JASA* 81.396 (1986), pp. 945–960.
- 51
- 52 [Hol92] J. H. Holland. *Adaptation in Natural and Ar-
tificial Systems*. https://mitpress.mit.edu/books/adaptation_natural_and_artificial_systems. Accessed: 2017-11-26. 1992.
- 53
- 54 [Hon+10] A. Honkela, T. Raiko, M. Kuusela, M.
Tornio, and J. Karhunen. "Approximate Riemannian
Conjugate Gradient Learning for Fixed-Form Varia-
tional Bayes". In: *JMLR* 11.Nov (2010), pp. 3235–
3268.
- 55
- 56 [Hop82] J. J. Hopfield. "Neural networks and physical
systems with emergent collective computational abili-
ties". In: *PNAS* 79.8 (1982), 2554–2558.
- 57
- 58 [Hor+05] E. Horvitz, J. Apacible, R. Sarin, and L. Liao.
"Prediction, Expectation, and Surprise: Methods, De-
signs, and Study of a Deployed Traffic Forecasting Ser-
vice". In: *UAI*. 2005.
- 59
- 60 [Hor61] P Horst. "Generalized canonical correlations
and their applications to experimental data". en. In:
J. Clin. Psychol. 17 (1961), pp. 331–347.
- 61
- 62 [Hos+20a] T. Hospedales, A. Antoniou, P. Micaelli,
and A. Storkey. "Meta-Learning in Neural Networks:
A Survey". In: (2020). arXiv: [2004 . 05439 \[cs.LG\]](https://arxiv.org/abs/2004.05439).
- 63
- 64 [Hos+20b] R. Hostettler, F. Tronarp, Á. F. García-
Fernández, and S. Särkkä. "Importance Densities for
Particle Filtering Using Iterated Conditional Expec-
tations". In: *IEEE Signal Process. Lett.* 27 (2020),
pp. 211–215.
- 65
- 66 [HOT06a] G. Hinton, S. Osindero, and Y. Teh. "A fast
learning algorithm for deep belief nets". In: *Neural
Computation* 18 (2006), pp. 1527–1554.
- 67
- 68 [HOT06b] G. E. Hinton, S. Osindero, and Y. W. Teh.
"A Fast Learning Algorithm for Deep Belief Nets". In:
Neural Computation 18 (2006), pp. 1527–1554.
- 69
- 70 [Hot36] H. Hotelling. "Relations Between Two Sets of
Variates". In: *Biometrika* 28.3/4 (1936), pp. 321–377.
- 71
- 72 [Hou+12] N. Houlsby, F. Huszar, Z. Ghahramani,
and J. M. Hernández-lobato. "Collaborative Gaussian
Processes for Preference Learning". In: *NIPS*. 2012,
pp. 2096–2104.

- 1 [Hou+19] N. Houldsby, A. Giurgiu, S. Jastrzebski, B.
2 Morrone, Q. De Laroussilhe, A. Gesmundo, M. At-
3 tariyan, and S. Gelly. “Parameter-efficient transfer
4 learning for NLP”. In: *International Conference on
Machine Learning*. PMLR. 2019, pp. 2790–2799.
- 5 [HOW11] P. Hall, J. T. Ormerod, and M. P. Wand.
6 “Theory of Gaussian Variational Approximation for
7 a Generalised Linear Mixed Model”. In: *Statistica
Sinica* 21 (2011), pp. 269–389.
- 8 [HP10] J. Hacker and P. Pierson. *Winner-Take-All
Politics: How Washington Made the Rich Richer —
and Turned Its Back on the Middle Class*. Simon &
Schuster, 2010.
- 9 [HPR19] C. Herzog né Hoffmann, E. Petersen, and P.
10 Rostalski. “Iterative Approximate Nonlinear Inference
11 via Gaussian Message Passing on Factor Graphs”. In:
12 *IEEE Control Systems Letters* 3.4 (2019), pp. 978–
983.
- 13 [HR17] C. Hoffmann and P. Rostalski. “Linear Optimal
14 Control on Factor Graphs — A Message Passing Per-
15 spective”. In: *Intl. Federation of Automatic Control*
50.1 (2017), pp. 6314–6319.
- 16 [HR20a] M. Hernan and J. Robins. *Causal Inference:
What If*. CRC Press, 2020.
- 17 [HR20b] M. Hernán and J. Robins. *Causal Inference:
What If*. Boca Raton: Chapman & Hall/CRC., 2020.
- 18 [HS05] H. Hoos and T. Stutzle. *Stochastic local search:
Foundations and applications*. Morgan Kauffman,
2005.
- 19 [HS06a] G. Hinton and R. Salakhutdinov. “Reducing
the dimensionality of data with neural networks”. In:
Science 313.5786 (2006), pp. 504–507.
- 20 [HS06b] G. E. Hinton and R. R. Salakhutdinov. “Re-
ducing the dimensionality of data with neural net-
works”. In: *science* 313.5786 (2006), pp. 504–507.
- 21 [HS09] M. Heaton and J. Scott. *Bayesian computa-
tion and the linear model*. Tech. rep. Duke, 2009.
- 22 [HS12] P. Hennig and C. Schuler. “Entropy search for
information-efficient global optimization”. In: *JMLR*
13 (2012), pp. 1809–1837.
- 23 [HS13] J. Y. Hsu and D. S. Small. “Calibrating Sen-
sitivity Analyses to Observed Covariates in Observa-
tional Studies”. In: *Biometrics* 69.4 (2013), pp. 803–
811. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/biom.12101>.
- 24 [HS18] D. Ha and J. Schmidhuber. “World Models”. In:
NIPS. 2018.
- 25 [HS21] J. Ho and T. Salimans. “Classifier-Free Diffu-
sion Guidance”. In: *NIPS Workshop on Deep Genera-
tive Models and Downstream Applications*. 2021.
- 26 [HS88] D. S. Hochbaum and D. B. Shmoys. “A poly-
nomial approximation scheme for scheduling on uniform
processors: Using the dual approximation approach”. In:
SICOMP. 1988.
- 27 [HS97] S Hochreiter and J Schmidhuber. “Flat minima”.
en. In: *Neural Comput.* 9.1 (1997), pp. 1–42.
- 28 [HSDK12] C. Hillar, J. Sohl-Dickstein, and K. Koepsell.
29 *Efficient and Optimal Binary Hopfield Associative
Memory Storage Using Minimum Probability Flow*.
Tech. rep. 2012. arXiv: [1204.2916](https://arxiv.org/abs/1204.2916).
- 30 [HSG06] J. D. Hol, T. B. Schon, and F. Gustafsson.
31 “On Resampling Algorithms for Particle Filters”. In:
32 *IEEE Nonlinear Statistical Signal Processing Work-
shop*. 2006, pp. 79–82.
- 33 [HSGF21] S. Hassan, S. Sarkka, and A. F. Garcia-
Fernandez. “Temporal Parallelization of Inference in
Hidden Markov Models”. In: *IEEE Trans. Signal Pro-
cessing* 69 (2021), pp. 4875–4887.
- 34 [Hsu+18] Y.-C. Hsu, Y.-C. Liu, A. Ramasamy, and
Z. Kira. “Re-evaluating Continual Learning Scenarios:
A Categorization and Case for Strong Baselines”. In:
NIPS Continual Learning Workshop. 2018.
- 35 [HT01] G. E. Hinton and Y. Teh. “Discovering multi-
ple constraints that are frequently approximately sat-
isfied”. In: *UAI*. 2001.
- 36 [HT09] H. Hoefling and R. Tibshirani. “Estimation
of Sparse Binary Pairwise Markov Networks using
Pseudo-likelihoods”. In: *JMLR* 10 (2009).
- 37 [HT15] J. H. Huggins and J. B. Tenenbaum. “Risk and
regret of hierarchical Bayesian learners”. In: *ICML*.
2015.
- 38 [HT17] T. J. Hastie and R. J. Tibshirani. *Generalized
additive models*. Routledge, 2017.
- 39 [HTF09] T. Hastie, R. Tibshirani, and J. Friedman.
The Elements of Statistical Learning. 2nd edition.
Springer, 2009.
- 40 [HTW15] T. Hastie, R. Tibshirani, and M. Wainwright.
*Statistical Learning with Sparsity: The Lasso and
Generalizations*. CRC Press, 2015.
- 41 [Hu+00] M. Hu, C. Ingram, M. Sirski, C. Pal, S. Swamy,
and C. Patten. *A Hierarchical HMM Implementation
for Vertebrate Gene Splice Site Prediction*. Tech. rep.
Dept. Computer Science, Univ. Waterloo, 2000.
- 42 [Hu+12] J. Hu, Y. Wang, E. Zhou, M. C. Fu, and S. I.
Marcus. “A Survey of Some Model-Based Methods
for Global Optimization”, en. In: *Optimization, Con-
trol, and Applications of Stochastic Systems. Systems
& Control: Foundations & Applications*. Birkhäuser,
Boston, 2012, pp. 157–179.
- 43 [Hu+17] W. Hu, C. J. Li, L. Li, and J.-G. Liu. “On the
diffusion approximation of nonconvex stochastic gradi-
ent descent”. In: (2017). arXiv: [1705.07562 \[stat.ML\]](https://arxiv.org/abs/1705.07562).
- 44 [Hu+18] W. Hu, G. Niu, I. Sato, and M. Sugiyama.
“Does Distributionally Robust Supervised Learning
Give Robust Classifiers?” In: *ICML*. 2018.
- 45 [Hua+17a] G. Huang, Y. Li, G. Pleiss, Z. Liu, J.
Hopcroft, and K. Weinberger. “Snapshot ensembles:
train 1, get *M* for free”. In: *ICLR*. 2017.
- 46 [Hua+17b] Y. Huang, Y. Zhang, N. Li, Z. Wu, and
J. A. Chambers. “A Novel Robust Student’s t-Based
Kalman Filter”. In: *IEEE Trans. Aerosp. Electron.
Syst.* 53.3 (2017), pp. 1545–1554.
- 47 [Hua+18a] C.-Z. A. Huang, A. Vaswani, J. Uszkoreit,
N. Shazeer, I. Simon, C. Hawthorne, A. M. Dai, M. D.
Hoffman, M. Dinculescu, and D. Eck. “Music Trans-
former”. In: (2018). arXiv: [1809.04281 \[cs.LG\]](https://arxiv.org/abs/1809.04281).
- 48 [Hua+18b] C.-W. Huang, D. Krueger, A. Lacoste,
and A. Courville. “Neural Autoregressive Flows”. In:
ICML. 2018.
- 49 [Hua+19] Y. Huang, Y. Zhang, Y. Zhao, and J. A.
Chambers. “A Novel Robust Gaussian-Student’s t
Mixture Distribution Based Kalman Filter”. In: *IEEE
Trans. Signal Process.* 67.13 (2019), pp. 3606–3620.
- 50 [Hug+19] J. H. Huggins, T. Campbell, M. Kasprzak,
and T. Broderick. “Scalable Gaussian Process Infer-
ence with Finite-data Mean and Variance Guarantees”.
In: *AISTATS*. 2019.
- 51 [Hug+20] J. Huggins, M. Kasprzak, T. Campbell, and
T. Broderick. “Validated Variational Inference via

- 1 Practical Posterior Error Bounds". In: *AISTATS*. Ed.
 2 by S. Chiappa and R. Calandra. Vol. 108. Proceedings
 3 of Machine Learning Research. PMLR, 2020, pp. 1792–
 1802.
- 4 [Hus17a] F. Huszár. *Is Maximum Likelihood Useful
 5 for Representation Learning?* 2017.
- 6 [Hus17b] F. Huszár. "Variational inference us-
 7 ing implicit distributions". In: *arXiv preprint
 arXiv:1702.08235* (2017).
- 8 [Hut89] M. F. Hutchinson. "A stochastic estimator
 9 of the trace of the influence matrix for Laplacian
 10 smoothing splines". In: *Communications in Statistics-
 Simulation and Computation* 18.3 (1989), pp. 1059–
 1076.
- 11 [HVD14] G. Hinton, O. Vinyals, and J. Dean. "Distill-
 12 ing the Knowledge in a Neural Network". In: *NIPS DL
 workshop*. 2014.
- 13 [HW13] A. Huang and M. P. Wand. "Simple Marginally
 14 Noninformative Prior Distributions for Covariance
 15 Matrices". en. In: *Bayesian Analysis* 8.2 (2013),
 pp. 439–452.
- 16 [HXW17] H. He, B. Xin, and D. Wipf. "From Bayesian
 17 Sparsity to Gated Recurrent Nets". In: *NIPS*. 2017.
- 18 [HY01] M. Hansen and B. Yu. "Model selection and the
 19 principle of minimum description length". In: *JASA*
 (2001).
- 20 [Hyv05] A. Hyvärinen. "Estimation of non-normalized
 21 statistical models by score matching". In: *JMLR* 6.Apr
 (2005), pp. 695–709.
- 22 [Hyv07a] A. Hyvärinen. "Connections between score
 23 matching, contrastive divergence, and pseudolikeli-
 24 hood for continuous-valued variables". In: *IEEE
 Transactions on neural networks* 18.5 (2007),
 pp. 1529–1531.
- 25 [Hyv07b] A. Hyvärinen. "Some extensions of score
 26 matching". In: *Computational statistics & data anal-
 27 ysis* 51.5 (2007), pp. 2499–2512.
- 28 [IB12] R. Iyer and J. Bilmes. "Algorithms for Approx-
 29 imate Minimization of the Difference Between Submod-
 30 ular Functions, with Applications". In: *Uncertainty in
 Artificial Intelligence (UAI)*. Catalina Island, USA:
 AUAI, 2012.
- 31 [IB13] R. Iyer and J. Bilmes. "Submodular Optimiza-
 32 tion with Submodular Cover and Submodular Knap-
 sack Constraints". In: *Neural Information Processing
 Society (NeurIPS, formerly NIPS)*. Lake Tahoe, CA,
 2013.
- 34 [IB15] R. K. Iyer and J. A. Bilmes. "Polyhedral as-
 35 pects of Submodularity, Convexity and Concavity". In:
Arxiv, CoRR abs/1506.07329 (2015).
- 37 [IB98] M. Isard and A. Blake. "CONDENSATION -
 38 conditional density propagation for visual tracking".
 In: *Intl. J. of Computer Vision* 29.1 (1998), pp. 5–
 18.
- 39 [IBK06] E. L. Ionides, C Bretó, and A. A. King. "Infer-
 40 ence for nonlinear dynamical systems". en. In: *PNAS*
 103.49 (2006), pp. 18438–18443.
- 41 [IFF00] S. Iwata, L. Fleischer, and S. Fujishige. "A com-
 42 binatorial strongly polynomial algorithm for minimiz-
 43 ing submodular functions". In: *Journal of the ACM*
 (2000).
- 44 [IFW05] A. T. Ihler, J. W. Fischer III, and A. S. Will-
 45 sky. "Loopy Belief Propagation: Convergence and Ef-
 46 fects of Message Errors". In: *JMLR* 6 (2005), pp. 905–
 936.
- [IJ01] H. Ishwaran and L. F. James. "Gibbs sampling
 methods for stick-breaking priors". In: *Journal of the
 American Statistical Association* 96.453 (2001),
 pp. 161–173.
- [IJB13] R. Iyer, S. Jegelka, and J. Bilmes. "Curvature
 and Optimal Algorithms for Learning and Minimizing
 Submodular Functions". In: *NIPS*. Lake Tahoe, CA,
 2013.
- [IKB21] A. Immer, M. Korzepa, and M. Bauer. "Im-
 proving predictions of Bayesian neural nets via local
 linearization". In: *AISTATS*. Ed. by A. Banerjee and
 K. Fukumizu. Vol. 130. Proceedings of Machine Learn-
 ing Research. PMLR, 2021, pp. 703–711.
- [IM17] J. Ingraham and D. Marks. "Bayesian Sparsity
 for Intractable Undirected Models". In: *ICML*. 2017.
- [Imb03] G. Imbens. "Sensitivity to Exogeneity Assump-
 tions in Program Evaluation". In: *The American Eco-
 nomic Review* (2003).
- [Imb19] G. W. Imbens. "Potential Outcome and Di-
 rected Acyclic Graph Approaches to Causality: Rele-
 vance for Empirical Practice in Economics". In: (2019).
 arXiv: 1907.07271 [stat.ME].
- [Imm+21] A. Immer, M. Bauer, V. Fortuin, G. Rätsch,
 and M. E. Khan. "Scalable Marginal Likelihood Es-
 timation for Model Selection in Deep Learning". In:
ICML. 2021.
- [IN09] S. Iwata and K. Nagano. "Submodular function
 minimization under covering constraints". In: *Proced-
 ings of the 50th Annual IEEE Symposium on Foun-
 dations of Computer Science (FOCS)*. 2009, pp. 671–
 680.
- [Ing20] M. Ingram. *Deterministic ADVI in JAX* (blog post). <https://martiningeram.github.io/deterministic-advi/>. 2020.
- [INK18] P. Izmailov, A. Novikov, and D. Kropotov.
 "Scalable Gaussian Processes with Billions of Induc-
 ing Inputs via Tensor Train Decomposition". In: *ICML*.
 2018.
- [Inn20] M. Innes. "Sense & Sensitivities: The Path
 to General-Purpose Algorithmic Differentiation". In:
Proceedings of Machine Learning and Systems. Ed.
 by I. Dhillon, D. Papailiopoulos, and V. Sze. Vol. 2.
 2020, pp. 58–69.
- [IO09] S. Iwata and J. B. Orlin. "A simple combina-
 torial algorithm for submodular function minimiza-
 tion". In: *Proceedings of the twentieth annual ACM-
 SIAM symposium on Discrete algorithms*. SIAM.
 2009, pp. 1230–1237.
- [IR00] D. R. Insua and F. Ruggeri. *Robust Bayesian
 Analysis*. Springer, 2000.
- [IR15] G. Imbens and D. Rubin. *Causal Inference in
 Statistics, Social and Biomedical Sciences: An Intro-
 duction*. Cambridge University Press, 2015.
- [IS15] S. Ioffe and C. Szegedy. "Batch Normalization:
 Accelerating Deep Network Training by Reducing In-
 ternal Covariate Shift". In: *ICML*. 2015, pp. 448–456.
- [Isa03] M. Isard. "PAMPAS: Real-Valued Graphical
 Models for Computer Vision". In: *CVPR*. Vol. 1. 2003,
 p. 613.
- [Isl+19] R. Islam, R. Seraj, S. Y. Arnob, and D. Precup.
 "Doubly Robust Off-Policy Actor-Critic Algorithms
 for Reinforcement Learning". In: *NeurIPS Workshop
 on Safety and Robustness in Decision Making*. 2019.

- 1 [Iso+17] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros.
“Image-to-Image Translation with Conditional Adversarial Networks”. In: *CVPR*. 2017.
- 2 [IX00] K Ito and K Xiong. “Gaussian filters for non-linear filtering problems”. In: *IEEE Trans. Automat. Contr.* 45.5 (2000), pp. 910–927.
- 3 [Iye+21] R. Iyer, N. Khargonkar, J. Bilmes, and H. Asnani. “Generalized Submodular Information Measures: Theoretical Properties, Examples, Optimization Algorithms, and Applications”. In: *IEEE Transactions on Information Theory* (2021).
- 4 [Iza+15] H. Izadinia, B. C. Russell, A. Farhadi, M. D. Hoffman, and A. Hertzmann. “Deep classifiers from image tags in the wild”. In: *Proceedings of the 2015 Workshop on Community-Organized Multi-modal Mining: Opportunities for Novel Solutions*. 2015, pp. 13–18.
- 5 [Izm+18] P. Izmailov, D. Podoprikhin, T. Garipov, D. Vetrov, and A. G. Wilson. “Averaging Weights Leads to Wider Optima and Better Generalization”. In: *UAI*. 2018.
- 6 [Izm+19] P. Izmailov, W. J. Maddox, P. Kirichenko, T. Garipov, D. Vetrov, and A. G. Wilson. “Subspace Inference for Bayesian Deep Learning”. In: *UAI*. 2019.
- 7 [Izm+21a] P. Izmailov, P. Nicholson, S. Lotfi, and A. G. Wilson. “Dangers of Bayesian Model Averaging under Covariate Shift”. In: *NIPS*. 2021.
- 8 [Izm+21b] P. Izmailov, S. Vikram, M. D. Hoffman, and A. G. Wilson. “What Are Bayesian Neural Network Posteriors Really Like?” In: *ICML*. 2021.
- 9 [Jaa01] T. Jaakkola. “Tutorial on variational approximation methods”. In: *Advanced mean field methods*. Ed. by M. Opper and D. Saad. MIT Press, 2001.
- 10 [Jac+21] M. Jacobs, M. F. Pradier, T. H. McCoy, R. H. Perlis, F. Doshi-Velez, and K. Z. Gajos. “How machine-learning recommendations influence clinician treatment selections: the example of antidepressant selection”. In: *Translational psychiatry* 11.1 (2021), pp. 1–9.
- 11 [Jad+17] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. “Reinforcement Learning with Unsupervised Auxiliary Tasks”. In: *ICLR*. 2017.
- 12 [Jak21] K. Jakkala. “Deep Gaussian Processes: A Survey”. In: (2021). arXiv: [2106.12135 \[cs.LG\]](#).
- 13 [Jan+17] P. A. Jang, A. Loeb, M. Davidow, and A. G. Wilson. “Scalable Levy Process Priors for Spectral Kernel Learning”. In: *Advances in Neural Information Processing Systems*. 2017.
- 14 [Jan18] E. Jang. *Normalizing Flows Tutorial*. 2018.
- 15 [Jan+19] M. Janner, J. Fu, M. Zhang, and S. Levine. “When to Trust Your Model: Model-Based Policy Optimization”. In: *NIPS*. 2019.
- 16 [Jas] Jason Antic and Jeremy Howard and Uri Manor. *Decrappification, DeOldification, and Super Resolution (Blog post)*.
- 17 [Jay03] E. T. Jaynes. *Probability theory: the logic of science*. Cambridge university press, 2003.
- 18 [Jay+20] S. M. Jayakumar, W. M. Czarnecki, J. Menick, J. Schwarz, J. Rae, S. Osindero, Y. W. Teh, T. Harley, and R. Pascanu. “Multiplicative Interactions and Where to Find Them”. In: *ICLR*. 2020.
- 19 [JB03] A. Jakulin and I. Bratko. “Analyzing Attribute Dependencies”. In: *Proc. 7th European Conf. on Principles and Practice of Knowledge Discovery in Databases*. 2003.
- 20 [JB16] S. Jegelka and J. Bilmes. “Graph cuts with interacting edge weights: examples, approximations, and algorithms”. In: *Mathematical Programming* (2016), pp. 1–42.
- 21 [JB65] C. Jacobi and C. W. Borchardt. “De investigando ordine systematis aequationum differentialium vulgarium cuiuscunque.” In: *Journal für die reine und angewandte Mathematik* 1865.64 (1865), pp. 297–320.
- 22 [JBB09] D. Jian, A. Barthels, and M. Beetz. “Adaptive Markov logic networks: Learning statistical relational models with dynamic parameters”. In: *9th European Conf. on AI*. 2009, 937–942.
- 23 [Jef04] R. Jeffrey. *Subjective Probability: The Real Thing*. Cambridge, 2004.
- 24 [Jen+17] R. Jenatton, C. Archambeau, J. González, and M. Seeger. “Bayesian Optimization with Tree-structured Dependencies”. In: *ICML*. 2017, pp. 1655–1664.
- 25 [Jeu+19] O. Jeunen, D. Mykhaylov, D. Rohde, F. Vasile, A. Gilotte, and M. Bompaire. “Learning from Bandit Feedback: An Overview of the State-of-the-art”. In: (2019). arXiv: [1909.08471 \[cs.IR\]](#).
- 26 [JG20] A. Jacovi and Y. Goldberg. “Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness?” In: *arXiv preprint arXiv:2004.09685* (2020).
- 27 [JG21] A. Jacovi and Y. Goldberg. “Aligning faithful interpretations with their social attribution”. In: *Transactions of the Association for Computational Linguistics* 9 (2021), pp. 294–310.
- 28 [JGH18] A. Jacot, F. Gabriel, and C. Hongler. “Neural Tangent Kernel: Convergence and Generalization in Neural Networks”. In: *NIPS*. 2018.
- 29 [JGP17] E. Jang, S. Gu, and B. Poole. “Categorical Reparameterization with Gumbel-Softmax”. In: *ICLR*. 2017.
- 30 [Jia+19] R. Jia, A. Raghunathan, K. Göksel, and P. Liang. “Certified Robustness to Adversarial Word Substitutions”. In: *EMNLP*. 2019.
- 31 [Jia+21] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. V. Le, Y. Sung, Z. Li, and T. Duerig. *Scaling Up Visual and Vision-Language Representation Learning With Noisy Text Supervision*. 2021. arXiv: [2102.05918 \[cs.CV\]](#).
- 32 [Jia21] H. Jiang. “Minimizing convex functions with integral minimizers”. In: *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*. SIAM. 2021, pp. 976–985.
- 33 [Jih+12] Jihong Min, J Kim, Seunghak Shin, and I. S. Kweon. “Efficient Data-Driven MCMC sampling for vision-based 6D SLAM”. In: *ICRA*. 2012, pp. 3025–3032.
- 34 [Jin11] Y. Jin. “Surrogate-assisted evolutionary computation: Recent advances and future challenges”. In: *Swarm and Evolutionary Computation* 1.2 (2011), pp. 61–70.
- 35 [Jit+16] W. Jitkrittum, Z. Szabó, K. P. Chwialkowski, and A. Gretton. “Interpretable Distribution Features with Maximum Testing Power”. In: *NIPS*. Curran Associates, Inc., 2016, pp. 181–189.
- 36 [JJ00] T. S. Jaakkola and M. I. Jordan. “Bayesian parameter estimation via variational methods”. In: *Statistics and Computing* 10 (2000), pp. 25–37.
- 37 [JKG18] H. Jiang, B. Kim, and M. Y. Guan. “To Trust Or Not To Trust A Classifier”. In: *NIPS*. 2018.

- 1
- 2 [JKJ12] K. Jiang, B. Kulis, and M. Jordan. “Small-
3 variance asymptotics for exponential family Dirichlet
4 process mixture models”. In: *Advances in Neural In-*
formation Processing Systems 25 (2012), pp. 3158–
3166.
- 5 [JKK95] C. S. Jensen, A. Kong, and U. Kjaerulff.
“Blocking-Gibbs Sampling in Very Large Probabilistic
6 Expert Systems”. In: *Intl. J. Human-Computer Studies*
7 (1995), pp. 647–666.
- 8 [JL15a] V. Jalali and D. Leake. “CBR meets big data:
A case study of large-scale adaptation rule generation”.
In: *International Conference on Case-Based Reasoning*. Springer. 2015, pp. 181–196.
- 9 [JL15b] V. Jalali and D. B. Leake. “Enhancing case-
based regression with automatically-generated ensem-
10 bles of adaptations”. In: *Journal of Intelligent Infor-*
mation Systems 46 (2015), pp. 237–258.
- 11 [JL16] N. Jiang and L. Li. “Doubly Robust Off-policy
Evaluation for Reinforcement Learning”. In: *ICML*.
2016, pp. 652–661.
- 12 [JM00] D. Jurafsky and J. H. Martin. *Speech and lan-*
guage processing: An Introduction to Natural Lan-
guage Processing, Computational Linguistics, and
13 Speech Recognition. Prentice-Hall, 2000.
- 14 [JM08] D. Jurafsky and J. H. Martin. *Speech and lan-*
guage processing: An Introduction to Natural Lan-
guage Processing, Computational Linguistics, and
15 Speech Recognition. 2nd edition. Prentice-Hall, 2008.
- 16 [JM18] A. Jolicoeur-Martineau. “The relativistic dis-
criminator: a key element missing from standard
GAN”. In: *arXiv preprint arXiv:1807.00734* (2018).
- 17 [JM70] D. H. Jacobson and D. Q. Mayne. *Differential
Dynamic Programming*. Elsevier Press, 1970.
- 18 [JMW06] J. K. Johnson, D. M. Malioutov, and A. S.
Willsky. “Walk-sum interpretation and analysis of
Gaussian belief propagation”. In: *NIPS*. 2006, pp. 579–
586.
- 19 [JNJ20] C. Jin, P. Netrapalli, and M. I. Jordan. “What
is local optimality in nonconvex-nonconcave minimax
optimization?” In: *Proceedings of the 34th Interna-*
tional Conference on Machine Learning- Volume 73.
2020.
- 20 [JO18] M. Jankowiak and F. Obermeyer. “Pathwise
Derivatives Beyond the Reparameterization Trick”. In:
21 *ICML*. 2018.
- 22 [JOA10] T. Jaksch, R. Ortner, and P. Auer. “Near-
optimal Regret Bounds for Reinforcement Learning”.
In: *JMLR* 11 (2010), pp. 1563–1600.
- 23 [JOA17] P. E. Jacob, J. O’Leary, and Y. F. Atchadé.
“Unbiased Markov Chain Monte Carlo with couplings”.
In: *arXiv preprint arXiv:1708.03625* (2017).
- 24 [Joh12] M. J. Johnson. “A Simple Explanation of A
Spectral Algorithm for Learning Hidden Markov Mod-
els”. In: (2012). *arXiv: 1204.2477 [stat.ME]*.
- 25 [Jon01] D. R. Jones. “A Taxonomy of Global Opti-
mization Methods Based on Response Surfaces”. In:
26 *J. Global Optimiz.* 21.4 (2001), pp. 345–383.
- 27 [Jor07] M. I. Jordan. *An Introduction to Probabilistic
Graphical Models*. In preparation. 2007.
- 28 [Jor11] M. I. Jordan. “The era of Big Data”. In: *ISBA
Bulletin*. Vol. 18. 2011, pp. 1–3.
- 29 [Jor+98] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola,
and L. K. Saul. “An introduction to variational meth-
ods for graphical models”. In: *Learning in Graphical
Models*. Ed. by M. Jordan. MIT Press, 1998.
- 30 [Jos+17] A. Joshi, S. Ghosh, M. Betke, S. Sclaroff, and
H. Pfister. “Personalizing gesture recognition using hi-
erarchical Bayesian neural networks”. In: *CVPR*. Hon-
olulu, HI: IEEE, July 2017.
- 31 [Jos20] C. Joshi. *Transformers are Graph Neural Net-
works*. Tech. rep. 2020.
- 32 [Jos+20] M. Joshi, D. Chen, Y. Liu, D. S. Weld, L.
Zettlemoyer, and O. Levy. “Spanbert: Improving pre-
training by representing and predicting spans”. In:
33 *Transactions of the Association for Computational
Linguistics* 8 (2020), pp. 64–77.
- 34 [Jos+22] L. V. Jospin, W. Buntine, F. Boussaid, H.
Laga, and M. Bennamoun. “Hands-on Bayesian Neu-
ral Networks – a Tutorial for Deep Learning Users”.
In: (2022).
- 35 [Jou+16] A. Joulin, L. Van Der Maaten, A. Jabri, and
N. Vasilache. “Learning visual features from large
weakly supervised data”. In: *European Conference on
Computer Vision*. Springer. 2016, pp. 67–84.
- 36 [JP95] R. Jirousek and S. Preucil. “On the effective im-
plementation of the iterative proportional fitting pro-
cedure”. In: *Computational Statistics & Data Analy-
sis* 19 (1995), pp. 177–189.
- 37 [JS93] M. Jerrum and A. Sinclair. “Polynomial-time ap-
proximation algorithms for the Ising model”. In: *SIAM
J. on Computing* 22 (1993), pp. 1087–1116.
- 38 [JS96] M. Jerrum and A. Sinclair. “The Markov chain
Monte Carlo method: an approach to approximate
counting and integration”. In: *Approximation Algo-
rithms for NP-hard problems*. Ed. by D. S. Hochbaum.
PWS Publishing, 1996.
- 39 [JSY19] P. Jaini, K. A. Selby, and Y. Yu. “Sum-of-
Squares Polynomial Flow”. In: *ICML*. 2019, pp. 3009–
3018.
- 40 [JU97] S. Julier and J. Uhlmann. “A New Extension of
the Kalman Filter to Nonlinear Systems”. In: *Proc. of
AeroSense: The 11th Intl. Symp. on Aerospace/De-
fence Sensing, Simulation and Controls*. 1997.
- 41 [JUDW00] S Julier, J Uhlmann, and H. F. Durrant-
Whyte. “A new method for the nonlinear transforma-
tion of means and covariances in filters and estima-
tors”. In: *IEEE Trans. Automat. Contr.* 45.3 (Mar.
2000), pp. 477–482.
- 42 [JW14] M. Johnson and A. Willsky. “Stochastic Varia-
tional Inference for Bayesian Time Series Models”. en.
In: *ICML*. 2014, pp. 1854–1862.
- 43 [JW19] S. Jain and B. C. Wallace. “Attention is not
explanation”. In: *arXiv preprint arXiv:1902.10186*
(2019).
- 44 [Kaa12] Kaare Brandt Petersen and Michael Syskind
Pedersen. *The Matrix Cookbook*. 2012.
- 45 [KAG19] A. Kirsch, J. van Amersfoort, and Y. Gal.
“BatchBALD: Efficient and Diverse Batch Acquisition
for Deep Bayesian Active Learning”. In: *NIPS*. 2019.
- 46 [KAH19] F. H. Kingma, P. Abbeel, and J. Ho. “Bit-
Swap: Recursive Bits-Back Coding for Lossless Com-
pression with Hierarchical Latent Variables”. In:
47 *ICML*. 2019.
- 48 [Kai58] H. Kaiser. “The varimax criterion for analytic
rotation in factor analysis”. In: *Psychometrika* 23.3
(1958).

- [Kak02] S. M. Kakade. “A Natural Policy Gradient”. In: *NIPS*. 2002, pp. 1531–1538.
- [Kal06] O. Kallenberg. *Foundations of modern probability*. Springer Science & Business Media, 2006.
- [Kal+11] M. Kalakrishnan, S. Chitta, E. A. Theodorou, P. Pastor, and S. Schaal. “STOMP: Stochastic Trajectory Optimization for Motion Planning”. In: *ICRA*. 2011, pp. 4569–4574.
- [Kal+18a] D. Kalashnikov et al. “QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation”. In: *CORL*. 2018.
- [Kal+18b] N. Kalchbrenner, E. Elsen, K. Simonyan, S. Noury, N. Casagrande, E. Lockhart, F. Stimberg, A. Oord, S. Dieleman, and K. Kavukcuoglu. “Efficient neural audio synthesis”. In: *International Conference on Machine Learning*. PMLR. 2018, pp. 2410–2419.
- [Kam16] E. Kamar. “Directions in Hybrid Intelligence: Complementing AI Systems with Human Intelligence.” In: *IJCAI*. 2016, pp. 4070–4073.
- [Kam+22] S. Kamthe, S. Takao, S. Mohamed, and M. P. Deisenroth. “Iterative state estimation in non-linear dynamical systems using approximate expectation propagation”. In: *Trans. on Machine Learning Research* (2022).
- [Kan+20] T. Kaneko, H. Kameoka, K. Tanaka, and N. Hojo. “CycleGAN-VC3: Examining and Improving CycleGAN-VCs for Mel-spectrogram Conversion”. In: *Interspeech conference proceedings* (2020).
- [Kan42] L. Kantorovich. “On the transfer of masses (in Russian)”. In: *Doklady Akademii Nauk* 37.2 (1942), pp. 227–229.
- [Kap+22] S. Kapoor, W. J. Maddox, P. Izmailov, and A. G. Wilson. “On Uncertainty, Tempering, and Data Augmentation in Bayesian Classification”. In: *arXiv preprint arXiv:2203.16481* (2022).
- [Kar+18] T. Karras, T. Aila, S. Laine, and J. Lehtinen. “Progressive Growing of GANs for Improved Quality, Stability, and Variation”. In: *ICLR*. 2018.
- [Kar+20a] A.-H. Karimi, G. Barthe, B. Balle, and I. Valera. *Model-Agnostic Counterfactual Explanations for Consequential Decisions*. 2020. arXiv: [1905.11190 \[cs.LG\]](https://arxiv.org/abs/1905.11190).
- [Kar+20b] A.-H. Karimi, G. Barthe, B. Schölkopf, and I. Valera. “A survey of algorithmic recourse: definitions, formulations, solutions, and prospects”. In: *arXiv preprint arXiv:2010.04050* (2020).
- [Kar+20c] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. “Analyzing and improving the image quality of stylegan”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2020, pp. 8110–8119.
- [Kar+21] T. Karras, M. Aittala, S. Laine, E. Hätkönen, J. Hellsten, J. Lehtinen, and T. Aila. “Alias-Free Generative Adversarial Networks”. In: *arXiv preprint arXiv:2106.12423* (2021).
- [Kat05] T. Katayama. *Subspace Methods for Systems Identification*. Springer Verlag, 2005.
- [Kat+06] H. G. Katzgraber, S. Trebst, D. A. Huse, and M. Troyer. “Feedback-optimized parallel tempering Monte Carlo”. In: *Journal of Statistical Mechanics: Theory and Experiment* 2006.03 (2006), P03018.
- [Kat+17] G. Katz, C. Barrett, D. L. Dill, K. Julian, and M. J. Kochenderfer. “Reluplex: An efficient SMT solver for verifying deep neural networks”. In: *International Conference on Computer Aided Verification*. Springer. 2017, pp. 97–117.
- [Kat+19] N. Kato, H. Osone, K. Oomori, C. W. Ooi, and Y. Ochiai. “GANs-Based Clothes Design: Pattern Maker Is All You Need to Design Clothing”. In: *Proceedings of the 10th Augmented Human International Conference 2019*. Association for Computing Machinery, 2019.
- [Kau+19] V. Kaushal, R. Iyer, S. Kothawade, R. Mahadev, K. Doctor, and G. Ramakrishnan. “Learning from less data: A unified data subset selection and active learning framework for computer vision”. In: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE. 2019, pp. 1289–1299.
- [Kau+20] H. Kaur, H. Nori, S. Jenkins, R. Caruana, H. Wallach, and J. Wortman Vaughan. “Interpreting interpretability: understanding data scientists’ use of interpretability tools for machine learning”. In: *Proceedings of the 2020 CHI conference on human factors in computing systems*. 2020, pp. 1–14.
- [Kau+21] D. Kaushik, A. Setlur, E. Hovy, and Z. C. Lipton. “Explaining The Efficacy of Counterfactually Augmented Data”. In: *ICLR*. 2021.
- [KB00] H. J. Kushner and A. S. Budhiraja. “A nonlinear filtering algorithm based on an approximation of the conditional distribution”. In: *IEEE Trans. Automat. Contr.* 45.3 (2000), pp. 580–585.
- [KB14a] D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [KB14b] K. Kirchhoff and J. Bilmes. “Submodularity for Data Selection in Machine Translation”. In: *Empirical Methods in Natural Language Processing (EMNLP)*. 2014.
- [KB15] D. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization”. In: *ICLR*. 2015.
- [KB16] T. Kim and Y. Bengio. “Deep directed generative models with energy-based probability estimation”. In: *arXiv preprint arXiv:1606.03439* (2016).
- [KHB19] F. Kunstner, L. Balles, and P. Hennig. “Limitations of the Empirical Fisher Approximation”. In: (2019). arXiv: [1905.12558 \[cs.LG\]](https://arxiv.org/abs/1905.12558).
- [KCC20] V. Kumar, A. Choudhary, and E. Cho. “Data Augmentation using Pre-trained Transformer Models”. In: *Proceedings of the 2nd Workshop on Lifelong Learning for Spoken Language Systems*. Suzhou, China: Association for Computational Linguistics, Dec. 2020, pp. 18–26.
- [KD18a] S. Kamthe and M. P. Deisenroth. “Data-Efficient Reinforcement Learning with Probabilistic Model Predictive Control”. In: *AISTATS*. 2018.
- [KD18b] D. P. Kingma and P. Dhariwal. “Glow: Generative Flow with Invertible 1×1 Convolutions”. In: *NIPS*. 2018.
- [Ke+19a] L. Ke, M. Barnes, W. Sun, G. Lee, S. Choudhury, and S. Srinivasa. “Imitation Learning as f -Divergence Minimization”. In: *arXiv preprint arXiv:1905.12888* (2019).
- [Ke+19b] L. Ke, S. Choudhury, M. Barnes, W. Sun, G. Lee, and S. Srinivasa. *Imitation Learning as f -Divergence Minimization*. arXiv:1905.12888. 2019.
- [Ke+21] A. Ke, W. Ellsworth, O. Banerjee, A. Y. Ng, and P. Rajpurkar. “CheXtransfer: performance and parameter efficiency of ImageNet models for chest X-Ray interpretation”. In: *Proceedings of the Conference on Health, Inference, and Learning*. 2021, pp. 116–124.

- 1
- 2 [Kei06] F. C. Keil. “Explanation and understanding”.
In: *Annu. Rev. Psychol.* 57 (2006), pp. 227–254.
- 3 [Kel+20] J. Kelly, J. Bettencourt, M. J. Johnson, and
D. Duvenaud. “Learning Differential Equations that
are Easy to Solve”. In: *Neural Information Processing
Systems*. 2020.
- 4 [Kem+06] C. Kemp, J. Tenenbaum, T. Y. T. Griffiths
and, and N. Ueda. “Learning systems of concepts with
an infinite relational model”. In: *AAAI*. 2006.
- 5 [Kem+10] C. Kemp, J. Tenenbaum, S. Niyogi, and T.
Griffiths. “A probabilistic model of theory formation”.
In: *Cognition* 114 (2010), pp. 165–196.
- 6 [Ken16] E. H. Kennedy. “Semiparametric theory and
empirical processes in causal inference”. In: *Statisti-
cal causal inferences and their applications in pub-
lic health research*. ICSA Book Ser. Stat. Springer,
[Cham], 2016, pp. 141–167.
- 7 [Ken17] E. H. Kennedy. *Semiparametric theory*. 2017.
arXiv: [1709.06418 \[stat.ME\]](https://arxiv.org/abs/1709.06418).
- 8 [Ken20] E. H. Kennedy. *Optimal doubly robust esti-
mation of heterogeneous causal effects*. 2020. arXiv:
[2004.14497 \[math.ST\]](https://arxiv.org/abs/2004.14497).
- 9 [Kes+17] N. S. Keskar, D. Mudigere, J. Nocedal, M.
Smelyanskiy, and P. T. P. Tang. “On Large-Batch
Training for Deep Learning: Generalization Gap and
Sharp Minima”. In: *ICLR*. 2017.
- 10 [KF09a] D. Koller and N. Friedman. *Probabilistic
Graphical Models: Principles and Techniques*. MIT
Press, 2009.
- 11 [KF09b] D. Krishnan and R. Fergus. “Fast Image De-
convolution using Hyper-Laplacian Priors”. In: *NIPS*.
2009, pp. 1033–1041.
- 12 [KFL01] F. Kschischang, B. Frey, and H.-A. Loeliger.
“Factor Graphs and the Sum-Product Algorithm”. In:
IEEE Trans Info. Theory (2001).
- 13 [KG05] A. Krause and C. Guestrin. “Near-optimal Non-
myopic Value of Information in Graphical Models”. In:
*Proc. of the 21st Annual Conf. on Uncertainty in
Artificial Intelligence (UAI 2005)*. AUAI Press, 2005,
pp. 324–331.
- 14 [KG17] A. Kendall and Y. Gal. “What Uncertainties
Do We Need in Bayesian Deep Learning for Com-
puter Vision?” In: *NIPS*. Curran Associates, Inc.,
2017, pp. 5574–5584.
- 15 [KGO12] H. J. Kappen, V. Gómez, and M. Opper. “Op-
timal control as a graphical model inference problem”.
In: *Mach. Learn.* 87.2 (2012), pp. 159–182.
- 16 [KGW11] M. Kalli, J. E. Griffin, and S. G. Walker.
“Slice sampling mixture models”. In: *Statistics and
computing* 21.1 (2011), pp. 93–105.
- 17 [KH22] L. Kurscheidt and M. Hein. “Lost in Translation:
Modern Image Classifiers still degrade even un-
der simple Translations”. In: *ICML 2022 Shift Happens
Workshop*. July 2022.
- 18 [Kha+10] M. E. Khan, B. Marlin, G. Bouchard, and
K. P. Murphy. “Variational bounds for mixed-data fac-
tor analysis”. In: *NIPS*. 2010.
- 19 [Kha+18] M. E. Khan, D. Nielsen, V. Tangkaratt, W.
Lin, Y. Gal, and A. Srivastava. “Fast and Scalable
Bayesian Deep Learning by Weight-Perturbation in
Adam”. In: *ICML*. 2018.
- 20 [Kha20] M. E. Khan. *Deep learning with Bayesian
principles*. NeurIPS tutorial. 2020.
- 21 [Kha+21] S. Khan, M. Naseer, M. Hayat, S. W. Zamir,
F. S. Khan, and M. Shah. “Transformers in Vision:
A Survey”. In: *ACM Computing Surveys December*
(2021).
- 22 [Khe+20] I. Khemakhem, R. Monti, D. Kingma, and
A. Hyvärinen. “ICE-BeeM: Identifiable Conditional
Energy-Based Deep Models Based on Nonlinear ICA”.
In: *Advances in Neural Information Processing Sys-
tems*. Ed. by H. Larochelle, M. Ranzato, R. Hadsell,
M. Balcan, and H. Lin. Vol. 33. Curran Associates,
Inc., 2020, pp. 12768–12778.
- 23 [KHH20] A. Kristiadi, M. Hein, and P. Hennig. “Being
Bayesian, Even Just a Bit, Fixes Overconfidence in
ReLU Networks”. In: *ICML*. 2020.
- 24 [KHL20] D. Kaushik, E. Hovy, and Z. C. Lipton.
*Learning the Difference that Makes a Difference with
Counterfactually-Augmented Data*. 2020. arXiv: [1909.12434 \[cs.CL\]](https://arxiv.org/abs/1909.12434).
- 25 [Kho+20] P. Khosla, P. Teterwak, C. Wang, A. Sarna,
Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Kr-
ishnam. “Supervised Contrastive Learning”. In: *ArXiv
abs/2004.11362* (2020).
- 26 [Kil+20] K. Killamsetty, D. Sivasubramanian, G. Ra-
makrishnan, and R. Iyer. “GLISTER: Generalization
based Data Subset Selection for Efficient and Ro-
bust Learning”. In: *arXiv preprint arXiv:2012.10630*
(2020).
- 27 [Kim+18a] B. Kim, M. Wattenberg, J. Gilmer, C. Cai,
J. Wexler, F. Viegas, and R. Sayres. “Interpretabil-
ity Beyond Feature Attribution: Quantitative Testing
with Concept Activation Vectors (TCAV)”. In: *ICML*.
2018.
- 28 [Kim+18b] B. Kim, M. Wattenberg, J. Gilmer, C. Cai,
J. Wexler, F. Viegas, et al. “Interpretability beyond
feature attribution: Quantitative testing with concept
activation vectors (tcav)”. In: *International confer-
ence on machine learning*. PMLR. 2018, pp. 2668–
2677.
- 29 [Kim+18c] Y. Kim, S. Wiseman, A. C. Miller, D. Son-
tag, and A. M. Rush. “Semi-Amortized Variational Au-
toencoders”. In: *ICML*. 2018.
- 30 [Kim+19] S. Kim, S.-G. Lee, J. Song, J. Kim, and S.
Yoon. “FloWaveNet : A Generative Flow for Raw Au-
dio”. In: *Proceedings of the 36th International Con-
ference on Machine Learning*. 2019, pp. 3370–3378.
- 31 [Kin+14a] D. P. Kingma, D. J. Rezende, S. Mohamed,
and M. Welling. “Semi-Supervised Learning with Deep
Generative Models”. In: *NIPS*. 2014.
- 32 [Kin+14b] D. P. Kingma, D. J. Rezende, S. Mohamed,
and M. Welling. *Semi-Supervised Learning with Deep
Generative Models*. 2014. arXiv: [1406.5298 \[cs.LG\]](https://arxiv.org/abs/1406.5298).
- 33 [Kin+16] D. P. Kingma, T. Salimans, R. Jozefowicz, X.
Chen, I. Sutskever, and M. Welling. “Improved Varia-
tional Inference with Inverse Autoregressive Flow”. In:
NIPS. 2016.
- 34 [Kin+19] P.-J. Kindermans, S. Hooker, J. Adebayo, M.
Alber, K. T. Schütt, S. Dähne, D. Erhan, and B.
Kim. “The (un) reliability of saliency methods”. In:
*Explainable AI: Interpreting, Explaining and Visu-
alizing Deep Learning*. Springer, 2019, pp. 267–280.
- 35 [Kin+21] D. P. Kingma, T. Salimans, B. Poole, and
J. Ho. “Variational Diffusion Models”. In: *NIPS*. July
2021.
- 36 [Kir+17] J. Kirkpatrick et al. “Overcoming catas-
trophic forgetting in neural networks”. en. In: *PNAS*
114.13 (2017), pp. 3521–3526.

