

Transfer Learning and Causality*

Sargur N. Srihari

University at Buffalo, The State University of New York
USA

*Based on work at MILA, e.g., “A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms,” Y. Bengio et.al. 2019

Plan of Discussion

- Transfer Learning
- Sparse change assumption due to intervention
- Advantage of Correct Causal Model
- Parameterization of Causal Structure
- Advantage of Correct Causal Structure
- Human-inspired biases

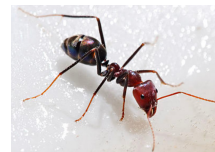
Simple example of transfer learning

- Visual classification

- Significantly more data in distribution sampled from cats and dogs



- Then learn to quickly generalize to ants and wasps



- Visual categories share

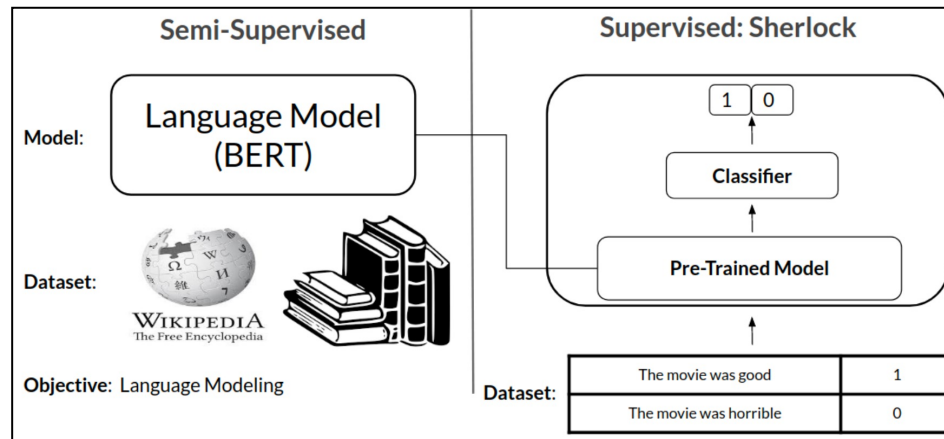
- Low-level notions of edges and visual shapes,
- Effects of geometric changes
- Changes in lighting

Practical Need for Transfer Learning

- Training a deep learning model requires
 - a lot of data and
 - more importantly a lot of time
- Useful to take advantage of pre-trained weights on huge datasets that took days/weeks to train, and leverage it towards our use case
- Depending on how much data we have at hand, here are the different ways to leverage this

