

Distributed Representation

Sargur N. Srihari
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 Distributed Representation

- Distributed representation and neural networks
- Non-distributed representations
- Similarity induced by distributed representations
- Dealing with the curse of dimensionality
- Capacity in terms of VC Dimension
- Experimental validation of concepts
- A distributed representation that disentangles causal factors

Distributed Representation of Concepts

- Distributed representations of concepts are representations composed of many elements that can be set separately from each other
- They are one of the most important tools for representation learning
- They are powerful because they can use n features with k values each to describe k^n different concepts
 - With $k=2$, 2^n concepts

Distributed rep. and Causal factors

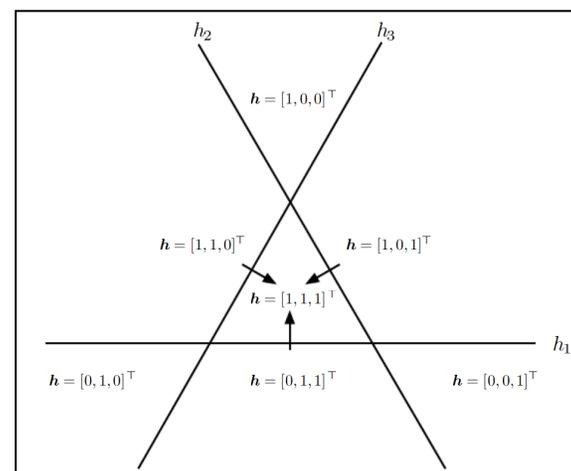
- Strategy of distributed representation used by
 - Neural networks with multiple hidden units
 - Probabilistic models with multiple latent variables
- Introduce an additional observation
 - Hidden units can learn to represent causal factors that explain the data
 - Distributed representation are natural for this approach
 - Each direction in representation space can correspond to the value of a different underlying configuration variable

Distributed Rep and Input Space

A learning algorithm based on a distributed representation breaks up the input space into regions

Three binary features h_1, h_2 and h_3 , each defined by thresholding the output of a learned linear transformation

Each feature divides R^2 into two half planes
 Let h_i^+ be the set of inputs for which $h_i=1$
 Each line represents the decision boundary
 For one h_i with the corresponding arrow pointing to the h_i^+ side of the boundary



Neural net uses distributed rep

- d inputs
- Extract n binary latent features (hidden units) by thresholding linear functions of input

$$d=2, n=3$$

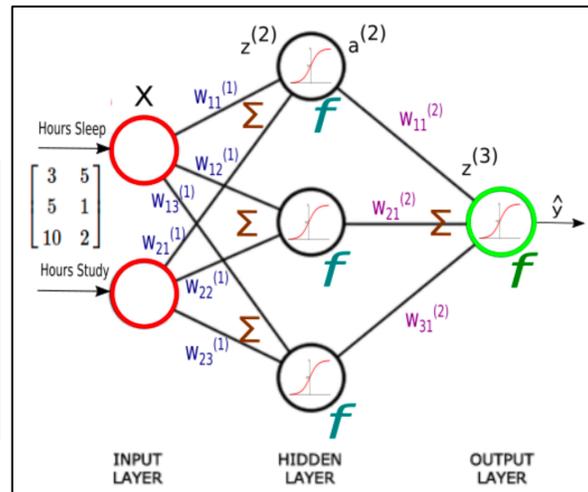
$$nd=6, n^d=9$$

$$z^{(2)} = XW^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = a^{(2)}W^{(2)}$$

$$\hat{y} = f(z^{(3)})$$



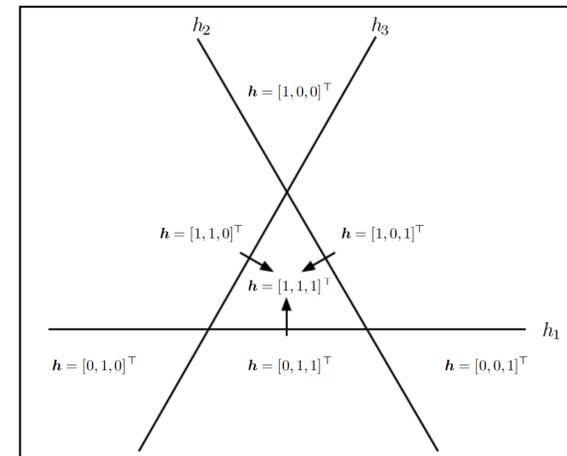
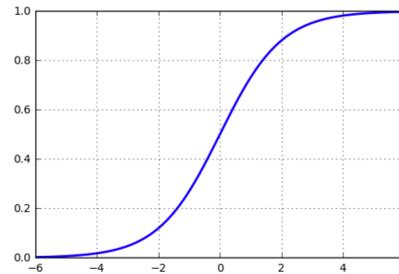
Output is test score given how many hours we sleep and how many hours we study. We pass in our input data (X) and get real outputs (\hat{y})