- [Kir+21] A. Kirsch, J. M. J. van Amersfoort, P. H. S. Torr, and Y. Gal. “On pitfalls in OoD detection: Predictive entropy considered harmful”. In: *ICML Workshop on Uncertainty in Deep Learning*. 2021.
- [Kit04] G. Kitagawa. “The two-filter formula for smoothing and an implementation of the Gaussian-sum smoother”. In: *Annals of the Institute of Statistical Mathematics* 46.4 (2004), pp. 605–623.
- [KIW20] P. Kirichenko, P. Izmailov, and A. G. Wilson. “Why Normalizing Flows Fail to Detect Out-of-Distribution Data”. In: (2020). arXiv: [2006.08545 \[stat.ML\]](#).
- [KJ12a] J. Z. Kolter and T. S. Jaakkola. “Approximate Inference in Additive Factorial HMMs with Application to Energy Disaggregation”. In: *AISTATS*. 2012.
- [KJ12b] B. Kulis and M. I. Jordan. “Revisiting K-Means: New Algorithms via Bayesian Nonparametrics”. In: *Proceedings of the 29th International Conference on International Conference on Machine Learning*. ICML’12. Omnipress, 2012, 1131–1138.
- [Kja90] U. Kjaerulff. *Triangulation of graphs – algorithms giving small total state space*. Tech. rep. R-90-09. Dept. of Math. and Comp. Sci., Aalborg Univ., Denmark, 1990.
- [Kja92] U. Kjaerulff. “Optimal decomposition of probabilistic networks by simulated annealing”. In: *Statistics and Computing*. Vol. 2. 1992, pp. 7–17.
- [KJD18] J. Knoblauch, J. Jewson, and T. Damoulas. “Doubly Robust Bayesian Inference for Non-Stationary Streaming Data with β -Divergences”. In: *NIPS*. 2018.
- [KJD19] J. Knoblauch, J. Jewson, and T. Damoulas. “Generalized Variational Inference: Three arguments for deriving new Posteriors”. In: (2019). arXiv: [1904.02063 \[stat.ML\]](#).
- [KJD21] J. Knoblauch, J. Jewson, and T. Damoulas. “An Optimization-centric View on Bayes’ Rule: Revisiting and Generalizing Variational Inference”. In: *JMLR* (2021).
- [KJM19] N. M. Kriege, F. D. Johansson, and C. Morris. “A Survey on Graph Kernels”. In: (2019). arXiv: [1903.11835 \[cs.LG\]](#).
- [KJV83] S. Kirkpatrick, C. G. Jr., and M. Vecchi. “Optimization by simulated annealing”. In: *Science* 220 (1983), pp. 671–680.
- [KK11] P. Krähenbühl and V. Koltun. “Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials”. In: *NIPS*. 2011.
- [KKG22] A. Kirsch, J. Kossen, and Y. Gal. “Marginal and Joint Cross-Entropies & Predictives for Online Bayesian Inference, Active Learning, and Active Sampling”. In: (May 2022). arXiv: [2205.08766 \[cs.LG\]](#).
- [KKH20] I. Khemakhem, D. P. Kingma, and A. Hyvärinen. “Variational Autoencoders and Nonlinear ICA: A Unifying Framework”. In: *AISTATS*. 2020.
- [KKL20] N. Kitaev, L. Kaiser, and A. Levskaya. “Reformer: The Efficient Transformer”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [KKR95] K. Kanazawa, D. Koller, and S. Russell. “Stochastic Simulation Algorithms for Dynamic Probabilistic Networks”. In: *UAI*. 1995.
- [KKS20] F. Kunstner, R. Kumar, and M. Schmidt. “Homeomorphic-Invariance of EM: Non-Asymptotic Convergence in KL Divergence for Exponential Families via Mirror Descent”. In: (2020). arXiv: [2011.01170 \[cs.LG\]](#).
- [KKT03] D. Kempe, J. Kleinberg, and É. Tardos. “Maximizing the spread of influence through a social network”. In: *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. 2003, pp. 137–146.
- [KL02] S. Kakade and J. Langford. “Approximately Optimal Approximate Reinforcement Learning”. In: *ICML*. ICML ’02. Morgan Kaufmann Publishers Inc., 2002, pp. 267–274.
- [KL09] H. Kawakatsu and A. Largey. “EM algorithms for ordered probit models with endogenous regressors”. In: *The Econometrics Journal* 12.1 (2009), pp. 164–186.
- [KL10] D. P. Kingma and Y. LeCun. “Regularized estimation of image statistics by score matching”. In: *Advances in neural information processing systems*. 2010, pp. 1126–1134.
- [KL17a] M. E. Khan and W. Lin. “Conjugate-Computation Variational Inference : Converting Variational Inference in Non-Conjugate Models to Inferences in Conjugate Models”. In: *AISTATS*. 2017.
- [KL17b] P. W. Koh and P. Liang. “Understanding black-box predictions via influence functions”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 1885–1894.
- [KL21] W. M. Kouw and M. Loog. “A review of domain adaptation without target labels”. en. In: *IEEE Trans. Pattern Anal. Mach. Intell.* (2021).
- [Kla+06] M. Klaas, M. Briers, N. de Freitas, A. Doucet, S. Maskell, and D. Lang. “Fast Particle Smoothing: If I Had a Million Particles”. In: *ICML*. 2006.
- [KLA19] T. Karras, S. Laine, and T. Aila. “A Style-Based Generator Architecture for Generative Adversarial Networks”. In: *CVPR*. 2019.
- [Kle17] G. A. Klein. *Sources of power: How people make decisions*. MIT press, 2017.
- [Kle18] J. Kleinberg. “Inherent Trade-Offs in Algorithmic Fairness”. In: *ACM International Conference on Measurement and Modeling of Computer Systems*. SIGMETRICS ’18. Irvine, CA, USA: Association for Computing Machinery, June 2018, p. 40.
- [KLM19] A. Kumar, P. Liang, and T. Ma. “Verified Uncertainty Calibration”. In: *NIPS*. 2019.
- [KMK00] J. Kwon and K. Murphy. *Modeling Freeway Traffic with Coupled HMMs*. Tech. rep. Univ. California, Berkeley, 2000.
- [KM08] U. Kjaerulff and A. Madsen. *Bayesian Networks and Influence Diagrams: A Guide to Construction and Analysis*. Springer, 2008.
- [KMN22] J. Kuan and J. Mueller. “Back to the Basics: Revisiting Out-of-Distribution Detection Baselines”. In: *ICML PODS Workshop*. July 2022.
- [KMB08] N. Kriegeskorte, M. Mur, and P. Bandettini. “Representational similarity analysis - connecting the branches of systems neuroscience”. en. In: *Front. Syst. Neurosci.* 2 (Nov. 2008), p. 4.
- [KML20] A. Kumar, T. Ma, and P. Liang. “Understanding Self-Training for Gradual Domain Adaptation”. In: *ICML*. Ed. by H. D. Iii and A. Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 5468–5479.

- 1
- 2 [KMR16] J. Kleinberg, S. Mullainathan, and M. Raghavan. “Inherent trade-offs in the fair determination
3 of risk scores”. In: *arXiv preprint arXiv:1609.05807* (2016).
- 4
- 5 [KMY04] D. Kersten, P. Mamassian, and A. Yuille.
6 “Object perception as Bayesian inference”. en. In:
7 *Annu. Rev. Psychol.* 55 (2004), pp. 271–304.
- 8
- 9 [KN09] J. Z. Kolter and A. Y. Ng. “Near-Bayesian Exploration
10 in Polynomial Time”. In: *ICML*. 2009.
- 11 [KN95] R. Kneser and H. Ney. “Improved backing-off
12 for n-gram language modeling”. In: *ICASSP*. Vol. 1.
13 1995, pp. 181–184.
- 14 [KN98] C.-J. Kim and C. Nelson. *State-Space Models with Regime-Switching: Classical and Gibbs-Sampling Approaches with Applications*. MIT Press, 1998.
- 15 [KNP11] K. Kersting, S. Natarajan, and D. Poole. *Statistical Relational AI: Logic, Probability and Computation*. Tech. rep. UBC, 2011.
- 16 [KNT20] I. Kostrikov, O. Nachum, and J. Tompson.
17 “Imitation Learning via Off-Policy Distribution Matching”. In: *ICLR*. 2020.
- 18 [Koe05] R. Koenker. *Quantile Regression*. en. Cambridge University Press, 2005.
- 19 [Koh+20a] P. W. Koh, T. Nguyen, Y. S. Tang, S. Mussmann, E. Pierson, B. Kim, and P. Liang. “Concept bottleneck models”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 5338–5348.
- 20 [Koh+20b] P. W. Koh et al. “WILDS: A Benchmark of in-the-Wild Distribution Shifts”. In: (Dec. 2020).
21 arXiv: 2012.07421 [cs.LG].
- 22 [Kol+20] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby. “Big
23 transfer (BiT): General visual representation learning”. In: *ECCV*. Springer. 2020, pp. 491–507.
- 24 [Kon+21] Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro. “DiffWave: A Versatile Diffusion Model
25 for Audio Synthesis”. In: *ICLR*. 2021.
- 26
- 27 [Koo03] G. Koop. *Bayesian econometrics*. Wiley, 2003.
- 28
- 29 [Kor+15] A. Korattikara, V. Rathod, K. Murphy, and M. Welling. “Bayesian Dark Knowledge”. In: *NIPS*.
30 2015.
- 31
- 32 [Kor+19] S. Kornblith, M. Norouzi, H. Lee, and G. Hin-
33 ton. “Similarity of neural network representations re-
34 visited”. In: *arXiv preprint arXiv:1905.00414* (2019).
- 35
- 36 [Kor+20] A. Korotin, V. Egiazarian, A. Asadulaev, A. Safin, and E. Burnaev. “Wasserstein-2 Generative Networks”. In: *International Conference on Learning Representations*. 2020.
- 37 [Kor+21] S. Kornblith, T. Chen, H. Lee, and M. Norouzi. “Why do better loss functions lead to less
38 transferable features?” In: *Advances in Neural Information Processing Systems*. Vol. 34. 2021.
- 39
- 40 [Kot+17] L. Kotthoff, C. Thornton, H. H. Hoos, F. Hutter, and K. Leyton-Brown. “Auto-WEKA 2.0: Auto-
41 matic model selection and hyperparameter optimization
42 in WEKA”. In: *JMLR* 18.25 (2017), pp. 1–5.
- 43
- 44 [Kot+22] S. Kothawade, V. Kaushal, G. Ramakrishnan, J. Bilmes, and R. Iyer. “PRISM: A Rich Class
45 of Parameterized Submodular Information Measures for Guided Subset Selection”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. 2022.
- 46
- 47
- [Koy+10] S. Koyama, L. C. Pérez-Bolde, C. R. Shalizi, and R. E. Kass. “Approximate Methods for State-Space Models”. en. In: *JASA* 105.489 (2010), pp. 170–180.
- [Koz92] J. R. Koza. *Genetic Programming*. <https://mitpress.mit.edu/books/genetic-programming>. Accessed: 2017-11-26. 1992.
- [KP20] A. Kumar and B. Poole. “On Implicit Regularization in β -VAEs”. In: *ICML*. 2020.
- [KPB19] I. Kobyzev, S. Prince, and M. A. Brubaker. “Normalizing Flows: An Introduction and Review of Current Methods”. In: (2019). arXiv: 1908.09257 [stat.ML].
- [KPHL17] M. J. Kusner, B. Paige, and J. M. Hernández-Lobato. “Grammar Variational Autoencoder”. In: *ICML*. 2017.
- [KPT21] J. S. Kim, G. Plumb, and A. Talwalkar. “Sanity Simulations for Saliency Methods”. In: *arXiv preprint arXiv:2105.06506* (2021).
- [KR21a] M. Khan and H. Rue. “The Bayesian Learning Rule”. In: (2021).
- [KR21b] S. Kong and D. Ramanan. “OpenGAN: Open-Set Recognition via Open Data Generation”. In: *ICCV*. 2021.
- [Kra+08] A. Krause, H. Brendan McMahan, C. Guestrin, and A. Gupta. “Robust Submodular Observation Selection”. In: *JMLR* 9 (2008), pp. 2761–2801.
- [Kri+05] B. Krishnapuram, L. Carin, M. Figueiredo, and A. Hartemink. “Learning sparse Bayesian classifiers: multi-class formulation, fast algorithms, and generalization bounds”. In: *IEEE Transaction on Pattern Analysis and Machine Intelligence* (2005).
- [Kri19] N. Kriegeskorte. *What’s the best measure of representational dissimilarity?* 2019.
- [KRL08] K. Kavukcuoglu, M. Ranzato, and Y. LeCun. “Fast Inference in Sparse Coding Algorithms with Applications to Object Recognition”. In: *NIPS workshop on optimization in machine learning*. 2008.
- [KRS14] B. Kim, C. Rudin, and J. A. Shah. “The bayesian case model: A generative approach for case-based reasoning and prototype classification”. In: *Advances in neural information processing systems*. 2014, pp. 1952–1960.
- [Kru13] J. K. Kruschke. “Bayesian estimation super-sedes the t test”. In: *J. Experimental Psychology: General* 142.2 (2013), pp. 573–603.
- [KS06] L. Kocsis and C. Szepesvári. “Bandit Based Monte-Carlo Planning”. In: *ECML*. 2006, pp. 282–293.
- [KS07] J. D. Y. Kang and J. L. Schafer. “Demystifying double robustness: a comparison of alternative strategies for estimating a population mean from incomplete data”. In: *Statist. Sci.* 22.4 (2007), pp. 523–539.
- [KS15] H. Kaya and A. A. Salah. “Adaptive Mixtures of Factor Analyzers”. In: (2015). arXiv: 1507.02801 [stat.ML].
- [KSB21] B. Kompa, J. Snoek, and A. Beam. “Empirical Frequentist Coverage of Deep Learning Uncertainty Quantification Procedures”. In: *Entropy* 23.12 (2021).
- [KSC98] S. Kim, N. Shephard, and S. Chib. “Stochastic Volatility: Likelihood Inference and Comparison with ARCH Models”. In: *Review of Economic Studies* 65.3 (1998), pp. 361–393.

- 1 [KSH12a] A. Krizhevsky, I. Sutskever, and G. Hinton.
2 “Imagenet classification with deep convolutional neural networks”. In: *NIPS*. 2012.
- 3 [KSH12b] A. Krizhevsky, I. Sutskever, and G. E. Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- 4 [KSL19] S. Kornblith, J. Shlens, and Q. V. Le. “Do better ImageNet models transfer better?” In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2019, pp. 2661–2671.
- 5 [KSL21] S. Kim, Q. Song, and F. Liang. “Stochastic gradient Langevin dynamics with adaptive drifts”. In: *J. Stat. Comput. Simul.* (2021), pp. 1–19.
- 6 [KSN99] N. L. Kleinman, J. C. Spall, and D. Q. Naiman. “Simulation-Based Optimization with Stochastic Approximation Using Common Random Numbers”. In: *Manage. Sci.* 45.11 (1999), pp. 1570–1578.
- 7 [KSS17] R. G. Krishnan, U. Shalit, and D. Sontag. “Structured Inference Networks for Nonlinear State Space Models”. In: *AAAI*. 2017.
- 8 [KT11] A. Kulesza and B. Taskar. “ k -DPPs: Fixed-size determinantal point processes”. In: *ICML*. 2011.
- 9 [KT+12] A. Kulesza, B. Taskar, et al. “Determinantal Point Processes for Machine Learning”. In: *Foundations and Trends in Machine Learning* 5.2–3 (2012), pp. 123–286.
- 10 [KTB11] D. P. Kroese, T. Taimre, and Z. I. Botev. *Handbook of Monte Carlo Methods*. en. 1 edition. Wiley, 2011.
- 11 [KTX20] R. Kohavi, D. Tang, and Y. Xu. *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing*. en. 1st ed. Cambridge University Press, 2020.
- 12 [KU21] S. Khan and J. Ugander. *Adaptive normalization for IPW estimation*. 2021. arXiv: [2106.07695 \[stat.ME\]](#).
- 13 [Kua+09] P. Kuan, G. Pan, J. A. Thomson, R. Stewart, and S. Keles. *A hierarchical semi-Markov model for detecting enrichment with application to ChIP-Seq experiments*. Tech. rep. U. Wisconsin, 2009.
- 14 [Kub04] M. Kubale. *Graph colorings*. Vol. 352. American Mathematical Society, 2004.
- 15 [Kuc+16] A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei. “Automatic Differentiation Variational Inference”. In: *JMLR* (2016).
- 16 [Kuh55] H. W. Kuhn. “The Hungarian method for the assignment problem”. In: *Naval Research Logistics Quarterly* 2 (1955), pp. 83–97.
- 17 [Kul+13] T. Kulesza, S. Stumpf, M. Burnett, S. Yang, I. Kwan, and W.-K. Wong. “Too much, too little, or just right? Ways explanations impact end users’ mental models”. In: *2013 IEEE Symposium on visual languages and human centric computing*. IEEE, 2013, pp. 3–10.
- 18 [Kul+19] M. Kull, M. Perello-Nieto, M. Kängsepp, T. S. Filho, H. Song, and P. Flach. “Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with Dirichlet calibration”. In: *NIPS*. 2019.
- 19 [Kum+19a] A. Kumar, J. Fu, M. Soh, G. Tucker, and S. Levine. “Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction”. In: *NeurIPS*. 2019, pp. 11761–11771.
- 20 [Kum+19b] M. Kumar, M. Babaeizadeh, D. Erhan, C. Finn, S. Levine, L. Dinh, and D. P. Kingma. “VideoFlow: A flow-based generative model for video”. In: *ICML Workshop on Invertible Neural Networks and Normalizing Flows* (2019).
- 21 [Kum+19c] R. Kumar, S. Ozair, A. Goyal, A. Courville, and Y. Bengio. “Maximum entropy generators for energy-based models”. In: *arXiv preprint arXiv:1901.08508* (2019).
- 22 [Kün+19] S. R. Künzel, J. S. Sekhon, P. J. Bickel, and B. Yu. “Metalearners for estimating heterogeneous treatment effects using machine learning”. In: *Proceedings of the National Academy of Sciences* 116.10 (2019), pp. 4156–4165. eprint: <https://www.pnas.org/content/116/10/4156.full.pdf>.
- 23 [Kun20] J. Kuntz. “Markov chains revisited”. In: (Jan. 2020). arXiv: [2001.02183 \[math.PR\]](#).
- 24 [Kur+19] T. Kurutach, I. Clavera, Y. Duan, A. Tamar, and P. Abbeel. “Model-Ensemble Trust-Region Policy Optimization”. In: *ICLR*. 2019.
- 25 [Kur+20] R. Kurle, B. Cseke, A. Klushyn, P. van der Smagt, and S. Günnemann. “Continual Learning with Bayesian Neural Networks for Non-Stationary Data”. In: *ICLR*. 2020.
- 26 [Kus+18] M. J. Kusner, J. R. Loftus, C. Russell, and R. Silva. *Counterfactual Fairness*. 2018. arXiv: [1703.06856 \[stat.ML\]](#).
- 27 [Kus64] H. J. Kushner. “A New Method of Locating the Maximum Point of an Arbitrary Multipeak Curve in the Presence of Noise”. In: *J. Basic Eng* 86.1 (1964), pp. 97–106.
- 28 [KVK10] A. Klami, S. Virtanen, and S. Kaski. “Bayesian exponential family projections for coupled data sources”. In: *UAI*. 2010.
- 29 [KW14] D. P. Kingma and M. Welling. “Auto-encoding variational Bayes”. In: *ICLR*. 2014.
- 30 [KW18] A. S. I. Kim and M. P. Wand. “On expectation propagation for generalised, linear and mixed models”. In: *Aust. N. Z. J. Stat.* 60.1 (2018), pp. 75–102.
- 31 [Kw19a] D. P. Kingma and M. Welling. “An Introduction to Variational Autoencoders”. In: *Foundations and Trends in Machine Learning* 12.4 (2019), pp. 307–392.
- 32 [Kw19b] M. J. Kochenderfer and T. A. Wheeler. *Algorithms for Optimization*. en. The MIT Press, 2019.
- 33 [KW70] G. S. Kimeldorf and G. Wahba. “A Correspondence Between Bayesian Estimation on Stochastic Processes and Smoothing by Splines”. en. In: *Ann. Math. Stat.* 41.2 (1970), pp. 495–502.
- 34 [Kw96] R. E. Kass and L. Wasserman. “The Selection of Prior Distributions by Formal Rules”. In: *JASA* 91.435 (1996), pp. 1343–1370.
- 35 [KwV06] K. Kurihara, M. Welling, and N. Vlassis. “Accelerated variational DP mixture models”. In: *NIPS*. 2006.
- 36 [KWW22] M. J. Kochenderfer, T. A. Wheeler, and K. Wray. *Algorithms for Decision Making*. The MIT Press, 2022.
- 37 [KY94] J. J. Kosowsky and A. L. Yuille. “The invisible hand algorithm: Solving the assignment problem with statistical physics”. In: *Neural networks* 7.3 (1994), pp. 477–490.
- 38 [Kyn+19] T. Kynkääniemi, T. Karras, S. Laine, J. Lehtinen, and T. Aila. “Improved Precision and Recall

- 1 Metric for Assessing Generative Models". In: *NeurIPS*.
 2 2019.
- 3 [KZ02] V. Kolmogorov and R. Zabih. "What energy
 4 functions can be minimized via graph cuts?" In: *Computer Vision—ECCV 2002* (2002), pp. 185–208.
- 5 [KZB19] A. Kolesnikov, X. Zhai, and L. Beyer. "Re-
 6 visiting self-supervised visual representation learning". In: *Proceedings of the IEEE/CVF conference*
 7 *on computer vision and pattern recognition*. 2019,
 8 pp. 1920–1929.
- 9 [LA87] P. J. M. Laarhoven and E. H. L. Aarts,
 10 eds. *Simulated Annealing: Theory and Applications*.
 Kluwer Academic Publishers, 1987.
- 11 [Lag+19] I. Lage, E. Chen, J. He, M. Narayanan, B.
 12 Kim, S. J. Gershman, and F. Doshi-Velez. "Human
 13 evaluation of models built for interpretability". In:
 14 *Proceedings of the AAAI Conference on Human
 Computation and Crowdsourcing*. Vol. 7. 1. 2019,
 pp. 59–67.
- 15 [Lak+17] B. M. Lake, T. D. Ullman, J. B. Tenenbaum,
 16 and S. J. Gershman. "Building Machines That Learn
 17 and Think Like People". en. In: *Behav. Brain Sci.*
 (2017), pp. 1–101.
- 18 [Lal+21] A. Lal, M. W. Lockhart, Y. Xu, and Z. Zu.
 19 "How Much Should We Trust Instrumental Variable
 Estimates in Political Science? Practical Advice based
 20 on Over 60 Replicated Studies". In: (2021).
- 21 [Lam92] D. Lambert. "Zero-Inflated Poisson Regres-
 22 sion, with an Application to Defects in Manufacturing". In: *Technometrics* 34.1 (1992), pp. 1–14.
- 22 [Lan95a] K. Lange. "A gradient algorithm locally equiv-
 23 alent to the em algorithm". en. In: *J. of Royal Stat.
 Soc. Series B* 57.2 (July 1995), pp. 425–437.
- 24 [Lan95b] K. Lange. "A QUASI-NEWTON ACCELER-
 25 ATION OF THE EM ALGORITHM". In: *Statistica
 Sinica* 5.1 (1995), pp. 1–18.
- 26 [Lao+20] J. Lao, C. Suter, I. Langmore, C. Chimisov,
 27 A. Saxena, P. Sountsov, D. Moore, R. A. Saurois,
 28 M. D. Hoffman, and J. V. Dillon. "tfpmcmc: Modern
 29 Markov Chain Monte Carlo Tools Built for Modern
 Hardware". In: *PROBPROG*. 2020.
- 30 [Lar+16] A. B. L. Larsen, S. K. Sonderby, H.
 31 Larochelle, and O. Winther. "Autoencoding beyond
 32 pixels using a learned similarity metric". In: *International conference on machine learning*. PMLR. 2016,
 pp. 1558–1566.
- 33 [Lar+21] B. W. Larsen, S. Fort, N. Becker, and S. Gan-
 34 guli. "How many degrees of freedom do we need to
 train deep networks: a loss landscape perspective". In:
 35 (2021). arXiv: [2107.05802 \[cs.LG\]](https://arxiv.org/abs/2107.05802).
- 36 [Las08] K. B. Laskey. "MEBN: A language for first-
 37 order Bayesian knowledge bases". In: *Artif. Intell.*
 172.2 (2008), pp. 140–178.
- 38 [Lau92] S. L. Lauritzen. "Propagation of probabilities,
 39 means and variances in mixed graphical association
 40 models". In: *JASA* 87.420 (1992), pp. 1098–1108.
- 41 [Lau95] S. L. Lauritzen. "The EM algorithm for graph-
 42 ical association models with missing data". In: *Com-
 putational Statistics and Data Analysis* 19 (1995),
 pp. 191–201.
- 43 [Law05] N. D. Lawrence. "Probabilistic non-linear prin-
 44 cipal component analysis with Gaussian process latent
 45 variable models". In: *JMLR* 6 (2005), pp. 1783–1816.
- 46
- 47 [Law+22] D. Lawson, A. Raventós, A. Warrington,
 and S. Linderman. "SIXO: Smoothing Inference with
 Twisted Objectives". In: (June 2022). arXiv: [2206.05952 \[cs.LG\]](https://arxiv.org/abs/2206.05952).
- [LB09] H. Lin and J. A. Bilmes. "How to Select a Good
 Training-data Subset for Transcription: Submodular
 Active Selection for Sequences". In: *Proc. Annual
 Conference of the International Speech Communica-
 tion Association (INTERSPEECH)*. Brighton, UK,
 2009.
- [LB10a] H. Lin and J. Bilmes. "Multi-document Sum-
 marization via Budgeted Maximization of Submodular
 Functions". In: *North American chapter of the Asso-
 ciation for Computational Linguistics/Human Lan-
 guage Technology Conference (NAACL/HLT-2010)*.
 Los Angeles, CA, 2010.
- [LB10b] H. Lin and J. A. Bilmes. "An Application
 of the Submodular Principal Partition to Training
 Data Subset Selection". In: *Neural Information Pro-
 cessing Society (NeurIPS, formerly NIPS) Work-
 shop*. NeurIPS (formerly NIPS) Workshop on Dis-
 crete Optimization in Machine Learning: Submodu-
 larity, Sparsity & Polyhedra (DISCML). Vancouver,
 Canada, 2010.
- [LB11] H. Lin and J. Bilmes. "A Class of Submodu-
 lar Functions for Document Summarization". In: *The
 49th Annual Meeting of the Association for Compu-
 tational Linguistics: Human Language Technologies
 (ACL/HLT-2011)*. (long paper). Portland, OR, 2011.
- [LB12] H. Lin and J. Bilmes. "Learning Mixtures of
 Submodular Shells with Application to Document
 Summarization". In: *Uncertainty in Artificial Intel-
 ligence (UAI)*. Catalina Island, USA: AUAI, 2012.
- [LB19] A. Levrat and V. Belle. "Learning Tractable
 Probabilistic Models in Open Worlds". In: (2019).
 arXiv: [1901.05847 \[cs.LG\]](https://arxiv.org/abs/1901.05847).
- [LBB20] Y. Liu, P.-L. Bacon, and E. Brunskill. "Un-
 derstanding the Curse of Horizon in Off-Policy Evalua-
 tion via Conditional Importance Sampling". In: *ICML*.
 2020.
- [LBH15] Y. LeCun, Y. Bengio, and G. Hinton. "Deep
 learning". en. In: *Nature* 521.7553 (May 2015),
 pp. 436–444.
- [LBH22] X. Lu, A. Boukouvalas, and J. Hensman. "Ad-
 ditive Gaussian Processes Revisited". In: *ICML*. June
 2022.
- [LBL16] H. Lakkaraju, S. H. Bach, and J. Leskovec.
 "Interpretable decision sets: A joint framework for de-
 scription and prediction". In: *Proceedings of the 22nd
 ACM SIGKDD international conference on knowl-
 edge discovery and data mining*. 2016, pp. 1675–1684.
- [LBM06] J. A. Lasserre, C. M. Bishop, and T. P. Minka.
 "Principled Hybrids of Generative and Discriminative
 Models". In: *CVPR*. Vol. 1. June 2006, pp. 87–94.
- [LBS01] T. Lefebvre, H. Bruyninckx, and J. D. Schut-
 ter. "Comment on "A New Method for the Nonlin-
 ear Transformation of Means and Covariances in Fil-
 ters and Estimators"". In: *IEEE Trans. on Automatic
 Control* 47.8 (2001), pp. 1406–1409.
- [LBS17] C. Lakshminarayanan, S. Bhattacharjee, and C.
 Szepesvári. "A Linearly Relaxed Approximate Linear
 Program for Markov Decision Processes". In: *IEEE
 Transactions on Automatic Control* 63.4 (2017),
 pp. 1185–1191.
- [LBW17] T. A. Le, A. G. Baydin, and F. Wood. "In-
 ference Compilation and Universal Probabilistic Pro-
 gramming". In: *AISTATS*. 2017.

- [LC02] J. Langford and R. Caruana. “(Not) bounding the true error”. In: *NIPS*. 2002.
- [LCG12] Y. Lou, R. Caruana, and J. Gehrke. “Intelligible models for classification and regression”. In: *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2012, pp. 150–158.
- [LCR21] T. Lesort, M. Caccia, and I. Rish. “Understanding Continual Learning Settings with Data Distribution Drift Analysis”. In: (2021). arXiv: [2104.01678 \[cs.LG\]](#).
- [LDZ11] Y Li, H Duan, and C. X. Zhai. “Cloudspeller: Spelling correction for search queries by using a unified hidden markov model with web-scale resources”. In: *SIGIR*. 2011.
- [LDZ12] Y. Li, H. Duan, and C. Zhai. “A Generalized Hidden Markov Model with Discriminative Training for Query Spelling Correction”. In: *SIGIR*. 2012, pp. 611–620.
- [Le+18] T. A. Le, M. Igl, T. Rainforth, T. Jin, and F. Wood. “Auto-Encoding Sequential Monte Carlo”. In: *ICLR*. 2018.
- [L'E18] P. L'Ecuyer. “Randomized Quasi-Monte Carlo: An Introduction for Practitioners”. In: *Monte Carlo and Quasi-Monte Carlo Methods*. Springer International Publishing, 2018, pp. 29–52.
- [Le+19] T. A. Le, A. R. Kosiorek, N Siddharth, Y. W. Teh, and F. Wood. “Revisiting Reweighted Wake-Sleep for Models with Stochastic Control Flow”. In: *UAI*. 2019.
- [LeC+89] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. “Backpropagation Applied to Handwritten Zip Code Recognition”. In: *Neural Computation* 1.4 (1989), pp. 541–551.
- [LeC+98] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. “Efficient BackProp”. en. In: *Neural Networks: Tricks of the Trade*. Lecture Notes in Computer Science. Springer, Berlin, Heidelberg, 1998, pp. 9–50.
- [Lee06] J. de Leeuw. “Principal Component Analysis of Binary Data by Iterated Singular Value Decomposition”. In: *Comput. Stat. Data Anal.* 50.1 (2006), pp. 21–39.
- [Lee+09] H. Lee, R. B. Grosse, R. Ranganath, and A. Ng. “Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations”. In: *ICML '09*. 2009.
- [Lee+10] J. Lee, V. Mirrokni, V. Nagarajan, and M. Sviridenko. “Maximizing Nonmonotone Submodular Functions under Matroid or Knapsack Constraints”. In: *SIAM Journal on Discrete Mathematics* 23.4 (2010), pp. 2053–2078.
- [Lee+18] J. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. “Deep Neural Networks as Gaussian Processes”. In: *ICLR*. 2018.
- [Lei+18] J. Lei, M. G'Sell, A. Rinaldo, R. J. Tibshirani, and L. Wasserman. “Distribution-Free Predictive Inference For Regression”. In: *JASA* (2018).
- [Lei+20] F. Leibfried, V. Dutordoir, S. T. John, and N. Durrande. “A Tutorial on Sparse Gaussian Processes and Variational Inference”. In: (2020). arXiv: [2012.13962 \[cs.LG\]](#).
- [Lem09] C. Lemieux. *Monte Carlo and Quasi-Monte Carlo Sampling*. Springer, New York, NY, 2009.
- [Léo14] C. Léonard. “A survey of the Schrödinger problem and some of its connections with optimal transport”. In: *Discrete & Continuous Dynamical Systems* 34.4 (2014), p. 1533.
- [Let+15a] B. Letham, C. Rudin, T. H. McCormick, and D. Madigan. “Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model”. In: *The Annals of Applied Statistics* 9.3 (2015).
- [Let+15b] B. Letham, C. Rudin, T. H. McCormick, and D. Madigan. “Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model”. In: *The Annals of Applied Statistics* 9.3 (2015), pp. 1350–1371.
- [Lev18] S. Levine. “Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review”. In: (2018). arXiv: [1805.00909 \[cs.LG\]](#).
- [Lev+20] S. Levine, A. Kumar, G. Tucker, and J. Fu. *Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems*. arXiv:2005.01643. 2020.
- [LF+21] L. Le Folgoc, V. Baltatzis, S. Desai, A. Devaraj, S. Ellis, O. E. Martinez Manzanera, A. Nair, H. Qiu, J. Schnabel, and B. Glocker. “Is MC Dropout Bayesian?” In: (2021). arXiv: [2110.04286 \[cs.LG\]](#).
- [LFG21] F. Lin, X. Fang, and Z. Gao. “Distributionally Robust Optimization: A review on theory and applications”. In: *Numer. Algebra Control Optim.* 12.1 (2021), pp. 159–212.
- [LG94] D. D. Lewis and W. A. Gale. “A sequential algorithm for training text classifiers”. In: *SIGIR94*. Springer, 1994, pp. 3–12.
- [LGMT11] F. Le Gland, V. Monbet, and V.-D. Tran. “Large Sample Asymptotics for the Ensemble Kalman Filter”. In: *Oxford Handbook of Nonlinear Filtering*. Ed. by D Crisan And. 2011.
- [LHF17] R. M. Levy, A. Haldane, and W. F. Flynn. “Potts Hamiltonian models of protein co-variation, free energy landscapes, and evolutionary fitness”. en. In: *Curr. Opin. Struct. Biol.* 43 (2017), pp. 55–62.
- [LHLT15] Y. Li, J. M. Hernandez-Lobato, and R. E. Turner. “Stochastic Expectation Propagation”. In: *NIPS*. 2015.
- [LHR20] A. Lucic, H. Hanned, and M. de Rijke. “Why does my model fail? contrastive local explanations for retail forecasting”. In: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 2020, pp. 90–98.
- [Li+10] L. Li, W. Chu, J. Langford, and R. E. Schapire. “A contextual-bandit approach to personalized news article recommendation”. In: *WWW*. 2010.
- [Li+16] C. Li, C. Chen, D. Carlson, and L. Carin. “Preconditioned Stochastic Gradient Langevin Dynamics for Deep Neural Networks”. In: *AAAI*. 2016.
- [Li+17a] A. Li, A. Jabri, A. Joulin, and L. van der Maaten. “Learning visual n-grams from web data”. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2017, pp. 4183–4192.
- [Li+17b] C.-L. Li, W.-C. Chang, Y. Cheng, Y. Yang, and B. Poczos. “Mmd gan: Towards deeper understanding of moment matching network”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 2203–2213.
- [Li+17c] L. Li, K. Jamieson, G. De Salvo, A. Rostamizadeh, and A. Talwalkar. “Hyperband: bandit-based configuration evaluation for hyperparameter optimization”. In: *ICLR*. 2017.

- 1
- 2 [Li+17d] O. Li, H. Liu, C. Chen, and C. Rudin. *Deep
Learning for Case-Based Reasoning through Proto-
types: A Neural Network that Explains Its Predictions*. 2017. arXiv: [1710.04806 \[cs.AI\]](https://arxiv.org/abs/1710.04806).
- 3
- 4 [Li+17e] T.-C. Li, J.-Y. Su, W. Liu, and J. M. Cor-
chado. “Approximate Gaussian conjugacy: parametric
recursive filtering under nonlinearity, multimodality,
uncertainty, and constraint, and beyond”. In: *Fron-
tiers of Information Technology & Electronic Engi-
neering* 18.12 (2017), pp. 1913–1939.
- 5
- 6 [Li18] Y. Li. “Deep Reinforcement Learning”. In:
7 (2018). arXiv: [1810.06339 \[cs.LG\]](https://arxiv.org/abs/1810.06339).
- 8
- 9 [Li+18a] C. Li, H. Farkhoor, R. Liu, and J. Yosinski.
10 “Measuring the Intrinsic Dimension of Objective Land-
scapes”. In: *ICLR*. 2018.
- 11
- 12 [Li+18b] C. Li, H. Farkhoor, R. Liu, and J. Yosinski.
13 “Measuring the Intrinsic Dimension of Objective Land-
scapes”. In: *ICLR*. 2018.
- 14
- 15 [Li+18c] X. Li, C. Li, J. Chi, J. Ouyang, and W. Wang.
16 “Black-box Expectation Propagation for Bayesian
Models”. In: *ICDM*. Proceedings. Society for Indus-
trial and Applied Mathematics, 2018, pp. 603–611.
- 17
- 18 [Li+19] J. Li, S. Qu, X. Li, J. Szwedley, J. Z. Kolter,
and F. Metze. “Adversarial Music: Real world Audio
Adversary against Wake-word Detection System”. In:
NIPS. Curran Associates, Inc., 2019, pp. 11908–11918.
- 19
- 20 [Li+20] C. Li, X. Gao, Y. Li, B. Peng, X. Li, Y. Zhang,
and J. Gao. “Optimus: Organizing Sentences via Pre-
trained Modeling of a Latent Space”. In: *EMNLP*.
2020.
- 21
- 22 [Li+21] Y. Li, R. Pogodin, D. J. Sutherland, and A.
23 Gretton. “Self-Supervised Learning with Kernel De-
pendence Maximization”. In: *NeurIPS*. 2021.
- 24
- 25 [Lia+07] L. Liao, D. J. Patterson, D. Fox, and H.
26 Kautz. “Learning and Inferring Transportation Rou-
tines”. In: *Artificial Intelligence* 171.5 (2007), pp. 311–
331.
- 27
- 28 [Lia+08] F. Liang, R. Paulo, G. Molina, M. Clyde, and
J. Berger. “Mixtures of g-priors for Bayesian Variable
Selection”. In: *JASA* 103.481 (2008), pp. 410–423.
- 29
- 30 [Lia+19] V. Liao, R. Ballamy, M. Muller, and H. Can-
dello. “Human-AI Collaboration: Towards Socially-
Guided Machine Learning”. In: *CHI Workshop on
Human-Centered Machine Learning Perspectives*.
2019.
- 31
- 32 [Lil+16] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess,
T. Erez, Y. Tassa, D. Silver, and D. Wierstra. “Con-
tinuous control with deep reinforcement learning”. In:
ICLR. 2016.
- 33
- 34 [Lin13a] W. Lin. “Agnostic notes on regression adjust-
ments to experimental data: Reexamining Freedman’s
critique”. In: *The Annals of Applied Statistics* 7.1
(2013), pp. 295–318.
- 35
- 36 [Lin13b] D. A. Linzer. “Dynamic Bayesian Forecast-
ing of Presidential Elections in the States”. In: *JASA*
108.501 (2013), pp. 124–134.
- 37
- 38 [Lin+17a] K. Lin, D. Li, X. He, Z. Zhang, and M.-T.
Sun. “Adversarial ranking for language generation”. In:
Advances in Neural Information Processing Systems.
2017, pp. 3155–3165.
- 39
- 40 [Lin+17b] T.-Y. Lin, P. Dollár, R. Girshick, K. He,
B. Hariharan, and S. Belongie. “Feature pyramid net-
works for object detection”. In: *Proceedings of the
IEEE conference on computer vision and pattern
recognition*. 2017, pp. 2117–2125.
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- [Lin+20] H. Lin, H. Chen, T. Zhang, C. Laroche, and K. Choromanski. “Demystifying Orthogonal Monte Carlo and Beyond”. In: (2020). arXiv: [2005.13590 \[cs.LG\]](https://arxiv.org/abs/2005.13590).
- [Lin+21a] T. Lin, Y. Wang, X. Liu, and X. Qiu. “A Survey of Transformers”. In: (2021). arXiv: [2106.04554 \[cs.LG\]](https://arxiv.org/abs/2106.04554).
- [Lin+21b] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön. *Machine Learning - A First Course for Engineers and Scientists*. 2021.
- [Lin+21c] J. Lindqvist, S. Särkkä, Á. F. García-
Fernández, M. Raitoharju, and L. Svensson. “Posterior
linearisation smoothing with robust iterations”. In:
(Dec. 2021). arXiv: [2112.03969 \[math.OC\]](https://arxiv.org/abs/2112.03969).
- [Lin+22] Z. Lin, J. Shi, D. Pathak, and D. Ramanan.
“The CLEAR Benchmark: Continual LEArning on
Real-World Imagery”. In: *NIPS Datasets and Bench-
marks Track*. Jan. 2022.
- [Lin56] D. Lindley. “On a measure of the information
provided by an experiment”. In: *The Annals of Math.
Stat.* (1956), 986–1005.
- [Lin88a] B. Lindsay. “Composite Likelihood Methods”.
In: *Contemporary Mathematics* 80.1 (1988), pp. 221–
239.
- [Lin88b] R. Linsker. “Self-organization in a perceptual
network”. In: *Computer* 21.3 (1988), pp. 105–117.
- [Lin88c] R. Linsker. “Self-organization in a perceptual
network”. In: *Computer* 21.3 (1988), pp. 105–117.
- [Lin92] L.-J. Lin. “Self-Improving Reactive Agents
Based on Reinforcement Learning, Planning and
Teaching”. In: *Mach. Learn.* 8.3-4 (1992), pp. 293–321.
- [Lip18] Z. C. Lipton. “The Mythos of Model Inter-
pretability: In machine learning, the concept of inter-
pretability is both important and slippery”. In: *Queue*
16.3 (2018), pp. 31–57.
- [Liu01] J. Liu. *Monte Carlo Strategies in Scientific
Computation*. Springer, 2001.
- [Liu+15] Z. Liu, P. Luo, X. Wang, and X. Tang. “Deep
Learning Face Attributes in the Wild”. In: *ICCV*.
2015.
- [Liu+18a] L. Liu, X. Liu, C.-J. Hsieh, and D. Tao.
“Stochastic Second-order Methods for Non-convex Opti-
mization with Inexact Hessian and Gradient”. In:
(2018). arXiv: [1809.09853 \[math.OC\]](https://arxiv.org/abs/1809.09853).
- [Liu+18b] Q. Liu, L. Li, Z. Tang, and D. Zhou. “Break-
ing the Curse of Horizon: Infinite-horizon Off-policy
Estimation”. In: *NeurIPS*. Curran Associates Inc.,
2018, pp. 5361–5371.
- [Liu+19a] H. Liu, Y.-S. Ong, Z. Yu, J. Cai, and X.
Shen. “Scalable Gaussian Process Classification with
Additive Noise for Various Likelihoods”. In: (2019).
arXiv: [1909.06541 \[stat.ML\]](https://arxiv.org/abs/1909.06541).
- [Liu+19b] R. Liu, J. Regier, N. Tripuraneni, M. I. Jor-
dan, and J. McAuliffe. “Rao-Blackwellized Stochastic
Gradients for Discrete Distributions”. In: *ICML*. 2019.
- [Liu+19c] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D.
Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoy-
anov. “Roberta: A robustly optimized bert pretrain-
ing approach”. In: *arXiv preprint arXiv:1907.11692*
(2019).
- [Liu+20a] F. Liu, W. Xu, J. Lu, G. Zhang, A. Gretton,
and D. J. Sutherland. “Learning Deep Kernels for Non-
Parametric Two-Sample Tests”. In: *ICML*. 2020.

