

Deep Learning Research: Learning Representations

Sargur N. Srihari
srihari@buffalo.edu

Topics

- Importance of Representation
- Representation learning methods
 - Autoencoders
 - Convolutional networks
 - Recurrent networks
- Transfer learning
- Disentangling variables
- Explainable artificial intelligence

Representation in daily life



Was Texas (MDCCCXLV) before California (MDCCCL)?

Ans: Use **MeDiCaL XaVler** to convert dates into Arabic numerals

How much is \$1 million in Indian Rupees?

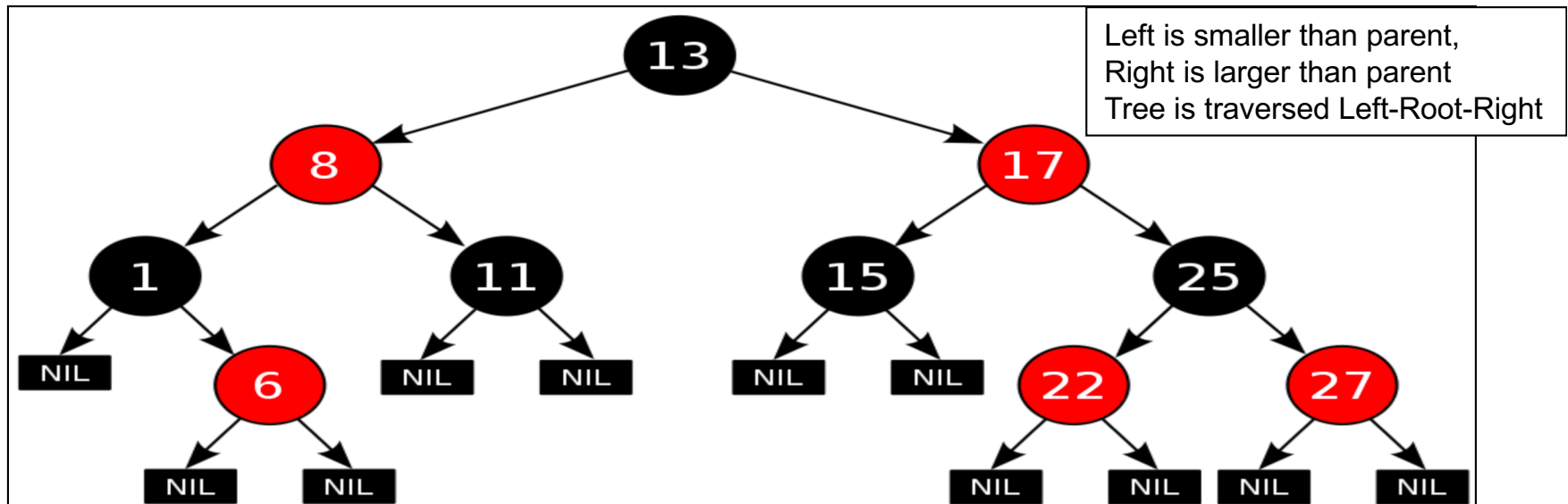
1 USD = 67.65 INR

1 Million = 1,000,000

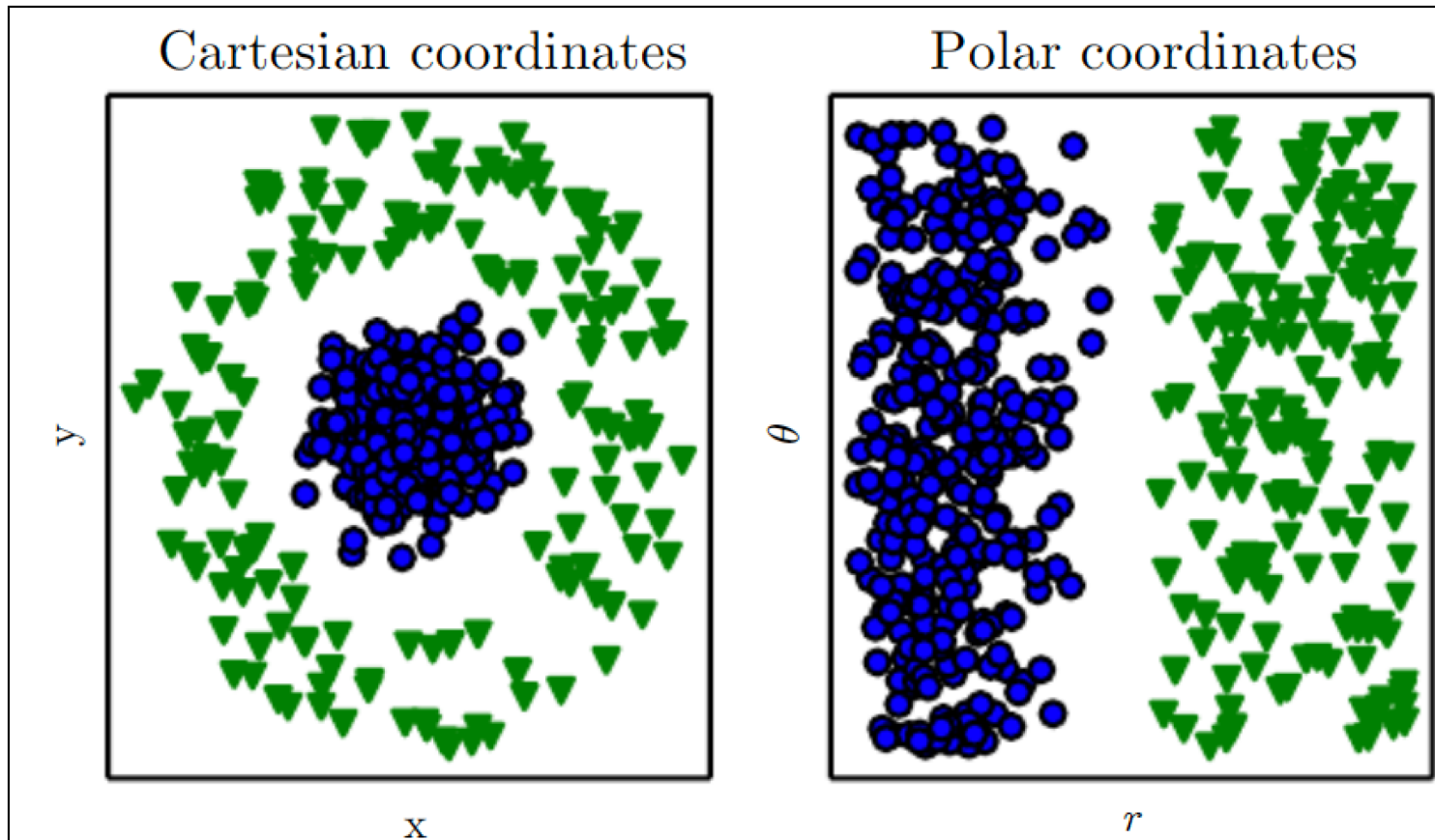
$67.65 \text{ INR} * 1,000,000 = 6,76,50,000 \text{ INR}$