Activation function f is sigmoid

$$z^{(2)} = XW^{(1)}$$

$$= \begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix} \begin{bmatrix} W_{11}^{(1)} & W_{12}^{(1)} & W_{13}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} & W_{23}^{(1)} \end{bmatrix}$$

$$= \begin{bmatrix} 3W_{11}^{(1)} + 5W_{21}^{(1)} & 3W_{12}^{(1)} + 5W_{22}^{(1)} & 3W_{13}^{(1)} + 5W_{23}^{(1)} \\ 5W_{11}^{(1)} + W_{21}^{(1)} & 5W_{12}^{(1)} + W_{22}^{(1)} & 5W_{13}^{(1)} + W_{23}^{(1)} \\ 10W_{11}^{(1)} + 2W_{21}^{(1)} & 10W_{12}^{(1)} + 2W_{22}^{(1)} & 10W_{13}^{(1)} + 2W_{23}^{(1)} \end{bmatrix}$$

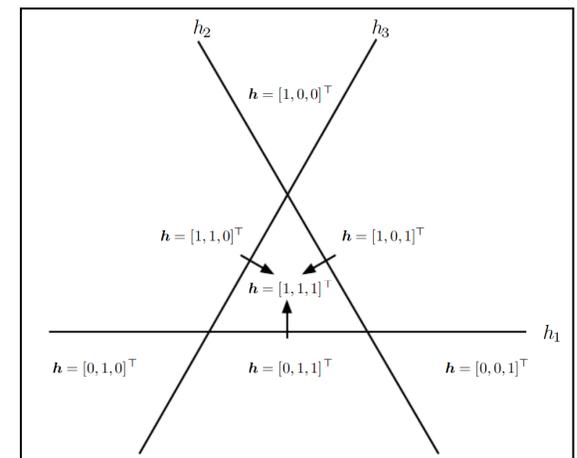


Partitioning of feature space by neural net

- d inputs
- Extract n binary latent features (hidden units) by thresholding linear functions of input
 - Each feature divides input space R^d into pair of half-spaces
 - By applying a general result concerning the intersection of hyperplanes (Zaslavsky, 1975), one can show
 - No. of regions distinguished

$$\sum_{j=0}^d \binom{n}{j} = n^d$$

$$d=2, n=3 \\ nd=6, n^d=9$$



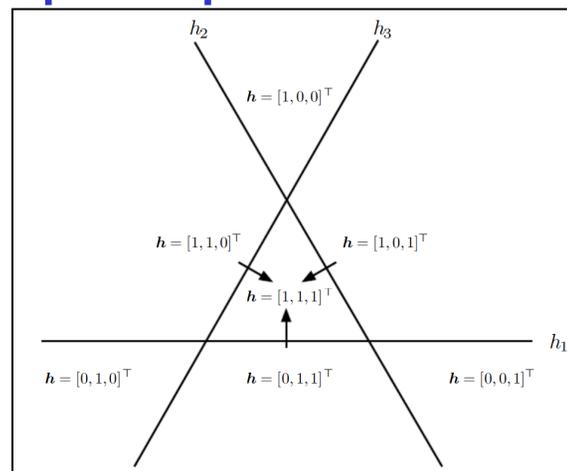
- Growth is exponential in input size d and polynomial in features n
- With $O(nd)$ parameters we can represent $O(n^d)$ regions in input space

Distributed/Non-distributed Reps.