- [Liu+20b] F. Liu, W. Xu, J. Lu, G. Zhang, A. Gretton, and D. J. Sutherland. “Learning deep kernels for non-parametric two-sample tests”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 6316–6326.
- [Liu+20c] H. Liu, Y.-S. Ong, X. Shen, and J. Cai. “When Gaussian Process Meets Big Data: A Review of Scalable GPs”. In: *IEEE Transactions on Neural Networks and Learning Systems* 31.1 (2020).
- [Liu+20d] J. Z. Liu, Z. Lin, S. Padhy, D. Tran, T. Bedrax-Weiss, and B. Lakshminarayanan. “Simple and Principled Uncertainty Estimation with Deterministic Deep Learning via Distance Awareness”. In: *NIPS*. 2020.
- [Liu+21a] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. “Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing”. In: (2021). arXiv: 2107.13586 [cs.CL].
- [Liu+21b] W. Liu, X. Wang, J. D. Owens, and Y. Li. “Energy-based Out-of-distribution Detection”. In: *NIPS*. 2021.
- [Liu+22] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. “A ConvNet for the 2020s”. In: (2022). arXiv: 2201.03545 [cs.CV].
- [LJ08] P. Liang and M. I. Jordan. “An Asymptotic Analysis of Generative, Discriminative, and Pseudo-likelihood Estimators”. In: *International Conference on Machine Learning (ICML)*. 2008.
- [Lju87] L. Ljung. *System Identification: Theory for the User*. Prentice Hall, 1987.
- [LJY19] S. Liu, Y. Jiang, and T. Yu. “Modelling RNA-Seq data with a zero-inflated mixture Poisson linear model”. en. In: *Genet. Epidemiol.* 43.7 (2019), pp. 786–799.
- [LK07] P. Liang and D. Klein. *Structured Bayesian Nonparametric Models with Variational Inference*. ACL Tutorial. 2007.
- [LK09] P. Liang and D. Klein. “Online EM for Unsupervised Models”. In: *NAACL*. 2009.
- [LKJ09] D. Lewandowski, D. Kurowicka, and H. Joe. “Generating random correlation matrices based on vines and extended onion method”. In: *J. Multivar. Anal.* 100.9 (2009), pp. 1989–2001.
- [LL02] P. Larrañaga and J. A. Lozano. *Estimation of Distribution Algorithms: A New Tool for Evolutionary Computation*. Kluwer Academic Publishers, 2002.
- [LL17] S. M. Lundberg and S.-I. Lee. “A unified approach to interpreting model predictions”. In: *NIPS*. 2017, pp. 4765–4774.
- [LLC20] M. Locher, K. B. Laskey, and P. C. G. Costa. “Design patterns for modeling first-order expressive Bayesian networks”. In: *Knowl. Eng. Rev.* 35 (2020).
- [LLJ16] Q. Liu, J. Lee, and M. Jordan. “A kernelized Stein discrepancy for goodness-of-fit tests”. In: *International conference on machine learning*. 2016, pp. 276–284.
- [LLN06] B. Lehmann, D. Lehmann, and N. Nisan. “Combinatorial auctions with decreasing marginal utilities”. In: *Games and Economic Behavior* 55.2 (2006), pp. 270–296.
- [Llo+14] J. R. Lloyd, D. Duvenaud, R. Grosse, J. B. Tenenbaum, and Z. Ghahramani. “Automatic Construction and Natural-Language Description of Non-parametric Regression Models”. In: *AAAI*. 2014.
- [LM03] T. S. Lee and D. Mumford. “Hierarchical Bayesian inference in the visual cortex.” In: *Journal of the Optical Society of America. A, Optics, image science, and vision* 20.7 (2003), pp. 1434–48.
- [LM11] H. Larochelle and I. Murray. “The neural autoregressive distribution estimator”. In: *AISTATS*. Vol. 15. 2011, pp. 29–37.
- [LM20] M. L. Leavitt and A. Morcos. “Towards falsifiable interpretability research”. In: *arXiv preprint arXiv:2010.12016* (2020).
- [LMP01] J. Lafferty, A. McCallum, and F. Pereira. “Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data”. In: *ICML*. 2001.
- [LMR15] F. Lavancier, J. Møller, and E. Rubak. “Determinantal point process models and statistical inference”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 77.4 (2015), pp. 853–877.
- [LMS16] B. Leimkuhler, C. Matthews, and G. Stoltz. “The computation of averages from equilibrium and nonequilibrium Langevin molecular dynamics”. In: *IMA J. Numer. Anal.* 36.1 (2016), pp. 13–79.
- [LN01] S. Lauritzen and D. Nilsson. “Representing and solving decision problems with limited information”. In: *Management Science* 47 (2001), pp. 1238–1251.
- [LN19] H. Lin and V. Ng. “Abstractive summarization: A survey of the state of the art”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 01. 2019, pp. 9815–9822.
- [LN96] o. Lee and J. A. Nelder. “Hierarchical Generalized Linear Models”. In: *J. of Royal Stat. Soc. Series B* 58.4 (1996), pp. 619–678.
- [Loc+18] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem. “Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations”. In: (2018). arXiv: 1811.12359 [cs.LG].
- [Loc+20a] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem. “A Sober Look at the Unsupervised Learning of Disentangled Representations and their Evaluation”. In: *Journal of Machine Learning Research* 21.209 (2020), pp. 1–62.
- [Loc+20b] F. Locatello, B. Poole, G. Rätsch, B. Schölkopf, O. Bachem, and M. Tschannen. “Weakly-supervised disentanglement without compromises”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 6348–6359.
- [Lod+02] H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins. “Text classification using string kernels”. en. In: *J. Mach. Learn. Res.* (2002).
- [Loe04] H Loeliger. “An introduction to factor graphs”. In: *IEEE Signal Process. Magazine* 21.1 (2004), pp. 28–41.
- [Loe+07] H Loeliger, J Dauwels, J Hu, S Korl, L Ping, and F. R. Kschischang. “The Factor Graph Approach to Model-Based Signal Processing”. In: *Proc. IEEE* 95.6 (2007), pp. 1295–1322.
- [Loe+16] H.-A. Loeliger, L. Bruderer, H. Malmberg, F. Wadehn, and N. Zalmai. “On Sparsity by NUV-EM, Gaussian Message Passing, and Kalman Smoothing”. In: *ITA Workshop*. 2016.
- [Loi+21] N. Loizou, S. Vaswani, I. Laradji, and S. Lacoste-Julien. “Stochastic Polyak step-size for SGD:

- 1 An adaptive learning rate for fast convergence". In: *AISTATS*. 2021.
- 2 [Lon+18] M. Long, Z. Cao, J. Wang, and M. I. Jordan.
3 "Conditional adversarial domain adaptation". In: *Neural Information Processing Systems* (2018).
- 4 [Lot+22] S. Lotfi, P. Izmailov, G. Benton, M. Gold-
5 blum, and A. G. Wilson. "Bayesian Model Selec-
6 tion, the Marginal Likelihood, and Generalization". In:
7 *ICML*. 2022.
- 8 [Lov+20] T. Lovett, M. Briers, M. Charalambides, R.
9 Jersakova, J. Lomax, and C. Holmes. "Inferring prox-
10 imity from Bluetooth Low Energy RSSI with Un-
11 scented Kalman Smoothers". In: (2020). arXiv: 2007.
12 05057 [eess.SP].
- 13 [Lov83] L. Lovász. "Submodular functions and convex-
14 ity". In: *Mathematical programming the state of the
art*. Springer, 1983, pp. 235–257.
- 15 [LÖW21] T. van de Laar, A. Özçelikkale, and H.
16 Wymeersch. "Application of the Free Energy Principle
17 to Estimation and Control". In: *IEEE Trans. Signal
Process.* 69 (2021), pp. 4234–4244.
- 18 [LP01] U. Lerner and R. Parr. "Inference in Hybrid Net-
19 works: Theoretical Limits and Practical Algorithms".
20 In: *UAI*. 2001.
- 21 [LP03] M. G. Lagoudakis and R. Parr. "Least-Squares
22 Policy Iteration". In: *JMLR* 4 (2003), pp. 1107–1149.
- 23 [LP06] N. Lartillot and H. Philippe. "Computing Bayes
24 factors using thermodynamic integration". en. In: *Sys-
tematic Biology* 55.2 (2006), pp. 195–207.
- 25 [LPB17] B. Lakshminarayanan, A. Pritzel, and C.
26 Blundell. "Simple and Scalable Predictive Uncertainty
27 Estimation using Deep Ensembles". In: *NIPS*. 2017.
- 28 [LPO17] D. Lopez-Paz and M. Oquab. "Revisiting clas-
29 sifier two-sample tests". In: *International Conference
on Learning Representations*. 2017.
- 30 [LPR17] D. Lopez-Paz and M. Ranzato. "Gradient
31 Episodic Memory for Continual Learning". In: *NIPS*.
32 2017.
- 33 [LR19] F. Lattimore and D. Rohde. "Replacing the do-
34 calculus with Bayes rule". In: (2019). arXiv: 1906.
35 07125 [stat.ML].
- 36 [LR85] T. L. Lai and H. Robbins. "Asymptotically ef-
37 ficient adaptive allocation rules". en. In: *Adv. Appl.
38 Math.* (1985).
- 39 [LR87] R. J. Little and D. B. Rubin. *Statistical Anal-
40 ysis with Missing Data*. Wiley and Son, 1987.
- 41 [LR95] C. Liu and D. Rubin. "ML Estimation of the T
42 distribution using EM and its extensions, ECM and
43 ECME". In: *Statistica Sinica* 5 (1995), pp. 19–39.
- 44 [LRC19] Y. Li, B. I. P. Rubinstein, and T. Cohn.
45 "Truth Inference at Scale: A Bayesian Model for Ad-
46 judicating Highly Redundant Crowd Annotations". In:
47 *WWW*. 2019.
- 48 [LS01] D. Lee and S. Seung. "Algorithms for non-
49 negative matrix factorization". In: *NIPS*. 2001.
- 50 [LS19] T. Lattimore and C. Szepesvari. *Bandit Algo-
51 rithms*. Cambridge, 2019.
- 52 [LS79] P. W. Lewis and G. S. Shedler. "Simulation
53 of nonhomogeneous Poisson processes by thinning".
54 In: *Naval research logistics quarterly* 26.3 (1979),
55 pp. 403–413.
- 56 [LS99] D. D. Lee and H. S. Seung. "Learning the parts
57 of objects by non-negative matrix factorization". In:
58 *Nature* 401.6755 (1999), pp. 788–791.
- 59 [LST15] B. M. Lake, R. Salakhutdinov, and J. B.
60 Tenenbaum. "Human-level concept learning through
61 probabilistic program induction". In: *Science* 350
62 (2015), pp. 1332–1338.
- 63 [LST21] N. Loo, S. Swaroop, and R. E. Turner. "Gen-
64 eralized Variational Continual Learning". In: *ICLR*.
65 2021.
- 66 [LST90] J. K. Lenstra, D. B. Shmoys, and É. Tar-
67 dós. "Approximation algorithms for scheduling unre-
68 lated parallel machines". In: *Mathematical program-
69 ming*. 1990.
- 70 [LSV09] J. Lee, M. Sviridenko, and J. Vondrák. "Sub-
71 modular maximization over multiple matroids via
72 generalized exchange properties". In: *Approximation,
73 Randomization, and Combinatorial Optimization.
74 Algorithms and Techniques* (2009), pp. 244–257.
- 75 [LSW15] Y. T. Lee, A. Sidford, and S. C.-w. Wong.
76 "A faster cutting plane method and its implications
77 for combinatorial and convex optimization". In: *2015
78 IEEE 56th Annual Symposium on Foundations of
79 Computer Science*. IEEE, 2015, pp. 1049–1065.
- 80 [LSZ15] Y. Li, K. Swersky, and R. Zemel. "Generative
81 Moment Matching Networks". In: *ICML*. 2015.
- 82 [LT16] M.-Y. Liu and O. Tuzel. "Coupled Generative
83 Adversarial Networks". In: *NIPS*. 2016, pp. 469–477.
- 84 [LTW15] Q. Li, C. Tai, and E. Weinan. "Stochastic mod-
85 ified equations and adaptive stochastic gradient algo-
86 rithms". In: *ICML*. 2015.
- 87 [Lu+21a] X. Lu, I. Osband, B. Van Roy, and Z. Wen.
88 "Evaluating Probabilistic Inference in Deep Learning:
89 Beyond Marginal Predictions". In: (2021). arXiv: 2107.
90 09224 [cs.LG].
- 91 [Lu+21b] X. Lu, B. Van Roy, V. Dwaracherla, M.
92 Ibrahim, I. Osband, and Z. Wen. "Reinforcement
93 Learning, Bit by Bit". In: (Mar. 2021). arXiv: 2103.
94 04047 [cs.LG].
- 95 [Lu+22] X. Lu, I. Osband, B. Van Roy, and Z. Wen.
96 "From Predictions to Decisions: The Importance of
97 Joint Predictive Distributions". In: (2022). arXiv: 2107.
98 09224 [cs.LG].
- 99 [Luc+18] A. Lucas, M. Iliadis, R. Molina, and A. K. Kat-
100 sagelos. "Using Deep Neural Networks for Inverse
101 Problems in Imaging: Beyond Analytical Methods". In:
102 *IEEE Signal Process. Mag.* 35.1 (2018), pp. 20–36.
- 103 [Luc+19] J. Lucas, G. Tucker, R. Grosse, and M.
104 Norouzi. "Don't blame the ELBO! A linear VAE per-
105 spective on posterior collapse". In: *NIPS*. 2019.
- 106 [Luh58] H. P. Luhn. "The automatic creation of litera-
107 ture abstracts". In: *IBM Journal of research and de-
108 velopment* 2.2 (1958), pp. 159–165.
- 109 [Lun+20] S. M. Lundberg, G. Erion, H. Chen, A. De-
110 Grave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb,
111 N. Bansal, and S.-I. Lee. "From local explanations to
112 global understanding with explainable AI for trees". In:
113 *Nature machine intelligence* 2.1 (2020), pp. 56–67.
- 114 [Luo+19] Y. Luo, H. Xu, Y. Li, Y. Tian, T. Darrell,
115 and T. Ma. "Algorithmic Framework for Model-based
116 Deep Reinforcement Learning with Theoretical Guar-
117 antees". In: *ICLR*. 2019.
- 118 [Lut16] J. Luttinen. "BayesPy: variational Bayesian in-
119 ference in Python". In: *JMLR* (2016).

- [LV06] F. Liese and I. Vajda. “On divergences and informations in statistics and information theory”. In: *IEEE Transactions on Information Theory* 52.10 (2006), pp. 4394–4412.
- [LW04] H. Lopes and M. West. “Bayesian model assessment in factor analysis”. In: *Statistica Sinica* 14 (2004), pp. 41–67.
- [LW16] C. Louizos and M. Welling. “Structured and Efficient Variational Deep Learning with Matrix Gaussian Posterioris”. In: *ICML*. 2016.
- [LW17] C. Louizos and M. Welling. “Multiplicative Normalizing Flows for Variational Bayesian Neural Networks”. In: *Proceedings of the 34th International Conference on Machine Learning*. 2017, pp. 2218–2227.
- [LWS18] Z. C. Lipton, Y.-X. Wang, and A. Smola. “Detecting and Correcting for Label Shift with Black Box Predictors”. In: *ICML*. 2018.
- [LY17] J. H. Lim and J. C. Ye. “Geometric gan”. In: *arXiv preprint arXiv:1705.02894* (2017).
- [Lyu11] S. Lyu. “Unifying non-maximum likelihood learning objectives with minimum KL contraction”. In: *Advances in Neural Information Processing Systems*. 2011, pp. 64–72.
- [Lyu12] S. Lyu. “Interpretation and generalization of score matching”. In: *arXiv preprint arXiv:1205.2629* (2012).
- [LZ20] B. Lim and S. Zohren. “Time Series Forecasting With Deep Learning: A Survey”. In: (2020). *arXiv: 2004.13408 [stat.ML]*.
- [MA10] I. Murray and R. P. Adams. “Slice sampling covariance hyperparameters of latent Gaussian models”. In: *NIPS*. 2010, pp. 1732–1740.
- [Maa+16] L. Maaløe, C. K. Sønderby, S. K. Sønderby, and O. Winther. “Auxiliary Deep Generative Models”. In: *ICML*. 2016.
- [Maa+19] L. Maaløe, M. Fraccaro, V. Liévin, and O. Winther. “BIVA: A Very Deep Hierarchy of Latent Variables for Generative Modeling”. In: *NIPS*. 2019.
- [Mac03] D. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- [Mac+11] J. H. Macke, L. Büsing, J. P. Cunningham, B. M. Y. Ece, K. V. Shenoy, and M. Sahani. “Empirical models of spiking in neural populations”. In: *NIPS*. 2011.
- [Mac+15] D. Maclaurin, D. Duvenaud, M. Johnson, and R. P. Adams. *Autograd: Reverse-mode differentiation of native Python*. 2015.
- [Mac+19] D. Maclaurin, A. Radul, M. J. Johnson, and D. Vytiniotis. “Dex: array programming with typed indices”. In: *NeurIPS workshop: Program Transformations for Machine Learning* (2019).
- [Mac75] O. Macchi. “The coincidence approach to stochastic point processes”. In: *Advances in Applied Probability* 7.1 (1975), pp. 83–122.
- [Mac92a] D. MacKay. “The evidence framework applied to classification networks”. In: *Neural Computation* 4.5 (1992), pp. 720–736.
- [Mac92b] D. J. C. MacKay. “A Practical Bayesian Framework for Backpropagation Networks”. In: *Neural Comput.* 4.3 (1992), pp. 448–472.
- [Mac95] D. MacKay. “Probable networks and plausible predictions — a review of practical Bayesian methods for supervised neural networks”. In: *Network: Computation in Neural Systems* 6.3 (1995), pp. 469–505.
- [Mac98] D. MacKay. “Introduction to Gaussian Processes”. In: *Neural Networks and Machine Learning*. Ed. by C. Bishop. 1998.
- [Mac99a] S. N. MacEachern. “Dependent nonparametric processes”. In: *ASA proceedings of the section on Bayesian statistical science*. Vol. 1. Alexandria, Virginia. Virginia: American Statistical Association; 1999. 1999, pp. 50–55.
- [Mac99b] D. MacKay. “Comparision of approximate methods for handling hyperparameters”. In: *Neural Computation* 11.5 (1999), pp. 1035–1068.
- [Mad+17] C. J. Maddison, D. Lawson, G. Tucker, N. Heess, M. Norouzi, A. Mnih, A. Doucet, and Y. W. Teh. “Filtering Variational Objectives”. In: *NIPS*. 2017.
- [Mad+18] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. “Towards Deep Learning Models Resistant to Adversarial Attacks”. In: *ICLR*. 2018.
- [Mad+19] W. J. Maddox, P. Izmailov, T. Garipov, D. P. Vetrov, and A. G. Wilson. “A Simple Baseline for Bayesian Uncertainty in Deep Learning”. In: *NIPS*. Curran Associates, Inc., 2019, pp. 13153–13164.
- [MAD20] W. R. Morningstar, A. A. Alemi, and J. V. Dillon. “PAC™-Bayes: Narrowing the Empirical Risk Gap in the Misspecified Bayesian Regime”. In: (2020). *arXiv: 2010.09629 [cs.LG]*.
- [Mah07] R. P. S. Mahler. *Statistical Multisource-Multitarget Information Fusion*. Artech House, Inc., 2007.
- [Mah13] R. Mahler. “Statistics 102 for Multisource-Multitarget Detection and Tracking”. In: *IEEE J. Sel. Top. Signal Process.* 7.3 (2013), pp. 376–389.
- [Mah+18] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. van der Maaten. “Exploring the limits of weakly supervised pretraining”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 181–196.
- [Mah+19] N. Maheswaranathan, A. H. Williams, M. D. Golub, S. Ganguli, and D. Sussillo. “Universality and individuality in neural dynamics across large populations of recurrent networks”. In: *Advances in Neural Information Processing Systems* 2019 (2019), p. 15629.
- [Mai+10] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. “Online learning for matrix factorization and sparse coding”. In: *JMLR* 11 (2010), pp. 19–60.
- [Mai13] J. Mairal. “Stochastic Majorization-minimization Algorithms for Large-scale Optimization”. In: *NIPS*. 2013, pp. 2283–2291.
- [Mai15] J. Mairal. “Incremental Majorization-Minimization Optimization with Application to Large-Scale Machine Learning”. In: *SIAM J. Optim.* 25.2 (2015), pp. 829–855.
- [Mai+22] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner. “Online continual learning in image classification: An empirical survey”. In: *Neurocomputing* 469 (2022), pp. 28–51.
- [Mak+15a] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. “Adversarial Autoencoders”. In: (2015). *arXiv: 1511.05644 [cs.LG]*.

- 1
- 2 [Mak+15b] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. “Adversarial autoencoders”. In: *arXiv preprint arXiv:1511.05644* (2015).
- 3
- 4 [Mak+20] A. Makkluva, A. Taghvaei, S. Oh, and J. Lee. “Optimal transport mapping via input convex neural networks”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 6672–6681.
- 5
- 6 [Mal+17] D. M. Malioutov, K. R. Varshney, A. Emad, and S. Dash. “Learning interpretable classification rules with boolean compressed sensing”. In: *Transparent Data Mining for Big and Small Data*. Springer, 2017, pp. 95–121.
- 7
- 8 [Man+17] B. Mann. *How many times should you shuffle a deck of cards?* Tech. rep. Dartmouth.
- 9
- 10 [Man+07] V. Mansinghka, D. Roy, R. Rifkin, and J. Tenenbaum. “AClass: An online algorithm for generative classification”. In: *AISTATS*. 2007.
- 11
- 12 [Man+19] D. J. Mankowitz, N. Levine, R. Jeong, Y. Shi, J. Kay, A. Abdolmaleki, J. T. Springenberg, T. Mann, T. Hester, and M. Riedmiller. “Robust Reinforcement Learning for Continuous Control with Model Misspecification”. In: (2019). arXiv: [1906.07516](#) [cs.LG].
- 13
- 14 [Man22] V. Manokhin. “Machine Learning for Probabilistic Prediction”. PhD thesis. Royal Holloway, University of London, June 2022.
- 15
- 16 [Man90] C. F. Manski. “Nonparametric Bounds on Treatment Effects”. In: *The American Economic Review* 80.2 (1990), pp. 319–323.
- 17
- 18 [Mao+17] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley. “Least Squares Generative Adversarial Networks”. In: *ICCV*. 2017.
- 19
- 20 [Mar03] B. Marlin. “Modeling User Rating Profiles for Collaborative Filtering”. In: *NIPS*. 2003.
- 21
- 22 [Mar+06] A. Margolin, I. Nemenman, K. Basso, C. Wiggins, G. Stolovitzky, and R. F. abd A. Califano. “ARACNE: An Algorithm for the Reconstruction of Gene Regulatory Networks in a Mammalian Cellular Context”. In: *BMC Bioinformatics* 7 (2006).
- 23
- 24 [Mar08] B. Marlin. “Missing Data Problems in Machine Learning”. PhD thesis. U. Toronto, 2008.
- 25
- 26 [Mar+10] B. M. Marlin, K. Swersky, B. Chen, and N. de Freitas. “Inductive Principles for Restricted Boltzmann Machine Learning”. In: *AISTATS*. 2010.
- 27
- 28 [Mar10a] J. Martens. “Deep learning via Hessian-free optimization”. In: *ICML*. 2010.
- 29
- 30 [Mar10b] J. Martens. “Learning the Linear Dynamical System with ASOS”. In: *ICML*. ICML’10. Omnipress, 2010, pp. 743–750.
- 31
- 32 [Mar16] J. Martens. “Second-order optimization for neural networks”. PhD thesis. Toronto, 2016.
- 33
- 34 [Mar18] O. Martin. *Bayesian analysis with Python*. Packt, 2018.
- 35
- 36 [Mar20] J. Martens. “New insights and perspectives on the natural gradient method”. In: *JMLR* (2020).
- 37
- 38 [Mas+20] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. van de Weijer. “Class-incremental learning: survey and performance evaluation on image classification”. In: (2020). arXiv: [2010.15277](#) [cs.LG].
- 39
- 40 [Mat14] C. Mattfeld. “Implementing spectral methods for hidden Markov models with real-valued emissions”. MA thesis. ETH Zurich, 2014.
- 41
- 42 [Mat+16] A. Matthews, J. Hensman, R. Turner, and Z. Ghahramani. “On Sparse Variational Methods and the Kullback-Leibler Divergence between Stochastic Processes”, en. In: *AISTATS*. 2016, pp. 231–239.
- 43
- 44 [Mav16] Mavrogonatou, L And Vyshevsky, “Sequential Importance Sampling for Online Bayesian Change-point Detection”. In: *22nd International Conference on Computational Statistics*. 2016.
- 45
- 46 [May+19] A. May, J. Zhang, T. Dao, and C. Ré. “On the downstream performance of compressed word embeddings”. In: *Advances in neural information processing systems* 32 (2019), p. 11782.
- 47
- 48 [May79] P. Maybeck. *Stochastic models, estimation, and control*. Academic Press, 1979.
- 49
- 50 [Maz+22] P. Mazzaglia, T. Verbelen, O. Çatal, and B. Dhoedt. “The Free Energy Principle for Perception and Action: A Deep Learning Perspective”, en. In: *Entropy* 24.2 (2022).
- 51
- 52 [MAZA18] X. B. P. and Marcin Andrychowicz, W. Zaremba, and P. Abbeel. “Sim-to-Real Transfer of Robotic Control with Dynamics Randomization”. In: *ICRA*. 2018, pp. 1–8.
- 53
- 54 [MB16] Y. Miao and P. Blunsom. “Language as a Latent Variable: Discrete Generative Models for Sentence Compression”. In: *EMNLP*. 2016.
- 55
- 56 [MB18] A. Mensch and M. Blondel. “Differentiable Dynamic Programming for Structured Prediction and Attention”. In: *ICML*. 2018.
- 57
- 58 [MB21] M. Y. Michelis and Q. Becker. “On Linear Interpolation in the Latent Space of Deep Generative Models”. In: *ICLR Workshop on Geometrical and Topological Representation Learning*. 2021.
- 59
- 60 [MB88] T. Mitchell and J. Beauchamp. “Bayesian Variable Selection in Linear Regression”. In: *JASA* 83 (1988), pp. 1023–1036.
- 61
- 62 [MBJ06] J. McAuliffe, D. Blei, and M. Jordan. “Nonparametric empirical Bayes for the Dirichlet process mixture model”. In: *Statistics and Computing* 16.1 (2006), pp. 5–14.
- 63
- 64 [MBJ20] T. M. Moerland, J. Broekens, and C. M. Jonker. “Model-based Reinforcement Learning: A Survey”, en. In: (2020). arXiv: [2006.16712](#) [cs.LG].
- 65
- 66 [MBK20] X. Meng, R. Bachmann, and M. E. Khan. “Training Binary Neural Networks using the Bayesian Learning Rule”, en. In: *ICML*. 2020.
- 67
- 68 [MBL20] B. Mirzasoleiman, J. Bilmes, and J. Leskovec. “Coresets for data-efficient training of machine learning models”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 6950–6960.
- 69
- 70 [MBW20] W. J. Maddox, G. Benton, and A. G. Wilson. “Rethinking Parameter Counting in Deep Models: Effective Dimensionality Revisited”, en. In: *arXiv preprint arXiv:2003.02139* (2020).
- 71
- 72 [MC03] P. Moscato and C. Cotta. “A Gentle Introduction to Memetic Algorithms”, en. In: *Handbook of Metaheuristics*. International Series in Operations Research & Management Science. Springer, Boston, MA, 2003, pp. 105–144.
- 73
- 74 [MC19] P. Moreno Comellas. “Vision as inverse graphics for detailed scene understanding”, en. PhD thesis. 2019.
- 75
- 76 [McA99] D. A. McAllester. “PAC-Bayesian model averaging”, en. In: *Proceedings of the twelfth annual conference on computational learning theory*. 1999.

- 1 [McC03] A. McCray. “An upper level ontology for the
2 biomedical domain”. In: *Comparative and Functional*
3 *Genomics* 4 (2003), pp. 80–84.
- 4 [McE20] R. McElreath. *Statistical Rethinking: A*
5 *Bayesian Course with Examples in R and Stan* (2nd
6 edition). en. Chapman and Hall/CRC, 2020.
- 7 [McG54] W. McGill. “Multivariate information trans-
mission”. In: *Psychometrika* 19 (1954), pp. 97–116.
- 8 [McM+13] H. B. McMahan, G. Holt, D. Sculley, M.
9 Young, D. Ebner, J. Grady, L. Nie, T. Phillips, E.
Davydov, D. Golovin, et al. “Ad click prediction: a
10 view from the trenches”. In: *KDD*. 2013, pp. 1222–
1230.
- 11 [Md+19] C. de Masson d’Autume, M. Rosca, J. Rae,
12 and S. Mohamed. “Training language GANs from
13 Scratch”. In: (2019). arXiv: [1905.09922 \[cs.CL\]](https://arxiv.org/abs/1905.09922).
- 14 [MD97] X. L. Meng and D. van Dyk. “The EM algo-
15 rithm — an old folk song sung to a fast new tune
(with Discussion)”. In: *J. Royal Stat. Soc. B* 59 (1997),
pp. 511–567.
- 16 [MDA15] D. Maclaurin, D. Duvenaud, and R. P.
17 Adams. “Gradient-based Hyperparameter Optimiza-
18 tion through Reversible Learning”. In: *ICML*. 2015.
- 19 [MDM19] S. Mahloujifar, D. I. Diochnos, and M. Mah-
20 moody. “The curse of concentration in robust learning:
Evasion and poisoning attacks from concentration of
21 measure”. In: *Proceedings of the AAAI Conference on
Artificial Intelligence*. Vol. 33. 2019, pp. 4536–4543.
- 22 [MDR94] S. Muggleton and L. De Raedt. “Inductive
23 logic programming: Theory and methods”. In: *The
Journal of Logic Programming* 19 (1994), pp. 629–
679.
- 24 [ME17] H. Mei and J. M. Eisner. “The neural hawkes
25 process: A neurally self-modulating multivariate point
process”. In: *Advances in Neural Information Pro-
26 cessing Systems*. 2017, pp. 6754–6764.
- 27 [Med+21] M. A. Medina, J. L. M. Olea, C. Rush, and
28 A. Velez. “On the Robustness to Misspecification of α -
29 Posteriors and Their Variational Approximations”. In:
(2021). arXiv: [2104.08324 \[stat.ML\]](https://arxiv.org/abs/2104.08324).
- 30 [Mee+18] J.-W. van de Meent, B. Paige, H. Yang, and
F. Wood. *An introduction to probabilistic program-
31 ming*. Foundations and Trends in Machine Learning,
2018.
- 32 [Mei18a] N. Meinshausen. “Causality from a distribu-
33 tional robustness point of view”. In: *IEEE Data Sci-
ence Workshop (DSW)*. 2018, pp. 6–10.
- 34 [Mei18b] N. Meinshausen. “CAUSALITY FROM A
35 DISTRIBUTIONAL ROBUSTNESS POINT OF
36 VIEW”. In: *2018 IEEE Data Science Workshop
(DSW)*. 2018, pp. 6–10.
- 37 [Mer] *Definition of interpret*. 2022. URL: <https://www.merriam-webster.com/dictionary/interpret>.
- 38 [Mer+00] R. van der Merwe, A. Doucet, N. de Fre-
39 itas, and E. Wan. “The Unscented Particle Filter”. In:
40 *NIPS-13*. 2000.
- 41 [Met16] C. Metz. *In Two Moves, AlphaGo and Lee
42 Sedol Redefined the Future*. 2016. URL: <https://www.wired.com/2016/03/two-moves-alpha-go-lee-sedol-redefined-future/> (visited on 01/07/2022).
- 43 [Met+16] L. Metz, B. Poole, D. Pfau, and J. Sohl-
44 Dickstein. “Unrolled Generative Adversarial Net-
45 works”. In: (2016).
- 46 [Met+17] L. Metz, J. Ibarz, N. Jaitly, and J. Davidson.
“Discrete Sequential Prediction of Continuous Actions
47 for Deep RL”. In: (2017). arXiv: [1705.05035 \[cs.LG\]](https://arxiv.org/abs/1705.05035).
- [Met+53] N. Metropolis, A. Rosenbluth, M. Rosen-
bluth, A. Teller, and E. Teller. “Equation of state
calculations by fast computing machines”. In: *J. of
Chemical Physics* 21 (1953), pp. 1087–1092.
- [Mey+18] F. Meyer, T. Kropfreiter, J. Williams, R.
Lau, F. Hlawatsch, P. Braca, and M. Win. “Message
passing algorithms for scalable multitarget tracking”.
In: *Proc. IEEE 106.2* (2018).
- [Mey+21] R. A. Meyer, C. Musco, C. Musco, and D. P.
Woodruff. “Hutch++: Optimal Stochastic Trace Es-
timation”. In: *SIAM Symposium on Simplicity in Al-
gorithms (SOSA21)*. 2021.
- [Mey22] S. Meyn. *Control Systems and Reinforcement
Learning*. Cambridge, 2022.
- [MFP00] A. McCallum, D. Freitag, and F. Pereira.
“Maximum Entropy Markov Models for Information
Extraction and Segmentation”. In: *ICML*. 2000.
- [MFR20] G. M. Martin, D. T. Frazier, and C. P. Robert.
“Computing Bayes: Bayesian Computation from 1763
to the 21st Century”. In: (2020). arXiv: [2004.06425 \[stat.CO\]](https://arxiv.org/abs/2004.06425).
- [MG15] J. Martens and R. Grosse. “Optimizing Neural
Networks with Kronecker-factored Approximate Cur-
vature”. In: *ICML*. 2015.
- [MGM06] I. Murray, Z. Ghahramani, and D. J. C.
MacKay. “MCMC for doubly-intractable distri-
butions”. In: *Proceedings of the 22nd Annual Confer-
ence on Uncertainty in Artificial Intelligence (UAI-
06)*. AUAI Press, 2006, pp. 359–366.
- [MGN18a] L. Mescheder, A. Geiger, and S. Nowozin.
“Which Training Methods for GANs do actually Con-
verge?” In: *ICML*. 2018.
- [MGN18b] L. Mescheder, A. Geiger, and S. Nowozin.
“Which training methods for GANs do actually con-
verge?” In: *International conference on machine
learning*. PMLR. 2018, pp. 3481–3490.
- [MGR18] H. Mania, A. Guy, and B. Recht. “Simple
random search of static linear policies is competitive
for reinforcement learning”. In: *NIPS*. Ed. by S Ben-
gio, H Wallach, H Larochelle, K Grauman, N Ces-
 Bianchi, and R Garnett. Curran Associates, Inc., 2018,
pp. 1800–1809.
- [MH12] R. Mazumder and T. Hastie. *The Graphi-
cal Lasso: New Insights and Alternatives*. Tech. rep.
Stanford Dept. Statistics, 2012.
- [MH14] J. W. Miller and M. T. Harrison. “Inconsis-
tency of Pitman-Yor process mixtures for the number
of components”. In: *JMLR* 15.1 (2014), pp. 3333–3370.
- [MH20] I. Mordatch and J. Hamrick. *ICML tu-
torial on model-based methods in reinforcement
learning*. <https://sites.google.com/corp/view/mbrl-tutorial>. 2020.
- [MHB17] S. Mandt, M. D. Hoffman, and D. M.
Blei. “Stochastic Gradient Descent As Approximate
Bayesian Inference”. In: *JMLR* 18.1 (2017), pp. 4873–
4907.
- [MHH14] F. Meyer, O. Hlinka, and F. Hlawatsch.
“Sigma point belief propagation”. In: *IEEE Signal
Processing Letters* 21.2 (2014), pp. 145–149.

- 1
- 2 [MHN13] A. L. Maas, A. Y. Hannun, and A. Y. Ng. “Rectifier Nonlinearities Improve Neural Network Acoustic Models”. In: *ICML*. Vol. 28. 2013.
- 3
- 4 [Mik+13] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. “Distributed representations of words and phrases and their compositionality”. In: *NIPS*. 2013, pp. 3111–3119.
- 5
- 6 [Mil+05] B. Milch, B. Marthi, S. Russell, D. Sontag, D. Ong, and A. Kolobov. “BLOG: Probabilistic Models with Unknown Objects”. In: *IJCAI*. 2005.
- 7
- 8 [Mil19] T. Miller. “Explanation in artificial intelligence: Insights from the social sciences”. In: *Artificial intelligence* 267 (2019), pp. 1–38.
- 9
- 10 [Mil+20] B. Millidge, A. Tschantz, A. K. Seth, and C. L. Buckley. “On the Relationship Between Active Inference and Control as Inference”. In: *International Workshop on Active Inference*. 2020.
- 11 [Mil+21] J. P. Miller, R. Taori, A. Raghunathan, S. Sagawa, P. W. Koh, V. Shankar, P. Liang, Y. Carmon, and L. Schmidt. “Accuracy on the Line: on the Strong Correlation Between Out-of-Distribution and In-Distribution Generalization”. In: *ICML*. Ed. by M. Meila and T. Zhang. Vol. 139. Proceedings of Machine Learning Research. PMLR, 2021, pp. 7721–7735.
- 12
- 13 [Min00a] T. Minka. *Bayesian linear regression*. Tech. rep. MIT, 2000.
- 14
- 15 [Min00b] T. Minka. *Bayesian model averaging is not model combination*. Tech. rep. MIT Media Lab, 2000.
- 16
- 17 [Min00c] T. Minka. *Estimating a Dirichlet distribution*. Tech. rep. MIT, 2000.
- 18
- 19 [Min01a] T. Minka. “A family of algorithms for approximate Bayesian inference”. PhD thesis. MIT, 2001.
- 20
- 21 [Min01b] T. Minka. “Expectation Propagation for approximate Bayesian inference”. In: *UAI*. 2001.
- 22
- 23 [Min04] T. Minka. *Power EP*. Tech. rep. MSR-TR-2004-149. 2004.
- 24
- 25 [Min05] T. Minka. *Divergence measures and message passing*. Tech. rep. MSR Cambridge, 2005.
- 26
- 27 [Min+18] T. Minka, J. Winn, J. Guiver, Y. Zaykov, D. Fabian, and J. Bronskill. *Infer.NET 0.3*. Microsoft Research Cambridge. 2018.
- 28
- 29 [Min78] M. Minoux. “Accelerated greedy algorithms for maximizing submodular set functions”. In: *Optimization Techniques*. Ed. by J. Stoer. Vol. 7. Lecture Notes in Control and Information Sciences. 10.1007/BFb0006528. Springer Berlin / Heidelberg, 1978, pp. 234–243.
- 30
- 31 [Mis+18] A. Mishkin, F. Kunstner, D. Nielsen, M. Schmidt, and M. E. Khan. “SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient”. In: *NIPS*. Curran Associates, Inc., 2018, pp. 6245–6255.
- 32
- 33 [Mit+19] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. “Model cards for model reporting”. In: *Proceedings of the conference on fairness, accountability, and transparency*. 2019, pp. 220–229.
- 34
- 35 [Mit+20] J. Mitrovic, B. McWilliams, J. Walker, L. Buesing, and C. Blundell. *Representation Learning via Invariant Causal Mechanisms*. 2020. arXiv: 2010.07922 [cs.LG].
- 36
- 37 [Mit97] T. Mitchell. *Machine Learning*. McGraw Hill, 1997.
- 38
- 39 [Miy+18a] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. “Spectral Normalization for Generative Adversarial Networks”. In: *ICLR*. 2018.
- 40
- 41 [Miy+18b] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. “Spectral Normalization for Generative Adversarial Networks”. In: *ICLR*. 2018.
- 42
- 43 [Miy+18c] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. “Spectral Normalization for Generative Adversarial Networks”. In: *International Conference on Learning Representations*. 2018.
- 44
- 45 [MJ97] M. Meila and M. Jordan. *Triangulation by continuous embedding*. Tech. rep. 1605. MIT AI Lab, 1997.
- 46
- 47 [MK05] J. Mooij and H. Kappen. “Sufficient conditions for convergence of loopy belief propagation”. In: *UAI*. 2005.
- 48 [MK07] G. J. McLachlan and T. Krishnan. *The EM Algorithm and Extensions (Second Edition)*. Wiley, 2007.
- 49 [MK18] T. Miyato and M. Koyama. “cGANs with Projection Discriminator”. In: *International Conference on Learning Representations*. 2018.
- 50 [MK19] J. Menick and N. Kalchbrenner. “Generating high fidelity images with subscale pixel networks and multidimensional upscaling”. In: *ICLR*. 2019.
- 51 [MK97] G. J. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. Wiley, 1997.
- 52 [MKG21] W. J. Ma, K. Kording, and D. Goldreich. *Bayesian models of perception and action*. MIT Press, 2021.
- 53 [MKH19] R. Müller, S. Kornblith, and G. E. Hinton. “When does label smoothing help?” In: *NIPS*. 2019, pp. 4694–4703.
- 54 [MKL11] O. Martin, R. Kumar, and J. Lao. *Bayesian Modeling and Computation in Python*. CRC Press, 2011.
- 55 [MKL21] O. A. Martin, R. Kumar, and J. Lao. *Bayesian Modeling and Computation in Python*. CRC Press, 2021.
- 56 [MKS21] K. Murphy, A. Kumar, and S. Serghiu. “Risk score learning for COVID-19 contact tracing apps”. In: *Machine Learning for Healthcare*. 2021.
- 57 [ML02] T. Minka and J. Lafferty. “Expectation-propagation for the Generative Aspect Model”. In: *UAI*. Morgan Kaufmann Publishers Inc., 2002, pp. 352–359.
- 58 [ML16] S. Mohamed and B. Lakshminarayanan. “Learning in Implicit Generative Models”. In: (2016). arXiv: 1610.03483 [stat.ML].
- 59 [MLN19] P. Michel, O. Levy, and G. Neubig. “Are Sixteen Heads Really Better than One?” In: *NIPS*. 2019.
- 60 [MLW19] V. Masrani, T. A. Le, and F. Wood. “The Thermodynamic Variational Objective”. In: *NIPS*. Curran Associates, Inc., 2019, pp. 11521–11530.
- 61 [MM01] T. K. Marks and J. R. Movellan. *Diffusion networks, products of experts, and factor analysis*. Tech. rep. University of California San Diego, 2001.
- 62 [MM90] D. Q. Mayne and H Michalska. “Receding horizon control of nonlinear systems”. In: *IEEE Trans. Automat. Contr.* 35.7 (1990), pp. 814–824.
- 63 [MMC98] R. J. McEliece, D. J. C. MacKay, and J. F. Cheng. “Turbo decoding as an instance of Pearl’s ‘belief propagation’ algorithm”. In: *IEEE J. on Selected Areas in Comm.* 16.2 (1998), pp. 140–152.