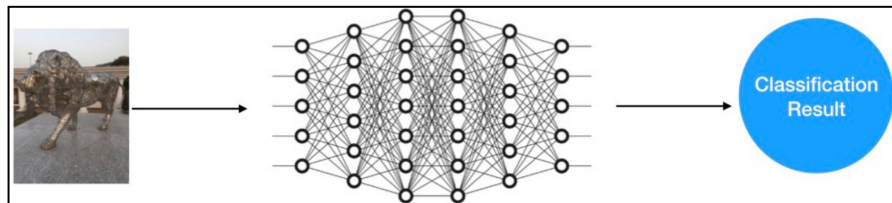
Transfer Learning Examples

Natural Language Processing

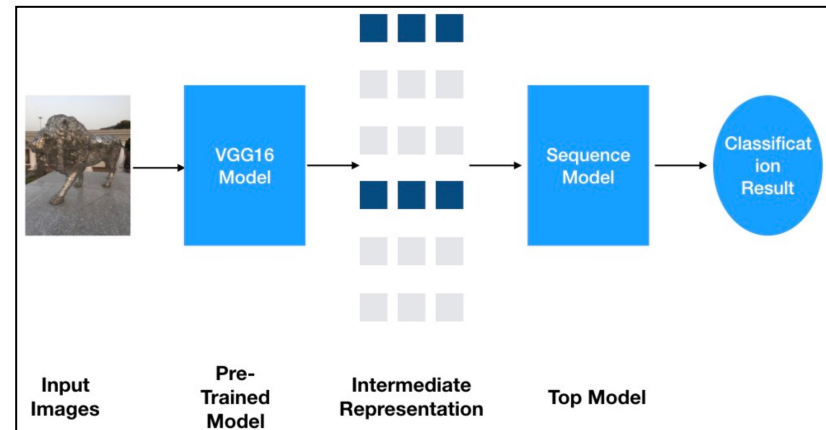


Computer Vision

Direct Learning



Transfer Learning



Fast transfer learning

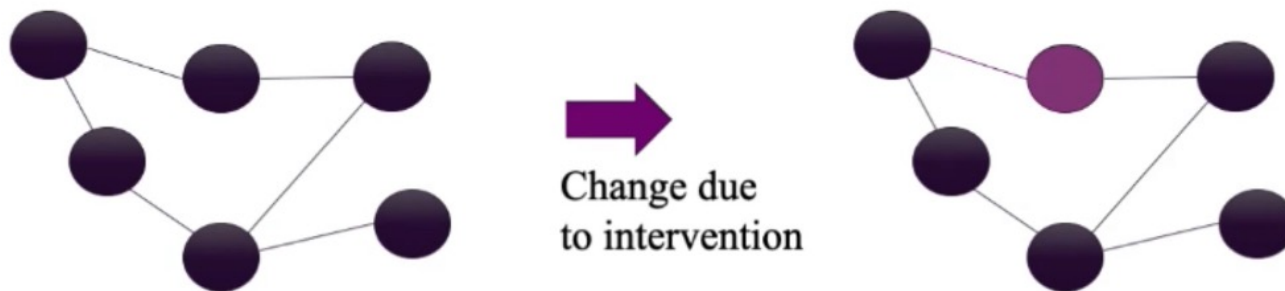
- Previous learning is a rich base to adapt to a new but related distribution
 - i.e., obtain good transfer
- Some new concept may have to be learned
 - but because most relevant concepts have already been captured, learning can be very fast on the transfer distribution

Weakness of Transfer Learning

- DL is weak to generalize beyond training distribution— which is what is needed in practice
- What has been learned in one setting needs to generalize well in other related distributions
 - Related distributions may involve the same concepts previously seen by learner
 - With changes due to actions of agents

Change due to intervention

- Changes are sparse when knowledge is represented in a modularized way, with few modules having changed
- Distributional change is due to actions by one or more agents
 - Such as interventions
 - a single causal variable is clamped to a particular value



Need for inferring causal structure

- A possible change to the joint distribution of the variables of interest due to an intervention, even one that has never been observed before
- Insufficient to learn joint distribution of observed variables
- Should also learn enough about the underlying high-level variables and their causal relations to be able to properly infer the effect of an intervention

Cause, Effect, Intervention

- Causal relations are inferred by intervention
 - A = Raining causes B = Open Umbrella (not vice-versa)
 - Changing probability of raining (because weather changed) doesn't change relation between A and B (captured by $P(B/A)$), but has an impact on $P(B)$
 - Conversely, an agent's intervention on B (Open umbrella) will have no effect on $P(A)$
- Asymmetry not visible from (A, B) training pairs
 - Until a change of distribution due to an intervention
 - Motivates setup where one learns from a set of distributions arising from not necessarily known interventions, not simply to capture a joint distribution but to discover some underlying causal structure

Changes in Distribution

WHAT CAUSES CHANGES IN DISTRIBUTION?

Hypothesis to replace iid assumption:

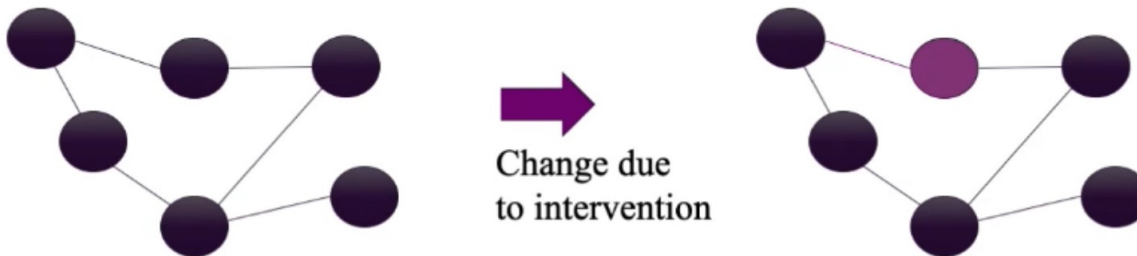
changes = consequence of an intervention on few causes or mechanisms

Underlying physics: actions are localized in space and time.