A learning algorithm based on a distributed representation breaks up the input space into regions

Three binary features h_1, h_2 and h_3 , each defined by thresholding the output of a learned linear transformation

Each feature divides R^2 into two half planes
 Let h_i^+ be the set of inputs for which $h_i=1$
 Each line represents the decision boundary
 For one h_i with the corresponding arrow pointing to the h_i^+ side of the boundary



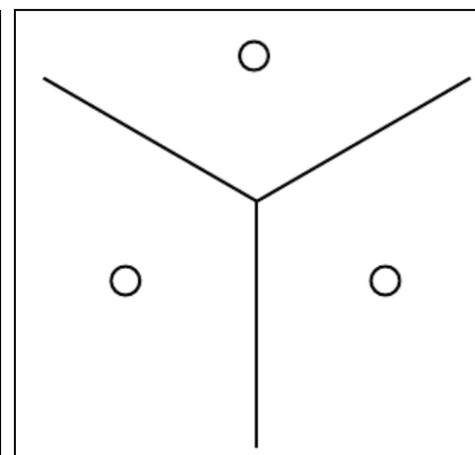
Nearest-neighbor is a non-distributed representation.

There is a separate set of parameters for each region

Advantage: given enough parameters it can fit the training set well without solving a difficult optimization problem.

Because it is straightforward to choose a different output independently for each region

Disadvantage: Generalizes only locally via the smoothness prior making it difficult to learn a complicated function with more peaks and troughs than the available no. of examples



Comparison of no of regions

1. Distributed representation

– Ex: extract n binary features (hidden units) by thresholding linear functions of input

- Each hidden unit divides input space R^d into a pair of half-spaces

– With n units no. of regions distinguished:

$$\sum_{j=0}^d \binom{n}{j} = 2^n$$

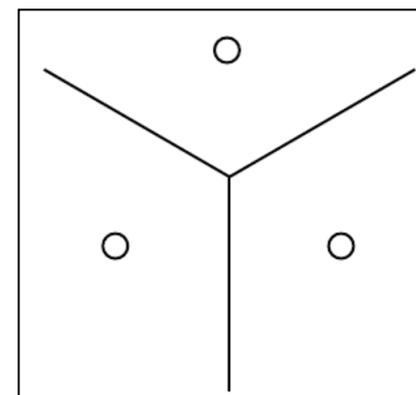
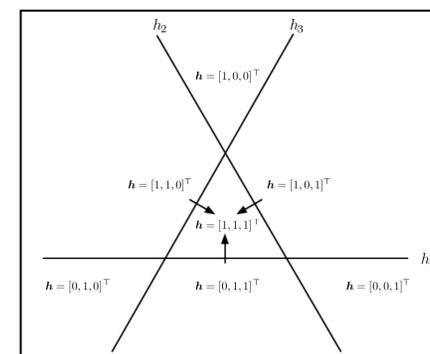
- With $O(nd)$ parameters we can represent $O(2^n)$ regions in input space

2. If we used a representation with a unique symbol for each region

– Specifying $O(2^n)$ regions would require $O(2^n)$ examples

- Fewer parameters means fewer samples needed

$d=2, n=3$



Non-distributed Representations

- Learning algorithms with non-distributed reps:
 - Clustering methods, including k -means
 - Each input point is assigned exactly one cluster
 - k nearest neighbor
 - One or a few prototypes are associated with a given input. For $k > 1$, multiple values describe each input, but not controlled separately
 - Decision trees
 - Only one leaf (and the nodes on the path from root to leaf) is activated when an input is given
 - GMMs
 - Templates are associated with a degree of activation. Each input represented by multiple values, but not controlled separately from each other
 - Kernel machines with a Gaussian kernel
 - Although the degree of activation of each support vector is continuous, same issue as with GMMs
 - Language or translation models based on n -grams
 - Set of contexts (sequences of symbols) partitioned according to a tree structure of suffixes. Separate parameters are estimated for each leaf of the tree

Linear relationship in non-distributed

- For some of these non-distributed algorithms, the output is not constant by parts but instead interpolates between neighboring regions.
- The relationship between the number of parameters (or examples) and the number of regions they can define remains linear.