- [MMP13] L. Malagò, M. Matteucci, and G. Pistoni. “Natural gradient, fitness modelling and model selection: A unifying perspective”. In: *IEEE Congress on Evolutionary Computation*. 2013.
- [MMP87] J. Marroquin, S. Mitter, and T. Poggio. “Probabilistic solution of ill-posed problems in computational vision”. In: *JASA* 82.297 (1987), pp. 76–89.
- [MMT17] C. J. Maddison, A. Mnih, and Y. W. Teh. “The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables”. In: *ICLR*. 2017.
- [MN89] P. McCullagh and J. Nelder. *Generalized linear models*. 2nd edition. Chapman and Hall, 1989.
- [MNG17a] L. Mescheder, S. Nowozin, and A. Geiger. “Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2391–2400.
- [MNG17b] L. Mescheder, S. Nowozin, and A. Geiger. “The numerics of gans”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 1825–1835.
- [Mni+15] V. Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature* 518.7540 (2015), pp. 529–533.
- [Mni+16] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. “Asynchronous Methods for Deep Reinforcement Learning”. In: *ICML*. 2016.
- [MO14] M. Mirza and S. Osindero. “Conditional generative adversarial nets”. In: *arXiv preprint arXiv:1411.1784* (2014).
- [Mob16] H. Mobahi. “Closed Form for Some Gaussian Convolutions”. In: (2016). arXiv: [1602.05610 \[math.CA\]](https://arxiv.org/abs/1602.05610).
- [Moc+96] J. Mockus, W. Eddy, A. Mockus, L. Mockus, and G. Reklaitis. *Bayesian Heuristic Approach to Discrete and Global Optimization: Algorithms, Visualization, Software, and Applications*. Kluwer, 1996.
- [Moh+19] S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih. “Monte Carlo Gradient Estimation in Machine Learning”. In: (2019). arXiv: [1906.10652 \[stat.ML\]](https://arxiv.org/abs/1906.10652).
- [Mon81] G. Monge. “Mémoire sur la théorie des déblais et des remblais”. In: *Histoire de l'Académie Royale des Sciences* (1781), pp. 666–704.
- [Mor+11] F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina, J. N. Onuchic, T. Hwa, and M. Weigt. “Direct-coupling analysis of residue coevolution captures native contacts across many protein families”. In: *Proc. Natl. Acad. Sci. U. S. A.* 108.49 (2011), E1293–301.
- [Mor+21a] W. Morningstar, C. Ham, A. Gallagher, B. Lakshminarayanan, A. Alemi, and J. Dillon. “Density of States Estimation for Out of Distribution Detection”. In: *AISTATS*. Ed. by A. Banerjee and K. Fukumizu. Vol. 130. Proceedings of Machine Learning Research. PMLR, 2021, pp. 3232–3240.
- [Mor+21b] W. Morningstar, S. Vikram, C. Ham, A. Gallagher, and J. Dillon. “Automatic Differentiation Variational Inference with Mixtures”. In: *AISTATS*. Ed. by A. Banerjee and K. Fukumizu. Vol. 130. Proceedings of Machine Learning Research. PMLR, 2021, pp. 3250–3258.
- [Mor63] T. Morimoto. “Markov Processes and the H-Theorem”. In: *J. Phys. Soc. Jpn.* 18.3 (1963), pp. 328–331.
- [MOT15] A. Mordvintsev, C. Olah, and M. Tyka. *Inceptionism: Going Deeper into Neural Networks*. <https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>. Accessed: NA-NA-NA. 2015.
- [Mov08] J. R. Movellan. “A minimum velocity approach to learning”. In: *unpublished draft*, Jan (2008).
- [MP01] K. Murphy and M. Paskin. “Linear time inference in hierarchical HMMs”. In: *NIPS*. 2001.
- [MP21] D. Mazza and M. Pagani. “Automatic Differentiation in PCF”. In: *Proc. ACM Program. Lang.* 5.POPL (2021).
- [MP95] D. MacKay and L. Peto. “A hierarchical dirichlet language model”. In: *Natural Language Engineering* 1.3 (1995), pp. 289–307.
- [MPS18] O. Mangoubi, N. S. Pillai, and A. Smith. “Does Hamiltonian Monte Carlo mix faster than a random walk on multimodal densities?”. In: (2018). arXiv: [1808.03230 \[math.PR\]](https://arxiv.org/abs/1808.03230).
- [MR09] A. Melkumyan and F. Ramos. “A Sparse Covariance Function for Exact Gaussian Process Inference in Large Datasets”. In: *IJCAI*. 2009, pp. 1936–1942.
- [MR10] B. Milch and S. Russell. “Extending Bayesian Networks to the Open-Universe Case”. In: *Heuristics, Probability and Causality. A Tribute to Judea Pearl*. Ed. by R. Dechter, H. Geffner, and J. Y. Halper. College Publications, 2010.
- [MRB18] A. S. Morcos, M. Raghu, and S. Bengio. “Insights on representational similarity in neural networks with canonical correlation”. In: *Advances in neural information processing systems* (2018).
- [MRW19] B. Mittelstadt, C. Russell, and S. Wachter. “Explaining explanations in AI”. In: *Proceedings of the conference on fairness, accountability, and transparency*. 2019, pp. 279–288.
- [MS96] V Matveev V and R Shrock. “Complex-temperature singularities in Potts models on the square lattice”. In: *Phys. Rev. E Stat. Phys. Plasmas Fluids Relat. Interdiscip. Topics* 54.6 (1996), pp. 6174–6185.
- [MS99] C. Manning and H. Schütze. *Foundations of statistical natural language processing*. MIT Press, 1999.
- [MSA18] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. “The M4 Competition: Results, findings, conclusion and way forward”. In: *Int. J. Forecast.* 34.4 (2018), pp. 802–808.
- [MT12] A. Mnih and Y. W. Teh. “A fast and simple algorithm for training neural probabilistic language models”. In: *ICML*. 2012, pp. 419–426.
- [MTM14] C. J. Maddison, D. Tarlow, and T. Minka. “A* Sampling”. In: *NIPS*. 2014.
- [MTS22] Y. Ma, D. Tsao, and H.-Y. Shum. “On the Principles of Parsimony and Self-Consistency for the Emergence of Intelligence”. In: (July 2022). arXiv: [2207.04630 \[cs.AI\]](https://arxiv.org/abs/2207.04630).
- [Mua+17] K. Muandet, K. Fukumizu, B. Sriperumbudur, and B. Schölkopf. “Kernel Mean Embedding of Distributions: A Review and Beyond”. In: *Foundations and Trends* 10.1–2 (2017), pp. 1–141.
- [Mua+20] K. Muandet, A. Mehrjou, S. K. Lee, and A. Raj. *Dual Instrumental Variable Regression*. 2020.

- 1
- 2 [Muk+18] S. Mukherjee, D. Shankar, A. Ghosh, N.
3 Tathawadekar, P. Kompalli, S. Sarawagi, and K.
4 Chaudhury. “ARMDN: Associative and Recurrent
5 Mixture Density Networks for eRetail Demand Fore-
6 casting”. In: (2018). arXiv: 1803.03800 [cs.LG].
- 7 [Müll+19a] T. Müller, B. McWilliams, F. Rousselle, M.
8 Gross, and J. Novák. “Neural Importance Sampling”.
9 In: SIGGRAPH. 2019.
- 10 [Müll+19b] T. Müller, B. McWilliams, F. Rousselle, M.
11 Gross, and J. Novák. “Neural importance sampling”.
12 In: ACM Transactions on Graphics 38.5 (2019), pp.
13 p. 145.
- 14 [Mun14] R. Munos. “From Bandits to Monte-Carlo
15 Tree Search: The Optimistic Principle Applied to Opti-
16 mization and Planning”. In: Foundations and Trends
17 in Machine Learning 7.1 (2014), pp. 1–129.
- 18 [Mun+16] R. Munos, T. Stepleton, A. Harutyunyan,
19 and M. G. Bellemare. “Safe and Efficient Off-Policy
20 Reinforcement Learning”. In: NIPS. 2016, pp. 1046–
21 1054.
- 22 [Mun57] J. Munkres. “Algorithms for the assignment
23 and transportation problems”. In: Journal of the society
24 for industrial and applied mathematics 5.1 (1957),
25 pp. 32–38.
- 26 [Mur00] K. Murphy. “Bayesian Map Learning in Dy-
27 namic Environments”. In: NIPS. Vol. 12. 2000.
- 28 [Mur02] K. Murphy. “Dynamic Bayesian Networks:
29 Representation, Inference and Learning”. PhD thesis.
30 Dept. Computer Science, UC Berkeley, 2002.
- 32 [Mur+19] W. J. Murdoch, C. Singh, K. Kumbier, R.
33 Abbasi-Asl, and B. Yu. “Definitions, methods, and ap-
34 plications in interpretable machine learning”. In: Pro-
35 ceedings of the National Academy of Sciences 116.44
36 (2019), pp. 22071–22080.
- 38 [Mur22] K. P. Murphy. *Probabilistic Machine Learn-
39 ing: An introduction*. MIT Press, 2022.
- 40 [MW07] J. Møller and R. P. Waagepetersen. “Modern
41 statistics for spatial point processes”. In: Scandinavian
42 Journal of Statistics 34.4 (2007), pp. 643–684.
- 44 [MW15] S. Morgan and C. Winship. *Counterfactu-
45 als and Causal Inference*. 2nd. Cambridge University
46 Press, 2015.
- 48 [MWJ99] K. Murphy, Y. Weiss, and M. Jordan. “Loopy
49 Belief Propagation for Approximate Inference: an Em-
50 pirical Study”. In: UAI. 1999.
- 52 [MYM18] J. Marino, Y. Yue, and S. Mandt. “Iterative
53 Amortized Inference”. In: ICML. 2018.
- 55 [Nac+17] O. Nachum, M. Norouzi, K. Xu, and D. Schu-
56 urmans. “Bridging the Gap Between Value and Pol-
57 icy Based Reinforcement Learning”. In: NIPS. 2017,
58 pp. 2772–2782.
- 60 [Nac+19a] O. Nachum, Y. Chow, B. Dai, and L.
61 Li. “DualDICE: Behavior-agnostic Estimation of Dis-
62 counted Stationary Distribution Corrections”. In:
63 NeurIPS. 2019, pp. 2315–2325.
- 65 [Nac+19b] O. Nachum, B. Dai, I. Kostrikov, Y. Chow,
66 L. Li, and D. Schuurmans. *Algae: Policy Gradient
67 from Arbitrary Experience*. CoRR abs/1912.02074.
68 2019.
- 70 [Nad+19] S. Naderi, K. He, R. Aghajani, S. Sclaroff,
71 and P. Felzenszwalb. “Generalized Majorization-
72 Minimization”. In: ICML. 2019.
- 74 [Nae+18] C. A. Naeseth, S. W. Linderman, R. Ran-
75 ganath, and D. M. Blei. “Variational Sequential Monte
76 Carlo”. In: AISTATS. 2018.
- 78 [Nal18] E. T. Nalisnick. “On Priors for Bayesian Neural
79 Networks”. PhD thesis. UC Irvine, 2018.
- 80 [Nal+19a] E. Nalisnick, A. Matsukawa, Y. W. Teh, D.
81 Gorur, and B. Lakshminarayanan. “Do Deep Gener-
82 ative Models Know What They Don’t Know?” In:
83 ICLR. 2019.
- 84 [Nal+19b] E. Nalisnick, A. Matsukawa, Y. W. Teh, D.
85 Gorur, and B. Lakshminarayanan. “Hybrid Models
86 with Deep and Invertible Features”. In: ICML. 2019,
87 pp. 4723–4732.
- 89 [Nau04] J. Naudts. “Estimators, escort probabilities
90 and ϕ -exponential families in statistical physics”. In:
91 J. of Inequalities in Pure and Applied Mathematics
92 5.4 (2004).
- 94 [NB05] M. Narasimhan and J. Bilmes. “A Submodular-
95 Supermodular Procedure with Applications to Dis-
96 criminative Structure Learning”. In: Uncertainty in
97 Artificial Intelligence: Proceedings of the Twentieth
98 Conference (UAI-2004). Edinburgh, Scotland: Mor-
99 gan Kaufmann Publishers, 2005.
- 100 [NB06] M. Narasimhan and J. Bilmes. *Learning
101 Graphical Models over partial k-trees*. Tech. rep.
102 UWETR-2006-0001. <https://vanavar.ece.uw.edu/techsite/papers/refer/UWETR-2006-0001.html>. Uni-
103 versity of Washington, Department of Electrical Engi-
104 neering, 2006.
- 105 [NBS18] B. Neyshabur, S. Bhojanapalli, and N. Srebro.
106 “A PAC-Bayesian Approach to Spectrally-Normalized
107 Margin Bounds for Neural Networks”. In: ICLR. 2018.
- 108 [NCH15] M. Naeini, G. Cooper, and M. Hauskrecht.
109 “Obtaining well calibrated probabilities using
110 Bayesian binning”. In: AAAI. 2015.
- 112 [NCL20] T. Nguyen, Z. Chen, and J. Lee. “Dataset
113 Meta-Learning from Kernel Ridge-Regression”. In: In-
114 ternational Conference on Learning Representations.
115 2020.
- 116 [NCT16a] S. Nowozin, B. Cseke, and R. Tomioka. “f-
117 GAN: Training Generative Neural Samplers using
118 Variational Divergence Minimization”. In: NIPS. Ed.
119 by D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon,
120 and R. Garnett. Curran Associates, Inc., 2016, pp. 271–
121 279.
- 123 [NCT16b] S. Nowozin, B. Cseke, and R. Tomioka. “f-
124 gan: Training generative neural samplers using vari-
125 ational divergence minimization”. In: NIPS. 2016,
126 pp. 271–279.
- 127 [NCT16c] S. Nowozin, B. Cseke, and R. Tomioka. “f-
128 gan: Training generative neural samplers using varia-
129 tional divergence minimization”. In: Advances in neu-
130 ral information processing systems. 2016, pp. 271–
131 279.
- 133 [ND20] O. Nachum and B. Dai. “Reinforcement Learn-
134 ing via Fenchel-Rockafellar Duality”. In: (2020). arXiv:
135 2001.01866 [cs.LG].
- 137 [ND21] A. Nichol and P. Dhariwal. “Improved Denois-
138 ing Diffusion Probabilistic Models”. In: ICML. Feb.
139 2021.
- 141 [NDL20] A. Nishimura, D. Dunson, and J. Lu. “Discon-
142 tinuous Hamiltonian Monte Carlo for discrete param-
143 eters and discontinuous likelihoods”. In: Biometrika
144 (2020).
- 146 [Nea00] R. Neal. “Markov Chain Sampling Methods
147 for Dirichlet Process Mixture Models”. In: JCGS 9.2
148 (2000), pp. 249–265.
- 149 [Nea01] R. M. Neal. “Annealed Importance Sampling”.
150 In: Statistics and Computing 11 (2001), pp. 125–139.

- 1
- [Nea03] R. Neal. "Slice sampling". In: *Annals of Statistics* 31.3 (2003), pp. 7–5–767.
- 2
- [Nea+08] R. Neal et al. "Computing likelihood functions for high-energy physics experiments when distributions are defined by simulators with nuisance parameters". In: (2008).
- 3
- [Nea10] R. Neal. "MCMC using Hamiltonian Dynamics". In: *Handbook of Markov Chain Monte Carlo*. Ed. by S. Brooks, A. Gelman, G. Jones, and X.-L. Meng. Chapman & Hall, 2010.
- 4
- [Nea12] R. C. Neath. "On Convergence Properties of the Monte Carlo EM Algorithm". In: *arXiv [math.ST]* (2012).
- 5
- [Nea20] B. Neal. *Introduction to Causal Inference from a Machine Learning Perspective*. 2020.
- 6
- [Nea92] R. Neal. "Connectionist learning of belief networks". In: *Artificial Intelligence* 56 (1992), pp. 71–113.
- 7
- [Nea93] R. M. Neal. *Probabilistic Inference using Markov Chain Monte Carlo Methods*. Tech. rep. CRG-TR-93-1. 144pp. Dept. of Computer Science, University of Toronto, 1993.
- 8
- [Nea96] R. Neal. *Bayesian learning for neural networks*. Springer, 1996.
- 9
- [Nef+02] A. Nefian, L. Liang, X. Pi, X. Liu, and K. Murphy. "Dynamic Bayesian Networks for Audio-Visual Speech Recognition". In: *J. Applied Signal Processing* (2002).
- 10
- [Neg+21] J. Negrea, J. Yang, H. Feng, D. M. Roy, and J. H. Huggins. "Statistical inference with stochastic gradient algorithms". 2021.
- 11
- [Nei+18] D. Neil, J. Briody, A. Lacoste, A. Sim, P. Creed, and A. Saffari. "Interpretable graph convolutional neural networks for inference on noisy knowledge graphs". In: *arXiv preprint arXiv:1812.00279* (2018).
- 12
- [Nel21] Nelson Elhage and Neel Nanda and Catherine Olsson and Tom Henighan and Nicholas Joseph and Ben Mann and Amanda Askell and Yuntao Bai and Anna Chen and Tom Conerly and Nova DasSarma and Dawn Drain and Deep Ganguli and Zac Hatfield-Dodds and Danny Hernandez and Andy Jones and Jackson Kernion and Liane Lovitt and Kamal Ndousse and Dario Amodei and Tom Brown and Jack Clark and Jared Kaplan and Sam McCandlish and Chris Olah. *A Mathematical Framework for Transformer Circuits*. Tech. rep. Anthropic, 2021.
- 13
- [Neu11] G. Neumann. "Variational Inference for Policy Search in Changing Situations". In: *ICML*. 2011, pp. 817–824.
- 14
- [Ney+17] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro. "Exploring generalization in deep learning". In: *NIPS*. 2017.
- 15
- [NF16] M. Noroozi and P. Favaro. "Unsupervised learning of visual representations by solving jigsaw puzzles". In: *European conference on computer vision*. Springer, 2016, pp. 69–84.
- 16
- [NG01] D. Nilsson and J. Goldberger. "Sequentially finding the N-Best List in Hidden Markov Models". In: *IJCAI*. 2001, pp. 1280–1285.
- 17
- [Ng+11] A. Ng et al. "Sparse autoencoder". In: *CS294A Lecture notes 72.2011* (2011), pp. 1–19.
- 18
- [Ngi+11] J. Ngiam, Z. Chen, P. W. Koh, and A. Y. Ng. "Learning deep energy models". In: *Proceedings of the 28th international conference on machine learning (ICML-11)*. 2011, pp. 1105–1112.
- 19
- [Ngu+18] J. Ngiam, D. Peng, V. Vasudevan, S. Kornblith, Q. V. Le, and R. Pang. "Domain adaptive transfer learning with specialist models". In: *arXiv preprint arXiv:1811.07056* (2018).
- 20
- [Ngu+16] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. *Synthesizing the preferred inputs for neurons in neural networks via deep generator networks*. 2016. arXiv: 1605.09304 [cs.NE].
- 21
- [Ngu+18] C. V. Nguyen, Y. Li, T. D. Bui, and R. E. Turner. "Variational Continual Learning". In: *ICLR*. 2018.
- 22
- [Ngu+19] T. T. Nguyen, Q. V. H. Nguyen, D. T. Nguyen, D. T. Nguyen, Thien Huynh-The, S. Nahavandi, T. T. Nguyen, Q.-V. Pham, and C. M. Nguyen. "Deep Learning for Deepfakes Creation and Detection: A Survey". In: (2019). arXiv: 1909.11573 [cs.CV].
- 23
- [Ngu+21] T. Nguyen, R. Novak, L. Xiao, and J. Lee. "Dataset Distillation with Infinitely Wide Convolutional Networks". In: *Advances in Neural Information Processing Systems*. Ed. by A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan. 2021.
- 24
- [NH98a] R. M. Neal and G. E. Hinton. "A new view of the EM algorithm that justifies incremental and other variants". In: *Learning in Graphical Models*. Ed. by M. Jordan. MIT Press, 1998.
- 25
- [NH98b] R. M. Neal and G. E. Hinton. "A View of the EM Algorithm that Justifies Incremental, Sparse, and other Variants". In: *Learning in Graphical Models*. Ed. by M. I. Jordan. Springer Netherlands, 1998, pp. 355–368.
- 26
- [NHL19] E. Nalisnick, J. M. Hernández-Lobato, and P. Smyth. "Dropout as a Structured Shrinkage Prior". In: *ICML*. 2019.
- 27
- [NHR99] A. Ng, D. Harada, and S. Russell. "Policy invariance under reward transformations: Theory and application to reward shaping". In: *ICML*. 1999.
- 28
- [NI92] H. Nagamochi and T. Ibaraki. "Computing edge-connectivity of multigraphs and capacitated graphs". In: *SIAM J. Discrete Math.* 5 (1992), pp. 54–66.
- 29
- [Nic+21] A. Nichol, P. Dhariwal, P. Ramesh Aditya and Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen. "GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models". In: (2021). arXiv: 2112.10741 [cs.CV].
- 30
- [Nig+00] K. Nigam, A. K. McCallum, S. Thrun, and T. Mitchell. "Text Classification from Labeled and Unlabeled Documents using EM". In: *MLJ* 39.2 (May 2000), pp. 103–134.
- 31
- [Nij+19] E. Nijkamp, M. Hill, S.-C. Zhu, and Y. N. Wu. "Learning non-convergent non-persistent short-run MCMC toward energy-based model". In: *NIPS*. 2019, pp. 5232–5242.
- 32
- [Nix+19] J. Nixon, M. Dusenberry, L. Zhang, G. Jerfel, and D. Tran. "Measuring Calibration in Deep Learning". In: (2019). arXiv: 1904.01685 [cs.LG].
- 33
- [NJ00] A. Y. Ng and M. Jordan. "PEGASUS: A policy search method for large MDPs and POMDPs". In: *UAI*. 2000.
- 34
- [NJB05] M. Narasimhan, N. Jojic, and J. A. Bilmes. "Q-clustering". In: *Advances in Neural Information Processing Systems* 18 (2005), pp. 979–986.
- 35
- [NK17] V. Nagarajan and J. Z. Kolter. "Gradient descent GAN optimization is locally stable". In: *Ad-*

- 1
2 *vances in neural information processing systems.*
3 2017, pp. 5585–5595.
- 4 [NKI10] K. Nagano, Y. Kawahara, and S. Iwata. “Minimum average cost clustering”. In: *Advances in Neural Information Processing Systems* 23 (2010), pp. 1759–1767.
- 5 [NLS15] C. A. Naesseth, F. Lindsten, and T. B. Schön. “Nested Sequential Monte Carlo Methods”. In: *ICML*. 2015.
- 6 [NLS19] C. A. Naesseth, F. Lindsten, and T. B. Schön. “Elements of Sequential Monte Carlo”. In: *Foundations and Trends in Machine Learning* (2019).
- 7 [NM12] A. Nenkova and K. McKeown. “A survey of text summarization techniques”. In: *Mining text data*. Springer, 2012, pp. 43–76.
- 8 [NMC05] A. Niculescu-Mizil and R. Caruana. “Predicting Good Probabilities with Supervised Learning”. In: *ICML*. 2005.
- 9 [NNP19] W. Nie, N. Narodytska, and A. Patel. “RelGAN: Relational Generative Adversarial Networks for Text Generation”. In: *International Conference on Learning Representations*. 2019.
- 10 [Noc+21] L. Noci, K. Roth, G. Bachmann, S. Nowozin, and T. Hofmann. “Disentangling the Roles of Curation, Data-Augmentation and the Prior in the Cold Posterior Effect”. In: *NIPS*. 2021.
- 11 [Noé+19] F. Noé, S. Olsson, J. Köhler, and H. Wu. “Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning”. In: *Science* 365 (2019).
- 12 [Nou+02] M. N. Nounou, B. R. Bakshi, P. K. Goel, and X. Shen. “Process modeling by Bayesian latent variable regression”. In: *Am. Inst. Chemical Engineers Journal* 48.8 (2002), pp. 1775–1793.
- 13 [Nov+19] R. Novak, L. Xiao, J. Lee, Y. Bahri, G. Yang, J. Hron, D. A. Abolafia, J. Pennington, and J. Sohl-Dickstein. “Bayesian Deep Convolutional Networks with Many Channels are Gaussian Processes”. In: *ICLR*. 2019.
- 14 [NP03] W. K. Newey and J. L. Powell. “Instrumental variable estimation of nonparametric models”. In: *Econometrica* 71.5 (2003), pp. 1565–1578.
- 15 [NR00a] A. Ng and S. Russell. “Algorithms for inverse reinforcement learning”. In: *ICML*. 2000.
- 16 [NR00b] A. Y. Ng and S. Russell. “Algorithms for Inverse Reinforcement Learning”. In: *in Proc. 17th International Conf. on Machine Learning*. Citeseer. 2000.
- 17 [NR94] M. Newton and A. Raftery. “Approximate Bayesian Inference with the Weighted Likelihood Bootstrap”. In: *J. of Royal Stat. Soc. Series B* 56.1 (1994), pp. 3–48.
- 18 [NS01] K. Nowicki and T. A. B. Snijders. “Estimation and Prediction for Stochastic Blockstructures”. In: *Journal of the American Statistical Association* 96.455 (2001), pp. 1077–1087.
- 19 [NS17] E. Nalisnick and P. Smyth. “Variational Reference Priors”. In: *ICLR Workshop*. 2017.
- 20 [NS18] E. Nalisnick and P. Smyth. “Learning Priors for Invariance”. In: *AISTATS*. 2018.
- 21 [NW06] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, 2006.
- 22 [NW20] X Nie and S Wager. “Quasi-oracle estimation of heterogeneous treatment effects”. In: *Biometrika* 108.2 (Sept. 2020), pp. 299–319. eprint: <https://academic.oup.com/biomet/article-pdf/108/2/299/37938939/asaa076.pdf>.
- 23 [NWF78] G. Nemhauser, L. Wolsey, and M. Fisher. “An analysis of approximations for maximizing submodular set functions—I”. In: *Mathematical Programming* 14.1 (1978), pp. 265–294.
- 24 [NWJ09] X. Nguyen, M. J. Wainwright, and M. I. Jordan. “On Surrogate Loss Functions and f-Divergences”. In: *Ann. Stat.* 37.2 (2009), pp. 876–904.
- 25 [NWJ+09] X. Nguyen, M. J. Wainwright, M. I. Jordan, et al. “On surrogate loss functions and f-divergences”. In: *The Annals of Statistics* 37.2 (2009), pp. 876–904.
- 26 [NWJ10a] X. Nguyen, M. J. Wainwright, and M. I. Jordan. “Estimating Divergence Functionals and the Likelihood Ratio by Convex Risk Minimization”. In: *IEEE Trans. Inf. Theory* 56.11 (2010), pp. 5847–5861.
- 27 [NWJ10b] X. Nguyen, M. J. Wainwright, and M. I. Jordan. “Estimating divergence functionals and the likelihood ratio by convex risk minimization”. In: *IEEE Transactions on Information Theory* 56.11 (2010), pp. 5847–5861.
- 28 [NYC15] A. Nguyen, J. Yosinski, and J. Clune. “Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images”. In: *CVPR*. 2015.
- 29 [OAA09] M. Opper and C. Archambeau. “The variational Gaussian approximation revisited”. In: *Neural Comput.* 21.3 (2009), pp. 786–792.
- 30 [OAC18] I. Osband, J. Aslanides, and A. Cassirer. “Randomized prior functions for deep reinforcement learning”. In: *NIPS*. 2018.
- 31 [Obe+19] F. Obermeyer, E. Bingham, M. Jankowiak, J. Chiu, N. Pradhan, A. Rush, and N. Goodman. “Tensor Variable Elimination for Plated Factor Graphs”. In: *ICML*. 2019.
- 32 [OCM21] L. A. Ortega, R. Cabañas, and A. R. Masegosa. “Diversity and Generalization in Neural Network Ensembles”. In: (2021). arXiv: [2110.13786](https://arxiv.org/abs/2110.13786) [cs.LG].
- 33 [ODK96] M. Ostendorf, V. Digalakis, and O. Kimball. “From HMMs to segment models: a unified view of stochastic modeling for speech recognition”. In: *IEEE Trans. on Speech and Audio Processing* 4.5 (1996), pp. 360–378.
- 34 [OED21] J. Ortiz, T. Evans, and A. J. Davison. “A visual introduction to Gaussian Belief Propagation”. In: *arXiv preprint arXiv:2107.02308* (2021).
- 35 [OF96] B. A. Olshausen and D. J. Field. “Emergence of simple cell receptive field properties by learning a sparse code for natural images”. In: *Nature* 381 (1996), pp. 607–609.
- 36 [O'H78] A. O'Hagan. “Curve Fitting and Optimal Design for Prediction”. In: *J. of Royal Stat. Soc. Series B* 40 (1978), pp. 1–42.
- 37 [OKK16] A. Van den Oord, N. Kalchbrenner, and K. Kavukcuoglu. “Pixel Recurrent Neural Networks”. In: *ICML*. 2016.
- 38 [Oll+17] Y. Ollivier, L. Arnold, A. Auger, and N. Hansen. “Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles”. In: *JMLR* 18 (2017), pp. 1–65.

- 1 [Oll18] Y. Ollivier. “Online natural gradient as a
2 Kalman filter”. en. In: *Electron. J. Stat.* 12.2 (2018),
3 pp. 2930–2961.
- 4 [OLV18a] A. van den Oord, Y. Li, and O. Vinyals.
5 “Representation Learning with Contrastive Predictive
Coding”. In: (2018). arXiv: [1807.03748 \[cs.LG\]](https://arxiv.org/abs/1807.03748).
- 6 [OLV18b] A. v. d. Oord, Y. Li, and O. Vinyals. “Representation learning with contrastive predictive coding”.
7 In: *arXiv preprint arXiv:1807.03748* (2018).
- 8 [OM96] P. V. Overschee and B. D. Moor. *Subspace
9 Identification for Linear Systems: Theory, Implementation,
10 Application, Applications*. Kluwer Academic Publishers, 1996.
- 11 [OMS17] C. Olah, A. Mordvintsev, and L. Schubert.
12 “Feature Visualization”. In: *Distill* 2.11 (2017).
- 13 [O’N09] B. O’Neill. “Exchangeability, Correlation, and
14 Bayes’ Effect”. In: *Int. Stat. Rev.* 77.2 (2009), pp. 241–
250.
- 15 [ONS18] V. M.-H. Ong, D. J. Nott, and M. S. Smith.
16 “Gaussian Variational Approximation With a Factor
Covariance Structure”. In: *J. Comput. Graph. Stat.*
17 27.3 (2018), pp. 465–478.
- 18 [oor+16] A. Van den oord, S. Dieleman, H. Zen, K. Si-
19 monyan, O. Vinyals, A. Graves, N. Kalchbrenner, A.
Senior, and K. Kavukcuoglu. “WaveNet: A Generative
20 Model for Raw Audio”. In: (2016). arXiv: [1609.03499 \[cs.SD\]](https://arxiv.org/abs/1609.03499).
- 21 [Oor+16] A. van den Oord, N. Kalchbrenner, O.
22 Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu.
“Conditional Image Generation with PixelCNN De-
coders”. In: (2016). arXiv: [1606.05328 \[cs.CV\]](https://arxiv.org/abs/1606.05328).
- 23 [Oor+18] A. van den Oord et al. “Parallel WaveNet:
24 Fast High-Fidelity Speech Synthesis”. In: *ICML*. Ed.
by J. Dy and A. Krause. Vol. 80. Proceedings of Ma-
chine Learning Research. PMLR, 2018, pp. 3918–3926.
- 25 [Oor+19] A. van den Oord, B. Poole, O. Vinyals, and A.
26 Razavi. “Fixing Posterior Collapse with delta-VAEs”.
In: *ICLR*. 2019.
- 27 [OOS17] A. Odena, C. Olah, and J. Shlens. “Condi-
28 tional image synthesis with auxiliary classifier gans”.
In: *International conference on machine learning*.
2017, pp. 2642–2651.
- 29 [Opt88] Optimal information processing and Bayes’
30 theorem. In: *The American Statistician* 42.4 (1988),
pp. 278–280.
- 31 [Oqu+14] M. Oquab, L. Bottou, I. Laptev, and J. Sivic.
“Learning and transferring mid-level image represen-
32 tations using convolutional neural networks”. In: *Pro-
ceedings of the IEEE conference on computer vision
and pattern recognition*. 2014, pp. 1717–1724.
- 33 [OR20] A. Owen and D. Rudolf. “A strong law of large
34 numbers for scrambled net integration”. In: (2020).
arXiv: [2002.07859 \[math.NA\]](https://arxiv.org/abs/2002.07859).
- 35 [Ore+20] B. N. Oreshkin, D. Carpov, N. Chapados,
36 and Y. Bengio. “N-BEATS: Neural basis expansion
analysis for interpretable time series forecasting”. In:
37 *ICLR*. 2020.
- 38 [ORW21] S. W. Ober, C. E. Rasmussen, and M. van
39 der Wilk. “The Promises and Pitfalls of Deep Kernel
Learning”. In: *ICML*. Feb. 2021.
- 40 [Osa+18] T. Osa, J. Pajarinen, G. Neumann, J. A. Bag-
41 nell, P. Abbeel, and J. Peters. “An Algorithmic Per-
42 spective on Imitation Learning”. In: *Foundations and
43 Trends in Robotics* 7.1–2 (2018), pp. 1–179.
- 44 [Osa+19a] K. Osawa, S. Swaroop, A. Jain, R. Eschen-
45 hagen, R. E. Turner, R. Yokota, and M. E. Khan.
“Practical Deep Learning with Bayesian Principles”.
In: *NIPS*. 2019.
- 46 [Osa+19b] K. Osawa, Y. Tsuji, Y. Ueno, A. Naruse,
R. Yokota, and S. Matsuoka. “Large-Scale Dis-
tributed Second-Order Optimization Using Kronecker-
Factored Approximate Curvature for Deep Convolu-
tional Neural Networks”. In: *CVPR*. 2019.
- [Osb16] I. Osband. “Risk versus Uncertainty in Deep
Learning: Bayes, Bootstrap and the Dangers of
Dropout”. In: *NIPS workshop on Bayesian deep
learning*. 2016.
- [Osb+21] I. Osband, Z. Wen, S. M. Asghari, V.
Dwarkerla, B. Hao, M. Ibrahim, D. Lawson, X.
Lu, B. O’Donoghue, and B. Van Roy. “The Neu-
ral Testbed: Evaluating Predictive Distributions”. In:
(2021). arXiv: [2110.04629 \[cs.LG\]](https://arxiv.org/abs/2110.04629).
- [Ose11] I. V. Oseledets. “Tensor-Train Decomposition”.
In: *SIAM J. Sci. Comput.* 33.5 (2011), pp. 2295–2317.
- [OT05] A. B. Owen and S. D. Tribble. “A quasi-Monte
Carlo Metropolis algorithm”. en. In: *PNAS* 102.25
(2005), pp. 8844–8849.
- [Ova+19] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D Scul-
ley, S. Nowozin, J. V. Dillon, B. Lakshminarayanan,
and J. Snoek. “Can You Trust Your Model’s Un-
certainty? Evaluating Predictive Uncertainty Under
Dataset Shift”. In: *NIPS*. 2019.
- [OVK17] A. van den Oord, O. Vinyals, and K.
Kavukcuoglu. “Neural Discrete Representation Learn-
ing”. In: *NIPS*. 2017.
- [Owe13] A. B. Owen. *Monte Carlo theory, methods
and examples*. 2013.
- [Owe17] A. B. Owen. “A randomized Halton algorithm
in R”. In: *arXiv [stat.CO]* (2017).
- [Oxl11] J. Oxley. *Matroid Theory: Second Edition*. Ox-
ford University Press, 2011.
- [Pac+14] J. Pacheco, S. Zuffi, M. Black, and E. Sud-
derth. “Preserving Modes and Messages via Diverse
Particle Selection”. en. In: *ICML*. 2014, pp. 1152–
1160.
- [Pai] *Explainable AI in Practice Falls Short of Trans-
parency Goals*. [https://partnershiponai.org/xai-in-
practice/](https://partnershiponai.org/xai-in-
practice/). Accessed: 2021-11-23.
- [Pai+14] T. L. Paine, P. Khorrami, W. Han, and
T. S. Huang. “An analysis of unsupervised pre-training
in light of recent advances”. In: *arXiv preprint
arXiv:1412.6597* (2014).
- [Pan+10] L. Paninski, Y. Ahmadian, D. G. Ferreira, S.
Koyama, K. Rahnama Rad, M. Vidne, J. Vogelstein,
and W. Wu. “A new look at state-space models for neu-
ral data”. en. In: *J. Comput. Neurosci.* 29.1–2 (2010),
pp. 107–126.
- [Pan+21] G. Pang, C. Shen, L. Cao, and A. Van Den
Hengel. “Deep Learning for Anomaly Detection: A Re-
view”. In: *ACM Comput. Surv.* 54.2 (2021), pp. 1–38.
- [Pap+17] N. Papernot, P. McDaniel, I. Goodfellow, S.
Jha, Z Berkay Celik, and A. Swami. “Practical Black-
Box Attacks against Deep Learning Systems using Ad-
versarial Examples”. In: *ACM Asia Conference on
Computer and Communications Security*. 2017.
- [Pap+19] G. Papamakarios, E. Nalisnick, D. J.
Rezende, S. Mohamed, and B. Lakshminarayanan.
“Normalizing Flows for Probabilistic Modeling and In-
ference”. In: (2019). arXiv: [1912.02762 \[stat.ML\]](https://arxiv.org/abs/1912.02762).

- 1
- 2 [Par+19] S Parameswaran, C Deledalle, L Denis, and
T. Q. Nguyen. “Accelerating GMM-Based Patch Priors for Image Restoration: Three Ingredients for a
100× Speed-Up”. In: *IEEE Trans. Image Process.*
28.2 (2019), pp. 687–698.
- 3
- 4
- 5 [Par81] G. Parisi. “Correlation functions and computer simulations”. In: *Nuclear Physics B* 180.3 (1981),
pp. 378–384.
- 6
- 7 [Pas+02] H. Pasula, B. Marthi, B. Milch, S. Russell,
and I. Shpitser. “Identity Uncertainty and Citation Matching”. In: *NIPS*. 2002.
- 8
- 9 [Pas+21a] A. Paszke, D. Johnson, D. Duvenaud, D.
Vytiniotis, A. Radul, M. Johnson, J. Ragan-Kelley,
and D. Maclaurin. “Getting to the Point: Index Sets
and Parallelism-Preserving Autodiff for Painful Array Programming”. In: *Proc. ACM Program. Lang.*
5.ICFP (2021).
- 10
- 11
- 12 [Pas+21b] A. Paszke, M. J. Johnson, R. Frostig, and
D. Maclaurin. “Parallelism-Preserving Automatic Differentiation for Second-Order Array Languages”. In:
Proceedings of the 9th ACM SIGPLAN International Workshop on Functional High-Performance and Numerical Computing. FHPNC 2021. Association for Computing Machinery, 2021, 13–23.
- 13
- 14 [Pat+16] D. Pathak, P. Krahenbuhl, J. Donahue, T.
Darrell, and A. A. Efros. “Context encoders: Feature learning by inpainting”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 2536–2544.
- 15
- 16 [PB14] R. Pascanu and Y. Bengio. “Revisiting Natural Gradient for Deep Networks”. In: *ICLR*. 2014.
- 17
- 18 [PBM16a] J. Peters, P. Bühlmann, and N. Meinshausen. “Causal inference by using invariant prediction: identification and confidence intervals”. In: *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 78.5 (2016), pp. 947–1012.
- 19
- 20 [PBM16b] J. Peters, P. Bühlmann, and N. Meinshausen. “Causal inference using invariant prediction: identification and confidence intervals”. In: *J. of Royal Stat. Soc. Series B* 78.5 (2016), pp. 947–1012.
- 21
- 22 [PC08] T. Park and G. Casella. “The Bayesian Lasso”. In: *JASA* 103.482 (2008), pp. 681–686.
- 23
- 24 [PC09] J. Paisley and L. Carin. “Nonparametric Factor Analysis with Beta Process Priors”. In: *ICML*. 2009.
- 25
- 26 [PC12] N. Pinto and D. D. Cox. “High-throughput derived biologically-inspired features for unconstrained face recognition”. In: *Image Vis. Comput.* 30.3 (2012), pp. 159–168.
- 27
- 28 [PD03] J. D. Park and A. Darwiche. “A Differential Semantics for Jontree Algorithms”. In: *NIPS*. MIT Press, 2003, pp. 801–808.
- 29
- 30 [PD11] H. Poon and P. Domingos. “Sum-Product Networks: A New Deep Architecture”. In: *UAI*. Java code at <http://alchemy.cs.washington.edu/spn/>. Short intro at <http://lessoned.blogspot.com/2011/10/intro-to-sum-product-networks.html>. 2011.
- 31
- 32 [PdC20] F.-P. Paty, A. d’Aspremont, and M. Cuturi. “Regularity as Regularization: Smooth and Strongly Convex Brenier Potentials in Optimal Transport”. In: *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*. Ed. by S. Chiappa and R. Calandra. Vol. 108. Proceedings of Machine Learning Research. PMLR, 2020, pp. 1222–1232.
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- [PDL+12] M. Parry, A. P. Dawid, S. Lauritzen, et al. “Proper local scoring rules”. In: *The Annals of Statistics* 40.1 (2012), pp. 561–592.
- [PE16] V. Papyan and M. Elad. “Multi-Scale Patch-Based Image Restoration”, en. In: *IEEE Trans. Image Process.* 25.1 (2016), pp. 249–261.
- [Pea09a] J. Pearl. *Causality*. 2nd. Cambridge University Press, 2009.
- [Pea09b] J. Pearl. *Causality: Models, Reasoning and Inference (Second Edition)*. Cambridge Univ. Press, 2009.
- [Pea09c] J. Pearl. “Causal inference in statistics: An overview”. In: *Stat. Surv.* 3.0 (2009), pp. 96–146.
- [Pea12] J. Pearl. “The Causal Foundations of Structural Equation Modeling”. In: *Handbook of structural equation modeling*. Ed. by R. H. Hoyle. Vol. 68. 2012.
- [Pea19] J. Pearl. “The Seven Tools of Causal Inference, with Reflections on Machine Learning”. In: *Comm. of the ACM* 62.3 (2019), pp. 54–60.
- [Pea36] K. Pearson. “Method of moments and method of maximum likelihood”. In: *Biometrika* 28.1/2 (1936), pp. 34–59.
- [Pea84] J. Pearl. *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Addison-Wesley Longman Publishing Co., Inc., 1984.
- [Pea88] J. Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann, 1988.
- [Pea94] B. A. Pearlmutter. “Fast Exact Multiplication by the Hessian”. In: *Neural Comput.* 6.1 (1994), pp. 147–160.
- [Peh+20] R. Peherz, S. Lang, A. Vergari, K. Stelzner, A. Molina, M. Trapp, G. Van den Broeck, K. Kersting, and Z. Ghahramani. “Einsum Networks: Fast and Scalable Learning of Tractable Probabilistic Circuits”. In: (2020). arXiv: [2004.06231 \[cs.LG\]](https://arxiv.org/abs/2004.06231).
- [Pel05] M. Pelikan. *Hierarchical Bayesian Optimization Algorithm: Toward a New Generation of Evolutionary Algorithms*. en. Softcover reprint of hardcover 1st ed. 2005 edition. Springer, 2005.
- [Pen13] J. Pena. “Reading dependencies from covariance graphs”. In: *Intl. J. of Approximate Reasoning* 54.1 (2013).
- [Per+18] E. Perez, F. Strub, H. de Vries, V. Dumoulin, and A. Courville. “FiLM: Visual Reasoning with a General Conditioning Layer”. In: *AAAI*. 2018.
- [Per+21] V. Perrone, M. Donini, M. B. Zafar, R. Schmucker, K. Kenthapadi, and C. Archambeau. “Fair Bayesian Optimization”. In: *AAAI/ACM Conference on AI, Ethics, and Society (AIES ’21)*. 2021.
- [Pes+21] H. Pesonen et al. “ABC of the Future”. In: (2021). arXiv: [2112.12841 \[stat.AP\]](https://arxiv.org/abs/2112.12841).
- [Pey20] G. Peyre. “Course notes on Optimization for Machine Learning”. 2020.
- [Pez+21] M. Pezeshki, S.-O. Kaba, Y. Bengio, A. Courville, D. Precup, and G. Lajoie. “Gradient Starvation: A Learning Proclivity in Neural Networks”. In: *NIPS*. 2021.
- [PF03] G. V. Puskorius and L. A. Feldkamp. “Parameter-based Kalman filter training: Theory and implementation”. In: *Kalman Filtering and Neural Networks*. John Wiley & Sons, Inc., 2003, pp. 23–67.