Extends the hypothesis of (informationally) Independent Mechanisms (*Scholkopf et al 2012*)

ICLR 2020: A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms, Bengio, Deleu, Rahaman, Ke, Lachapelle, Bilaniuk, Goyal, Pal

→ **local inference or adaptation in the right model**



Ke et al 2019, 2020; Brouillard et al NeurIPS 2020



Assumption of distribution changes

- ML methods exploit assumptions of data distribution
 - Or else *no-free lunch theorem* states that we cannot generalize
- Here we assume not just data distribution but how it changes
 - That when knowledge is appropriately represented these changes will be small
 - That only a few ground truth mechanisms change

Few updates needed for transfer

- Assume ground-truth data generative process is composed of independent mechanisms
- Few ground truth mechanisms and parameters change to go from training to transfer distribution
 - Model requires few updates, a few examples
- Gradient of unchanged parameters would be 0
- Leads to fast adaptation

Meta Learning Objective

- Measures speed of adaptation
 - i.e., a form of regret
 - In order to optimize the way in which knowledge should be represented, factorized and structured
- A stronger signal can be obtained when there are more non-stationarities
 - i.e., many changes in the distribution

From nuisance to training signal

- Nuisance in ML:
 - changes in distribution due to non-stationarity, uncontrolled intervention
- Converting to training signal
 - Factorize knowledge into components and mechanisms that match assumption of small change
- End-to-end optimization
 - Fast transfer and robustness to distributional changes

Factorization and Independence

- If data is generated from composing independent causal mechanisms
 - Then there exists a good factorization
 - Like factorizing code into independent modules



- At each time step few variables are changed
 - So small change in representation

Which is Cause and Which is effect

- Problem of determining whether variable A causes variable B or vice versa
- Training distribution $P(A,B)$
 - Lots of examples (a,b) are available from training distribution
 - Both $A \rightarrow B$ and $B \rightarrow A$ perform equally well in terms of ordinary generalization
- Transfer distribution
 - Very few samples available

Advantage of Correct Causal Model

- A and B are discrete with N possible values
- Consider two parametrizations ($A \rightarrow B$ model and $B \rightarrow A$ model) to estimate their joint:
 - $P_{A \rightarrow B}(A;B) = P_{A \rightarrow B}(A)P_{A \rightarrow B}(B | A)$
 - $P_{B \rightarrow A}(A;B) = P_{B \rightarrow A}(B)P_{B \rightarrow A}(A | B)$
 - Both models have $O(N^2)$ parameters
 - Maximum likelihood estimation leads to indistinguishable test set performance (where the test set is sampled from the training distribution)

Comparing ($A \rightarrow B$ vs $B \rightarrow A$)

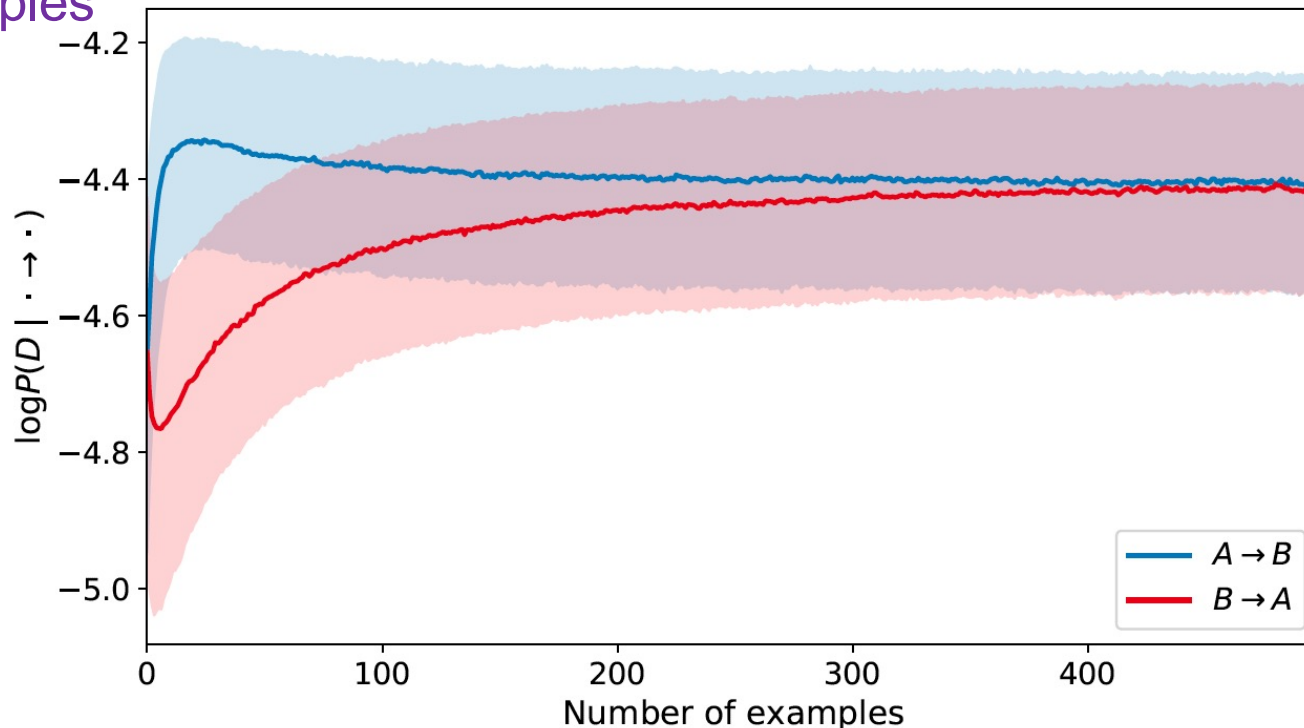
- How fast models adapt on a transfer distribution after being trained on training distribution
 - Simple SGD is used
- Let $A \rightarrow B$ be the correct causal model.
- Consider change between the two distributions is a random change in parameters of the true $P(A)$ for the cause A (because this will have an impact on the effect B , which can be picked up and reveal the causal direction)