Representation in Computer Science

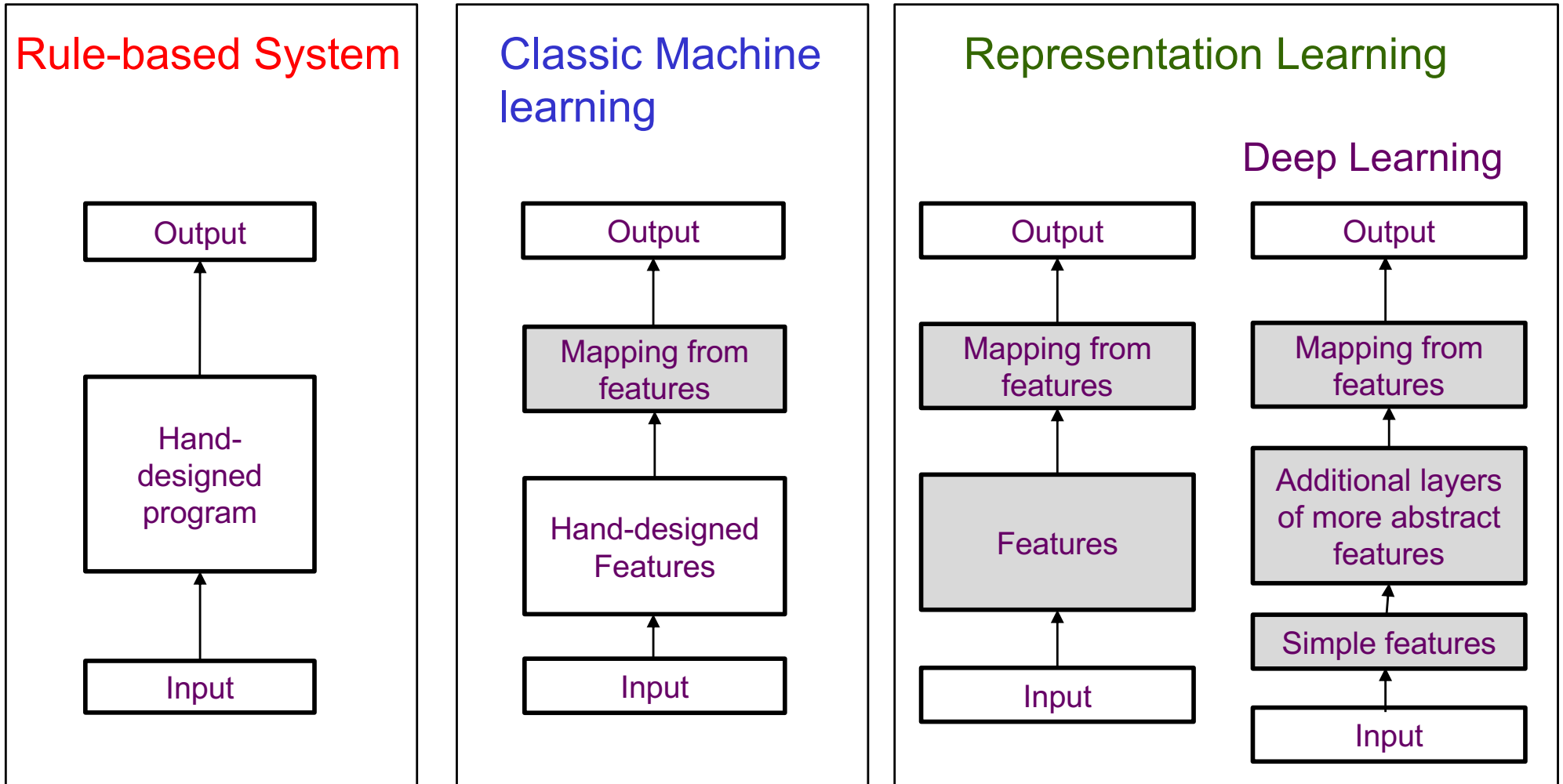
- Inserting a number in a sorted list is $O(n)$
 - But only $O(\log n)$ if list represented as a red-black tree