- [PF91] G. V. Puskorius and L. A. Feldkamp. “Decoupled extended Kalman filter training of feedforward layered networks”. In: *International Joint Conference on Neural Networks*. Vol. i. 1991, 771–777 vol.1.
- [PFW21] R. Prado, M. Ferreira, and M. West. *Time Series: Modelling, Computation and Inference (2nd ed)*. CRC Press, 2021.
- [PG98] M. Popescu and P. D. Gader. “Image content retrieval from image databases using feature integration by Choquet integral”. In: *Storage and Retrieval for Image and Video Databases VII*. Ed. by M. M. Yeung, B.-L. Yeo, and C. A. Bouman. Vol. 3656. International Society for Optics and Photonics. SPIE, 1998, pp. 552 – 560.
- [PGCP00] M Pelikan, D. E. Goldberg, and E Cantú-Paz. “Linkage problem, distribution estimation, and Bayesian networks”. en. In: *Evol. Comput.* 8.3 (2000), pp. 311–340.
- [PGJ16] J. Pearl, M. Glymour, and N. Jewell. *Causal inference in statistics: a primer*. Wiley, 2016.
- [PHL12] M. Pelikan, M. Hausschild, and F. Lobo. *Introduction to estimation of distribution algorithms*. Tech. rep. U. Missouri, 2012.
- [PHR18] E. Petersen, C. Hoffmann, and P. Rostalski. “On Approximate Nonlinear Gaussian Message Passing On Factor Graphs”. In: *IEEE Statistical Signal Processing Workshop (SSP)*. 2018.
- [Phu+18] M. Phuong, M. Welling, N. Kushman, R. Tomioka, and S. Nowozin. “The Mutual Autoencoder: Controlling Information in Latent Code Representations”. In: *Arxiv* (2018).
- [Pir+13] M. Pirotta, M. Restelli, A. Pecorino, and D. Calandriello. “Safe Policy Iteration”. In: *ICML*. 3. 2013, pp. 307–317.
- [PJD21] Y. Petetin, Y. Janati, and F. Desbouvries. “Structured Variational Bayesian Inference for Gaussian State-Space Models With Regime Switching”. In: *IEEE Signal Process. Lett.* 28 (2021), pp. 1953–1957.
- [PJS17] J. Peters, D. Janzing, and B. Schölkopf. *Elements of Causal Inference: Foundations and Learning Algorithms (Adaptive Computation and Machine Learning series)*. The MIT Press, 2017.
- [PKP21] A. Plaat, W. Kosters, and M. Preuss. “High-Accuracy Model-Based Reinforcement Learning, a Survey”. In: (2021). arXiv: [2107.08241 \[cs.LG\]](#).
- [PL03] M. A. Paskin and G. D. Lawrence. *Junction Tree Algorithms for Solving Sparse Linear Systems*. Tech. rep. UCB/CSD-03-1271. UC Berkeley, 2003.
- [Pla00] J. Platt. “Probabilities for SV machines”. In: *Advances in Large Margin Classifiers*. Ed. by A. Smola, P. Bartlett, B. Schoelkopf, and D. Schuurmans. MIT Press, 2000.
- [Pla18] E. Plaut. “From Principal Subspaces to Principal Components with Linear Autoencoders”. In: *ArXiv* abs/1804.10253 (2018).
- [Ple+18] G. Pleiss, J. R. Gardner, K. Q. Weinberger, and A. G. Wilson. “Constant-Time Predictive Distributions for Gaussian Processes”. In: *International Conference on Machine Learning*. 2018.
- [Plu+20] G. Plumb, M. Al-Shedivat, Á. A. Cabrera, A. Perer, E. Xing, and A. Talwalkar. “Regularizing black-box models for improved interpretability”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 33.
- [PM18a] N. Papernot and P. McDaniel. *Deep k-Nearest Neighbors: Towards Confident, Interpretable and Robust Deep Learning*. 2018. arXiv: [1803.04765 \[cs.LG\]](#).
- [PM18b] J. Pearl and D. Mackenzie. *The book of why: the new science of cause and effect*. 2018.
- [PMT18] G. Plumb, D. Molitor, and A. Talwalkar. “Supervised Local Modeling for Interpretability”. In: *CoRR* abs/1807.02910 (2018). arXiv: [1807.02910](#).
- [Pol+19] A. A. Pol, V. Berger, G. Cerminali, C. Germain, and M. Pierini. “Anomaly Detection With Conditional Variational Autoencoders”. In: *IEEE International Conference on Machine Learning and Applications*. 2019.
- [Pom89] D. Pomerleau. “ALVINN: An Autonomous Land Vehicle in a Neural Network”. In: *NIPS*. 1989, pp. 305–313.
- [Poo+12] D. Poole, D. Buchman, S. Natarajan, and K. Kersting. “Aggregation and Population Growth: The Relational Logistic Regression and Markov Logic Cases”. In: *Statistical Relational AI workshop*. 2012.
- [Poo+19a] B. Poole, S. Ozair, A. van den Oord, A. A. Alemi, and G. Tucker. “On Variational Bounds of Mutual Information”. In: *ICML*. 2019.
- [Poo+19b] B. Poole, S. Ozair, A. van den Oord, A. A. Alemi, and G. Tucker. “On variational lower bounds of mutual information”. In: *ICML*. 2019.
- [Pou04] M. Pourahmadi. *Simultaneous Modelling of Covariance Matrices: GLM, Bayesian and Nonparametric Perspectives*. Tech. rep. Northern Illinois University, 2004.
- [Roy+20] R. Poyiadzi, K. Sokol, R. Santos-Rodriguez, T. De Bie, and P. Flach. “FACE: Feasible and actionable counterfactual explanations”. In: *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. 2020, pp. 344–350.
- [PPC09] G. Petris, S. Petrone, and P. Campagnoli. *Dynamic linear models with R*. Springer, 2009.
- [PPG91] C. S. Pomerleau, O. F. Pomerleau, and A. W. Garcia. “Biobehavioral research on nicotine use in women”. In: *British Journal of Addiction* 86.5 (1991), pp. 527–531.
- [PPM17] G. Papamakarios, T. Pavlakou, and I. Murray. “Masked Autoregressive Flow for Density Estimation”. In: *NIPS*. 2017.
- [PPS18] T. Pierrot, N. Perrin, and O. Sigaud. “First-order and second-order variants of the gradient descent in a unified framework”. In: (2018). arXiv: [1810.08102 \[cs.LG\]](#).
- [PR03] O. Papaspiliopoulos and G. O. Roberts. “Non-Centered Parameterisations for Hierarchical Models and Data Augmentation”. In: *Bayesian Statistics 7* (2003), pp. 307–326.
- [Pra+18] S. Prabhhumoye, Y. Tsvetkov, R. Salakhutdinov, and A. W. Black. “Style Transfer Through Back-Translation”. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2018, pp. 866–876.
- [Pre05] S. J. Press. *Applied multivariate analysis, using Bayesian and frequentist methods of inference*. Second edition. Dover, 2005.
- [Pre+17a] V. Premachandran, D. Tarlow, A. L. Yuille, and D. Batra. “Empirical Minimum Bayes Risk Prediction”. en. In: *IEEE PAMI* 39.1 (Jan. 2017), pp. 75–86.

- 1
- 2 [Pre+17b] O. Press, A. Bar, B. Bogin, J. Berant, and L.
3 Wolf. "Language generation with recurrent generative
4 adversarial networks without pre-training". In: *arXiv
5 preprint arXiv:1706.01399* (2017).
- 6 [Pre+88] W. Press, W. Vetterling, S. Teukolosky, and
7 B. Flannery. *Numerical Recipes in C: The Art of
8 Scientific Computing*. Second. Cambridge University
9 Press, 1988.
- 10 [PRG17] M. Probst, F. Rothlauf, and J. Grahl. "Scal-
11 ability of using Restricted Boltzmann Machines for
12 combinatorial optimization". In: *Eur. J. Oper. Res.*
13 256.2 (2017), pp. 368–383.
- 14 [Pri58] R. Price. "A useful theorem for nonlinear de-
15 vices having Gaussian inputs". In: *IRE Trans. Info.
16 Theory* 4.2 (1958), pp. 69–72.
- 17 [PS07] J. Peters and S. Schaal. "Reinforcement Learn-
18 ing by Reward-Weighted Regression for Operational
19 Space Control". In: *ICML*. 2007, pp. 745–750.
- 20 [PS08a] B. A. Pearlmutter and J. M. Siskind. "Reverse-
21 Mode AD in a Functional Framework: Lambda the
22 Ultimate Backpropagator". In: *ACM Trans. Program.
Lang. Syst.* 30.2 (2008).
- 23 [PS08b] J. Peters and S. Schaal. "Reinforcement Learn-
24 ing of Motor Skills with Policy Gradients". In: *Neural
Networks* 21.4 (2008), pp. 682–697.
- 25 [PS12] N. G. Polson and J. G. Scott. "On the Half-
26 Cauchy Prior for a Global Scale Parameter". en. In:
27 *Bayesian Anal.* 7.4 (2012), pp. 887–902.
- 28 [PS17] N. G. Polson and V. Sokolov. "Deep Learning:
29 A Bayesian Perspective". en. In: *Bayesian Anal.* 12.4
30 (2017), pp. 1275–1304.
- 31 [PSCP06] M. Pelikan, K. Sastry, and E. Cantú-Paz.
32 "Scalable Optimization via Probabilistic Modeling:
33 From Algorithms to Applications (Studies in Compu-
34 tational Intelligence)". Springer-Verlag New York, Inc.,
2006.
- 35 [PSD00] J. K. Pritchard, M. Stephens, and P. Donnelly.
36 "Inference of population structure using multilocus
37 genotype data". In: *Genetics* 155.2 (2000), pp. 945–
959.
- 38 [PSDG14] B. Poole, J. Sohl-Dickstein, and S. Ganguli.
39 "Analyzing noise in autoencoders and deep networks".
40 In: *arXiv preprint arXiv:1406.1831* (2014).
- 41 [PSM19] G. Papamakarios, D. Sterratt, and I. Mur-
42 ray. "Sequential Neural Likelihood: Fast Likelihood-
43 free Inference with Autoregressive Flows". In: *AIS-
44 TATS*. 2019.
- 45 [PSS00] D. Precup, R. S. Sutton, and S. P. Singh. "El-
igibility Traces for Off-Policy Policy Evaluation". In:
46 *ICML*. ICML '00. Morgan Kaufmann Publishers Inc.,
2000, pp. 759–766.
- 47 [PT13] S. Patterson and Y. W. Teh. "Stochastic Gradi-
48 ent Riemannian Langevin Dynamics on the Probabil-
49 ity Simplex". In: *NIPS*. 2013.
- 50 [PT87] C. Papadimitriou and J. Tsitsiklis. "The com-
plexity of Markov decision processes". In: *Mathemat-
ics of Operations Research* 12.3 (1987), pp. 441–450.
- 51 [PT94] P. Paatero and U. Tapper. "Positive Matrix Fac-
torization: A Non-negative Factor Model with Opti-
mal Utilization of Error Estimates of Data Values". In:
52 *Environmetrics* 5 (1994), pp. 111–126.
- 53 [PTD20] A. Prabhu, P. H. S. Torr, and P. K. Dokan-
nia. "GDumb: A simple approach that questions our
54 progress in continual learning". In: *ECCV*. Lecture
notes in computer science. Springer International Pub-
55 lishing, 2020, pp. 524–540.
- 56 [Put94] M. L. Puterman. *Markov Decision Processes:
Discrete Stochastic Dynamic Programming*. Wiley,
1994.
- 57 [PVC19] R. Prenger, R. Valle, and B. Catanzaro.
"WaveGLOW: A flow-based generative network for
58 speech synthesis". In: *Proceedings of the 2019 IEEE
International Conference on Acoustics, Speech and
Signal Processing*. IEEE, 2019, pp. 3617–3621.
- 59 [PW05] S. Parise and M. Welling. "Learning in Markov
Random Fields: An Empirical Study". In: *Joint Sta-
tistical Meeting*. 2005.
- 60 [PY10] G. Papandreou and A. L. Yuille. "Gaussian
sampling by local perturbations". In: *NIPS*. 2010.
- 61 [PY11] G. Papandreou and A. L. Yuille. "Perturb-and-
MAP random fields: Using discrete optimization to
learn and sample from energy models". In: *ICCV*. Nov.
2011, pp. 193–200.
- 62 [PY14] G. Papandreou and A. Yuille. "Perturb-and-
MAP Random Fields: Reducing Random Sampling
to Optimization, with Applications in Computer Vi-
sion". In: *Advanced Structured Prediction*. Ed. by S.
Nowozin, P. Gehler, J. Jancsary, C. Lampert. MIT
Press, 2014.
- 63 [PY97] J. Pitman and M. Yor. "The two-parameter
Poisson-Dirichlet distribution derived from a stable
subordinator". In: *The Annals of Probability* (1997),
pp. 855–900.
- 64 [QC+06] J. Quiñonero-Candela, C. E. Rasmussen, F.
Sinz, O. Bousquet, and B. Schölkopf. "Evaluating Pre-
dictive Uncertainty Challenge". In: *Machine Learn-
ing Challenges: Evaluating Predictive Uncertainty,
Visual Object Classification, and Recognising Tech-
nical Entailment*. Lecture Notes in Computer Science.
Springer Berlin Heidelberg, 2006, pp. 1–27.
- 65 [QC+08] J. Quiñonero-Candela, M. Sugiyama, A.
Schwaighofer, and N. D. Lawrence, eds. *Dataset Shift
in Machine Learning*. en. The MIT Press, 2008.
- 66 [QCR05] J. Quiñonero-Candela and C. Rasmussen. "A
unifying view of sparse approximate Gaussian process
regression". In: *JMLR* 6.3 (2005), pp. 1939–1959.
- 67 [Qin+20] C. Qin, Y. Wu, J. T. Springenberg, A. Brock,
J. Donahue, T. P. Lillicrap, and P. Kohli. "Train-
ing Generative Adversarial Networks by Solving Or-
dinary Differential Equations". In: *arXiv preprint
arXiv:2010.15040* (2020).
- 68 [Qu+21] H. Qu, H. Rahmani, L. Xu, B. Williams, and
J. Liu. "Recent Advances of Continual Learning in
Computer Vision: An Overview". In: (2021). *arXiv:
2109.11369 [cs.CV]*.
- 69 [Qua+07] A. Quattoni, S. Wang, L.-P. Morency, M.
Collins, and T. Darrell. "Hidden conditional random
fields". In: *IEEE PAMI* 29.10 (2007), pp. 1848–1852.
- 70 [Que98] M. Queyranne. "Minimizing symmetric sub-
modular functions". In: *Math. Programming* 82
(1998), pp. 3–12.
- 71 [QZW19] Y. Qiu, L. Zhang, and X. Wang. "Unbi-
ased Contrastive Divergence Algorithm for Training
Energy-Based Latent Variable Models". In: *ICLR*.
2019.
- 72 [RA13] O. Rippel and R. P. Adams. "High-dimensional
probability estimation with deep density models". In:
73 *ArXiv Preprint arXiv:1302.5125* (2013).

- [Rab89] L. R. Rabiner. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition". In: *Proc. of the IEEE* 77.2 (1989), pp. 257–286.
- [Rad+18] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. *Improving Language Understanding by Generative Pre-Training*. Tech. rep. OpenAI, 2018.
- [Rad+19] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. *Language Models are Unsupervised Multitask Learners*. Tech. rep. OpenAI, 2019.
- [Rad+21] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. "Learning transferable visual models from natural language supervision". In: *arXiv preprint arXiv:2103.00020* (2021).
- [Raf+20a] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". In: *JMLR* (2020).
- [Raf+20b] C. Raffel, N. M. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer". In: *ArXiv abs/1910.10683* (2020).
- [RAG04] B. Ristic, S. Arulampalam, and N. Gordon. *Beyond the Kalman Filter: Particle Filters for Tracking Applications*. Artech House Radar Library, 2004.
- [Rag+17] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein. "Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability". In: *Advances in Neural Information Processing Systems*. 2017, pp. 6076–6085.
- [Rag+19] M. Raghu, C. Zhang, J. Kleinberg, and S. Bengio. "Transfusion: Understanding transfer learning for medical imaging". In: *NIPS*. 2019, pp. 3347–3357.
- [Rag+21] M. Raghu, T. Unterthiner, S. Kornblith, C. Zhang, and A. Dosovitskiy. "Do Vision Transformers See Like Convolutional Neural Networks?". In: *NIPS*. 2021.
- [Rai+18a] T. Rainforth, A. R. Kosiorek, T. A. Le, C. J. Maddison, M. Igl, F. Wood, and Y. W. Teh. "Tighter Variational Bounds are Not Necessarily Better". In: *ICML*. 2018.
- [Rai+18b] M. Raitoharju, L. Svensson, Á. F. García-Fernández, and R. Piché. "Damped Posterior Linearization Filter". In: *IEEE Signal Process. Lett.* 25.4 (2018).
- [Rai+20] T. Rainforth, A. Golinski, F. Wood, and S. Zaidi. "Target-Aware Bayesian Inference: How to Beat Optimal Conventional Estimators". In: *JMLR* 21.88 (2020), pp. 1–54.
- [Rai68] H. Raiffa. *Decision Analysis*. Addison Wesley, 1968.
- [Rak+08] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. "SimpleMKL". In: *JMLR* 9 (2008), pp. 2491–2521.
- [Ram+21a] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. "Zero-Shot Text-to-Image Generation". In: (2021). arXiv: 2102.12092 [cs.CV].
- [Ram+21b] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. "Zero-shot text-to-image generation". In: *International Conference on Machine Learning*. PMLR. 2021, pp. 8821–8831.
- [Ram+22] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. "Hierarchical Text-Conditional Image Generation with CLIP Latents". In: (Apr. 2022). arXiv: 2204.06125 [cs.CV].
- [Ran+06] M. Ranzato, C. S. Poultney, S. Chopra, and Y. LeCun. "Efficient Learning of Sparse Representations with an Energy-Based Model". In: *NIPS*. 2006.
- [Ran16] R. Ranganath. "Hierarchical Variational Models". In: *ICML*. 2016.
- [Ran+18] S. S. Rangapuram, M. W. Seeger, J. Gasthaus, L. Stella, Y. Wang, and T. Januszkiewicz. "Deep State Space Models for Time Series Forecasting". In: *NIPS*. Curran Associates, Inc., 2018, pp. 7796–7805.
- [Rao10] A. V. Rao. "A Survey of Numerical Methods for Optimal Control". In: *Adv. Astronaut. Sci.* 135.1 (2010).
- [Rao99] R. P. Rao. "An optimal estimation approach to visual perception and learning". en. In: *Vision Res.* 39.11 (1999), pp. 1963–1989.
- [Ras00] C. Rasmussen. "The Infinite Gaussian Mixture Model". In: *NIPS*. 2000.
- [Ras+15] A. Rasmus, H. Valpola, M. Honkala, M. Berglund, and T. Raiko. *Semi-Supervised Learning with Ladder Networks*. 2015. arXiv: 1507.02672 [cs.NE].
- [Rat+09] M. Ratnayake, O. Stegle, K. Sharp, and J. Winn. "Inference algorithms and learning theory for Bayesian sparse factor analysis". In: *Proc. Intl. Workshop on Statistical-Mechanical Informatics*. 2009.
- [Rav+18] S. Ravuri, S. Mohamed, M. Rosca, and O. Vinyals. "Learning Implicit Generative Models with the Method of Learned Moments". In: *International Conference on Machine Learning*. 2018, pp. 4314–4323.
- [RB16] G. P. Rigby BR. "The Efficacy of Equine-Assisted Activities and Therapies on Improving Physical Function". In: *J Altern Complement Med.* (2016).
- [RBB18a] H. Ritter, A. Botev, and D. Barber. "A Scalable Laplace Approximation for Neural Networks". In: *ICLR*. 2018.
- [RBB18b] H. Ritter, A. Botev, and D. Barber. "Online Structured Laplace Approximations for Overcoming Catastrophic Forgetting". In: *NIPS*. Curran Associates, Inc., 2018, pp. 3738–3748.
- [RBS16] L. J. Ratliff, S. A. Burden, and S. S. Sastry. "On the characterization of local Nash equilibria in continuous games". In: *IEEE transactions on automatic control* 61.8 (2016), pp. 2301–2307.
- [RBS84] J. Ramsay, J. ten Berge, and G. Styan. "Matrix correlation". In: *Psychometrika* 49.3 (1984), pp. 403–423.
- [RC04] C. Robert and G. Casella. *Monte Carlo Statistical Methods*. 2nd edition. Springer, 2004.
- [RC+18] M. Rojas-Carulla, B. Schölkopf, R. Turner, and J. Peters. "Invariant Models for Causal Transfer Learning". In: *Journal of Machine Learning Research* 19.36 (2018), pp. 1–34.
- [RD06] M. Richardson and P. Domingos. "Markov logic networks". In: *Machine Learning* 62 (2006), pp. 107–136.
- [RDV18] A. S. Ross and F. Doshi-Velez. "Improving the adversarial robustness and interpretability of deep neural networks by regularizing their input gradients".

- 1 In: *Thirty-second AAAI conference on artificial intelligence*. 2018.
- 2 [RE76] P Robert and Y Escoufier. “A unifying tool for
3 linear multivariate statistical methods: The RV- coefficient”. In: *J. R. Stat. Soc. Ser. C Appl. Stat.* 25.3
4 (1976), p. 257.
- 5 [Rea+19] E. Real, A. Aggarwal, Y. Huang, and Q. V.
6 Le. “Regularized Evolution for Image Classifier Archi-
7 tecture Search”. In: *AAAI*. 2019.
- 8 [Rec19] B. Recht. “A Tour of Reinforcement Learning:
9 The View from Continuous Control”. In: *Annual Re-
10 view of Control, Robotics, and Autonomous Systems*
11 2 (2019), pp. 253–279.
- 12 [Ree+17] S. Reed, A. van den Oord, N. Kalchbrenner,
13 S. G. Colmenarejo, Z. Wang, D. Belov, and N. de Fre-
14 icas. “Parallel Multiscale Autoregressive Density Es-
15 timation”. In: (2017). arXiv: [1703.03664 \[cs.CV\]](https://arxiv.org/abs/1703.03664).
- 16 [Rei+10] J. Reisinger, A. Waters, B. Silverthorn, and
17 R. Mooney. “Spherical topic models”. In: *ICML*. 2010.
- 18 [Rei13] S. Reich. “A Nonparametric Ensemble Trans-
19 form Method for Bayesian Inference”. In: *SIAM J. Sci.
Comput.* 35.4 (2013), A2013–A2024.
- 20 [Rei16] P. C. Reiss. “Just How Sensitive are Instrumen-
21 tal Variable Estimates?” In: *Foundations and Trends
in Accounting* 10.2-4 (2016).
- 22 [Rei22] Reinforcement Learning: Theory and Algo-
23 rithms. *Alekh Agarwal; Nan Jiang; Sham M. Kakade;*
24 *Wen Sun*. 2022.
- 25 [Rei+22] P. Reizinger, L. Gresele, J. Brady, J. von
26 Kügelgen, D. Zietlow, B. Scholkopf, G. Martius, W.
27 Brendel, and M. Besserve. “Embrace the Gap: VAEs
28 Perform Independent Mechanism Analysis”. In: 2022.
- 29 [Ren+19] J. Ren, P. J. Liu, E. Fertig, J. Snoek, R.
30 Poplin, M. A. DePristo, J. V. Dillon, and B. Lakshmi-
31 narayanan. “Likelihood Ratios for Out-of-Distribution
32 Detection”. In: *NIPS*. 2019.
- 33 [Rén61] A. Rényi. “On Measures of Entropy and Infor-
34 mation”. en. In: *Proceedings of the Fourth Berkeley
35 Symposium on Mathematical Statistics and Probabil-
36 ity, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California, 1961.
- 37 [RFB15] O. Ronneberger, P. Fischer, and T. Brox. “U-
38 Net: Convolutional Networks for Biomedical Image
39 Segmentation”. In: *MICCAI (Intl. Conf. on Medical
40 Image Computing and Computer Assisted Interven-
41 tions)*. 2015.
- 42 [RG17] M Roth and F Gustafsson. “Computation and
43 visualization of posterior densities in scalar nonlinear
44 and non-Gaussian Bayesian filtering and smoothing
45 problems”. In: *ICASSP*. 2017, pp. 4686–4690.
- 46 [RGB11] S. Ross, G. J. Gordon, and D. Bagnell. “A
47 Reduction of Imitation Learning and Structured Pre-
48 diction to No-Regret Online Learning”. In: *AISTATS*.
49 2011, pp. 627–635.
- 50 [RGB14] R. Ranganath, S. Gerrish, and D. M. Blei.
51 “Black Box Variational Inference”. In: *AISTATS*.
52 2014.
- 53 [RGL19] S. Rabanser, S. Günnemann, and Z. C. Lip-
54 ton. “Failing Loudly: An Empirical Study of Methods
55 for Detecting Dataset Shift”. In: *NIPS*. 2019.
- 56 [RH05] H. Rue and L. Held. *Gaussian Markov Ran-
57 dom Fields: Theory and Applications*. Vol. 104. Mono-
58 graphs on Statistics and Applied Probability. London:
59 Chapman & Hall, 2005.
- 60 [RHDV17] A. S. Ross, M. C. Hughes, and F. Doshi-
61 Velez. “Right for the right reasons: Training differen-
62 tiable models by constraining their explanations”. In:
63 *IJCAI* (2017).
- 64 [RHG16] D. Ritchie, P. Horsfall, and N. D. Good-
65 man. “Deep Amortized Inference for Probabilistic Pro-
66 grams”. In: (2016). arXiv: [1610.05735 \[cs.AI\]](https://arxiv.org/abs/1610.05735).
- 67 [RHS22] S. Rissanen, M. Heinonen, and A. Solin. “Gen-
68 erative Modelling With Inverse Heat Dissipation”. In:
69 (June 2022). arXiv: [2206.13397 \[cs.CV\]](https://arxiv.org/abs/2206.13397).
- 70 [RHW86a] D. Rumelhart, G. Hinton, and R. Williams.
71 “Learning internal representations by error propa-
72 gation”. In: *Parallel Distributed Processing: Explora-
73 tions in the Microstructure of Cognition*. Ed. by
74 D. Rumelhart, J. McClelland, and the PDD Research
75 Group. MIT Press, 1986.
- 76 [RHW86b] D. Rumelhart, G. E. Hinton, and
77 R. J. Williams. “Learning representations by back-
78 propagating errors”. In: *Nature* 323 (1986), pp. 533–
79 536.
- 80 [Ric03] T. Richardson. “Markov properties for acyclic
81 directed mixed graphs”. In: *Scandinavian J. of Statis-
82 tics* 30 (2003), pp. 145–157.
- 83 [Ric95] J. Rice. *Mathematical statistics and data anal-
84 ysis*. 2nd edition. Duxbury, 1995.
- 85 [Rif+11] S. Rifai, P. Vincent, X. Muller, X. Glorot, and
86 Y. Bengio. “Contractive auto-encoders: Explicit invar-
87 iance during feature extraction”. In: *ICML*. 2011.
- 88 [Rip05] B. D. Ripley. *Spatial statistics*. Vol. 575. John
89 Wiley & Sons, 2005.
- 90 [Rip77] B. D. Ripley. “Modelling spatial patterns”. In:
91 *Journal of the Royal Statistical Society: Series B
(Methodological)* 39.2 (1977), pp. 172–192.
- 92 [Ris+08] I. Rish, G. Grabarnik, G. Cecchi, F. Pereira,
93 and G. Gordon. “Closed-form supervised dimension-
94 ality reduction with generalized linear models”. In:
95 *ICML*. 2008.
- 96 [Riv87] R. L. Rivest. “Learning decision lists”. In: *Ma-
97 chine learning* 2.3 (1987), pp. 229–246.
- 98 [RK04] R. Rubinstein and D. Kroese. *The Cross-
99 Entropy Method: A Unified Approach to Combinato-
100 rial Optimization, Monte-Carlo Simulation, and Ma-
101 chine Learning*. Springer-Verlag, 2004.
- 102 [RL17] S. Ravi and H. Larochelle. “Optimization as a
103 Model for Few-Shot Learning”. In: *ICLR*. 2017.
- 104 [RM15] D. J. Rezende and S. Mohamed. “Variational
105 Inference with Normalizing Flows”. In: *ICML*. 2015.
- 106 [RMB08] N. L. Roux, P.-A. Manzagol, and Y. Bengio.
107 “Topomoumoute Online Natural Gradient Algorithm”.
108 In: *NIPS*. 2008, pp. 849–856.
- 109 [RMC09] H. Rue, S. Martino, and N. Chopin. “Appox-
110 imate Bayesian Inference for Latent Gaussian Models
111 Using Integrated Nested Laplace Approximations”. In:
112 *J. of Royal Stat. Soc. Series B* 71 (2009), pp. 319–
113 392.
- 114 [RMC15] A. Radford, L. Metz, and S. Chintala. “Un-
115 supervised Representation Learning with Deep Convo-
116 lutional Generative Adversarial Networks”. In: *arXiv*
117 (2015).
- 118 [RMC16a] A. Radford, L. Metz, and S. Chintala. “Un-
119 supervised Representation Learning with Deep Convo-
120 lutional Generative Adversarial Networks”. In: *ICLR*.
121 2016.

- 1 [RMC16b] A. Radford, L. Metz, and S. Chintala. “Un-
2 supervised Representation Learning with Deep Convo-
3 lutional Generative Adversarial Networks”. In: *CoRR*
4 abs/1511.06434 (2016).
- 5 [RMK21] G. Roeder, L. Metz, and D. P. Kingma. “On
6 Linear Identifiability of Learned Representations”. In:
7 *ArXiv* abs/2007.00810 (2021).
- 8 [RMW14a] D. Rezende, S. Mohamed, and D. Wierstra.
9 “Stochastic backpropagation and approximate infer-
10 ence in deep generative models”. In: *ICML*. 2014.
- 11 [RMW14b] D. J. Rezende, S. Mohamed, and D. Wier-
12 stra. “Stochastic Backpropagation and Approximate
13 Inference in Deep Generative Models”. In: *ICML*. Ed.
14 by E. P. Xing and T. Jebara. Vol. 32. Proceedings of
15 Machine Learning Research. PMLR, 2014, pp. 1278–
16 1286.
- 17 [RN02] S. Russell and P. Norvig. *Artificial Intelli-
18 gence: A Modern Approach*. 2nd edition. Prentice
19 Hall, 2002.
- 20 [RN10] S. Russell and P. Norvig. *Artificial Intelli-
21 gence: A Modern Approach*. 3rd edition. Prentice Hall,
22 2010.
- 23 [RN19] S. Russell and P. Norvig. *Artificial Intelli-
24 gence: A Modern Approach*. 4th edition. Prentice Hall,
25 2019.
- 26 [RN94] G. A. Rummery and M Niranjan. *On-Line Q-
27 Learning Using Connectionist Systems*. Tech. rep.
28 Cambridge Univ. Engineering Dept., 1994.
- 29 [RN95] S. Russell and P. Norvig. *Artificial Intelli-
30 gence: A Modern Approach*. Prentice Hall, 1995.
- 31 [RNA22] L. Regenwetter, A. H. Nobari, and F. Ahmed.
32 “Deep Generative Models in Engineering Design: A Re-
33 view”. In: *J. Mech. Des.* (2022).
- 34 [Rob07] C. P. Robert. *The Bayesian Choice: From
35 Decision-Theoretic Foundations to Computational
36 Implementation*. en. 2nd edition. Springer Verlag,
37 New York, 2007.
- 38 [Rob+13] S Roberts, M Osborne, M Ebden, S Reece,
39 N Gibson, and S Aigrain. “Gaussian processes for
40 time-series modelling”. en. In: *Philos. Trans. A Math.
41 Phys. Eng. Sci.* 371.1984 (2013), p. 20110550.
- 42 [Rob+18] C. P. Robert, V. Elvira, N. Tawn, and C. Wu.
43 “Accelerating MCMC Algorithms”. In: (2018). arXiv:
44 1804.02719 [stat.CO].
- 45 [Rob+21] J. Robinson, C.-Y. Chuang, S. Sra, and S.
46 Jegelka. “Contrastive Learning with Hard Negative
47 Samples”. In: *ArXiv* abs/2010.04592 (2021).
- 48 [Rob63] L. G. Roberts. “Machine Perception of Three-
49 Dimensional Solids”. In: *Outstanding Dissertations in
50 the Computer Sciences*. 1963.
- 51 [Rob86] J. Robins. “A new approach to causal infer-
52 ence in mortality studies with a sustained exposure pe-
53 riod—application to control of the healthy worker sur-
54 vivor effect”. In: *Mathematical Modelling* 7.9 (1986),
55 pp. 1393–1512.
- 56 [Rob95a] C. Robert. “Simulation of truncated normal
57 distributions”. In: *Statistics and computing* 5 (1995),
58 pp. 121–125.
- 59 [Rob95b] A. Robins. “Catastrophic Forgetting, Re-
60 hearsal and Pseudorehearsal”. In: *Conn. Sci.* 7.2
61 (1995), pp. 123–146.
- 62 [Rod14] J. Rodu. “Spectral estimation of hidden
63 Markov models”. PhD thesis. U. Penn, 2014.
- 64 [RÖG13] M. Roth, E. Özkan, and F. Gustafsson. “A
65 Student’s t filter for heavy tailed process and mea-
66 surement noise”. In: *ICASSP*. 2013, pp. 5770–5774.
- 67 [Roh21] D. Rohde. “Causal Inference, is just Inference:
68 A beautifully simple idea that not everyone accepts”.
69 In: *I (Still) Can’t Believe It’s Not Better! NeurIPS
70 2021 Workshop*. 2021.
- 71 [Ros10] P. Rosenbaum. *Design of Observational Stud-
72 ies*. 2010.
- 73 [Ros+21] M. Rosca, Y. Wu, B. Dherin, and D. G. Bar-
74 rett. “Discretization Drift in Two-Player Games”. In:
75 (2021).
- 76 [Ros+22] C. Rosato, L. Devlin, V. Beraud, P. Horridge,
77 T. B. Schön, and S. Maskell. “Efficient Learning of the
78 Parameters of Non-Linear Models Using Differentiable
79 Resampling in Particle Filters”. In: *IEEE Trans. Sig-
80 nal Process.* 70 (2022), pp. 3676–3692.
- 81 [Ros22] C. Ross. *AI gone astray: How subtle shifts in
82 patient data send popular algorithms reeling, under-
83 mining patient safety*. en. <https://www.statnews.com/2022/02/28/sepsis-hospital-algorithms-data-shift/>. Accessed: 2022-3-2. 2022.
- 84 [Rot+17] M. Roth, G. Hendeby, C. Fritzsche, and F.
85 Gustafsson. “The Ensemble Kalman filter: a signal
86 processing perspective”. In: *EURASIP J. Adv. Signal
87 Processing* 2017.1 (2017), p. 56.
- 88 [Rot+18] W. Roth, R. Peharz, S. Tschiatschek, and F.
89 Pernkopf. “Hybrid generative-discriminative training
90 of Gaussian mixture models”. In: *Pattern Recognit.
91 Lett.* 112 (Sept. 2018), pp. 131–137.
- 92 [Rot96] D. Roth. “On the hardness of approximate
93 reasoning”. In: *Artificial Intelligence* 82.1-2 (1996),
94 pp. 273–302.
- 95 [ROV19] A. Razavi, A. van den Oord, and O. Vinyals.
96 “Generating diverse high resolution images with VA-
97 VAE-2”. In: *NIPS*. 2019.
- 98 [Row97] S. Roweis. “EM algorithms for PCA and
99 SPCA”. In: *NIPS*. 1997.
- 100 [Roy+21] N. Roy et al. “From Machine Learning to
101 Robotics: Challenges and Opportunities for Embod-
102 ited Intelligence”. In: (Oct. 2021). arXiv: 2110.15245
103 [cs.RO].
- 104 [RPC19] Y. Romano, E. Patterson, and E. J. Candès.
105 “Conformalized Quantile Regression”. In: *NIPS*. 2019.
- 106 [RPH21] A. Robey, G. J. Pappas, and H. Hassani.
107 *Model-Based Domain Generalization*. 2021. arXiv:
108 2102.11436 [stat.ML].
- 109 [RR01a] A. Rao and K. Rose. “Deterministically An-
110 nealed Design of Hidden Markov Model Speech Recog-
111 nizers”. In: *IEEE Trans. on Speech and Audio Proc.*
112 9.2 (2001), pp. 111–126.
- 113 [RR01b] G. Roberts and J. Rosenthal. “Optimal scal-
114 ing for various Metropolis-Hastings algorithms”. In:
115 *Statistical Science* 16 (2001), pp. 351–367.
- 116 [RR08] A. Rahimi and B. Recht. “Random Features for
117 Large-Scale Kernel Machines”. In: *NIPS*. Curran As-
118 sociates, Inc., 2008, pp. 1177–1184.
- 119 [RR09] A. Rahimi and B. Recht. “Weighted Sums of
120 Random Kitchen Sinks: Replacing minimization with
121 randomization in learning”. In: *NIPS*. Curran As-
122 sociates, Inc., 2009, pp. 1313–1320.
- 123 [RR11] T. S. Richardson and J. M. Robins. “Single
124 World Intervention Graphs: A Primer”. In: *Second
125 UAI workshop on causal structure learning*. 2011.

- 1
- 2 [RR13] T. S. Richardson and J. M. Robins. "Single
3 World Intervention Graphs (SWIGs): A Unification
4 of the Counterfactual and Graphical Approaches to
Causality". 2013.
- 5 [RR14] D. Russo and B. V. Roy. "Learning to Optimize
6 via Posterior Sampling". In: *Math. Oper. Res.*
39.4 (2014), pp. 1221–1243.
- 7 [RR83] P. R. Rosenbaum and D. B. Rubin. "Assessing
8 Sensitivity to an Unobserved Binary Covariate in
an Observational Study with Binary Outcome". In:
*Journal of the Royal Statistical Society. Series B
(Methodological)* 45.2 (1983), pp. 212–218.
- 10 [RRR21] E. Rosenfeld, P. Ravikumar, and A. Risteski.
11 "The Risks of Invariant Risk Minimization". In: *ICML*.
2021.
- 12 [RRS00] J. M. Robins, A. Rotnitzky, and D. O. Scharf-
13 stein. "Sensitivity analysis for selection bias and un-
measured confounding in missing data and causal in-
ference models". In: *Statistical models in epidemiology,
the environment, and clinical trials*. Springer, 2000,
pp. 1–94.
- 16 [RS07] M. Raphan and E. P. Simoncelli. "Learning
17 to be Bayesian without supervision". In: *Advances
in neural information processing systems*. 2007,
pp. 1145–1152.
- 19 [RS11] M. Raphan and E. P. Simoncelli. "Least squares
estimation without priors or supervision". In: *Neural
computation* 23.2 (2011), pp. 374–420.
- 21 [RS20] A. Rotnitzky and E. Smucler. "Efficient Adjustment
Sets for Population Average Causal Treatment
Effect Estimation in Graphical Models". In: *J. Mach.
Learn. Res.* 21 (2020), pp. 188–1.
- 24 [RS97a] G. O. Roberts and S. K. Sahu. "Updating
25 Schemes, Correlation Structure, Blocking and Paramet-
erization for the Gibbs Sampler". In: *J. of Royal Stat.
Soc. Series B* 59.2 (1997), pp. 291–317.
- 26 [RS97b] G. O. Roberts and S. K. Sahu. "Updating
27 schemes, correlation structure, blocking and paramet-
erization for the Gibbs sampler". In: *J. of Royal Stat.
Soc. Series B* 59.2 (1997), pp. 291–317.
- 29 [RSC20] Y. Romano, M. Sesia, and E. J. Candès.
30 "Classification with Valid and Adaptive Coverage". In:
NIPS. 2020.
- 31 [RSG16a] M. T. Ribeiro, S. Singh, and C. Guestrin. "
32 Why should i trust you?" Explaining the predictions
of any classifier". In: *Proceedings of the 22nd ACM
SIGKDD international conference on knowledge dis-
covery and data mining*. 2016, pp. 1135–1144.
- 35 [RSG16b] M. T. Ribeiro, S. Singh, and C. Guestrin.
36 "Model-agnostic interpretability of machine learning".
In: *arXiv preprint arXiv:1606.05386* (2016).
- 37 [RSG17] S. Rabanser, O. Shchur, and S. Günnemann.
38 "Introduction to Tensor Decompositions and their Ap-
plications in Machine Learning". In: (2017). arXiv:
1711.10781 [stat.ML].
- 40 [RT16] S. Reid and R. Tibshirani. "Sparse regression
and marginal testing using cluster prototypes". In:
Biostatistics 17.2 (2016), pp. 364–376.
- 42 [RT82] D. B. Rubin and D. T. Thayer. "EM algo-
rithms for ML factor analysis". In: *Psychometrika* 47.1
(1982), pp. 69–76.
- 44 [RT96] G. O. Roberts and R. L. Tweedie. "Exponen-
45 tial convergence of Langevin distributions and their
discrete approximations". In: *Bernoulli* 2.4 (1996),
pp. 341–363.
- 47 [RTS18] C. Riquelme, G. Tucker, and J. Snoek. "Deep
Bayesian Bandits Showdown: An Empirical Compari-
son of Bayesian Deep Networks for Thompson Sam-
pling". In: *ICLR*. 2018.
- [RTS65] H. E. Rauch, F. Tung, and C. T. Striebel.
"Maximum likelihood estimates of linear dynamic sys-
tems". In: *AIAA Journal* 3.8 (1965), pp. 1445–1450.
- [Rub+20] Y. Rubanova, D. Dohan, K. Swersky, and K.
Murphy. "Amortized Bayesian Optimization over Dis-
crete Spaces". In: *UAI*. 2020.
- [Rub74] D. B. Rubin. "Estimating causal effects of
treatments in randomized and nonrandomized stud-
ies". In: *J. Educ. Psychol.* 66.5 (1974), pp. 688–701.
- [Rub84] D. B. Rubin. "Bayesianly Justifiable and Rel-
evant Frequency Calculations for the Applied Statisti-
cian". In: *Ann. Stat.* 12.4 (1984), pp. 1151–1172.
- [Rub97] R. Y. Rubinstein. "Optimization of computer
simulation models with rare events". In: *Eur. J. Oper.
Res.* 99.1 (1997), pp. 89–112.
- [Rud19] C. Rudin. *Stop Explaining Black Box Ma-
chine Learning Models for High Stakes Decisions
and Use Interpretable Models Instead*. 2019. arXiv:
1811.10154 [stat.ML].
- [Ruf+21] L. Ruff, J. R. Kauffmann, R. A. Vander-
meulen, G. Montavon, W. Samek, M. Kloft, T. G. Diet-
terich, and K.-R. Müller. "A Unifying Review of Deep
and Shallow Anomaly Detection". In: *Proc. IEEE*
109.5 (2021), pp. 756–795.
- [Rus15] S. Russell. "Unifying Logic and Probability".
In: *Commun. ACM* 58.7 (2015), pp. 88–97.
- [Rus+16] A. A. Rusu, N. C. Rabinowitz, G. Desjardins,
H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu,
and R. Hadsell. "Progressive Neural Networks". In:
(2016). arXiv: 1606.04671 [cs.LG].
- [Rus+18] D. J. Russo, B. Van Roy, A. Kazerouni,
I. Osband, and Z. Wen. "A Tutorial on Thompson
Sampling". In: *Foundations and Trends in Machine
Learning* 11.1 (2018), pp. 1–96.
- [Rus+95] S. Russell, J. Binder, D. Koller, and K.
Kanazawa. "Local learning in probabilistic networks
with hidden variables". In: *IJCAI*. 1995.
- [RV19] S. Ravuri and O. Vinyals. "Classification ac-
curacy score for conditional generative models". In:
Advances in Neural Information Processing Systems.
2019, pp. 12268–12279.
- [RW06] C. E. Rasmussen and C. K. I. Williams. *Gau-
sian Processes for Machine Learning*. MIT Press,
2006.
- [RW11] M. D. Reid and R. C. Williamson. "Informa-
tion, Divergence and Risk for Binary Experiments". In:
Journal of Machine Learning Research 12.3 (2011).
- [RW15] D. Rosebaum and Y. Weiss. "The Return of
the Gating Network: Combining Generative Models
and Discriminative Training in Natural Image Priors".
In: *NIPS*. 2015, pp. 2665–2673.
- [RW18] E. Richardson and Y. Weiss. "On GANs and
GMMs". In: *NIPS*. 2018.
- [RWD17] G. Roeder, Y. Wu, and D. Duvenaud. "Stick-
ing the Landing: An Asymptotically Zero-Variance
Gradient Estimator for Variational Inference". In:
NIPS. 2017.
- [RY21] D. Roberts and S. Yaida. *The Principles of
Deep Learning Theory: An Effective Theory Ap-
proach to Understanding Neural Network*. 2021.