Distributed Rep induces Similarity

- Concept distinguishing distributed vs. symbolic
 - In symbolic: cat and dog are as far from each other as any other two symbols
 - In distributed: generalization due to shared attributes between concepts
 - Things about cats generalize to cats
 - Entries “has_fur” or “no_of_legs” that have the same value
- Distributed reps induce a rich similarity space
 - Semantically close concepts are close in distance
 - Property absent in purely symbolic representations

Statistical Advantage of Distributed Rep

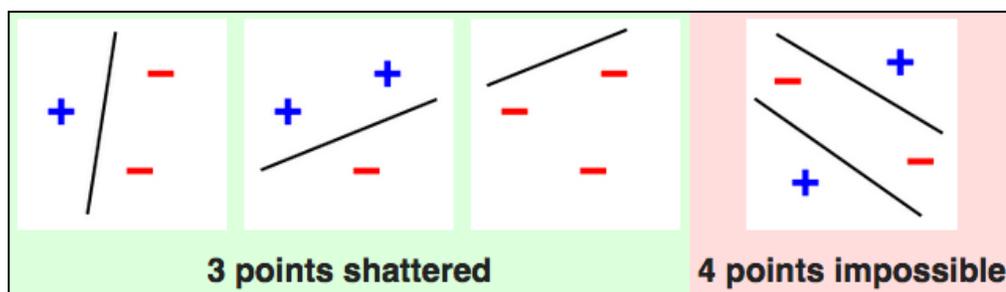
- Q: What is the advantage of using a distributed representation?
- A: When a complicated structure is to be compactly represented by a small no. of parameters

Symbolic Algos and Dimensionality

- Symbolic algorithms assume smoothness
 - i.e., if $u \approx v$, learned f has property $f(u) \approx f(v)$
 - if $f(x) \approx y$, choose \hat{f} that changes little over input $x + \epsilon$
- Suffers from curse of dimensionality
 - Needs no of examples as large as no of regions
 - Each region is a category or symbol:
 - A separate degree of freedom for each symbol (or region)
 - Can learn an arbitrary mapping from symbol to value
 - However no generalization to new symbols for new regions
- If lucky, regularity in f , besides smoothness
 - E.g., CNN with max-pool invariant to translation not corresponding to smooth transformations

Capacity of Distributed Rep

- Distributed representation can generalize well
 - because capacity remains limited despite being able to encode so many regions
- VC dimension of a model f : max. no. of points that can be arranged so that f shatters them



- VC Dimension of Linear Classifier in R^d is: $d+1$
- VC Dimension of Distributed Rep is only: $O(w \log w)$
 - where w is the no. of weights

Explaining low VC Dimension

- VCD of Distributed Rep is only: $O(w \log w)$
- Limitation because:
 - while we can assign many unique codes to representation space, we cannot use all of code space, nor learn arbitrary functions mapping from the representation space h to the output y using a linear classifier
- A distributed representation to learn h combined with a linear classifier to determine y
 - Expresses belief that classes are linearly separable of the underlying causal factors captured by h

Justifying linear classification

- We typically want to learn concepts such as
 - set of all green objects or set of all cars
- But not categories that require nonlinear XOR logic
 - set of all red cars and green trucks as one class and



- set of all green cars and red trucks as another class



Experimental validation of concepts

- Ideas discussed so far experimentally validated
- Hidden units in a deep CNN trained on ImageNet and Places benchmark datasets
 - Learn features that are interpretable
 - Correspond to a label that humans would assign

Hidden units learn interpretable features

- Not always that hidden units learn something with a linguistic name
- But emerge near top levels of deep networks
- They tend to be concepts we could learn about each without having to see configurations of all the others

Ex: a generative model for faces

- A generative model that has learned a distributed representation with directions capturing underlying factors of variation
- Images of faces:
 - Disentangles concepts of
 - Person is male or female
 - Wearing glasses



Learnt causal factors

- One direction in representation space is gender, another is whether wearing glasses
- Features discovered automatically and not discovered a priori
- No need to have labels for hidden units



Concept of
man with
glasses



Concept
of man without
glasses



Concept
of woman
without
glasses

We can learn about the distinction between male and female, or about the presence or absence of glasses, without having to characterize all of the configurations of the $n - 1$ other features by examples covering all of these combinations of values. This form of statistical separability is what allows one to generalize to new configurations of a person's features that have never been seen during training.

Disentangling gender and glasses

- Begin with representation of concept of a man with glasses, then subtract the vector representing concept of a man without glasses, and finally add vector representing concept of a woman without glasses, we obtain vector representing concept of a woman with glasses.
- Generative model correctly decodes all these representation vectors to images that may be recognized as belonging to the correct class

