

# Providing Clues to Discover Underlying Causes

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in Representation Learning

1. Greedy Layer-Wise Unsupervised Pretraining
2. Transfer Learning and Domain Adaptation
3. Semi-supervised Disentangling of Causal Factors
4. Distributed Representation
5. Exponential Gains from depth
6. Providing Clues to Discover Underlying Causes

# Topics in Clues to Causes

- Ideal representation
- Clues to find factors of variation
- Generic regularization strategies
- Summary of representation learning

# What is an ideal representation?

- Q: What makes a representation better than another?
- A: Ideal representation disentangles underlying *causal factors of variation* that generated the data, especially those relevant to application
- Most strategies introduce clues that help the learning to find underlying factors of variation

# Clues to find factors of variation

- Supervised learning provides a strong clue
  - Label  $y$  presented with each  $x$  specifies one factor of variation directly
- To make use of unlabeled data, representation learning makes use of other less direct hints about underlying factors
  - From designers of learning algo to guide learner
- Goal of deep learning: find regularization strategies applicable to many AI tasks that people solve

# Generic Regularization Strategies

- To encourage learning algorithms to discover features that correspond to factors of variation

- |                                                             |                                       |
|-------------------------------------------------------------|---------------------------------------|
| 1. Smoothness                                               | 6. Shared factors across tasks        |
| 2. Linearity                                                | 7. Manifolds                          |
| 3. Multiple explanatory factors                             | 8. Natural Clustering                 |
| 4. Causal factors                                           | 9. Temporal/Spatial Coherence         |
| 5. Depth (hierarchical organization of explanatory factors) | 10. Sparsity                          |
|                                                             | 11. Simplicity of factor dependencies |

# 1. Smoothness

- Assume  $f(\mathbf{x} + \varepsilon \mathbf{d}) \approx f(\mathbf{x})$  for unit  $\mathbf{d}$  and small  $\varepsilon$
- Allows generalizing from training points to nearby points
- Many ML algorithms use this idea
- But it is insufficient to overcome curse of dimensionality

## 2. Linearity

- Assume linear relationships between variables
- Allows predictions far from observed data
  - But sometimes leads to extreme predictions

# 3. Multiple Explanatory Factors

- Motivated by assumption that data is generated by multiple underlying explanatory factors
- Most tasks can be solved easily given the state of each of these factors
- Learning  $p(\mathbf{x})$  requires some of the same features that are useful for modeling  $p(y|\mathbf{x})$  because both refer to the same underlying explanatory factors
  - This view motivates the use of distributed representations, with separate directions in representation space corresponding to separate factors of variation

## 4. Causal Factors

- Model treats factors of variation described by  $h$  as the causes of  $x$  and not vice versa
- This is useful for
  - Semi-supervised learning
  - Learned model more robust when distribution over underlying causes changes
  - When model used for a new task

## 5. Depth

- By using hierarchy of explanatory factors, high-level concepts defined in terms of simple concepts
- Another viewpoint: task should be accomplished via a multi-step program
  - With each step referring back to the output of the processing accomplished via previous steps

## 6. Shared factors across new tasks

- When we have many tasks corresponding to:
  - Different  $y_i$  variables sharing the same input  $\mathbf{x}$ , or
  - When each task is associated with a subset or a function  $f^{(i)}(\mathbf{x})$  of a global input  $\mathbf{x}$
- The assumption is that each each  $y_i$  is associated with a different subset from a common pool of relevant factors  $\mathbf{h}$
- Because these subsets overlap, learning all the  $P(y_i|\mathbf{x})$  via a shared intermediate representation  $P(\mathbf{h}|\mathbf{x})$  allows sharing of statistical strength between the tasks

# 7. Manifolds

- Probability mass concentrates
  - In locally connected regions
  - Occupy a tiny volume
- In continuous cases, regions characterized by low-dimensional manifolds
- Many ML algorithms behave sensibly only on this volume
- Many ML algorithms, especially autoencoders attempt to explicitly learn the structure of manifolds

## 8. Natural Clustering

- Many ML algorithms assume that each connected manifold in the input space may be assigned to a single class
  - But class data may lie on may lie on disconnected manifolds
- This motivates many learning algorithms
  - Tangent propagation, double backprop, manifold tangent classifier, adversarial training

# 9. Temporal and Spatial Coherence

- Most important explanatory factors are assumed to change slowly with time
- It is easier to predict the true underlying explanatory factors than to predict raw observations such as pixel values

# 10. Sparsity

- Most features not relevant to most inputs
  - No need to use a feature that detects elephant trunks when representing an image of a cat
- Reasonable to impose a prior that any feature that can be interpreted as “present” or “absent” should be absent most of the time

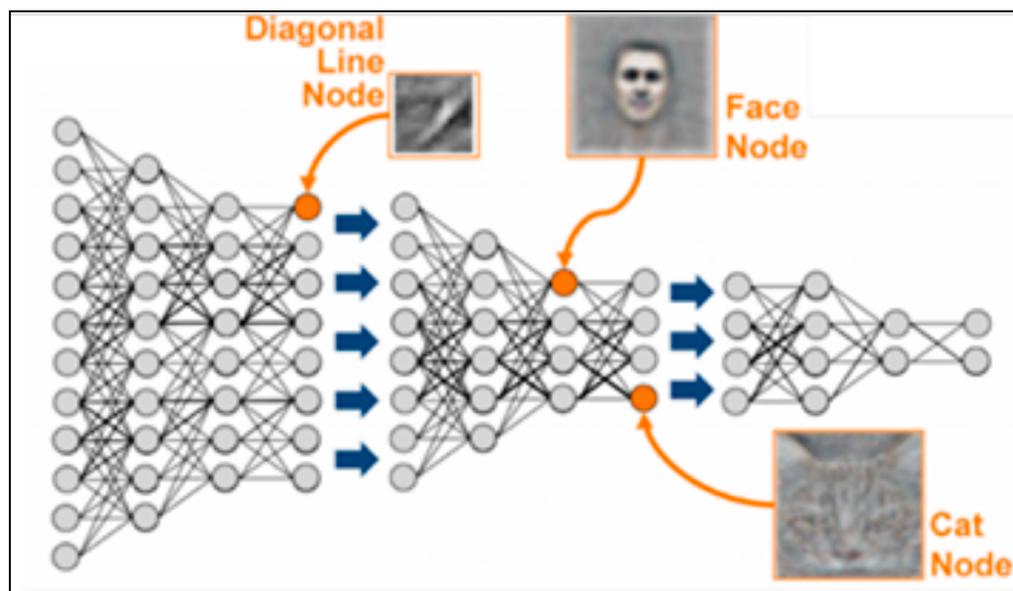
# 11. Simplicity of factor dependencies

- Factors related through simple dependencies
  - Simplest is marginal independence  $P(\mathbf{h}) = \prod_i P(h_i)$
  - Linear dependency also reasonable
    - Many laws of physics
    - Assumed when plugging a linear predictor of factorized prior

# Summary of representation learning

- Representation learning ties together all the many forms of deep learning
  - Feedforward/recurrent networks, autoencoders and deep PGMs all learn and exploit representations
- Learning the best possible representation remains an exciting avenue of research

# A representation



# Qualities of Representation

- Representations should have four qualities:
  1. Be as informative as the data
  2. Be invariant to nuisance factors, or “noise”
  3. Be disentangled
  4. Be simple and easy to work with
    - For representation optimized to task, find tradeoff between accuracy and complexity