- 1 [Ryc+19] B. Rychalska, D. Basaj, A. Gosiewska, and
2 P. Biecek. “Models in the Wild: On Corruption
3 Robustness of Neural NLP Systems”. In: *International Conference on Neural Information Processing (ICONIP)*. Springer International Publishing, 2019, pp. 235–247.
- 4 [Ryu+20] M. Ryu, Y. Chow, R. Anderson, C. Tjandraatmadja, and C. Boutilier. “CAQL: Continuous Action Q-Learning”. In: *ICLR*. 2020.
- 5 [RZL17] P. Ramachandran, B. Zoph, and Q. V. Le.
6 “Searching for Activation Functions”. In: (2017). arXiv: 1710.05941 [cs.NE].
- 7 [SA19] F. Schafer and A. Anandkumar. “Competitive gradient descent”. In: *NIPS*. 2019, pp. 7625–7635.
- 8 [Sac+05] K. Sachs, O. Perez, D. Pe’er, D. Lauffenburger, and G. Nolan. “Causal Protein-Signaling Networks Derived from Multiparameter Single-Cell Data”. In: *Science* 308 (2005).
- 9 [SAC17] J. Schulman, P. Abbeel, and X. Chen. *Equivalence Between Policy Gradients and Soft Q-Learning*. arXiv:1704.06440. 2017.
- 10 [Sag+20] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. “Distributionally Robust Neural Networks for Group Shifts: On the Importance of Regularization for Worst-Case Generalization”. In: *ICLR*. 2020.
- 11 [Sah+21] C. Saharia, W. Chan, H. Chang, C. A. Lee, J. Ho, T. Salimans, D. J. Fleet, and M. Norouzi. “Palette: Image-to-Image Diffusion Models”. In: (Nov. 2021). arXiv: 2111.05826 [cs.CV].
- 12 [Sah+22] C. Saharia et al. “Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding”. In: (May 2022). arXiv: 2205.11487 [cs.CV].
- 13 [Sai+20] M. Saito, S. Saito, M. Koyama, and S. Kobayashi. “Train Sparsely, Generate Densely: Memory-Efficient Unsupervised Training of High-Resolution Temporal GAN”. In: *International Journal of Computer Vision* 128 (2020), pp. 2586–2606.
- 14 [Saj+18] M. S. Sajjadi, O. Bachem, M. Lucic, O. Bousquet, and S. Gelly. “Assessing generative models via precision and recall”. In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems*. 2018, pp. 5234–5243.
- 15 [Sal16] T. Salimans. “A Structured Variational Autoencoder for Learning Deep Hierarchies of Sparse Features”. In: (2016). arXiv: 1602.08734 [stat.ML].
- 16 [Sal+16] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. “Improved Techniques for Training GANs”. In: (2016). arXiv: 1606.03498 [cs.LG].
- 17 [Sal+17a] M. Salehi, A. Karbasi, D. Scheinost, and R. T. Constable. “A Submodular Approach to Create Individualized Parcellations of the Human Brain”. In: *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017*. Ed. by M. Descoteaux, L. Maier-Hein, A. Franz, P. Jannin, D. L. Collins, and S. Duchesne. Cham: Springer International Publishing, 2017, pp. 478–485.
- 18 [Sal+17b] T. Salimans, J. Ho, X. Chen, and I. Sutskever. “Evolution Strategies as a Scalable Alternative to Reinforcement Learning”. In: (2017). arXiv: 1703.03864 [stat.ML].
- 19 [Sal+17c] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma. “PixelCNN++: Improving the Pixel-CNN with Discretized Logistic Mixture Likelihood and Other Modifications”. In: *ICLR*. 2017.
- 20 [Salinas+19a] D. Salinas, M. Bohlke-Schneider, L. Callot, R. Medico, and J. Gasthaus. “High-Dimensional Multivariate Forecasting with Low-Rank Gaussian Copula Processes”. In: *NIPS*. 2019.
- 21 [Salinas+19b] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski. “DeepAR: Probabilistic forecasting with autoregressive recurrent networks”. In: *International Journal of Forecasting* (2019).
- 22 [Salman+20] H. Salman, A. Ilyas, L. Engstrom, A. Kapoor, and A. Madry. “Do Adversarially Robust ImageNet Models Transfer Better?” In: *arXiv preprint arXiv:2007.08489* (2020).
- 23 [Salehi+21] M. Salehi, H. Mirzaei, D. Hendrycks, Y. Li, M. H. Rohban, and M. Sabokrou. “A Unified Survey on Anomaly, Novelty, Open-Set, and Out-of-Distribution Detection: Solutions and Future Challenges”. In: (2021). arXiv: 2110.14051 [cs.CV].
- 24 [Sam68] F. Sampson. “A Novitiate in a Period of Change: An Experimental and Case Study of Social Relationships”. PhD thesis. Cornell, 1968.
- 25 [Sam74] P. A. Samuelson. “Complementarity: An essay on the 40th anniversary of the Hicks-Allen revolution in demand theory”. In: *Journal of Economic literature* 12.4 (1974), pp. 1255–1289.
- 26 [San17] R. Santana. “Gray-box optimization and factorized distribution algorithms: where two worlds collide”. In: (2017). arXiv: 1707.03093 [cs.NE].
- 27 [San+17] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap. “A simple neural network module for relational reasoning”. In: *Advances in neural information processing systems*. 2017, pp. 4967–4976.
- 28 [Sarkka08] S. Sarkka. “Unscented Rauch–Tung–Striebel Smoother”. In: *IEEE Trans. Automat. Contr.* 53.3 (Apr. 2008), pp. 845–849.
- 29 [Sarkka13] S. Sarkka. *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013.
- 30 [Sarin18] H. Sarin. “Playing a game of GANstruction”. In: *The Gradient* (2018).
- 31 [Sayres+19] R. Sayres, S. Xu, T Saensuksopa, M. Le, and D. R. Webster. “Assistance from a deep learning system improves diabetic retinopathy assessment in optometrists”. In: *Investigative Ophthalmology & Visual Science* 60.9 (2019), pp. 1433–1433.
- 32 [SB01] A. J. Smola and P. L. Bartlett. “Sparse Greedy Gaussian Process Regression”. In: *NIPS*. Ed. by T. K. Leen, T. G. Dietterich, and V. Tresp. MIT Press, 2001, pp. 619–625.
- 33 [SB12] J. Staines and D. Barber. “Variational Optimization”. In: (2012). arXiv: 1212.4507 [stat.ML].
- 34 [SB13] J. Staines and D. Barber. “Optimization by Variational Bounding”. In: *European Symposium on ANNs*. elen.ucl.ac.be, 2013.
- 35 [SB18] R. Sutton and A. Barto. *Reinforcement learning: an introduction* (2nd edn). MIT Press, 2018.
- 36 [SBG07] S. Siddiqi, B. Boots, and G. Gordon. “A constraint generation approach to learning stable linear dynamical systems”. In: *NIPS*. 2007.
- 37 [SBP17] Y. Sun, P. Babu, and D. P. Palomar. “Majorization-Minimization Algorithms in Signal Processing, Communications, and Machine Learning”. In: *IEEE Trans. Signal Process.* 65.3 (2017), pp. 794–816.

- 1
- 2 [SC13] C. Schäfer and N. Chopin. “Sequential Monte
3 Carlo on large binary sampling spaces”. In: *Stat. Comput.* 23.2 (2013), pp. 163–184.
- 4 [SC86] R. Smith and P. Cheeseman. “On the Repre-
5 sentation and Estimation of Spatial Uncertainty”. In: *Intl.
J. Robotics Research* 5.4 (1986), pp. 56–68.
- 6 [SC90] R. Schwarz and Y. Chow. “The n-best algo-
7 rithm: an efficient and exact procedure for finding the
8 n most likely hypotheses”. In: *ICASSP*. 1990.
- 9 [Sca21] S. Scardapane. *Lecture 8: Beyond single-task
supervised learning*. 2021.
- 10 [Sch00] A. Schrijver. “A combinatorial algorithm min-
11 imizing submodular functions in strongly polynomial
12 time”. In: *Journal of Combinatorial Theory, Series
B* 80.2 (2000), pp. 346–355.
- 13 [Sch02] N. N. Schraudolph. “Fast Curvature Matrix-
Vector Products for Second-Order Gradient Descent”.
14 In: *Neural Computation* 14 (2002).
- 15 [Sch04] A. Schrijver. *Combinatorial Optimization*.
16 Springer, 2004.
- 17 [Sch+12a] B. Schoelkopf, D. Janzing, J. Peters, E.
Sgouritsa, K. Zhang, and J. Mooij. “On Causal and
18 Anticausal Learning”. In: *ICML*. 2012.
- 19 [Sch+12b] B. Schölkopf, D. Janzing, J. Peters, E.
Sgouritsa, K. Zhang, and J. Mooij. “On causal and
20 anticausal learning”. In: *Proceedings of the 29th Inter-
national Conference on International Conference
on Machine Learning*. 2012, pp. 459–466.
- 22 [Sch14] J. Schmidhuber. *Deep Learning in Neural Net-
works: An Overview*. Tech. rep. 2014.
- 23
- 24 [Sch+15a] J. Schulman, N. Heess, T. Weber, and P.
Abbeel. “Gradient Estimation Using Stochastic Com-
25 putation Graphs”. In: *NIPS*. 2015.
- 26 [Sch+15b] J. Schulman, S. Levine, P. Moritz, M. I. Jor-
27 dan, and P. Abbeel. “Trust Region Policy Optimiza-
tion”. In: *ICML*. 2015.
- 28 [Sch+16a] T. Schaul, J. Quan, I. Antonoglou, and
29 D. Silver. “Prioritized Experience Replay”. In: *ICLR*.
2016.
- 30 [Sch+16b] J. Schulman, P. Moritz, S. Levine, M. Jor-
31 dan, and P. Abbeel. “High-Dimensional Continuous
Control Using Generalized Advantage Estimation”. In:
32 *ICLR*. 2016.
- 33 [Sch+17] J. Schulman, F. Wolski, P. Dhariwal, A. Rad-
34 ford, and O. Klimov. “Proximal Policy Optimiza-
35 tion Algorithms”. In: (2017). arXiv: [1707.06347 \[cs.LG\]](#).
- 36 [Sch+18] J. Schwarz, J. Luketina, W. M. Czarnecki, A.
Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R.
Hadsell. “Progress & Compress: A scalable framework
for continual learning”. In: *ICML*. 2018.
- 37 [Sch19] B. Schölkopf. “Causality for Machine Learning”.
38 In: (2019). arXiv: [1911.10500 \[cs.LG\]](#).
- 39 [Sch20] J. Schmidhuber. *Planning & Reinforcement
Learning with Recurrent World Models and Artificial
Curiosity*. 2020.
- 40 [Sch+20] J. Schrittwieser et al. “Mastering Atari, Go,
41 Chess and Shogi by Planning with a Learned Model”.
In: *Nature* (2020).
- 42
- 43 [Sch+21a] D. O. Scharfstein, R. Nabi, E. H. Kennedy,
44 M.-Y. Huang, M. Bonvini, and M. Smid. *Semipara-
45 metric Sensitivity Analysis: Unmeasured Confounding
in Observational Studies*. 2021. arXiv: [2104.08300
\[stat.ME\]](#).
- 46
- 47
- [Sch+21b] B. Schölkopf, F. Locatello, S. Bauer, N. R.
Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio.
“Toward Causal Representation Learning”. In: *Proc.
IEEE* 109.5 (2021), pp. 612–634.
- [Sch+21c] B. Schölkopf, F. Locatello, S. Bauer, N. R.
Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio. “To-
wards Causal Representation Learning”. In: *CoRR*
abs/2102.11107 (2021). arXiv: [2102.11107](#).
- [Sch78] G. Schwarz. “Estimating the dimension of a
model”. In: *Annals of Statistics* 6.2 (1978), pp. 461–
464.
- [Sco02] S Scott. “Bayesian methods for hidden Markov
models: Recursive computing in the 21st century.” In:
JASA (2002).
- [Sco09] S. Scott. “Data augmentation, frequentist esti-
mation, and the Bayesian analysis of multinomial logit
models”. In: *Statistical Papers* (2009).
- [Sco10] S. Scott. “A modern Bayesian look at the multi-
armed bandit”. In: *Applied Stochastic Models in Busi-
ness and Industry* 26 (2010), pp. 639–658.
- [SCPD22] R. Sanchez-Cauce, I. Paris, and F. J. Diez.
“Sum-Product Networks: A Survey”. en. In: *IEEE
PAMI* 44.7 (July 2022), pp. 3821–3839.
- [SCS19] A. Subbaswamy, B. Chen, and S. Saria. *A Uni-
versal Hierarchy of Shift-Stable Distributions and
the Tradeoff Between Stability and Performance*.
2019. arXiv: [1905.11374 \[stat.ML\]](#).
- [SCS22] A. Subbaswamy, B. Chen, and S. Saria. “A
unifying causal framework for analyzing dataset shift-
stable learning algorithms”. en. In: *Journal of Causal
Inference* 10.1 (Jan. 2022), pp. 64–89.
- [SD12] J. Sohl-Dickstein. “The Natural Gradient by
Analogy to Signal Whitening, and Recipes and Tricks
for its Use”. In: (2012). arXiv: [1205.1828 \[cs.LG\]](#).
- [SD+15a] J. Sohl-Dickstein, E. Weiss, N. Mah-
eswaranathan, and S. Ganguli. “Deep Unsupervised
Learning using Nonequilibrium Thermodynamics”. In:
48 *ICML*. 2015, pp. 2256–2265.
- [SD+15b] J. Sohl-Dickstein, E. A. Weiss, N. Mah-
eswaranathan, and S. Ganguli. “Deep Unsupervised
Learning using Nonequilibrium Thermodynamics”. In:
49 *ICML*. 2015.
- [SDBD11] J. Sohl-Dickstein, P. Battaglino, and M. R.
DeWeese. “Minimum probability flow learning”. In:
50 *Proceedings of the 28th International Conference
on International Conference on Machine Learning*.
2011, pp. 905–912.
- [SE19] Y. Song and S. Ermon. “Generative Modeling
by Estimating Gradients of the Data Distribution”. In:
51 *NIPS*. 2019, pp. 11895–11907.
- [SE20a] J. Song and S. Ermon. “Multi-label Con-
trastive Predictive Coding”. In: *NIPS*. 2020.
- [SE20b] Y. Song and S. Ermon. “Improved Techniques
for Training Score-Based Generative Models”. In:
52 *NIPS*. 2020.
- [See+17] M. Seeger, S. Rangapuram, Y. Wang, D. Salin-
53 as, J. Gasthaus, T. Januschowski, and V. Flunkert.
“Approximate Bayesian Inference in Linear State
Space Models for Intermittent Demand Forecasting at
Scale”. In: (2017). arXiv: [1709.07638 \[stat.ML\]](#).
- [Sej20] T. J. Sejnowski. “The unreasonable effective-
54 ness of deep learning in artificial intelligence”. en. In:
55 *PNAS* 117.48 (Dec. 2020), pp. 30033–30038.

- 1 [Sel+17] R. R. Selvaraju, M. Cogswell, A. Das, R.
2 Vedantam, D. Parikh, and D. Batra. “Grad-cam: Vi-
3 sual explanations from deep networks via gradient-
4 based localization”. In: *Proceedings of the IEEE in-
5 ternational conference on computer vision*. 2017,
6 pp. 618–626.
- 7 [Sel+19] A. D. Selbst, D. Boyd, S. A. Friedler, S.
8 Venkatasubramanian, and J. Vertesi. “Fairness and
9 Abstraction in Sociotechnical Systems”. In: *Pro-
10 ceedings of the Conference on Fairness, Accountabil-
11 ity, and Transparency*. FAT* ’19, Atlanta, GA, USA: As-
12 sociation for Computing Machinery, 2019, 59–68.
- 13 [Sen+08] P. Sen, G. Namata, M. Bilgic, L. Getoor,
14 B. Galligher, and T. Eliassi-Rad. “Collective Classi-
15 fication in Network Data”. en. In: *AI Magazine* 29.3
16 (2008), pp. 93–93.
- 17 [Ser+20] J. Serrà, D. Álvarez, V. Gómez, O. Sli-
18 zovskaia, J. F. Núñez, and J. Luque. “Input complex-
19 ity and out-of-distribution detection with likelihood-
20 based generative models”. In: *ICLR*. 2020.
- 21 [Set12] B. Settles. “Active learning”. In: *Synthesis Lectures on Artificial Intelligence and Machine Learn-
22 ing* 6 (2012), 1–114.
- 23 [Set94] J. Sethuraman. “A constructive definition of
24 Dirichlet priors”. In: *Statistica Sinica* (1994), pp. 639–
25 650.
- 26 [SF08] Z. Svitkina and L. Fleischer. “Submodular ap-
27 proximation: Sampling-based algorithms and lower
28 bounds”. In: *FOCS*. 2008.
- 29 [SF20] K. Sokol and P. Flach. “Explainability fact
30 sheets: a framework for systematic assessment of ex-
31 plainable approaches”. In: *Proceedings of the 2020
32 Conference on Fairness, Accountability, and Trans-
33 parency*. 2020, pp. 56–67.
- 34 [SFB18] S. A. Sisson, Y. Fan, and M. A. Beau-
35 mont. “Overview of ABC”. In: *Handbook of approx-
36 imate Bayesian computation*. Chapman and Hall/
37 CRC, 2018, pp. 3–54.
- 38 [SG05] E. Snelson and Z. Ghahramani. “Compact Ap-
39 proximations to Bayesian Predictive Distributions”. In:
40 *ICML*. 2005.
- 41 [SG06a] E. Snelson and Z. Ghahramani. “Sparse Gaus-
42 sian processes using pseudo-inputs”. In: *NIPS*. 2006.
- 43 [SG06b] E. Snelson and Z. Ghahramani. “Sparse Gaus-
44 sian Processes using Pseudo-inputs”. In: *Advances in
45 Neural Information Processing Systems*. Ed. by Y.
46 Weiss, B. Schölkopf, and J. Platt. Vol. 18. MIT Press,
47 2006.
- 48 [SG07] M. Steyvers and T. Griffiths. “Probabilistic
49 topic models”. In: *Latent Semantic Analysis: A Road
50 to Meaning*. Ed. by T. Landauer, D. McNamara, S.
51 Dennis, and W. Kintsch. Laurence Erlbaum, 2007.
- 52 [SG09] R. Silva and Z. Ghahramani. “The Hidden Life
53 of Latent Variables: Bayesian Learning with Mixed
54 Graph Models”. In: *JMLR* 10 (2009), pp. 1187–1238.
- 55 [SGF21] S. Särkkä and Á. F. García-Fernández.
56 “Temporal Parallelization of Bayesian Filters and
57 Smoothers”. In: *IEEE Trans. Automat. Contr.* 66.1
58 (2021).
- 59 [SGS16] A. Sharghi, B. Gong, and M. Shah. “Query-
60 focused extractive video summarization”. In: *Euro-
61 pean Conference on Computer Vision*. Springer.
62 2016, pp. 3–19.
- 63 [SH07] R. Salakhutdinov and G. Hinton. “Using Deep
64 Belief Nets to Learn Covariance Kernels for Gaussian
65 Processes”. In: *NIPS*. 2007.
- 66 [SH09] R. Salakhutdinov and G. Hinton. “Deep Boltz-
67 mann Machines”. In: *AISTATS*. Vol. 5. 2009, pp. 448–
68 455.
- 69 [SH10] R. Salakhutdinov and G. Hinton. “Replicated
70 Softmax: an Undirected Topic Model”. In: *NIPS*. 2010.
- 71 [SH22] T. Salimans and J. Ho. “Progressive Distillation
72 for Fast Sampling of Diffusion Models”. In: *ICLR*. Feb.
73 2022.
- 74 [SHA15] M. A. Skoglund, G. Hendeby, and D. Axe-
75 hill. “Extended Kalman filter modifications based on
76 an optimization view point”. In: *2015 18th Interna-
77 tional Conference on Information Fusion (Fusion)*.
78 July 2015, pp. 1856–1861.
- 79 [Sha+16] B. Shahriari, K. Swersky, Z. Wang, R. P.
80 Adams, and N de Freitas. “Taking the Human Out
81 of the Loop: A Review of Bayesian Optimization”. In:
82 *Proc. IEEE* 104.1 (2016), pp. 148–175.
- 83 [Sha16] L. S. Shapley. *17. A value for n-person games*.
84 Princeton University Press, 2016.
- 85 [Sha+16] M. Sharif, S. Bhagavatula, L. Bauer, and
86 M. K. Reiter. “Accessorize to a Crime: Real and
87 Stealthy Attacks on State-of-the-Art Face Recogni-
88 tion”. In: *Proceedings of the 2016 ACM SIGSAC
89 Conference on Computer and Communications Se-
90 curity*. ACM, 2016, pp. 1528–1540.
- 91 [Sha+19] A. Shaikhha, A. Fitzgibbon, D. Vytiniotis,
92 and S. Peyton Jones. “Efficient differentiable program-
93 ming in a functional array-processing language”. In:
94 *Proceedings of the ACM on Programming Languages*
95 3.ICFP (2019), pp. 1–30.
- 96 [Sha+20] H. Shah, K. Tamuly, A. Raghunathan, P.
97 Jain, and P. Netrapalli. “The Pitfalls of Simplicity
98 Bias in Neural Networks”. In: *NIPS*. 2020.
- 99 [Sha22] C. Shalizi. *Advanced Data Analysis from an
100 Elementary Point of View*. Cambridge University
101 Press, 2022.
- 102 [Sha48] C. Shannon. “A mathematical theory of com-
103 munication”. In: *Bell Systems Tech. Journal* 27
104 (1948), pp. 379–423.
- 105 [Sha98] R. Shachter. “Bayes-Ball: The Rational Pas-
106 time (for determining Irrelevance and Requisite Infor-
107 mation in Belief Networks and Influence Diagrams)”.
108 In: *UAI*. 1998.
- 109 [She+11] C. Shen, X. Li, L. Li, and M. C. Were. “Sensi-
110 tivity analysis for causal inference using inverse proba-
111 bility weighting”. In: *Biometrical Journal* 53.5 (2011),
112 pp. 822–837. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/bimj.201100042>.
- 113 [She+17] T. Shen, T. Lei, R. Barzilay, and T. Jaakkola.
114 “Style transfer from non-parallel text by cross-
115 alignment”. In: *Advances in neural information pro-
116 cessing systems* 30 (2017), pp. 6830–6841.
- 117 [She+20] T. Shen, J. Mueller, R. Barzilay, and T.
118 Jaakkola. “Educating Text Autoencoders: Latent Rep-
119 resentation Guidance via Denoising”. In: *ICML*. 2020.
- 120 [She+21] Z. Shen, J. Liu, Y. He, X. Zhang, R. Xu, H.
121 Yu, and P. Cui. “Towards Out-Of-Distribution Gen-
122 eralization: A Survey”. In: (2021). arXiv: [2108.13624 \[cs.LG\]](https://arxiv.org/abs/2108.13624).
- 123 [SHF15] R. Steorts, R. Hall, and S. Fienberg. “A
124 Bayesian Approach to Graphical Record Linkage and
125 De-duplication”. In: *JASA* (2015).

- 1
- 2 [Shi00] B. Shipley. *Cause and Correlation in Biology: A User's Guide to Path Analysis, Structural Equations and Causal Inference*. Cambridge, 2000.
- 3
- 4 [Shi+21] C. Shi, D. Sridhar, V. Misra, and D. M. Blei. “On the Assumptions of Synthetic Control Methods”. In: (2021). arXiv: 2112.05671 [stat.ME].
- 5
- 6 [SHM14] D. Soudry, I. Hubara, and R. Meir. “Expectation backpropagation: Parameter-free training of multilayer neural networks with continuous or discrete weights”. In: *NIPS*. 2014.
- 7
- 8
- 9 [Shr+16] A. Shrikumar, P. Greenside, A. Shcherbina, and A. Kundaje. “Not just a black box: Learning important features through propagating activation differences”. In: *arXiv preprint arXiv:1605.01713* (2016).
- 10
- 11
- 12 [SHS] D. Stutz, M. Hein, and B. Schiele. “Confidence-calibrated adversarial training: Generalizing to unseen attacks”. In: () .
- 13
- 14 [SHS01] B. Schölkopf, R. Herbrich, and A. J. Smola. “A Generalized Representer Theorem”. In: *COLT. COLT '01/EuroCOLT '01*. Springer-Verlag, 2001, pp. 416–426.
- 15
- 16
- 17 [Shu+19a] R. Shu, L. Cui, S. Wang, D. Lee, and H. Liu. “defend: Explainable fake news detection”. In: *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 2019, pp. 395–405.
- 18
- 19
- 20 [Shu+19b] R. Shu, Y. Chen, A. Kumar, S. Ermon, and B. Poole. “Weakly supervised disentanglement with guarantees”. In: *arXiv preprint arXiv:1910.09772* (2019).
- 21
- 22
- 23 [SI00] M. Sato and S. Ishii. “On-line EM algorithm for the normalized Gaussian network”. In: *Neural Computation* 12 (2000), pp. 407–432.
- 24
- 25 [Sil+14] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller. “Deterministic Policy Gradient Algorithms”. In: *ICML. ICML'14. JMLR.org*, 2014, pp. I–387–I–395.
- 26
- 27 [Sil+16] D. Silver et al. “Mastering the game of Go with deep neural networks and tree search”. en. In: *Nature* 529.7587 (2016), pp. 484–489.
- 28
- 29
- 30 [Sil18] D. Silver. *Lecture 9L Exploration and Exploitation*. 2018.
- 31 [Sil+18] D. Silver et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play”. en. In: *Science* 362.6419 (2018), pp. 1140–1144.
- 32
- 33
- 34 [Sil85] B. W. Silverman. “Some Aspects of the Spline Smoothing Approach to Non-Parametric Regression Curve Fitting”. In: *J. R. Stat. Soc. Series B Stat. Methodol.* 47.1 (1985), pp. 1–52.
- 35
- 36
- 37 [Sim06] D. Simon. *Optimal State Estimation: Kalman, H Infinity, and Nonlinear Approaches*. Wiley, 2006.
- 38
- 39 [Sin+00] S. Singh, T. Jaakkola, M. L. Littman, and C. Szepesvári. “Convergence Results for Single-Step On-Policy Reinforcement-Learning Algorithms”. In: *MLJ* 38.3 (2000), pp. 287–308.
- 40
- 41
- 42 [Sin67] R. Sinkhorn. “Diagonal Equivalence to Matrices with Prescribed Row and Column Sums”. In: *The American Mathematical Monthly* 74.4 (1967), pp. 402–405.
- 43
- 44 [SJ08] S. Shirdhonkar and D. W. Jacobs. “Approximate earth mover's distance in linear time”. In: *2008 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2008, pp. 1–8.
- 45
- 46
- 47
- [SJ15] A. Swaminathan and T. Joachims. “Batch Learning from Logged Bandit Feedback through Counterfactual Risk Minimization”. In: *JMLR* 16.1 (2015), pp. 1731–1755.
- [SJ95] L. Saul and M. Jordan. “Exploiting tractable substructures in intractable networks”. In: *NIPS*. Vol. 8. 1995.
- [SJ99] L. Saul and M. Jordan. “Mixed memory Markov models: Decomposing complex stochastic processes as mixture of simpler ones”. In: *Machine Learning* 37.1 (1999), pp. 75–87.
- [JJ96] L. Saul, T. Jaakkola, and M. Jordan. “Mean Field Theory for Sigmoid Belief Networks”. In: *JAIR* 4 (1996), pp. 61–76.
- [SJR04] S. Singh, M. James, and M. Rudary. “Predictive state representations: A new theory for modeling dynamical systems”. In: *UAI*. 2004.
- [SK19] C. Shorten and T. M. Khoshgoftaar. “A survey on Image Data Augmentation for Deep Learning”. en. In: *Journal of Big Data* 6.1 (2019), pp. 1–48.
- [SK20] S. Singh and S. Krishnan. “Filter Response Normalization Layer: Eliminating Batch Dependence in the Training of Deep Neural Networks”. In: *CVPR*. 2020.
- [SK89] R. Shachter and C. R. Kenley. “Gaussian Influence Diagrams”. In: *Management Science* 35.5 (1989), pp. 527–550.
- [Ski89] J. Skilling. “The eigenvalues of mega-dimensional matrices”. In: *Maximum Entropy and Bayesian Methods*. Springer, 1989, pp. 455–466.
- [SKM18] S. Schwöbel, S. Kiebel, and D. Marković. “Active Inference, Belief Propagation, and the Bethe Approximation”. en. In: *Neural Comput.* 30.9 (2018), pp. 2530–2567.
- [SKM21] M. Shanahan, C. Kaplanić, and J. Mitrović. “Encoders and Ensembles for Task-Free Continual Learning”. In: (2021). arXiv: 2105.13327 [cs.LG].
- [SKP15] F. Schroff, D. Kalenichenko, and J. Philbin. “Facenet: A unified embedding for face recognition and clustering”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 815–823.
- [SKTF18] H. Shao, A. Kumar, and P. Thomas Fletcher. “The Riemannian Geometry of Deep Generative Models”. In: *CVPR*. 2018, pp. 315–323.
- [SKW15] T. Salimans, D. Kingma, and M. Welling. “Markov Chain Monte Carlo and Variational Inference: Bridging the Gap”. In: *ICML*. 2015, pp. 1218–1226.
- [SL08] A. L. Strehl and M. L. Littman. “An Analysis of Model-based Interval Estimation for Markov Decision Processes”. In: *J. of Comp. and Sys. Sci.* 74.8 (2008), pp. 1309–1331.
- [SL18] S. L. Smith and Q. V. Le. “A Bayesian Perspective on Generalization and Stochastic Gradient Descent”. In: *ICLR*. 2018.
- [SL90] D. J. Spiegelhalter and S. L. Lauritzen. “Sequential updating of conditional probabilities on directed graphical structures”. In: *Networks* 20 (1990).
- [Sla+20] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju. “Fooling lime and shap: Adversarial attacks on post hoc explanation methods”. In: *Proceed-*

- ings of the AAAI/ACM Conference on AI, Ethics, and Society.* 2020, pp. 180–186.

[SLG17] A. Sharghi, J. S. Laurel, and B. Gong. “Query-focused video summarization: Dataset, evaluation, and a memory network based approach”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 4788–4797.

[Sli19] A. Slivkins. “Introduction to Multi-Armed Bandits”. In: *Foundations and Trends in Machine Learning* (2019).

[SLL09] A. L. Strehl, L. Li, and M. L. Littman. “Reinforcement Learning in Finite MDPs: PAC Analysis”. In: *JMLR* 10 (2009), pp. 2413–2444.

[SLM92] B. Selman, H. Levesque, and D. Mitchell. “A New Method for Solving Hard Satisfiability Problems”. In: *Proceedings of the Tenth National Conference on Artificial Intelligence*. AAAI’92. AAAI Press, 1992, pp. 440–446.

[SLW19] M. Sadinle, J. Lei, and L. Wasserman. “Least Ambiguous Set-Valued Classifiers With Bounded Error Levels”. In: *JASA* 114.525 (2019), pp. 223–234.

[SM07] C. Sutton and A. McCallum. “Improved Dynamic Schedules for Belief Propagation”. In: *UAI*. 2007.

[SM12] Y Saika and K Morimoto. “Generalized MAP estimation via parameter scheduling and maximizer of the posterior marginal estimate for image reconstruction using multiple halftone images”. In: *12th International Conference on Control, Automation and Systems*. 2012, pp. 1285–1289.

[SMB10] H. Schulz, A. Müller, and S. Behnke. “Investigating convergence of restricted Boltzmann machine learning”. In: *NIPS 2010 Workshop on Deep Learning and Unsupervised Feature Learning*. Vol. 1. 2. 2010, pp. 6–1.

[SME21] J. Song, C. Meng, and S. Ermon. “Denoising Diffusion Implicit Models”. In: *ICLR*. 2021.

[SMH07] R. R. Salakhutdinov, A. Mnih, and G. E. Hinton. “Restricted Boltzmann machines for collaborative filtering”. In: *ICML*. Vol. 24. 2007, pp. 791–798.

[Smi+00] G. Smith, J. F. G. de Freitas, T. Robinson, and M. Niranjani. “Speech Modelling Using Subspace and EM Techniques”. In: *NIPS*. MIT Press, 2000, pp. 796–802.

[Smi+06] V. Smith, J. Yu, T. Smulders, A. Hartemink, and E. Jarvis. “Computational Inference of Neural Information Flow Networks”. In: *PLOS Computational Biology* 2 (2006), pp. 1436–1439.

[Smi11] N. Smith. *Linguistic structure prediction*. Morgan Claypool, 2011.

[Smi+17] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg. *SmoothGrad: removing noise by adding noise*. 2017. arXiv: 1706.03825 [cs.LG].

[Smo86] P. Smolensky. “Information processing in dynamical systems: foundations of harmony theory”. In: *Parallel Distributed Processing: Explorations in the Microstructure of Cognition. Volume 1*. Ed. by D. Rumelhart and J. McClelland. McGraw-Hill, 1986.

[SMS17] M. Saito, E. Matsumoto, and S. Saito. “Temporal Generative Adversarial Nets with Singular Value Clipping”. In: *ICCV*. 2017.

[SMT18] M. R. U. Saputra, A. Markham, and N. Trigoni. “Visual SLAM and Structure from Motion in Dynamic Environments: A Survey”. In: *ACM Computing Surveys* 51.2 (2018), pp. 1–36.

[Smy20] S. Smyl. “A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting”. In: *Int. J. Forecast.* 36.1 (2020), pp. 75–85.

[Sn+16] C. K. Sønderby, T. Raiko, L. Maaløe, S. R. K. Sønderby, and O. Winther. “Ladder Variational Autoencoders”. In: *NIPS*. Curran Associates, Inc., 2016, pp. 3738–3746.

[SMN16] M. Suzuki, K. Nakayama, and Y. Matsuo. “Joint Multimodal Learning with Deep Generative Models”. In: (2016). arXiv: 1611.01891 [stat.ML].

[SOB12] R. Snyder, J. K. Ord, and A. Beaumont. “Forecasting the intermittent demand for slow-moving inventories: A modelling approach”. In: *Int. J. Forecast.* 28.2 (2012), pp. 485–496.

[Soh16] K. Sohn. “Improved deep metric learning with multi-class n-pair loss objective”. In: *Advances in Neural Information Processing Systems*. 2016, pp. 1857–1865.

[Søn+16] C. Sønderby, T. Raiko, L. Maaløe, S. Sønderby, and O. Winther. “How to Train Deep Variational Autoencoders and Probabilistic Ladder Networks”. In: *ICML*. 2016.

[Son+19] Y. Song, S. Garg, J. Shi, and S. Ermon. “Sliced Score Matching: A Scalable Approach to Density and Score Estimation”. In: *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2019, Tel Aviv, Israel, July 22–25, 2019*, p. 204.

[Son+21] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. “Score-Based Generative Modeling through Stochastic Differential Equations”. In: *ICLR*. 2021.

[Son98] E. D. Sontag. *Mathematical Control Theory: Deterministic Finite Dimensional Systems*. 2nd. Vol. 6. Texts in Applied Mathematics. Springer, 1998.

[SOS92] H. Seung, M. Opper, and H. Sompolinsky. “Query by committee”. In: *5th Annual Workshop on Computational Learning Theory*. 1992, 287–294.

[SPD92] S. Shah, F. Palmieri, and M. Datum. “Optimal filtering algorithms for fast learning in feedforward neural networks”. In: *Neural Netw.* 5.5 (1992), pp. 779–787.

[Spi71] M. Spivak. *Calculus On Manifolds: A Modern Approach To Classical Theorems Of Advanced Calculus*. Westview Press; 5th edition, 1971.

[SPL20] M. B. Sariyildiz, J. Perez, and D. Larlus. “Learning visual representations with caption annotations”. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*. Springer, 2020, pp. 153–170.

[Spr+14] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller. “Striving for simplicity: The all convolutional net”. In: *arXiv preprint arXiv:1412.6806* (2014).

[Spr+16] J. T. Springenberg, A. Klein, S. Falkner, and F. Hutter. “Bayesian Optimization with Robust Bayesian Neural Networks”. In: *NIPS*. 2016, pp. 4141–4149.

[SPW18] M. H. S. Segler, M. Preuss, and M. P. Waller. “Planning chemical syntheses with deep neural networks and symbolic AI”. In: *Nature* 555.7698 (2018), pp. 604–610.

[SPZ09] P. Schniter, L. C. Potter, and J. Ziniel. “Fast Bayesian Matching Pursuit: Model Uncertainty and

- 1 Parameter Estimation for Sparse Linear Models". In: *IEEE Trans. on Signal Processing* (2009).
- 2 [SR+14] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. "CNN features off-the-shelf: an astounding baseline for recognition". In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 2014, pp. 806–813.
- 3 [SRG03] R. Salakhutdinov, S. T. Roweis, and Z. Ghahramani. "Optimization with EM and Expectation-Conjugate-Gradient". In: *ICML*. 2003.
- 4 [Sri+09] B. K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. R. G. Lanckriet. "On integral probability metrics, φ -divergences and binary classification". In: (2009). arXiv: [0901.2698 \[cs.IT\]](#).
- 5 [Sri+10] N. Srinivas, A. Krause, S. Kakade, and M. Seeger. "Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design". In: *ICML*. 2010, pp. 1015–1022.
- 6 [Sri+14a] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". In: *JMLR* (2014).
- 7 [Sri+14b] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting". In: *J. Mach. Learn. Res.* 15 (2014), pp. 1929–1958.
- 8 [Sri+17] A. Srivastava, L. Valkov, C. Russell, M. U. Gutmann, and C. Sutton. "Veegan: Reducing mode collapse in gans using implicit variational learning". In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 2017, pp. 3310–3320.
- 9 [SRS10] P. Schnitzspan, S. Roth, and B. Schiele. "Automatic discovery of meaningful object parts with latent CRFs". In: *CVPR*. 2010.
- 10 [SS17a] A. Srivastava and C. Sutton. "Autoencoding Variational Inference For Topic Models". In: *ICLR*. 2017.
- 11 [SS17b] A. Srivastava and C. Sutton. "Autoencoding Variational Inference For Topic Models". In: *ICLR*. 2017.
- 12 [SS18a] O. Sener and S. Savarese. "Active Learning for Convolutional Neural Networks: A Core-Set Approach". In: *International Conference on Learning Representations*. 2018.
- 13 [SS18b] A. Subbaswamy and S. Saria. "Counterfactual Normalization: Proactively Addressing Dataset Shift and Improving Reliability Using Causal Mechanisms". In: *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018. 2018.
- 14 [SS19] S. Sarkka and A. Solin. *Applied stochastic differential equations*. en. Cambridge University Press, 2019.
- 15 [SS20a] S. Sarkka and L. Svensson. "Levenberg-Marquardt and Line-Search Extended Kalman Smoothers". In: *ICASSP*. Barcelona, Spain: IEEE, May 2020.
- 16 [SS20b] K. E. Smith and A. O. Smith. "Conditional GAN for timeseries generation". In: *arXiv preprint arXiv:2006.16477* (2020).
- 17 [SS21] I. Sucholutsky and M. Schonlau. "Soft-Label Dataset Distillation and Text Dataset Distillation". In: *2021 International Joint Conference on Neural Networks (IJCNN)*. 2021, pp. 1–8.
- 18 [SS82] R. H. Shumway and D. S. Stoffer. "An approach to time series smoothing and forecasting using the em algorithm". en. In: *J. Time Ser. Anal.* 3.4 (July 1982), pp. 253–264.
- 19 [SSA14] K. Swersky, J. Snoek, and R. P. Adams. "Freeze-Thaw Bayesian Optimization". In: (2014). arXiv: [1406.3896 \[stat.ML\]](#).
- 20 [SSA18] K. Shmelkov, C. Schmid, and K. Alahari. "How good is my GAN?" In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 213–229.
- 21 [SSB17] S. Semeniuta, A. Severyn, and E. Barth. "A Hybrid Convolutional Variational Autoencoder for Text Generation". In: (2017). arXiv: [1702.02390 \[cs.CL\]](#).
- 22 [SSE18] Y. Song, J. Song, and S. Ermon. "Accelerating Natural Gradient with Higher-Order Invariance". In: *ICML*. 2018.
- 23 [SSF16] M. W. Seeger, D. Salinas, and V. Flunkert. "Bayesian Intermittent Demand Forecasting for Large Inventories". In: *NIPS*. 2016, pp. 4646–4654.
- 24 [SSG18a] S. Semeniuta, A. Severyn, and S. Gelly. "On Accurate Evaluation of GANs for Language Generation". In: *arXiv preprint arXiv:1806.04936* (2018).
- 25 [SSG18b] S. Semeniuta, A. Severyn, and S. Gelly. "On Accurate Evaluation of GANs for Language Generation". In: (2018). arXiv: [1806.04936 \[cs.CL\]](#).
- 26 [SSG19] R. Singh, M. Sahani, and A. Gretton. "Kernel Instrumental Variable Regression". In: *Advances in Neural Information Processing Systems*. 2019, pp. 4593–4605.
- 27 [SSH13] S. Sarkka, A. Solin, and J. Hartikainen. "Spatio-Temporal Learning via Infinite-Dimensional Bayesian Filtering and Smoothing: A look at Gaussian process regression through Kalman filtering". In: *IEEE Signal Processing Magazine* (2013).
- 28 [SSJ12] R. Sipos, P. Shivaswamy, and T. Joachims. "Large-margin learning of submodular summarization models". In: *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*. 2012, pp. 224–233.
- 29 [SSK12] M. Sugiyama, T. Suzuki, and T. Kanamori. *Density Ratio Estimation in Machine Learning*. en. Cambridge University Press, 2012.
- 30 [SSM18] S. Santurkar, L. Schmidt, and A. Madry. "A classification-based study of covariate shift in gan distributions". In: *International Conference on Machine Learning*. PMLR. 2018, pp. 4480–4489.
- 31 [SSZ17] J. Snell, K. Swersky, and R. Zemel. "Prototypical networks for few-shot learning". In: *NIPS*. 2017, pp. 4077–4087.
- 32 [Sta] *Scientific Explanation*. <https://plato.stanford.edu/entries/scientific-explanation/#ConcOpenIssueFuture>. Accessed: 2021-11-23.
- 33 [Sta07] K. O. Stanley. "Compositional pattern producing networks: A novel abstraction of development". In: *Genet. Program. Evolvable Mach.* 8.2 (2007), pp. 131–162.
- 34 [Sta+17] N. Stallard, F. Miller, S. Day, S. W. Hee, J. Madan, S. Zohar, and M. Posch. "Determination of the optimal sample size for a clinical trial accounting for the population size". en. In: *Biom. J.* 59.4 (2017), pp. 609–625.
- 35 [Sta+19] K. O. Stanley, J. Clune, J. Lehman, and R. Miikkulainen. "Designing neural networks through

- 1 “neuroevolution”. In: *Nature Machine Intelligence* 1.1
 2 (2019).
- 3 [Ste81] C. M. Stein. “Estimation of the mean of a multi-
 4 variate normal distribution”. In: *The annals of Statistics* (1981), pp. 1135–1151.
- 5 [Sto09] A. J. Storkey. “When Training and Test Sets
 6 are Different: Characterising Learning Transfer”. In:
 7 *Dataset Shift in Machine Learning*. 2009.
- 8 [Sto17] J. Stoehr. “A review on statistical inference
 9 methods for discrete Markov random fields”. In: (2017).
 arXiv: [1704.03331 \[stat.ME\]](https://arxiv.org/abs/1704.03331).
- 10 [STR10] Y. Saatchi, R. Turner, and C. E. Rasmussen.
 11 “Gaussian Process Change Point Models”. In: *ICML*.
 unknown, 2010, pp. 927–934.
- 12 [Str+17] H. Strobelt, S. Gehrmann, H. Pfister, and
 13 A. M. Rush. “Lstmyis: A tool for visual analysis of
 14 hidden state dynamics in recurrent neural networks”.
 15 In: *IEEE transactions on visualization and computer
 graphics* 24.1 (2017), pp. 667–676.
- 16 [Str19] M. Streeter. “Bayes Optimal Early Stopping
 17 Policies for Black-Box Optimization”. In: (2019).
 arXiv: [1902.08285 \[cs.LG\]](https://arxiv.org/abs/1902.08285).
- 18 [Stu+22] D. Stutz, Krishnamurthy, Dvijotham, A. T.
 19 Cemgil, and A. Doucet. “Learning Optimal Conformal
 Classifiers”. In: *ICLR*. 2022.
- 20 [STY17] M. Sundararajan, A. Taly, and Q. Yan. *Axiomatic Attribution for Deep Networks*. 2017. arXiv:
 1703.01365 [cs.LG].
- 21 [Suc+20] F. P. Such, A. Rawal, J. Lehman, K. Stanley,
 22 and J. Clune. “Generative teaching networks: Accelerating
 23 neural architecture search by learning to generate synthetic
 24 training data”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 9206–
 25 9216.
- 26 [Sud+03] E. Sudderth, A. Ihler, W. Freeman, and
 27 A. Willsky. “Nonparametric Belief Propagation”. In:
 28 *CVPR*. 2003.
- 29 [Sud06] E. Sudderth. “Graphical Models for Visual Ob-
 30 ject Recognition and Tracking”. PhD thesis. MIT,
 2006.
- 31 [Sud+10] E. Sudderth, A. Ihler, M. Isard, W. Freeman,
 32 and A. Willsky. “Nonparametric Belief Propagation”.
 In: *Comm. of the ACM* 53.10 (2010).
- 33 [Sug+13] M. Sugiyama, T. Kanamori, T. Suzuki, M. C.
 34 du Plessis, S. Liu, and I. Takeuchi. “Density-difference
 estimation”. en. In: *Neural Comput.* 25.10 (2013), pp. 2734–2775.
- 35 [Sun+09] L. Sun, S. Ji, S. Yu, and J. Ye. “On the Equiva-
 36 lence Between Canonical Correlation Analysis and
 Orthonormalized Partial Least Squares”. In: *IJCAI*.
 2009.
- 37 [Sun+17] C. Sun, A. Shrivastava, S. Singh, and A.
 38 Gupta. “Revisiting unreasonable effectiveness of data
 39 in deep learning era”. In: *Proceedings of the IEEE
 international conference on computer vision*. 2017,
 pp. 843–852.
- 40 [Sun+18] S. Sun, G. Zhang, C. Wang, W. Zeng, J. Li,
 41 and R. Grosse. “Differentiable Compositional Kernel
 Learning for Gaussian Processes”. In: *ICML*. 2018.
- 42 [Sun+19a] S. Sun, G. Zhang, J. Shi, and R. Grosse.
 43 “Functional variational bayesian neural networks”. In:
 44 *arXiv preprint arXiv:1903.05779* (2019).
- 45 [Sun+19b] S. Sun, Z. Cao, H. Zhu, and J. Zhao. “A
 46 Survey of Optimization Methods from a Machine
 47 Learning Perspective”. In: (2019). arXiv: [1906.06821 \[cs.LG\]](https://arxiv.org/abs/1906.06821).
- [Sun+19c] M. Sundararajan, J. Xu, A. Taly, R. Sayres,
 and A. Najmi. “Exploring Principled Visualizations
 for Deep Network Attributions.” In: *IUI Workshops*. Vol. 4. 2019.
- [Sun+20] Y. Sun, X. Wang, Z. Liu, J. Miller, A.
 Efros, and M. Hardt. “Test-Time Training with Self-
 Supervision for Generalization under Distribution
 Shifts”. In: *ICML*. Vol. 119. Proceedings of Machine
 Learning Research. PMLR, 2020, pp. 9229–9248.
- [Sun+22] T. Sun, M. Segu, J. Postels, Y. Wang, L. Van
 Cool, B. Schiele, F. Tombari, and F. Yu. “SHIFT: A
 Synthetic Driving Dataset for Continuous Multi-Task
 Domain Adaptation”. In: *CVPR*. June 2022.
- [Sut+17] D. J. Sutherland, H.-Y. Tung, H. Strathmann,
 S. De, A. Ramdas, A. Smola, and A. Gretton. “Gener-
 ative Models and Model Criticism via Optimized Max-
 imum Mean Discrepancy”. In: *ICLR*. 2017.
- [Sut19] R. Sutton. *The Bitter Lesson*. 2019.
- [Sut88] R. Sutton. “Learning to predict by the meth-
 ods of temporal differences”. In: *Machine Learning* 3.1
 (1988), pp. 9–44.
- [Sut90] R. S. Sutton. “Integrated Architectures for
 Learning, Planning, and Reacting Based on Approx-
 imating Dynamic Programming”. In: *ICML*. Ed. by
 B. Porter and R. Mooney. Morgan Kaufmann, 1990,
 pp. 216–224.
- [Sut96] R. S. Sutton. “Generalization in Reinforce-
 ment Learning: Successful Examples Using Sparse
 Coarse Coding”. In: *NIPS*. Ed. by D. S. Touretzky,
 M. C. Mozer, and M. E. Hasselmo. MIT Press, 1996,
 pp. 1038–1044.
- [Sut+99] R. Sutton, D. McAllester, S. Singh, and Y.
 Mansour. “Policy Gradient Methods for Reinforce-
 ment Learning with Function Approximation”. In:
NIPS. 1999.
- [SV08] G. Shafer and V. Vovk. “A Tutorial on Confor-
 mal Prediction”. In: *JMLR* 9.Mar (2008), pp. 371–421.
- [SV98] M. Studenty and J. Vejnarova. “The multi-
 information function as a tool for measuring stochastic
 dependence”. In: *Learning in graphical models*. Ed. by
 M. Jordan. MIT Press, 1998, pp. 261–297.
- [SVE04] A. Smola, S. V. N. Vishwanathan, and E.
 Eskin. “Laplace Propagation”. In: *NIPS*. MIT Press,
 2004, pp. 441–448.
- [Svi04] M. Sviridenko. “A note on maximizing a sub-
 modular set function subject to a knapsack constraint”.
 In: *Operations Research Letters* 32.1 (2004), pp. 41–
 43.
- [SVK19] J. Su, D. V. Vargas, and S. Kouichi. “One
 pixel attack for fooling deep neural networks”. In:
IEEE Trans. Evol. Comput. 23.5 (2019).
- [SVZ13] K. Simonyan, A. Vedaldi, and A. Zisserman.
 “Deep inside convolutional networks: Visualising im-
 age classification models and saliency maps”. In: *arXiv
 preprint arXiv:1312.6034* (2013).
- [SW06] J. E. Smith and R. L. Winkler. “The Opti-
 mizer’s Curse: Skepticism and Postdecision Surprise
 in Decision Analysis”. In: *Manage. Sci.* 52.3 (2006),
 pp. 311–322.
- [SW13] G. J. Sussman and J. Wisdom. *Functional Dif-
 ferential Geometry*. Functional Differential Geometry.
 MIT Press, 2013.