Adaptation to transfer distribution

As more transfer distribution examples are seen by the learner (horizontal axis), in terms of the log-likelihood on the transfer distribution (on a large test set from the transfer distribution, tested after each update of the parameters)

$N = 10$. Curves are the median over 10,000 runs.

Correct causal model adapts faster. Most informative part of trajectory is first 10-20 samples



Human-inspired biases

- Inductive biases for out of distribution (OOD) generalization
- Conscious processing helps humans deal with OOD settings
 - In novel situations, we call upon conscious attention to combine on-the-fly the appropriate pieces of knowledge, to reason with them and imagine solutions
 - Not follow habitual routines, think hard to solve problems

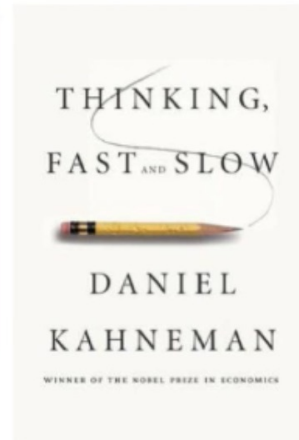


System 1 vs System 2 Cognition

2 systems (and categories of cognitive tasks):

System 1

- Intuitive, fast, **UNCONSCIOUS**, 1-step parallel, non-linguistic, habitual
- Implicit knowledge
- Current DL



System 2

- Slow, logical, **sequential**, **CONSCIOUS**, linguistic, algorithmic, planning, reasoning
- Explicit knowledge
- DL 2.0

Manipulates high-level / semantic concepts, which can be recombined combinatorially



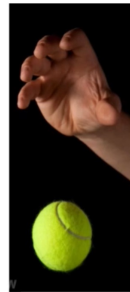
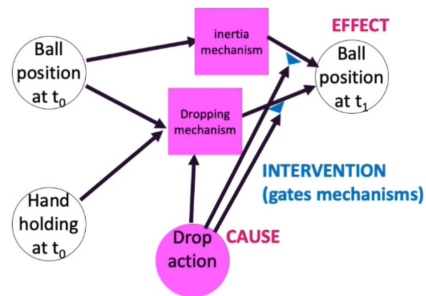
Reusable Causal Mechanisms

SPARSE DEPENDENCIES BETWEEN ABSTRACT VARIABLES

Also consistent with Baar's Global Workspace Theory (1997) of conscious processing.

Linguistic example:

"if I drop the ball, it will fall on the ground"



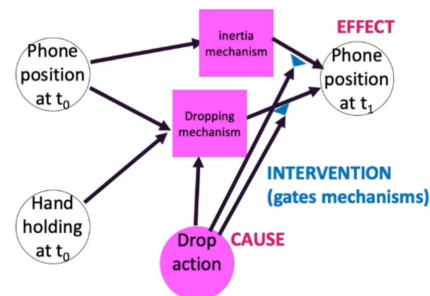
An abstract outcome can be predicted accurately from very few conditioning abstract variables

REUSABLE CAUSAL MECHANISMS

Linguistic example:

COUNTERFACTUAL

"if I had dropped the phone, it would have fallen on the ground"



The same mechanism can be reused on many instance tuples