- 1
- 2 [SW20] V. G. Satorras and M. Welling. “Neural Enhanced Belief Propagation on Factor Graphs”. In: 3 (2020). arXiv: 2003.01998 [cs.LG].
- 4 [SW87a] M. Shewry and H. Wynn. “Maximum entropy sampling”. In: *J. Applied Statistics* 14 (1987), 5 165–170.
- 6 [SW87b] R. Swendsen and J.-S. Wang. “Nonuniversal critical dynamics in Monte Carlo simulations”. In: 7 *Physical Review Letters* 58 (1987), pp. 86–88.
- 8 [SW89] S. Singhal and L. Wu. “Training Multilayer Perceptrons with the Extended Kalman Algorithm”. In: 9 *NIPS*. Vol. 1. 1989.
- 10 [Swe+10] K. Swersky, B. Chen, B. Marlin, and N. de 11 Freitas. “A Tutorial on Stochastic Approximation Algorithms for Training Restricted Boltzmann Machines 12 and Deep Belief Nets”. In: *Information Theory and Applications (ITA) Workshop*. 2010.
- 14 [Swe+13] K. Swersky, D. Duvenaud, J. Snoek, F. Hutter, and M. A. Osborne. “Raiders of the Lost Architecture: Kernels for Bayesian Optimization in Conditional Parameter Spaces”. In: *NIPS BayesOpt workshop*. 2013.
- 17 [SWL03] M. Seeger, C. K. I. Williams, and N. D. Lawrence. “Fast Forward Selection to Speed Up Sparse Gaussian Process Regression”. In: *AISTATS*. 2003.
- 19 [SWW08] E. Sudderth, M. Wainwright, and A. Willsky. “Loop series and Bethe variational bounds in attractive graphical models”. In: *NIPS* (2008).
- 21 [SYD19] R. Sen, H.-F. Yu, and I. Dhillon. “Think Globally, Act Locally: A Deep Neural Network Approach to High-Dimensional Time Series Forecasting”. In: *NIPS*. 23 2019.
- 24 [SZ22] R. Schwartz-Ziv. “Information Flow in Deep Neural Networks”. PhD thesis. 2022.
- 25 [SZ+22] R. Schwartz-Ziv, M. Goldblum, H. Souris, S. Kapoor, C. Zhu, Y. LeCun, and A. G. Wilson. “Pre-Train Your Loss: Easy Bayesian Transfer Learning 27 with Informative Priors”. In: (May 2022). arXiv: 2205. 28 10279 [cs.LG].
- 29 [Sze10] C. Szepesvari. *Algorithms for Reinforcement Learning*. Morgan Claypool, 2010.
- 30 [Sze+14] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. “Intriguing properties of neural networks”. In: *ICLR*. 32 2014.
- 33 [Sze+15] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and 34 A. Rabinovich. “Going Deeper with Convolutions”. In: *CVPR*. 2015.
- 35 [TAH20] D. Teney, E. Abbasnejad, and A. van den Hengel. “Learning What Makes a Difference from Counterfactual Examples and Gradient Supervision”. In: 37 *CoRR* abs/2004.09034 (2020). arXiv: 2004.09034.
- 38 [Tan+16] X. Tan, S. A. Naqvi, K. A. Heller, and V. A. Rao. “Content-based Modeling of Reciprocal Relationships using Hawkes and Gaussian Processes.” In: *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*. 2016.
- 41 [TB16] P. S. Thomas and E. Brunskill. “Data-Efficient Off-Policy Policy Evaluation for Reinforcement Learning”. In: *ICML*. 2016, pp. 2139–2148.
- 44 [TB22] A. Tiulpin and M. B. Blaschko. “Greedy Bayesian Posterior Approximation with Deep Ensembles”. In: *Trans. on Machine Learning Research* 45 (2022).
- 46 [TB99] M. Tipping and C. Bishop. “Probabilistic principal component analysis”. In: *J. of Royal Stat. Soc. Series B* 21.3 (1999), pp. 611–622.
- [TBA19] N. Tremblay, S. Barthélémy, and P.-O. Amblard. “Determinantal Point Processes for Coresets”. In: *J. Mach. Learn. Res.* 20 (2019), pp. 168–1.
- [TBB19] J. Townsend, T. Bird, and D. Barber. “Practical Lossless Compression with Latent Variables using Bits Back Coding”. In: *ICLR*. 2019.
- [TBF06] S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics*. MIT Press, 2006.
- [TBS13] R. Turner, S. Bottone, and C. Stanek. “Online variational approximations to non-exponential family change point models: With application to radar tracking”. In: *NIPS*. 2013.
- [TCG21] Y. Tian, X. Chen, and S. Ganguli. “Understanding self-supervised learning dynamics without contrastive pairs”. In: *International Conference on Machine Learning*. PMLR. 2021, pp. 10268–10278.
- [Teh06] Y. W. Teh. “A hierarchical Bayesian language model based on Pitman-Yor processes”. In: *Proc. of the Assoc. for Computational Linguistics*. 2006, 985–992.
- [Teh+06a] Y.-W. Teh, M. Jordan, M. Beal, and D. Blei. “Hierarchical Dirichlet processes”. In: *JASA* 101.476 (2006), pp. 1566–1581.
- [Teh+06b] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. “Hierarchical Dirichlet Processes”. In: *JASA* 101.476 (2006), pp. 1566–1581.
- [Teh+20] N. Tehrani, N. S. Arora, Y. L. Li, K. D. Shah, D. Noursi, M. Tingley, N. Torabi, S. Masouleh, E. Lippert, and E. Meijer. “Bean Machine: A Declarative Probabilistic Programming Language For Efficient Programmable Inference”. In: *Proceedings of the 10th International Conference on Probabilistic Graphical Models*. Ed. by M. Jaeger and T. D. Nielsen. Vol. 138. Proceedings of Machine Learning Research. PMLR, 2020, pp. 485–496.
- [Ten+20] I. Tenney, J. Wexler, J. Bastings, T. Bolukbasi, A. Coenen, S. Gehrmann, E. Jiang, M. Pushkarna, C. Radefbaugh, E. Reif, et al. “The language interpretability tool: Extensible, interactive visualizations and analysis for NLP models”. In: *arXiv preprint arXiv:2008.05122* (2020).
- [TF03] M. Tipping and A. Faul. “Fast marginal likelihood maximisation for sparse Bayesian models”. In: *AI/Stats*. 2003.
- [TG09] Y. W. Teh and D. Gorur. “Indian buffet processes with power-law behavior”. In: *NIPS*. 2009, pp. 1838–1846.
- [TG18] L. Tu and K. Gimpel. “Learning Approximate Inference Networks for Structured Prediction”. In: *ICLR*. 2018.
- [TGFS18] F. Tronarp, Á. F. García-Fernández, and S. Särkkä. “Iterative Filtering and Smoothing in Nonlinear and Non-Gaussian Systems Using Conditional Moments”. In: *IEEE Signal Process. Lett.* 25.3 (2018), pp. 408–412.
- [TGK03] B. Taskar, C. Guestrin, and D. Koller. “Max-Margin Markov Networks”. In: *NIPS*. 2003.
- [TH09] T. Tielemans and G. Hinton. “Using Fast Weights to Improve Persistent Contrastive Divergence”. In: *ICML*. 2009, pp. 1033–1040.
- [Tho+19] V. Thomas, F. Pedregosa, B. van Merriënboer, P.-A. Mangazol, Y. Bengio, and N. Le Roux.

- 1 "Information matrices and generalization". In: (2019).
 2 arXiv: [1906.07774 \[cs.LG\]](https://arxiv.org/abs/1906.07774).
- 3 [Tho33] W. R. Thompson. "On the Likelihood that
 4 One Unknown Probability Exceeds Another in View of
 5 the Evidence of Two Samples". In: *Biometrika* 25.3/4
 6 (1933), pp. 285–294.
- 7 [Thr+04] S. Thrun, M. Montemerlo, D. Koller, B. Weg-
 8 breit, J. Nieto, and E. Nebot. "FastSLAM: An efficient
 9 solution to the simultaneous localization and mapping
 10 problem with unknown data association". In: *JMLR*
 11 2004 (2004).
- 12 [Thr98] S. Thrun. "Lifelong learning algorithms". In:
 13 *Learning to learn*. Ed. by S. Thrun and L. Pratt.
 14 Kluwer, 1998, pp. 181–209.
- 15 [Thu+21] S. Thulasidasan, S. Thapa, S. Dhaubhadel,
 16 G. Chennupati, T. Bhattacharya, and J. Bilmes. "An
 17 Effective Baseline for Robustness to Distributional
 18 Shift". In: (2021). arXiv: [2105.07107 \[cs.LG\]](https://arxiv.org/abs/2105.07107).
- 19 [Tia+20] Y. Tian, C. Sun, B. Poole, D. Krishnan, C.
 20 Schmid, and P. Isola. "What makes for good views
 21 for contrastive learning". In: *ArXiv* abs/2005.10243
 22 (2020).
- 23 [Tib96] R. Tibshirani. "Regression shrinkage and selec-
 24 tion via the lasso". In: *J. Royal. Statist. Soc B* 58.1
 25 (1996), pp. 267–288.
- 26 [Tie08] T. Tielemans. "Training restricted Boltzmann
 27 machines using approximations to the likelihood gra-
 28 dient". In: *ICML*. ACM New York, NY, USA, 2008,
 29 pp. 1064–1071.
- 30 [Tip01] M. Tipping. "Sparse Bayesian learning and
 31 the relevance vector machine". In: *JMLR* 1 (2001),
 32 pp. 211–244.
- 33 [Tip98] M. Tipping. "Probabilistic visualization of
 34 high-dimensional binary data". In: *NIPS*. 1998.
- 35 [Tit09] M. K. Titsias. "Variational Learning of Induc-
 36 ing Variables in Sparse Gaussian Processes". In: *AIS-
 37 TATS*. 2009.
- 38 [TK86] L. Tierney and J. Kadane. "Accurate approxi-
 39 mations for posterior moments and marginal densities". In: *JASA* 81.393 (1986).
- 40 [TKI20] Y. Tian, D. Krishnan, and P. Isola. "Con-
 41 trastive Multiview Coding". In: *ECCV*. 2020.
- 42 [TKW08] Y. W. Teh, K. Kurihara, and M. Welling.
 43 "Collapsed variational inference for HDP". In: *Ad-
 44 vances in neural information processing systems*.
 45 2008, pp. 1481–1488.
- 46 [TL05] E. Todorov and W. Li. "A Generalized Iterative
 47 LQG Method for Locally-optimal Feedback Control of
 48 Constrained Nonlinear Stochastic Systems". In: *ACC*.
 49 2005, pp. 300–306.
- 50 [TL18a] S. J. Taylor and B. Letham. "Forecasting at
 51 scale". en. In: *The American Statistician* 72.1 (2018),
 52 pp. 37–45.
- 53 [TL18b] G. Tucker and D. Lawson. "Doubly Reparam-
 54 eterized Gradient Estimators for Monte Carlo Objec-
 55 tives". In: *1st Symposium on Advances in Approximate
 56 Bayesian Inference*. 2018.
- 57 [TL19] M. Tan and Q. Le. "Efficientnet: Rethink-
 58 ing model scaling for convolutional neural networks".
 59 In: *International Conference on Machine Learning*.
 60 PMLR, 2019, pp. 6105–6114.
- 61 [TLG14] M. Titsias and M. Lázaro-Gredilla. "Doubly
 62 Stochastic Variational Bayes for non-Conjugate Infer-
 63 ence". In: *ICML*. 2014, pp. 1971–1979.
- 64 [TMD12] A. Talhouk, K. Murphy, and A. Doucet. "Effi-
 65 cient Bayesian Inference for Multivariate Probit Mod-
 66 els with Sparse Inverse Correlation Matrices". In: *J.
 67 Comp. Graph. Statist.* 21.3 (2012), pp. 739–757.
- 68 [TMK04] G. Theocharous, K. Murphy, and L. Kael-
 69 bling. "Representing hierarchical POMDPs as DBNs
 70 for multi-scale robot localization". In: *ICRA*. 2004.
- 71 [TN13] L. S. L. Tan and D. J. Nott. "Variational In-
 72 ference for Generalized Linear Mixed Models Using
 73 Partially Noncentered Parametrizations". In: *Stat. Sci.*
 74 (2013).
- 75 [TN18] L. S. L. Tan and D. J. Nott. "Gaussian varia-
 76 tional approximation with sparse precision matrices".
 77 In: *Stat. Comput.* 28.2 (2018), pp. 259–275.
- 78 [TND21] M.-N. Tran, T.-N. Nguyen, and V.-H. Dao.
 79 "A practical tutorial on Variational Bayes". In: (2021).
 80 arXiv: [2103.01327 \[stat.CO\]](https://arxiv.org/abs/2103.01327).
- 81 [TOB16] L. Theis, A. van den Oord, and M. Bethge.
 82 "A note on the evaluation of generative models". In:
 83 *ICLR*. 2016.
- 84 [Tol22] S. Toledo. "UltimateKalman: Flexible Kalman
 85 Filtering and Smoothing Using Orthogonal Trans-
 86 formations". In: (July 2022). arXiv: [2207 . 13526 \[math.NA\]](https://arxiv.org/abs/2207.13526).
- 87 [Tom+20] R. Tomsett, D. Harborne, S. Chakraborty,
 88 P. Gurram, and A. Preece. "Sanity checks for saliency
 89 metrics". In: *Proceedings of the AAAI conference on
 90 artificial intelligence*. Vol. 34. 04. 2020, pp. 6021–
 91 6029.
- 92 [Tom22] J. M. Tomczak. *Deep Generative Modeling*.
 93 en. 1st ed. Springer, 2022.
- 94 [Tou09] M. Toussaint. "Robot Rrajectory Optimiza-
 95 tion using Approximate Inference". In: *ICML*. 2009,
 96 pp. 1049–1056.
- 97 [Tou14] M. Toussaint. *Bandits, Global Optimization,
 98 Active Learning, and Bayesian RL – understanding
 99 the common ground*. Autonomous Learning Summer
 100 School, 2014.
- 101 [Tou+19] J Toubeau, J Bottreau, F Vallée, and Z De
 102 Grève. "Deep Learning-Based Multivariate Probabilis-
 103 tic Forecasting for Short-Term Scheduling in Power
 104 Markets". In: *IEEE Trans. Power Syst.* 34.2 (2019),
 105 pp. 1203–1215.
- 106 [TOV18] C. Truong, L. Oudre, and N. Vayatis. "Selec-
 107 tive review of offline change point detection methods".
 108 In: (2018). arXiv: [1801.00718 \[cs.CE\]](https://arxiv.org/abs/1801.00718).
- 109 [TP97] S. Thrun and L. Pratt, eds. *Learning to learn*.
 110 Kluwer, 1997.
- 111 [TPB00] N. Tishby, F. C. Pereira, and W. Bialek.
 112 "The information bottleneck method". In: *ArXiv
 113 physics/0004057* (2000).
- 114 [TPB99] N. Tishby, F. Pereira, and W. Bialek. "The
 115 Information Bottleneck method". In: *The 37th an-
 116 nual Allerton Conf. on Communication, Control,
 117 and Computing*. 1999, pp. 368–377.
- 118 [TR19] M. K. Titsias and F. Ruiz. "Unbiased Im-
 119 plicit Variational Inference". In: *AISTATS*. Ed. by K.
 120 Chaudhuri and M. Sugiyama. Vol. 89. Proceedings of
 121 Machine Learning Research. PMLR, 2019, pp. 167–
 122 176.
- 123 [TR97] J. Tsitsiklis and B. V. Roy. "An analysis of
 124 temporal-difference learning with function approxima-
 125 tion". In: *IEEE Trans. on Automatic Control* 42.5
 126 (1997), pp. 674–690.

- 1
- 2 [Tra+19] D. Tran, K. Vafa, K. K. Agrawal, L. Dinh,
and B. Poole. “Discrete Flows: Invertible Generative
Models of Discrete Data”. In: *Advances in Neural Information Processing Systems*. 2019.
- 3
- 4 [Tra+20a] L. Tran, B. S. Veeling, K. Roth, J.
Swiatkowski, J. V. Dillon, J. Snoek, S. Mandt, T.
Salimans, S. Nowozin, and R. Jenatton. “Hydra: Pre-
serving Ensemble Diversity for Model Distillation”. In:
(2020). arXiv: 2001.04694 [cs, LG].
- 5
- 6 [Tra+20b] M.-N. Tran, N Nguyen, D Nott, and R Kohn.
“Bayesian Deep Net GLM and GLMM”. In: *J. Comput.
Graph. Stat.* 29.1 (2020), pp. 97–113.
- 7
- 8 [TRB16] D. Tran, R. Ranganath, and D. M. Blei. “The
Variational Gaussian Process”. In: *ICLR*. 2016.
- 9
- 10 [Tri21] K. Triantafyllopoulos. *Bayesian Inference of
State Space Models: Kalman Filtering and Beyond*.
en. 1st ed. Springer, 2021.
- 11
- 12 [TS06] M. Toussaint and A. Storkey. “Probabilistic
inference for solving discrete and continuous state
Markov Decision Processes”. In: *ICML*. 2006, pp. 945–
952.
- 13
- 14 [Tsa+18] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn,
M.-H. Yang, and M. Chandraker. “Learning to adapt
structured output space for semantic segmentation”.
In: *Proceedings of the IEEE conference on computer
vision and pattern recognition*. 2018, pp. 7472–7481.
- 15
- 16 [Tsa+19] Y.-H. H. Tsai, S. Bai, M. Yamada, L.-P.
Morency, and R. Salakhutdinov. “Transformer Dissec-
tion: An Unified Understanding for Transformer’s At-
tention via the Lens of Kernel”. In: *EMNLP*. Associa-
tion for Computational Linguistics, 2019, pp. 4344–
4353.
- 17
- 18 [Tsa+21] Y.-H. H. Tsai, S. Bai, L.-P. Morency, and R.
Salakhutdinov. “A Note on Connecting Barlow Twins
with Negative-Sample-Free Contrastive Learning”. In:
ArXiv abs/2104.13712 (2021).
- 19
- 20 [Tsa88] C. Tsallis. “Possible generalization of
Boltzmann-Gibbs statistics”. In: *J. of Statistical
Physics* 52 (1988), pp. 479–487.
- 21
- 22 [Tsc+14] S. Tschwartschek, R. Iyer, H. Wei, and J.
Bilmes. “Learning Mixtures of Submodular Functions
for Image Collection Summarization”. In: *NIPS*. Mon-
tréal, Canada, 2014.
- 23
- 24 [Tsc+19] M. Tschannen, J. Djolonga, P. K. Ruben-
stein, S. Gelly, and M. Lucic. “On Mutual Information
Maximization for Representation Learning”. In: *arXiv
preprint arXiv:1907.13625* (2019).
- 25
- 26 [Tsi+17] P. A. Tsividis, T. Pouncy, J. L. Xu, J. B.
Tenenbaum, and S. J. Gershman. “Human Learning in
Atari”. en. In: *AAAI Spring Symposium Series*. 2017.
- 27
- 28 [Tso+05] I. Tsochantaridis, T. Joachims, T. Hofmann,
and Y. Altun. “Large Margin Methods for Structured
and Interdependent Output Variables”. In: *JMLR* 6
(2005), pp. 1453–1484.
- 29
- 30 [TT13] E. G. Tabak and C. V. Turner. “A family of non-
parametric density estimation algorithms”. In: *Com-
munications on Pure and Applied Mathematics* 66.2
(2013), pp. 145–164.
- 31
- 32 [TT17] B. Trippe and R. Turner. “Overpruning in Vari-
ational Bayesian Neural Networks”. In: *NIPS Work-
shop on Advances in Approximate Bayesian Infer-
ence*. 2017.
- 33
- 34 [Tuc+19] G. Tucker, D. Lawson, B. Dai, and R. Ran-
ganath. “Revisiting auxiliary latent variables in gener-
ative models”. In: (2019).
- 35
- 36 [Uch17] T. Taketomi, H. Uchiyama, and S. Ikeda. “Vi-
sual SLAM algorithms: a survey from 2010 to 2016”.
en. In: *IPSJ Transactions on Computer Vision and
Applications* 9.1 (2017), p. 16.
- 37
- 38 [Tul+18] S. Tulyakov, M.-Y. Liu, X. Yang, and J.
Kautz. “Mocogan: Decomposing motion and content
for video generation”. In: *Proceedings of the IEEE
conference on computer vision and pattern recogni-
tion*. 2018, pp. 1526–1535.
- 39
- 40 [Tur+08] R. Turner, P. Berkes, M. Sahani, and D.
Mackay. *Counterexamples to variational free energy
compactness folk theorems*. Tech. rep. U. Cambridge,
2008.
- 41
- 42 [TVE10] E. G. Tabak and E. Vanden-Eijnden. “Den-
sity estimation by dual ascent of the log-likelihood”.
In: *Communications in Mathematical Sciences* 8.1
(2010), pp. 217–233.
- 43
- 44 [TW16] J. M. Tomczak and M. Welling. “Improv-
ing variational auto-encoders using Householder flow”.
In: *NeurIPS Workshop on Bayesian Deep Learning*
(2016).
- 45
- 46 [TX00] J. B. Tenenbaum and F. Xu. “Word learning as
Bayesian inference”. In: *Proc. 22nd Annual Conf. of
the Cognitive Science Society*. 2000.
- 47
- 48 [TZ02] Z. Tu and S. Zhu. “Image Segmentation by
Data-Driven Markov Chain Monte Carlo”. In: *IEEE
PAMI* 24.5 (2002), pp. 657–673.
- 49
- 50 [Tze+17] E. Tzeng, J. Hoffman, K. Saenko, and T. Dar-
rell. “Adversarial discriminative domain adaptation”.
In: *Proceedings of the IEEE conference on computer
vision and pattern recognition*. 2017, pp. 7167–7176.
- 51
- 52 [UCS17] S. Ubaru, J. Chen, and Y. Saad. “Fast Esti-
mation of $\text{tr}(f(A))$ via Stochastic Lanczos Quadra-
ture”. In: *SIAM J. Matrix Anal. Appl.* 38.4 (2017),
pp. 1075–1099.
- 53
- 54 [Ude+16] M. Udell, C. Horn, R. Zadeh, and S. Boyd.
“Generalized Low Rank Models”. In: *Foundations and
Trends in Machine Learning* 9.1 (2016), pp. 1–118.
- 55
- 56 [UHJ20] M. Uehara, J. Huang, and N. Jiang. “Minimax
Weight and Q-Function Learning for Off-Policy Eval-
uation”. In: *ICML*. 2020.
- 57
- 58 [UML13] B. Uria, I. Murray, and H. Larochelle.
“RNADe: The real-valued neural autoregressive
density-estimator”. In: *NIPS*. 2013.
- 59
- 60 [UML14] B. Uria, I. Murray, and H. Larochelle. “A
Deep and Tractable Density Estimator”. In: *ICML*.
2014.
- 61
- 62 [UN98] N. Ueda and R. Nakano. “Deterministic anneal-
ing EM algorithm”. In: *Neural Networks* 11 (1998),
pp. 271–282.
- 63
- 64 [UR16] B. Ustun and C. Rudin. “Supersparse linear in-
teger models for optimized medical scoring systems”.
In: *Machine Learning* 102.3 (2016), pp. 349–391.
- 65
- 66 [Uri+16] B. Uria, M.-A. Côté, K. Gregor, I. Murray,
and H. Larochelle. “Neural Autoregressive Distribu-
tion Estimation”. In: *JMLR* (2016).
- 67
- 68 [UTR14] B. Ustun, S. Tracà, and C. Rudin. *Super-
sparse Linear Integer Models for Interpretable Clas-
sification*. 2014. arXiv: 1306.6677 [stat, ML].
- 69
- 70 [Uur+17] V. Uurtio, J. M. Monteiro, J. Kandola, J.
Shawe-Taylor, D. Fernandez-Reyes, and J. Rousu. “A
Tutorial on Canonical Correlation Methods”. In: *ACM
Computing Surveys* (2017).

- [UVL16] D. Ulyanov, A. Vedaldi, and V. Lempitsky. “Instance Normalization: The Missing Ingredient for Fast Stylization”. In: (2016). arXiv: 1607.08022 [cs.CV].
- [UVL18] D. Ulyanov, A. Vedaldi, and V. Lempitsky. “Deep Image Prior”. In: CVPR. 2018.
- [Vaa00] A. W. Van der Vaart. *Asymptotic statistics*. Vol. 3. Cambridge university press, 2000.
- [Val00] H. Valpola. “Bayesian Ensemble Learning for Nonlinear Factor Analysis”. PhD thesis. Helsinki University of Technology, 2000.
- [Van10] J. Vanhatalo. “Speeding up the inference in Gaussian process models”. PhD thesis. Helsinki Univ. Technology, 2010.
- [van+18] H. van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. *Deep Reinforcement Learning and the Deadly Triad*. arXiv:1812.02648. 2018.
- [Vas+17a] A. B. Vasudevan, M. Gygli, A. Volokitin, and L. Van Gool. “Query-adaptive video summarization via quality-aware relevance estimation”. In: *Proceedings of the 25th ACM international conference on Multimedia*. 2017, pp. 582–590.
- [Vas+17b] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. “Attention is all you need”. In: NIPS. 2017, pp. 5998–6008.
- [Vas+17c] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. “Attention Is All You Need”. In: NIPS. 2017.
- [Vaz+22] S. Vaze, K. Han, A. Vedaldi, and A. Zisserman. “Open-Set Recognition: A Good Closed-Set Classifier is All You Need”. In: ICLR. 2022.
- [VBW15] S. S. Villar, J. Bowden, and J. Wason. “Multi-armed Bandit Models for the Optimal Design of Clinical Trials: Benefits and Challenges”. en. In: *Stat. Sci.* 30.2 (2015), pp. 199–215.
- [VDMW03] R. Van Der Merwe and E. Wan. “Sigma-Point Kalman Filters for probabilistic inference in dynamic state-space models”. In: *Workshop on Advances in ML*. 2003.
- [Ved+18] R. Vedantam, I. Fischer, J. Huang, and K. Murphy. “Generative Models of Visually Grounded Imagination”. In: ICLR. 2018.
- [Veh+19] A. Vehtari, A. Gelman, D. Simpson, B. Carpenter, and P.-C. Bürkner. “Rank-normalization, folding, and localization: An improved \hat{R} for assessing convergence of MCMC”. In: (2019). arXiv: 1903.08008 [stat.CO].
- [Vei+21] V. Veitch, A. D’Amour, S. Yadlowsky, and J. Eisenstein. “Counterfactual Invariance to Spurious Correlations: Why and How to Pass Stress Tests”. In: *Advances in Neural Information Processing Systems*. 2021.
- [Vel+17] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio. “Graph attention networks”. In: *arXiv preprint arXiv:1710.10903* (2017).
- [Ver18] R. Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. en. 1 edition. Cambridge University Press, 2018.
- [Ver+19] A. Vergari, A. Molina, R. Peherz, Z. Ghahramani, K. Kersting, and I. Valera. “Automatic Bayesian Density Analysis”. In: AAAI. 2019.
- [VF+80] B. C. Van Fraassen et al. *The scientific image*. Oxford University Press, 1980.
- [VGG17] A. Vehtari, A. Gelman, and J. Gabry. “Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC”. In: *Stat. Comput.* 27.5 (2017), pp. 1413–1432.
- [VGS05] V. Vovk, A. Gammerman, and G. Shafer. *Algorithmic Learning in a Random World*. en. 2005th ed. Springer, 2005.
- [Vid99] P. Vidoni. “Exponential Family State Space Models Based on a Conjugate Latent Process”. In: *J. R. Stat. Soc. Series B Stat. Methodol.* 61.1 (1999), pp. 213–221.
- [Vil08] C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008.
- [Vil+19] R. Villegas, A. Pathak, H. Kannan, D. Erhan, Q. V. Le, and H. Lee. “High Fidelity Video Prediction with Large Stochastic Recurrent Neural Networks”. In: NIPS. 2019.
- [Vin+08] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. “Extracting and composing robust features with denoising autoencoders”. In: *Proceedings of the 25th international conference on Machine learning*. 2008, pp. 1096–1103.
- [Vin+10a] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. “Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion”. In: *Journal of machine learning research* 11.Dec (2010), pp. 3371–3408.
- [Vin+10b] M. Vinyals, J. Cerquides, J. Rodriguez-Aguilar, and A. Farinelli. “Worst-case bounds on the quality of max-product fixed-points”. In: NIPS. 2010.
- [Vin11] P. Vincent. “A connection between score matching and denoising autoencoders”. In: *Neural computation* 23.7 (2011), pp. 1661–1674.
- [Vir10] S. Virtanen. “Bayesian exponential family projections”. MA thesis. Aalto University, 2010.
- [Vis+06] S. V. N. Vishwanathan, N. N. Schraudolph, M. W. Schmidt, and K. P. Murphy. “Accelerated training of conditional random fields with stochastic gradient methods”. In: ICML. ACM Press, 2006.
- [Vis+10] S. V. N. Vishwanathan, N. N. Schraudolph, R. Kondor, and K. M. Borgward. “Graph Kernels”. In: JMLR 11 (2010), pp. 1201–1242.
- [Vit67] A. Viterbi. “Error bounds for convolutional codes and an asymptotically optimal decoding algorithm”. In: *IEEE Trans. on Information Theory* 13.2 (1967), pp. 260–269.
- [VJMP22] V. Voleti, A. Jolicoeur-Martineau, and C. Pal. “MCVD: Masked Conditional Video Diffusion for Prediction, Generation, and Interpolation”. In: (May 2022). arXiv: 2205.09853 [cs.CV].
- [VJV09] J. Vanhatalo, P. Jylänki, and A. Vehtari. “Gaussian process regression with Student-t likelihood”. In: NIPS. Vol. 22. 2009.
- [VK20a] A. Vahdat and J. Kautz. “NVAE: A Deep Hierarchical Variational Autoencoder”. In: NIPS. 2020.
- [VK20b] A. Vahdat and J. Kautz. “Nvae: A deep hierarchical variational autoencoder”. In: *arXiv preprint arXiv:2007.03898* (2020).
- [VLT21] G. M. van de Ven, Z. Li, and A. S. Tolias. “Class-Incremental Learning with Generative Classi-

- 1 fiers". In: *CVPR workshop on Continual Learning in*
2 *Computer Vision (CLVision)*. 2021.
- 3 [Vo+15] B.-N. Vo, M. Mallick, Y. Bar-Shalom, S.
4 Coraluppi, R. Osborne, R. Mahler, B. t Vo, and J.
5 Webster. *Multitarget tracking*. John Wiley and Sons,
6 2015.
- 7 [Von13] P. Vontobel. "The Bethe permanent of a non-
8 negative matrix". In: *IEEE Trans. Info. Theory* 59.3
9 (2013).
- 10 [Vov13] V. Vovk. "Kernel Ridge Regression". In: *Em-
11 pirical Inference: Festschrift in Honor of Vladimir
12 N. Vapnik*. Ed. by B. Schölkopf, Z. Luo, and V. Vovk.
13 Springer Berlin Heidelberg, 2013, pp. 105–116.
- 14 [VPV10] J. Vanhatalo, V. Pietiläinen, and A. Vehtari.
15 "Approximate inference for disease mapping with
16 sparse Gaussian processes". In: *Statistics in Medicine*
17 29.15 (2010), pp. 1580–1607.
- 18 [vR11] M. van der Laan and S. Rose. *Targeted Learn-
19 ing: Causal Inference for Observational and Experi-
20 mental Data*. Jan. 2011.
- 21 [VR18] S. Verma and J. Rubin. "Fairness Defini-
22 tions Explained". In: *2018 IEEE/ACM International
23 Workshop on Software Fairness (FairWare)*. May
24 2018, pp. 1–7.
- 25 [VRF11] C. Varin, N. Reid, and D. Firth. "An overview
26 of composite likelihood methods". In: *Stat. Sin.* 21.1
27 (2011), pp. 5–42.
- 28 [VT11] S. I. VanderWeele TJ. "A new criterion for con-
29 founder selection". In: *Biometrics* (2011).
- 30 [VT18] G. M. van de Ven and A. S. Tolias. "Three sce-
31 narios for continual learning". In: *NeurIPS Continual
32 Learning workshop*. 2018.
- 33 [Vyt+19] D. Vytiniotis, D. Belov, R. Wei, G. Plotkin,
34 and M. Abadi. "The differentiable curry". In: *NeurIPS
35 2019 Workshop Program Transformations*. 2019.
- 36 [VZ20] V. Veitch and A. Zaveri. "Sense and Sensitivity
37 Analysis: Simple Post-Hoc Analysis of Bias Due to Un-
38 Observed Confounding". In: *Advances in Neural Infor-
39 mation Processing Systems*. Ed. by H. Larochelle, M.
40 Ranzato, R. Hadsell, M. F. Balcan, and H. Lin. Vol. 33.
41 Curran Associates, Inc., 2020, pp. 10999–11009.
- 42 [WA13] A. Wilson and R. Adams. "Gaussian process
43 kernels for pattern discovery and extrapolation". In:
44 *International conference on machine learning*. 2013,
45 pp. 1067–1075.
- 46 [Wal+09] H. Wallach, I. Murray, R. Salakhutdinov, and
47 D. Mimno. "Evaluation Methods for Topic Models". In:
48 *ICML*. 2009.
- 49 [Wal+20] M. Walmsley et al. "Galaxy Zoo: probabilistic
50 metric morphology through Bayesian CNNs and active
51 learning". In: *Monthly Notices Royal Astronomical So-
52 ciety* 491.2 (2020), pp. 1554–1574.
- 53 [Wan+16] Z. Wang, T. Schaul, M. Hessel, H. van Has-
54 selst, M. Lanctot, and N. de Freitas. "Dueling Network
55 Architectures for Deep Reinforcement Learning". In:
56 *ICML*. 2016.
- 57 [Wan17] M. P. Wand. "Fast Approximate Inference
58 for Arbitrarily Large Semiparametric Regression Mod-
59 els via Message Passing". In: *JASA* 112.517 (2017),
60 pp. 137–168.
- 61 [Wan+17a] T. Wang, C. Rudin, F. Doshi-Velez, Y. Liu,
62 E. Klampfl, and P. MacNeille. "A bayesian framework
63 for learning rule sets for interpretable classification".
64 In: *The Journal of Machine Learning Research* 18.1
65 (2017), pp. 2357–2393.
- 66 [Wan+17b] X. Wang, T. Li, S. Sun, and J. M. Cor-
67 chado. "A Survey of Recent Advances in Particle Fil-
68ters and Remaining Challenges for Multitarget Track-
69ing". en. In: *Sensors* 17.12 (2017).
- 70 [Wan+17c] Y. Wang et al. "Tacotron: Towards End-to-
71End Speech Synthesis". In: *Interspeech*. 2017.
- 72 [Wan+18] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, G. Liu,
73 A. Tao, J. Kautz, and B. Catanzaro. "Video-to-video
74 synthesis". In: *Proceedings of the 32nd International
75 Conference on Neural Information Processing Sys-
76tems*. 2018, pp. 1152–1164.
- 77 [Wan+19a] K. Wang, G. Pleiss, J. Gardner, S. Tyree,
78 K. Q. Weinberger, and A. G. Wilson. "Exact Gaussian
79 Processes on a Million Data Points". In: *NIPS*. 2019,
80 pp. 14622–14632.
- 81 [Wan+19b] S. Wang, W. Bai, C. Lavania, and J.
82 Bilmes. "Fixing Mini-batch Sequences with Hierar-
83chical Robust Partitioning". In: *Proceedings of the
84 Twenty-Second International Conference on Artifi-
85cial Intelligence and Statistics*. Ed. by K. Chaudhuri
86 and M. Sugiyama. Vol. 89. Proceedings of Machine
87 Learning Research. PMLR, 2019, pp. 3352–3361.
- 88 [Wan+19c] Y. Wang, A. Smola, D. C. Maddix, J.
89 Gasthaus, D. Foster, and T. Januschowski. "Deep Fac-
90tors for Forecasting". In: *ICML*. 2019.
- 91 [Wan+20a] D. Wang, E. Shelhamer, S. Liu, B. Ol-
92shaugen, and T. Darrell. "Tent: Fully Test-Time Adap-
93tation by Entropy Minimization". In: *ICLR*. 2020.
- 94 [Wan+20b] H. Wang, X. Wu, Z. Huang, and E. P. Xing.
95 "High-frequency component helps explain the general-
96ization of convolutional neural networks". In: *Pro-
97ceedings of the IEEE/CVF Conference on Computer Vi-
98sion and Pattern Recognition*. 2020, pp. 8684–8694.
- 99 [Wan+20c] T. Wang, J.-Y. Zhu, A. Torralba, and A. A.
100 Efros. *Dataset Distillation*. 2020. arXiv: [1811.10959](https://arxiv.org/abs/1811.10959) [cs.LG].
- 101 [Wan+20d] Z. Wang, S. Cheng, L. Yueru, J. Zhu,
102 and B. Zhang. "A Wasserstein Minimum Velocity Ap-
103proach to Learning Unnormalized Models". In: *Pro-
104ceedings of the Twenty Third International Confer-
105ence on Artificial Intelligence and Statistics*. Ed. by
106 S. Chiappa and R. Calandra. Vol. 108. Proceedings of
107 Machine Learning Research. PMLR, 2020, pp. 3728–
108 3738.
- 109 [Wan+21] J. Wang, C. Lan, C. Liu, Y. Ouyang, W.
110 Zeng, and T. Qin. "Generalizing to Unseen Domains:
111 A Survey on Domain Generalization". In: *IJCAI*. 2021.
- 112 [Wan+22] H. Wang, Y. Yang, D. Pati, and A.
113 Bhattacharya. "Structured Variational Inference in
114 Bayesian State-Space Models". In: (2022).
- 115 [Was06] L. Wasserman. *All of Nonparametric Statistics*. Springer, 2006.
- 116 [Wat10] S. Watanabe. "Asymptotic Equivalence of
117 Bayes Cross Validation and Widely Applicable Infor-
118mation Criterion in Singular Learning Theory". In:
119 *JMLR* 11 (2010), pp. 3571–3594.
- 120 [Wat13] S. Watanabe. "A Widely Applicable Bayesian
121 Information Criterion". In: *JMLR* 14 (2013), pp. 867–
122 897.
- 123 [Wat60] S. Watanabe. "Information theoretical analy-
124sis of multivariate correlation". In: *IBM J. of Research
125 and Development* 4 (1960), pp. 66–82.
- 126 [WB05] J. Winn and C. Bishop. "Variational Message
127 Passing". In: *JMLR* 6 (2005), pp. 661–694.

- [WB12] C. Wang and D. M. Blei. “Truncation-free online variational inference for Bayesian nonparametric models”. In: *Advances in neural information processing systems*. 2012, pp. 413–421.
- [WB20] T. Wang and J. Ba. “Exploring Model-based Planning with Policy Networks”. In: *ICLR*. 2020.
- [WBC21] Y. Wang, D. M. Blei, and J. P. Cunningham. “Posterior Collapse and Latent Variable Non-identifiability”. In: *NIPS*. 2021.
- [WCS08] M. Welling, C. Chemudugunta, and N. Sutera. “Deterministic Latent Variable Models and their Pitfalls”. In: *ICDM*. 2008.
- [WD92] C. Watkins and P. Dayan. “Q-learning”. In: *Machine Learning* 8.3 (1992), pp. 279–292.
- [WDN15] A. G. Wilson, C. Dann, and H. Nickisch. “Thoughts on Massively Scalable Gaussian Processes”. In: *arXiv preprint arXiv:1511.01870* (2015). <https://arxiv.org/abs/1511.01870>.
- [Web+17] T. Weber et al. “Imagination-Augmented Agents for Deep Reinforcement Learning”. In: *NIPS*. 2017.
- [Wei00] Y. Weiss. “Correctness of local probability propagation in graphical models with loops”. In: *Neural Computation* 12 (2000), pp. 1–41.
- [Wei+13] K. Wei, Y. Liu, K. Kirchhoff, and J. Bilmes. “Using Document Summarization Techniques for Speech Data Subset Selection.” In: *HLT-NAACL*. 2013, pp. 721–726.
- [Wei+14] K. Wei, Y. Liu, K. Kirchhoff, and J. Bilmes. “Unsupervised Submodular Subset Selection for Speech Data”. In: *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*. Florence, Italy, 2014.
- [Wei+15a] K. Wei, R. Iyer, S. Wang, W. Bai, and J. Bilmes. “How to Intelligently Distribute Training Data to Multiple Compute Nodes: Distributed Machine Learning via Submodular Partitioning”. In: *Neural Information Processing Society (NeurIPS, formerly NIPS) Workshop. LearningSys Workshop*, <http://learningsys.org>. Montreal, Canada, 2015.
- [Wei+15b] K. Wei, R. Iyer, S. Wang, W. Bai, and J. Bilmes. “Mixed Robust/Average Submodular Partitioning: Fast Algorithms, Guarantees, and Applications”. In: *Neural Information Processing Society (NeurIPS, formerly NIPS)*. Montreal, Canada, 2015.
- [Wel11] M. Welling. “Bayesian Learning via Stochastic Gradient Langevin Dynamics”. In: *ICML*. 2011.
- [Wen+17] R. Wen, K. Torkkola, B. Narayanaswamy, and D. Madeka. “A Multi-Horizon Quantile Recurrent Forecaster”. In: *NIPS Time Series Workshop*. 2017.
- [Wen+19a] L. Wenliang, D. Sutherland, H. Strathmann, and A. Gretton. “Learning deep kernels for exponential family densities”. In: *International Conference on Machine Learning*. 2019, pp. 6737–6746.
- [Wen+19b] F. Wenzel, T. Galy-Fajou, C. Donner, M. Kloft, and M. Opper. “Efficient Gaussian Process Classification Using Polya-Gamma Data Augmentation”. In: *AAAI*. 2019.
- [Wen+20a] C. Wendler, A. Amrollahi, B. Seifert, A. Krause, and M. Püschel. “Learning set functions that are sparse in non-orthogonal Fourier bases”. In: *arXiv preprint arXiv:2010.00439* (2020).
- [Wen+20b] F. Wenzel, K. Roth, B. S. Veeling, J. Świątkowski, L. Tran, S. Mandt, J. Snoek, T. Salimans, R. Jenatton, and S. Nowozin. “How Good is the Bayes Posterior in Deep Neural Networks Really?” In: *ICML*. 2020.
- [Wen+20c] F. Wenzel, J. Snoek, D. Tran, and R. Jenatton. “Hyperparameter Ensembles for Robustness and Uncertainty Quantification”. In: *NIPS*. 2020.
- [Wen21] L. Weng. “What are diffusion models?” In: *lilianweng.github.io/lil-log* (2021).
- [Wes03] M. West. “Bayesian Factor Regression Models in the “Large p, Small n” Paradigm”. In: *Bayesian Statistics 7* (2003).
- [Wes87] M. West. “On scale mixtures of normal distributions”. In: *Biometrika* 74 (1987), pp. 646–648.
- [WF01a] Y. Weiss and W. T. Freeman. “Correctness of belief propagation in Gaussian graphical models of arbitrary topology”. In: *Neural Computation* 13.10 (2001), pp. 2173–2200.
- [WF01b] Y. Weiss and W. T. Freeman. “On the optimality of solutions of the max-product belief propagation algorithm in arbitrary graphs”. In: *IEEE Trans. Information Theory, Special Issue on Codes on Graphs and Iterative Algorithms* 47.2 (2001), pp. 723–735.
- [WF14] Z. Wang and N. de Freitas. “Theoretical Analysis of Bayesian Optimisation with Unknown Gaussian Process Hyper-Parameters”. In: (2014). *arXiv: 1406.7758 [stat.ML]*.
- [WF16] Z. Wang and N. de Freitas. “Theoretical Analysis of Bayesian Optimisation with Unknown Gaussian Process Hyper-Parameters”. In: *BayesOpt Workshop*. 2016.
- [WG18] M. Wu and N. Goodman. “Multimodal Generative Models for Scalable Weakly-Supervised Learning”. In: *NIPS*. 2018.
- [WGY21] G. Weiss, Y. Goldberg, and E. Yahav. “Thinking Like Transformers”. In: *ICML*. 2021.
- [WH02] M. Welling and G. E. Hinton. “A new learning algorithm for mean field Boltzmann machines”. In: *International Conference on Artificial Neural Networks*. Springer, 2002, pp. 351–357.
- [WH18] Y. Wu and K. He. “Group Normalization”. In: *ECCV*. 2018.
- [WH97] M. West and J. Harrison. *Bayesian forecasting and dynamic models*. Springer, 1997.
- [WHD18] J. T. Wilson, F. Hutter, and M. P. Deisenroth. “Maximizing acquisition functions for Bayesian optimization”. In: *NIPS*. 2018.
- [Whi16] T. White. “Sampling Generative Networks”. In: *arXiv* (2016).
- [Whi88] P. Whittle. “Restless bandits: activity allocation in a changing world”. In: *J. Appl. Probab.* 25.A (1988), pp. 287–298.
- [WHT19] Y. Wang, H. He, and X. Tan. “Truly Proximal Policy Optimization”. In: *UAI*. 2019.
- [WI20] A. G. Wilson and P. Izmailov. “Bayesian Deep Learning and a Probabilistic Perspective of Generalization”. In: *NIPS*. 2020.
- [WIB15] K. Wei, R. Iyer, and J. Bilmes. “Submodularity in Data Subset Selection and Active Learning”. In: *Proceedings of the 32nd international conference on Machine learning*. Lille, France, 2015.

- 1 [Wie+14] D Wierstra, T Schaul, J Peters, and J Schmidhuber. “Natural Evolution Strategies”. In: *JMLR* 15.1 (2014), pp. 949–980.
- 2 [Wik21] Wikipedia contributors. *CliffsNotes — Wikipedia, The Free Encyclopedia*. [Online; accessed 29-December-2021]. 2021.
- 3 [Wil+14] A. G. Wilson, E. Gilboa, A. Nehorai, and J. P. Cunningham. “Fast kernel learning for multidimensional pattern extrapolation”. In: *Advances in Neural Information Processing Systems*. 2014, pp. 3626–3634.
- 4 [Wil14] A. G. Wilson. “Covariance kernels for fast automatic pattern discovery and extrapolation with Gaussian processes”. PhD thesis. University of Cambridge, 2014.
- 5 [Wil+16] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. “Deep Kernel Learning”. en. In: *AISTATS*. 2016, pp. 370–378.
- 6 [Wil+17] A. G. Wills, J. Hendriks, C. Renton, and B. Ninness. “A Bayesian Filtering Algorithm for Gaussian Mixture Models”. In: (2017). arXiv: [1705.05495 \[stat.ML\]](#).
- 7 [Wil20] A. G. Wilson. “The Case for Bayesian Deep Learning”. In: (2020). arXiv: [2001.10995 \[cs.LG\]](#).
- 8 [Wil+20a] J. T. Wilson, V. Borovitskiy, A. Terenin, P. Mostovsky, and M. P. Deisenroth. “Efficiently Sampling Functions from Gaussian Process Posteriors”. In: *ICML*. 2020.
- 9 [Wil+20b] A. B. Wiltschko, B. Sanchez-Lengeling, B. Lee, E. Reif, J. Wei, K. J. McCloskey, L. Colwell, W. Qian, and Y. Wang. “Evaluating Attribution for Graph Neural Networks”. In: (2020).
- 10 [Wil69] A. G. Wilson. “The use of entropy maximising models, in the theory of trip distribution, mode split and route split”. In: *Journal of transport economics and policy* (1969), pp. 108–126.
- 11 [Wil92] R. J. Williams. “Simple statistical gradient-following algorithms for connectionist reinforcement learning”. In: *MLJ* 8.3-4 (1992), pp. 229–256.
- 12 [Wil98] C. Williams. “Computation with infinite networks”. In: *Neural Computation* 10.5 (1998), pp. 1203–1216.
- 13 [Win] J. Winn. *VIBES*.
- 14 [Win+19] J. Winn, C. Bishop, T. Diethe, J. Guiver, and Y. Zaykov. *Model-based Machine Learning*. 2019.
- 15 [WIP20] J. Watson, A. Imohosen, and J. Peters. “Active Inference or Control as Inference? A Unifying View”. In: *International Workshop on Active Inference*. 2020.
- 16 [Wit] DEEPFAKES: PREPARE NOW (PERSPECTIVES FROM BRAZIL). <https://lab.witness.org/brazil-deepfakes-prepare-now/>. Accessed: 2021-08-18.
- 17 [Wiy+19] R. R. Wiyatno, A. Xu, O. Dia, and A. de Berker. “Adversarial Examples in Modern Machine Learning: A Review”. In: (2019). arXiv: [1911.05268 \[cs.LG\]](#).
- 18 [WJ08] M. J. Wainwright and M. I. Jordan. “Graphical models, exponential families, and variational inference”. In: *Foundations and Trends in Machine Learning* 1–2 (2008), pp. 1–305.
- 19 [WJ21] Y. Wang and M. I. Jordan. *Desiderata for Representation Learning: A Causal Perspective*. 2021. arXiv: [2109.03795 \[stat.ML\]](#).
- 20 [WJW03] M. J. Wainwright, T. S. Jaakkola, and A. S. Willsky. “Tree-based reparameterization framework for analysis of sum-product and related algorithms”. In: *IEEE Trans. on Information Theory* 49.5 (2003), pp. 1120–1146.
- 21 [WK18] E. Wong and Z. Kolter. “Provable defenses against adversarial examples via the convex outer adversarial polytope”. In: *International Conference on Machine Learning*. PMLR. 2018, pp. 5286–5295.
- 22 [WK96] G. Widmer and M. Kubat. “Learning in the presence of concept drift and hidden contexts”. In: *Mach. Learn.* 23.1 (1996), pp. 69–101.
- 23 [WKS21] V. Wild, M. Kanagawa, and D. Sejdinovic. “Connections and Equivalences between the Nyström Method and Sparse Variational Gaussian Processes”. In: (2021). arXiv: [2106.01121 \[stat.ML\]](#).
- 24 [WL14a] N. Whiteley and A. Lee. “Twisted particle filters”. en. In: *Annals of Statistics* 42.1 (Feb. 2014), pp. 115–141.
- 25 [WL14b] J. L. Williams and R. A. Lau. “Approximate evaluation of marginal association probabilities with belief propagation”. In: *IEEE Trans. Aerosp. Electron. Syst.* 50.4 (2014).
- 26 [WL19] A. Wehenkel and G. Louppe. “Unconstrained Monotonic Neural Networks”. In: *NIPS*. 2019.
- 27 [WLL16] W. Wang, H. Lee, and K. Livescu. “Deep Variational Canonical Correlation Analysis”. In: *arXiv* (2016).
- 28 [WLZ19] D. Widmann, F. Lindsten, and D. Zachariah. “Calibration tests in multi-class classification: A unifying framework”. In: *NIPS*. Curran Associates, Inc., 2019, pp. 12236–12246.
- 29 [WM12] K. Wakabayashi and T. Miura. “Forward-Backward Activation Algorithm for Hierarchical Hidden Markov Models”. In: *NIPS*. 2012.
- 30 [WMR17] S. Wachter, B. Mittelstadt, and C. Russell. “Counterfactual explanations without opening the black box: Automated decisions and the GDPR”. In: *Harv. JL & Tech.* 31 (2017), p. 841.
- 31 [WMR18] S. Wachter, B. Mittelstadt, and C. Russell. *Counterfactual Explanations without Opening the Black Box: Automated Decisions and the GDPR*. 2018. arXiv: [1711.00399 \[cs.AI\]](#).
- 32 [WN01] E. A. Wan and A. T. Nelson. “Dual EKF Methods”. In: *Kalman Filtering and Neural Networks*. Ed. by S. Haykin. Wiley, 2001.
- 33 [WN07] D. Wipf and S. Nagarajan. “A new view of automatic relevancy determination”. In: *NIPS*. 2007.
- 34 [WN10] D. Wipf and S. Nagarajan. “Iterative Reweighted ℓ_1 and ℓ_2 Methods for Finding Sparse Solutions”. In: *J. of Selected Topics in Signal Processing (Special Issue on Compressive Sensing)* 4.2 (2010).
- 35 [WN15] A. G. Wilson and H. Nickisch. “Kernel Interpolation for Scalable Structured Gaussian Processes (KISS-GP)”. In: *ICML*. ICML’15. JMLR.org, 2015, pp. 1775–1784.
- 36 [WN18] C. K. I. Williams and C. Nash. “Autoencoders and Probabilistic Inference with Missing Data: An Exact Solution for The Factor Analysis Case”. In: (2018). arXiv: [1801.03851 \[cs.LG\]](#).
- 37 [WO12] M. Wiering and M. van Otterlo, eds. *Reinforcement learning: State-of-the-art*. Springer, 2012.

- [Woł+21] M. Wołczyk, M. Zając, R. Pascanu, Ł. Kuściński, and P. Miłoś. “Continual World: A Robotic Benchmark For Continual Reinforcement Learning”. In: *NIPS*. 2021.
- [Wol76] P. Wolfe. “Finding the nearest point in a polytope”. In: *Mathematical Programming* 11 (1976), pp. 128–149.
- [Wol92] D. Wolpert. “Stacked Generalization”. In: *Neural Networks* 5.2 (1992), pp. 241–259.
- [Woo+09] F. Wood, C. Archambeau, J. Gasthaus, L. James, and Y. W. Teh. “A Stochastic Memoizer for Sequence Data”. In: *ICML*. 2009.
- [Woo+11] F. Wood, J. Gasthaus, C. Archambeau, L. James, and Y. W. Teh. “The sequence memoizer”. In: *Comm. of the ACM* 54.2 (2011), pp. 91–98.
- [Woo+19] B. Woodworth, S. Gunasekar, P. Savarese, E. Moroshko, I. Golan, J. Lee, D. Soudry, and N. Srebro. “Kernel and Rich Regimes in Overparametrized Models”. In: (2019). arXiv: [1906.05827 \[cs.LG\]](https://arxiv.org/abs/1906.05827).
- [WP18] H. Wang and H. Poon. “Deep Probabilistic Logic: A Unifying Framework for Indirect Supervision”. In: *EMNLP*. 2018.
- [WP19] S. Wiegrefe and Y. Pinter. “Attention is not explanation”. In: *arXiv preprint arXiv:1908.04626* (2019).
- [WR15] F. Wang and C. Rudin. “Falling Rule Lists”. In: *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*. Ed. by G. Lebanon and S. V. N. Vishwanathan. Vol. 38. Proceedings of Machine Learning Research. San Diego, California, USA: PMLR, 2015, pp. 1013–1022.
- [WRN10] D. Wipf, B. Rao, and S. Nagarajan. “Latent Variable Bayesian Models for Promoting Sparsity”. In: *IEEE Transactions on Information Theory* (2010).
- [WRZH04] M. Welling, M. Rosen-Zvi, and G. Hinton. “Exponential family harmoniums with an application to information retrieval”. In: *NIPS-14*. 2004.
- [WS01] C. K. I. Williams and M. Seeger. “Using the Nyström Method to Speed Up Kernel Machines”. In: *NIPS*. Ed. by T. K. Leen, T. G. Dietterich, and V. Tresp. MIT Press, 2001, pp. 682–688.
- [WS05] M. Welling and C. Sutton. “Learning in Markov Random Fields with Contrastive Free Energies”. In: *Tenth International Workshop on Artificial Intelligence and Statistics (AISTATS)*. 2005.
- [WS93] D. Wolpert and C. Strauss. “What Bayes has to say about the evidence procedure”. In: *Proc. Workshop on Maximum Entropy and Bayesian methods*. 1993.
- [WSG21] C. Wang, S. Sun, and R. Grosse. “Beyond Marginal Uncertainty: How Accurately can Bayesian Regression Models Estimate Posterior Predictive Correlations?” In: *AISTATS*. Ed. by A. Banerjee and K. Fukumizu. Vol. 130. Proceedings of Machine Learning Research. PMLR, 2021, pp. 2476–2484.
- [WSN00] B. Williams, T. Santner, and W. Notz. “Sequential design of computer experiments to minimize integrated response functions”. In: *Statistica Sinica* 10 (2000), pp. 1133–1152.
- [WSS21] W. J. Wilkinson, S. Särkkä, and A. Solin. “Bayes-Newton Methods for Approximate Bayesian Inference with PSD Guarantees”. In: (2021). arXiv: [2111.01721 \[stat.ML\]](https://arxiv.org/abs/2111.01721).
- [WT01] M. Welling and Y.-W. Teh. “Belief Optimization for Binary Networks: a Stable Alternative to Loopy Belief Propagation”. In: *UAI*. 2001.
- [WT19] R. Wen and K. Torkkola. “Deep Generative Quantile-Copula Models for Probabilistic Forecasting”. In: *ICML*. 2019.
- [WT90] G. Wei and M. Tanner. “A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithms”. In: *JASA* 85.411 (1990), pp. 699–704.
- [WTB20] Y. Wen, D. Tran, and J. Ba. “BatchEnsemble: an Alternative Approach to Efficient Ensemble and Lifelong Learning”. In: *ICLR*. 2020.
- [WTN19] Y. Wu, G. Tucker, and O. Nachum. *Behavior Regularized Offline Reinforcement Learning*. arXiv: [1911.11361](https://arxiv.org/abs/1911.11361). 2019.
- [Wu+06] Y. Wu, D. Hu, M. Wu, and X. Hu. “A Numerical Integration Perspective on Gaussian Filters”. In: *IEEE Trans. Signal Process.* 54.8 (2006), pp. 2910–2921.
- [Wu+17] Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba. “Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation”. In: *NIPS*. 2017.
- [Wu+19a] A. Wu, S. Nowozin, E. Meeds, R. E. Turner, J. M. Hernández-Lobato, and A. L. Gaunt. “Fixing Variational Bayes: Deterministic Variational Inference for Bayesian Neural Networks”. In: *ICLR*. 2019.
- [Wu+19b] M. Wu, S. Parbhoo, M. Hughes, R. Kindle, L. Celi, M. Zazzi, V. Roth, and F. Doshi-Velez. “Regional tree regularization for interpretability in black box models”. In: *AAAI* (2019).
- [Wu+21] T. Wu, M. T. Ribeiro, J. Heer, and D. S. Weld. “Polyjuice: Generating Counterfactuals for Explaining, Evaluating, and Improving Models”. In: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2021.
- [Wüt+16] M. Wüthrich, S. Trimpe, C. Garcia Cifuentes, D. Kappler, and S. Schaal. “A new perspective and extension of the Gaussian Filter”. In: *The International Journal of Robotics Research* 35.14 (2016), pp. 1731–1749.
- [WW12] Y. Wu and D. P. Wipf. “Dual-Space Analysis of the Sparse Linear Model”. In: *NIPS*. 2012.
- [WY02] D. Wilkinson and S. Yeung. “Conditional simulation from highly structured Gaussian systems with application to blocking-MCMC for the Bayesian analysis of very large linear models”. In: *Statistics and Computing* 12 (2002), pp. 287–300.
- [WYG14] J. Wen, C.-N. Yu, and R. Greiner. “Robust Learning under Uncertain Test Distributions: Relating Covariate Shift to Model Misspecification”. In: *ICML*. Vol. 32. Proceedings of Machine Learning Research. PMLR, 2014, pp. 631–639.
- [WZ19] J. Wei and K. Zou. “EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, Nov. 2019, pp. 6382–6388.
- [WZR20] S. Wu, H. R. Zhang, and C. Ré. “Understanding and Improving Information Transfer in Multi-Task Learning”. In: *International Conference on Learning Representations*. 2020.

- 1
- 2 [XC19] Z. Xia and A. Chakrabarti. “Training Image
Estimators without Image Ground-Truth”. In: *NIPS*.
3 2019.
- 4 [XD18] J. Xu and G. Durrett. “Spherical Latent Spaces
for Stable Variational Autoencoders”. In: *EMNLP*.
5 2018.
- 6 [Xia+21] K. Xiao, L. Engstrom, A. Ilyas, and A. Madry.
“Noise or Signal: The Role of Image Backgrounds in
Object Recognition”. In: *ICLR*. 2021.
- 7 [Xie+16] J. Xie, Y. Lu, S.-C. Zhu, and Y. N. Wu. “A
Theory of Generative ConvNet”. In: *ICML*. 2016.
- 8 [Xie+18] J. Xie, Y. Lu, R. Gao, and Y. N. Wu. “Co-
operative Learning of Energy-Based Model and La-
tent Variable Model via MCMC Teaching”. In: *AAAI*.
9 Vol. 1. 6. 2018, p. 7.
- 10 [Xie+22] S. M. Xie, A. Raghunathan, P. Liang, and T.
Ma. “An Explanation of In-context Learning as Im-
plicit Bayesian Inference”. In: *ICLR*. 2022.
- 11 [XJ96] L. Xu and M. I. Jordan. “On Convergence Prop-
erties of the EM Algorithm for Gaussian Mixtures”. In:
Neural Computation 8 (1996), pp. 129–151.
- 12 [XKV22] Z. Xiao, K. Kreis, and A. Vahdat. “Tackling
the Generative Learning Trilemma with Denoising Di-
fusion GANs”. In: *ICLR*. 2022.
- 13 [Xu+06] Z. Xu, V. Tresp, K. Yu, and H.-P. Kriegel. “In-
finite hidden relational models”. In: *UAI*. 2006.
- 14 [Xu+07] Z. Xu, V. Tresp, S. Yu, K. Yu, and H.-P.
Kriegel. “Fast Inference in Infinite Hidden Relational
Models”. In: *Workshop on Mining and Learning with
Graphs*. 2007.
- 15 [Xu+15] J. Xu, L. Mukherjee, Y. Li, J. Warner, J. M.
Rehg, and V. Singh. “Gaze-enabled egocentric video
summarization via constrained submodular maximiza-
tion”. In: *Proceedings of the IEEE conference on com-
puter vision and pattern recognition*. 2015, pp. 2235–
2244.
- 16 [Xu18] J. Xu. “Distance-based Protein Folding Pow-
ered by Deep Learning”. In: (2018). arXiv: [1811.03481](https://arxiv.org/abs/1811.03481)
[q-bio.BM].
- 17 [Xu+19] M. Xu, M. Quiroz, R. Kohn, and S. A. Sisson.
“Variance reduction properties of the reparameteriza-
tion trick”. In: *AISTATS*. Ed. by K. Chaudhuri and M.
Sugiyama. Vol. 89. *Proceedings of Machine Learning*
Research. PMLR, 2019, pp. 2711–2720.
- 18 [Yad+18] S. Yadlowsky, H. Namkoong, S. Basu, J.
Duchi, and L. Tian. “Bounds on the conditional and
average treatment effect with unobserved confounding
factors”. In: *arXiv e-prints*, arXiv:1808.09521 (Aug.
2018), arXiv:1808.09521. arXiv: [1808.09521](https://arxiv.org/abs/1808.09521) [stat.ME].
- 19 [Yad+21] S. Yadlowsky, S. Fleming, N. Shah, E. Brun-
skill, and S. Wager. *Evaluating Treatment Prioritiza-
tion Rules via Rank-Weighted Average Treatment
Effects*. 2021. arXiv: [2111.07966](https://arxiv.org/abs/2111.07966) [stat.ME].
- 20 [Yan+17] S. Yang, L. Xie, X. Chen, X. Lou, X. Zhu,
D. Huang, and H. Li. “Statistical parametric speech
synthesis using generative adversarial networks under
a multi-task learning framework”. In: *IEEE Auto-
matic Speech Recognition and Understanding Work-
shop (ASRU)*. 2017, pp. 685–691.
- 21 [Yan19] G. Yang. “Scaling Limits of Wide Neural Net-
works with Weight Sharing: Gaussian Process Behav-
ior, Gradient Independence, and Neural Tangent Ker-
nel Derivation”. In: (2019). arXiv: [1902.04760](https://arxiv.org/abs/1902.04760) [cs.NE].
- 22 [Yan+21] J. Yang, K. Zhou, Y. Li, and Z. Liu. “Gen-
eralized OOD Detection: A Survey”. In: (2021).
- 23 [Yan74] H. Yanai. “Unification of various techniques of
multivariate analysis by means of generalized coeffi-
cient of determination (GCD)”. In: *J. Behaviormet-
rics* 1.1 (1974), pp. 45–54.
- 24 [Yan81] M. Yannakakis. “Computing the minimum fill-
in is NP-complete”. In: *SIAM J. Alg. Discrete Meth-
ods* 2 (1981), pp. 77–79.
- 25 [Yao+18a] Y. Yao, A. Vehtari, D. Simpson, and A.
Gelman. “Using Stacking to Average Bayesian Predic-
tive Distributions (with Discussion)”. en. In: *Bayesian
Analysis* 13.3 (2018), pp. 917–1007.
- 26 [Yao+18b] Y. Yao, A. Vehtari, D. Simpson, and A. Gel-
man. “Yes, but Did It Work?: Evaluating Variational
Inference”. In: *ICML*. Vol. 80. *Proceedings of Machine
Learning Research*. PMLR, 2018, pp. 5581–5590.
- 27 [YBM20] Y. Yang, R. Bamler, and S. Mandt. *Improv-
ing Inference for Neural Image Compression*. 2020.
arXiv: [2006.04240](https://arxiv.org/abs/2006.04240) [eess.IV].
- 28 [YBW15] F. Yang, S. Balakrishnan, and M. J. Wain-
wright. “Statistical and Computational Guarantees for
the Baum-Welch Algorithm”. In: (2015). arXiv: [1512.08269](https://arxiv.org/abs/1512.08269) [stat.ML].
- 29 [Yed11] J. S. Yedidia. “Message-Passing Algorithms for
Inference and Optimization”. In: *J. Stat. Phys.* 145.4
(2011), pp. 860–890.
- 30 [Yeh+18] C.-K. Yeh, J. S. Kim, I. E. H. Yen, and P.
Ravikumar. *Representer Point Selection for Explain-
ing Deep Neural Networks*. 2018. arXiv: [1811.09720](https://arxiv.org/abs/1811.09720) [cs.LG].
- 31 [Yeh+19a] C.-K. Yeh, C.-Y. Hsieh, A. Suggala, D. I.
Inouye, and P. K. Ravikumar. “On the (in) fidelity
and sensitivity of explanations”. In: *Advances in
Neural Information Processing Systems* 32 (2019),
pp. 10967–10978.
- 32 [Yeh+19b] C.-K. Yeh, B. Kim, S. O. Arik, C.-L. Li, T.
Pfister, and P. K. Ravikumar. “On completeness-aware
concept-based explanations in deep neural networks”.
In: *arXiv preprint arXiv:1910.07969* (2019).
- 33 [Yeu+17] S. Yeung, A. Kannan, Y. Dauphin, and L.
Fei-Fei. “Tackling Over-pruning in Variational Au-
toencoders”. In: *ICML Workshop on “Principled Ap-
proaches to Deep Learning”*. 2017.
- 34 [Yeu91a] R. W. Yeung. “A new outlook on Shannon’s
information measures”. In: *IEEE Trans. Inf. Theory*
37.3 (1991), pp. 466–474.
- 35 [Yeu91b] R. W. Yeung. “A new outlook on Shannon’s
information measures”. In: *IEEE Trans. on Informa-
tion Theory* 37 (1991), pp. 466–474.
- 36 [YFW00] J. Yedidia, W. T. Freeman, and Y. Weiss.
“Generalized Belief Propagation”. In: *NIPS*. 2000.
- 37 [YH21] G. Yang and E. J. Hu. “Feature Learning in
Infinite-Width Neural Networks”. In: *ICML*. 2021.
- 38 [Yin+19a] D. Yin, R. G. Lopes, J. Shlens, E. D. Cubuk,
and J. Gilmer. “A Fourier Perspective on Model Ro-
bustness in Computer Vision”. In: *NIPS*. 2019.
- 39 [Yin+19b] P. Yin, J. Lyu, S. Zhang, S. Osher, Y. Qi,
and J. Xin. “Understanding Straight-Through Estima-
tor in Training Activation Quantized Neural Nets”. In:
ICLR. 2019.
- 40 [Yin+19c] R. Ying, D. Bourgeois, J. You, M. Zit-
nik, and J. Leskovec. “Gnnexplainer: Generating ex-
planations for graph neural networks”. In: *Advances
in neural information processing systems* 32 (2019),
p. 9240.
- 41
- 42
- 43
- 44
- 45
- 46
- 47

- [Yin+20] D. Yin, M. Farajtabar, A. Li, N. Levine, and A. Mott. “Optimization and Generalization of Regularization-Based Continual Learning: a Loss Approximation Viewpoint”. In: (2020). arXiv: [2006.10974 \[cs.LG\]](#).
- [YK06] A. Yuille and D. Kersten. “Vision as Bayesian inference: analysis by synthesis?” en. In: *Trends Cogn. Sci.* 10.7 (2006), pp. 301–308.
- [YK19] M. Yang and B. Kim. *Benchmarking Attribution Methods with Relative Feature Importance*. 2019. arXiv: [1907.09701 \[cs.LG\]](#).
- [YMT22] Y. Yang, S. Mandt, and L. Theis. “An Introduction to Neural Data Compression”. In: (2022). arXiv: [2202.06533 \[cs.LG\]](#).
- [Yoo+18] K. Yoon, R. Liao, Y. Xiong, L. Zhang, E. Feitaya, R. Urtasun, R. Zemel, and X. Pitkow. “Inference in Probabilistic Graphical Models by Graph Neural Networks”. In: *ICLR Workshop*. 2018.
- [You19] A. Young. “Consistency without inference: Instrumental variables in practical application”. In: (2019).
- [You89] L. Younes. “Parameter estimation for imperfectly observed Gibbsian fields”. In: *Probab. Theory and Related Fields* 82 (1989), pp. 625–645.
- [You99] L. Younes. “On the convergence of Markovian stochastic algorithms with rapidly decreasing ergodicity rates”. In: *Stochastics: An International Journal of Probability and Stochastic Processes* 65.3-4 (1999), pp. 177–228.
- [Yu+06] S. Yu, K. Yu, V. Tresp, K. H-P., and M. Wu. “Supervised probabilistic principal component analysis”. In: *KDD*. 2006.
- [Yu10] S.-Z. Yu. “Hidden Semi-Markov Models”. In: *Artificial Intelligence J.* 174.2 (2010).
- [Yu+16] F. X. X. Yu, A. T. Suresh, K. M. Choromanski, D. N. Holtmann-Rice, and S. Kumar. “Orthogonal Random Features”. In: *NIPS*. Curran Associates, Inc., 2016, pp. 1975–1983.
- [Yu+17] L. Yu, W. Zhang, J. Wang, and Y. Yu. “Seqgan: Sequence generative adversarial nets with policy gradient”. In: *Thirty-first AAAI conference on artificial intelligence*. 2017.
- [Yu+18] Y. Yu et al. “Dynamic Control Flow in Large-Scale Machine Learning”. In: *Proceedings of the Thirteenth EuroSys Conference*. EuroSys ’18. Association for Computing Machinery, 2018.
- [Yu+20] L. Yu, Y. Song, J. Song, and S. Ermon. “Training Deep Energy-Based Models with f-Divergence Minimization”. In: *arXiv preprint arXiv:2003.03463* (2020).
- [Yu+21] J. Yu, X. Li, J. Y. Koh, H. Zhang, R. Pang, J. Qin, A. Ku, Y. Xu, J. Baldridge, and Y. Wu. “Vector-quantized Image Modeling with Improved VQGAN”. In: (2021). arXiv: [2110.04627 \[cs.CV\]](#).
- [Yu+22] J. Yu et al. “Scaling Autoregressive Models for Content-Rich Text-to-Image Generation”. 2022.
- [Yua+19] X. Yuan, P. He, Q. Zhu, and X. Li. “Adversarial Examples: Attacks and Defenses for Deep Learning”. en. In: *IEEE Trans. Neural Networks and Learning Systems* 30.9 (2019), pp. 2805–2824.
- [Yui01] A. Yuille. “CCCP algorithms to minimize the Bethe and Kikuchi free energies: convergent alternatives to belief propagation”. In: *Neural Computation* 14 (2001), pp. 1691–1722.
- [YW04] C. Yanover and Y. Weiss. “Finding the M Most Probable Configurations in Arbitrary Graphical Models”. In: *NIPS*. 2004.
- [YWX17] J.-g. Yao, X. Wan, and J. Xiao. “Recent advances in document summarization”. In: *Knowledge and Information Systems* 53.2 (2017), pp. 297–336.
- [YZ19] G. Yaroslavtsev and S. Zhou. “Approximate F_2 -Sketching of Valuation Functions”. In: *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2019)*. Ed. by D. Achlioptas and L. A. Végh. Vol. 145. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2019, 69:1–69:21.
- [YZ22] Y. Yang and P. Zhai. “Click-through rate prediction in online advertising: A literature review”. In: *Inf. Process. Manag.* 59.2 (2022), p. 102853.
- [ZA12] J. Zou and R. Adams. “Priors for Diversity in Generative Latent Variable Models”. In: *NIPS*. 2012.
- [Zaf+22] M. Zaffran, A. Dieuleveut, O. Féron, Y. Goude, and J. Josse. “Adaptive Conformal Predictions for Time Series”. In: (2022). arXiv: [2202.07282 \[stat.ML\]](#).
- [Zai+20] S. Zaidi, A. Zela, T. Elskens, C. Holmes, F. Hutter, and Y. W. Teh. “Neural Ensemble Search for Performant and Calibrated Predictions”. In: (2020). arXiv: [2006.08573 \[cs.LG\]](#).
- [Zan21] N. Zanichelli. *IAML Distill Blog: Transformers in Vision*. 2021.
- [ZB18] T. Zhou and J. Bilmes. “Minimax curriculum learning: Machine teaching with desirable difficulties and scheduled diversity”. In: *International Conference on Learning Representations*. 2018.
- [ZB21] B. Zhao and H. Bilen. “Dataset Condensation with Differentiable Siamese Augmentation”. In: *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*. Ed. by M. Meila and T. Z. 0001. Vol. 139. Proceedings of Machine Learning Research. PMLR, 2021, pp. 12674–12685.
- [Zbo+21] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. “Barlow twins: Self-supervised learning via redundancy reduction”. In: *arXiv preprint arXiv:2103.03230* (2021).
- [ZDK15] J. Zhang, J. Djolonga, and A. Krause. “Higher-order inference for multi-class log-supermodular models”. In: *Proceedings of the IEEE International Conference on Computer Vision*. 2015, pp. 1859–1867.
- [ZE01a] B. Zadrozny and C. Elkan. “Obtaining calibrated probability estimates from decision trees and naive Bayesian classifiers”. In: *ICML*. 2001.
- [ZE01b] B. Zadrozny and C. Elkan. “Transforming classifier scores into accurate multiclass probability estimates”. In: *KDD*. 2001.
- [Zec+18a] J. R. Zech, M. A. Badgeley, M. Liu, A. B. Costa, J. J. Titano, and E. K. Oermann. “Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study”. en. In: *PLoS Med.* 15.11 (Nov. 2018), e1002683.
- [Zec+18b] J. R. Zech, M. A. Badgeley, M. Liu, A. B. Costa, J. J. Titano, and E. K. Oermann. “Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study”. In: *PLoS medicine* 15.11 (2018), e1002683.

- 1
- 2 [Zel76] A. Zellner. “Bayesian and non-Bayesian analysis
3 of the regression model with multivariate Student-t
error terms”. In: *JASA* 71.354 (1976), pp. 400–405.
- 4 [Zel86] A. Zellner. “On assessing prior distributions
5 and Bayesian regression analysis with g-prior distributions”. In: *Bayesian inference and decision techniques, Studies of Bayesian and Econometrics and Statistics volume 6*. North Holland, 1986.
- 6 [Zen+18] C. Zeno, I. Golan, E. Hoffer, and D. Soudry.
7 “Task Agnostic Continual Learning Using Online
8 Variational Bayes”. In: (2018). arXiv: [1803 . 10123 \[stat.ML\]](https://arxiv.org/abs/1803.10123).
- 9 [Zen+21] C. Zeno, I. Golan, E. Hoffer, and D. Soudry.
10 “Task-Agnostic Continual Learning Using Online Variational Bayes With Fixed-Point Updates”. In: *Neural Comput.* 33.11 (2021), pp. 3139–3177.
- 11 [Zer+19] J. Zerilli, A. Knott, J. Maclaurin, and C. Gavaghan. “Transparency in algorithmic and human
12 decision-making: is there a double standard?” In: *Philosophy & Technology* 32.4 (2019), pp. 661–683.
- 13 [ZF14] M. D. Zeiler and R. Fergus. “Visualizing and
14 understanding convolutional networks”. In: *European conference on computer vision*. Springer, 2014,
15 pp. 818–833.
- 16 [ZFV20] G. Zeni, M. Fontana, and S. Vantini. “Conformal Prediction: a Unified Review of Theory and New
17 Challenges”. In: (2020). arXiv: [2005 . 07972 \[cs.LG\]](https://arxiv.org/abs/2005.07972).
- 18 [ZG21] D. Zou and Q. Gu. “On the Convergence of
19 Hamiltonian Monte Carlo with Stochastic Gradients”.
20 In: *ICML*. Ed. by M. Meila and T. Zhang, Vol. 139. Proceedings of Machine Learning Research. PMLR,
21 pp. 13012–13022.
- 22 [ZGR21] L. Zhang, M. Goldstein, and R. Ranganath.
23 “Understanding Failures in Out-of-Distribution Detection
24 with Deep Generative Models”. In: *ICML*. Ed.
25 by M. Meila and T. Zhang, Vol. 139. Proceedings of
26 Machine Learning Research. PMLR, 2021, pp. 12427–
27 12436.
- 28 [Zha+13a] K. Zhang, B. Schölkopf, K. Muandet, and
29 Z. Wang. “Domain Adaptation under Target and
30 Conditional Shift”. In: *Proceedings of the 30th International Conference on Machine Learning*. 2013,
31 pp. 819–827.
- 32 [Zha+13b] K. Zhang, B. Schölkopf, K. Muandet, and
33 Z. Wang. “Domain Adaptation under Target and Conditional Shift”. In: *ICML*. Vol. 28. 2013.
- 34 [Zha+16] S. Zhai, Y. Cheng, R. Feris, and Z. Zhang.
35 “Generative adversarial networks as variational training
36 of energy based models”. In: *arXiv preprint arXiv:1611.01799* (2016).
- 37 [Zha+17] C. Zhang, S. Bengio, M. Hardt, B. Recht, and
38 O. Vinyals. “Understanding deep learning requires rethinking generalization”. In: *ICLR*. 2017.
- 39 [Zha+18] G. Zhang, S. Sun, D. Duvenaud, and R. Grosse. “Noisy Natural Gradient as Variational Inference”. In: *ICML*. 2018.
- 40 [Zha+19a] X. Zhai, J. Puigcerver, A. Kolesnikov, P.
41 Ruyssen, C. Riquelme, M. Lucic, J. Djolonga, A. S.
42 Pinto, M. Neumann, A. Dosovitskiy, et al. “A large-
43 scale study of representation learning with the visual
44 task adaptation benchmark”. In: *arXiv preprint arXiv:1910.04867* (2019).
- 45 [Zha+19b] C. Zhang, J. Butepage, H. Kjellstrom, and
46 S. Mandt. “Advances in Variational Inference”. In:
47 *IEEE PAMI* (2019), pp. 2008–2026.
- [Zha+19c] H. Zhang, I. Goodfellow, D. Metaxas, and
A. Odena. “Self-attention generative adversarial networks”. In: *International conference on machine learning*. PMLR, 2019, pp. 7354–7363.
- [Zha+19d] L. Zhao, K. Korovina, W. Si, and M. Cheung. “Approximate inference with Graph Neural Networks”. 2019.
- [Zha+20a] A. Zhang, Z. Lipton, M. Li, and A. Smola. *Drive into deep learning*. 2020.
- [Zha+20b] H. Zhang, A. Li, J. Guo, and Y. Guo.
“Hybrid models for open set recognition”. In: *European Conference on Computer Vision*. Springer, 2020,
pp. 102–117.
- [Zha+20c] R. Zhang, B. Dai, L. Li, and D. Schuurmans.
“GenDICE: Generalized Offline Estimation of Stationary Values”. In: *ICLR*. 2020.
- [Zha+20d] R. Zhang, C. Li, J. Zhang, C. Chen, and
A. G. Wilson. “Cyclical stochastic gradient MCMC for
Bayesian deep learning”. In: *ICLR*. 2020.
- [Zha+20e] X. Zhang, Y. Li, Z. Zhang, and Z.-L. Zhang.
“*f*-GAIL: Learning *f*-Divergence for Generative Adversarial Imitation Learning”. In: *Neural Information Processing Systems* (2020).
- [Zha+20f] Y. Zhang, X. Chen, Y. Yang, A. Ramamurthy, B. Li, Y. Qi, and L. Song. “Efficient Probabilistic Logic Reasoning with Graph Neural Networks”. In: *ICLR*. 2020.
- [Zha+21a] X. Zhai, A. Kolesnikov, N. Houlsby, and
L. Beyer. “Scaling vision transformers”. In: *arXiv preprint arXiv:2106.04560* (2021).
- [Zha+21b] H. Zhang, J. Y. Koh, J. Baldridge, H. Lee,
and Y. Yang. “Cross-Modal Contrastive Learning for
Text-to-Image Generation”. In: *CVPR*. 2021.
- [Zhe+15] S. Zheng, S. Jayasumana, B. Romera-Paredes,
V. Vineet, Z. Su, D. Du, C. Huang, and P. Torr.
“Conditional Random Fields as Recurrent Neural Networks”. In: *ICCV*. 2015.
- [Zho+19a] S. Zhou, M. Gordon, R. Krishna, A. Narcomay, L. Fei-Fei, and M. Bernstein. “HYPE: A Benchmark for Human eYe Perceptual Evaluation of Generative Models”. In: *NIPS*. Curran Associates, Inc., 2019, pp. 3444–3456.
- [Zho+19b] S. Zhou, M. L. Gordon, R. Krishna, A. Narcomay, L. Fei-Fei, and M. S. Bernstein. “HYPE: a benchmark for human eye perceptual evaluation of generative models”. In: *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 2019, pp. 3449–3461.
- [Zho20] G. Zhou. “Mixed Hamiltonian Monte Carlo for
Mixed Discrete and Continuous Variable”. In: (2020).
arXiv: [1909 . 04852 \[stat.CO\]](https://arxiv.org/abs/1909.04852).
- [ZHT06] H. Zou, T. Hastie, and R. Tibshirani. “Sparse principal component analysis”. In: *JCGS* 15.2 (2006),
pp. 262–286.
- [Zhu+17] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros.
“Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks”. In: *ICCV*. 2017.
- [Zhu+18] X. Zhu, A. Singla, S. Zilles, and A. N. Rafferty. “An overview of machine teaching”. In: *arXiv preprint arXiv:1801.05927* (2018).
- [Zhu+21] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He. “A Comprehensive Survey on Transfer Learning”. In: *Proc. IEEE* 109.1 (2021).

- 1 [Zie+08] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and
2 A. K. Dey. “Maximum Entropy Inverse Reinforcement
3 Learning”. In: *AAAI*. 2008, pp. 1433–1438.
- 4 [ZIE16] R. Zhang, P. Isola, and A. A. Efros. “Colorful
5 image colorization”. In: *European conference on com-
puter vision*. Springer. 2016, pp. 649–666.
- 6 [ZIE17] R. Zhang, P. Isola, and A. A. Efros. “Split-
7 brain autoencoders: Unsupervised learning by cross-
8 channel prediction”. In: *Proceedings of the IEEE Con-
ference on Computer Vision and Pattern Recogni-
tion*. 2017, pp. 1058–1067.
- 9 [Zin+20] L. Zintgraf, K. Shiarlis, M. Igl, S. Schulze, Y.
10 Gal, K. Hofmann, and S. Whiteson. “VariBAD: A Very
11 Good Method for Bayes-Adaptive Deep RL via Meta-
12 Learning”. In: *ICLR*. 2020.
- 13 [Ziy+19] L. Ziyin, Z. Wang, P. P. Liang, R. Salakhut-
14 dinov, L.-P. Morency, and M. Ueda. “Deep gamblers:
15 Learning to abstain with portfolio theory”. In: *NIPS*.
16 June 2019.
- 17 [ZL21] A. Zhou and S. Levine. “Training on Test Data
18 with Bayesian Adaptation for Covariate Shift”. In:
19 *NIPS*. 2021.
- 20 [ZLF21] M. Zhang, S. Levine, and C. Finn. “MEMO:
21 Test Time Robustness via Adaptation and Augmen-
22 tation”. In: (2021). arXiv: [2110.09506 \[cs.LG\]](https://arxiv.org/abs/2110.09506).
- 23 [ZLG20] R. Zivan, O. Lev, and R. Galiki. “Beyond
24 Trees: Analysis and Convergence of Belief Propagation
25 in Graphs with Multiple Cycles”. In: *AAAI*. 2020.
- 26 [ZMB21] B. Zhao, K. R. Mopuri, and H. Bilen.
27 “Dataset Condensation with Gradient Matching”. In:
28 *International Conference on Learning Representa-
tions*. 2021.
- 29 [ZMG19] G. Zhang, J. Martens, and R. B. Grosse.
30 “Fast Convergence of Natural Gradient Descent for
31 Over-Parameterized Neural Networks”. In: *NIPS*.
32 2019, pp. 8082–8093.
- 33 [ZML16] J. J. Zhao, M. Mathieu, and Y. LeCun.
34 “Energy-based Generative Adversarial Network”. In:
35 (2016).
- 36 [Zob09] O. Zobay. “Mean field inference for the Dirich-
37 let process mixture model”. In: *Electronic J. of Statis-
38 tics* 3 (2009), pp. 507–545.
- 39 [Zoe07] O. Zoeter. “Bayesian generalized linear models
40 in a terabyte world”. In: *Proc. 5th International Sym-
41 posium on image and Signal Processing and Analy-
42 sis*. 2007.
- 43 [Zon+18] B. Zong, Q. Song, M. R. Min, W. Cheng, C.
44 Lumezanu, D. Cho, and H. Chen. “Deep Autoencoding
45 Gaussian Mixture Model for Unsupervised Anomaly
46 Detection”. In: *ICLR*. 2018.
- 47
- [ZP00] G. Zweig and M. Padmanabhan. “Exact alpha-
beta computation in logarithmic space with applica-
tion to map word graph construction”. In: *ICSLP*.
2000.
- [ZP96] N. Zhang and D. Poole. “Exploiting causal in-
dependence in Bayesian network inference”. In: *JAIR*
(1996), pp. 301–328.
- [ZR19a] Z. Ziegler and A. Rush. “Latent Normalizing
Flows for Discrete Sequences”. In: *Proceedings of the
36th International Conference on Machine Learning*.
2019, pp. 7673–7682.
- [ZR19b] Z. M. Ziegler and A. M. Rush. “Latent Normal-
izing Flows for Discrete Sequences”. In: *ICML*. 2019.
- [ZSB19] Q. Zhao, D. S. Small, and B. B. Bhattacharya.
“Sensitivity analysis for inverse probability weighting
estimators via the percentile bootstrap”. In: *Journal
of the Royal Statistical Society: Series B (Statistical
Methodology)* 81.4 (2019), pp. 735–761. eprint: <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/rssb.12327>.
- [ZSE19] S. Zhao, J. Song, and S. Ermon. “InfoVAE: In-
formation Maximizing Variational Autoencoders”. In:
AAAI. 2019.
- [ZW11] D. Zoran and Y. Weiss. “From learning models
of natural image patches to whole image restoration”.
In: *ICCV*. 2011.
- [ZW12] D. Zoran and Y. Weiss. “Natural Images, Gaus-
sian Mixtures and Dead Leaves”. In: *NIPS*. 2012,
pp. 1736–1744.
- [ZWM97] C. S. Zhu, N. Y. Wu, and D. Mumford. “Min-
imax Entropy Principle and Its Application to Texture
Modeling”. In: *Neural Computation* 9.8 (1997).
- [ZY08] J.-H. Zhao and P. L. H. Yu. “Fast ML Estima-
tion for the Mixture of Factor Analyzers via an ECM
Algorithm”. In: *IEEE Trans. on Neural Networks*
19.11 (2008).
- [ZY21] Y. Zhang and Q. Yang. “A Survey on Multi-
Task Learning”. In: *IEEE Trans. Knowl. Data Eng.*
(2021).
- [ZY97] Z. Zhang and R. W. Yeung. “A non-Shannon-
type conditional inequality of information quantities”.
In: *IEEE Transactions on Information Theory* 43.6
(1997), pp. 1982–1986.
- [ZY98] Z. Zhang and R. W. Yeung. “On characteriza-
tion of entropy function via information inequalities”.
In: *IEEE Transactions on Information Theory* 44.4
(1998), pp. 1440–1452.