Representation in Machine Learning



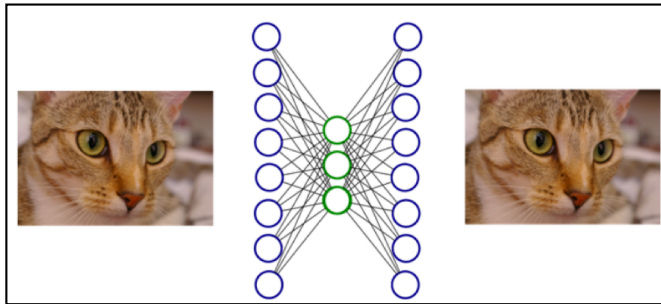
Deep Learning finds representations



■ Shaded boxes indicate components that can learn from data

Autoencoder Representation

- A network trained to copy its input to output

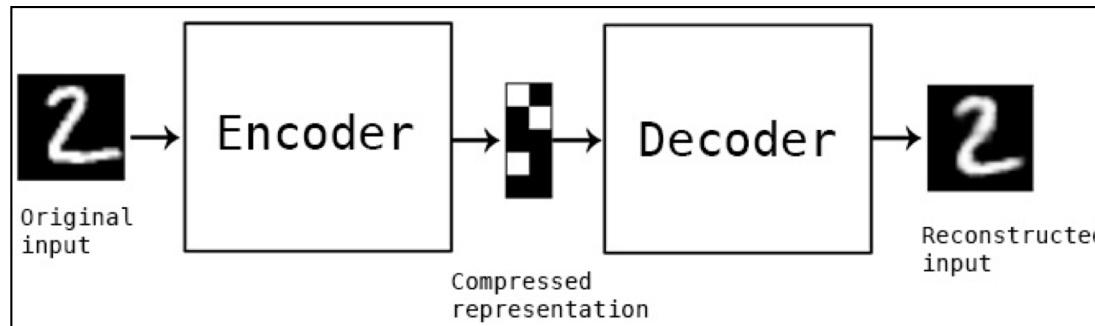


Encoder:

Converts input into a representation

Decoder:

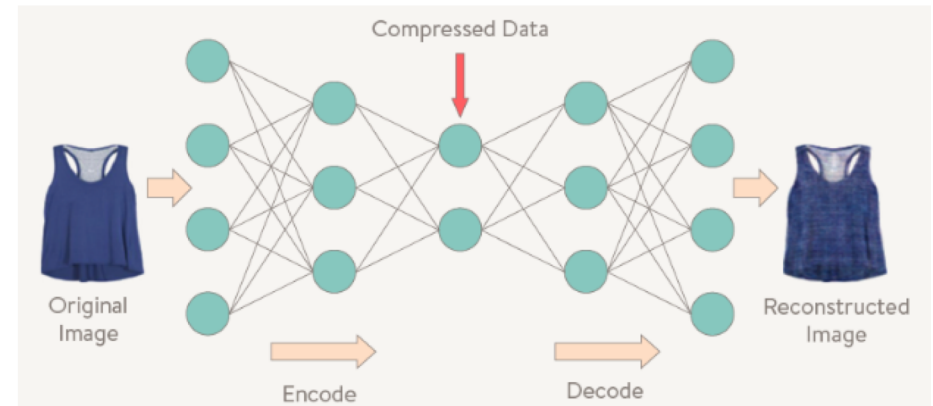
Converts representation back to input



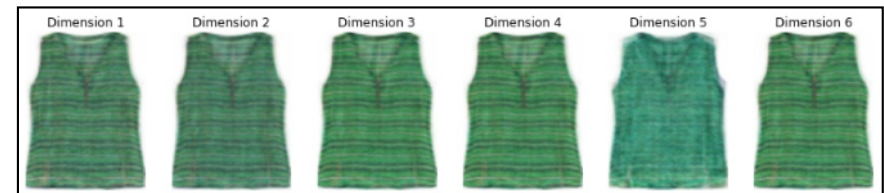
- Hidden layer describes code to represent input
- Autoencoders are data-specific
 - Differs from MP3 or JPEG compression
 - Autoencoders are lossy
 - Autoencoders are learnt

Autoencoder Example: Deepstyle

- Random points in encoded space are new images of styles not seen
 - Using one input image and changing values along different dimensions of feature space you can see how the generated image changes (patterning, color texture) in style space

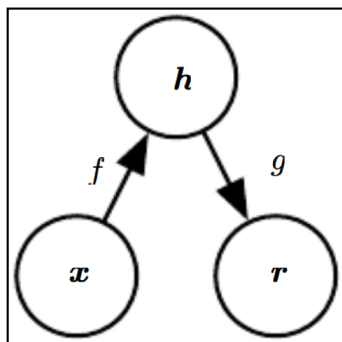


New designs from representation

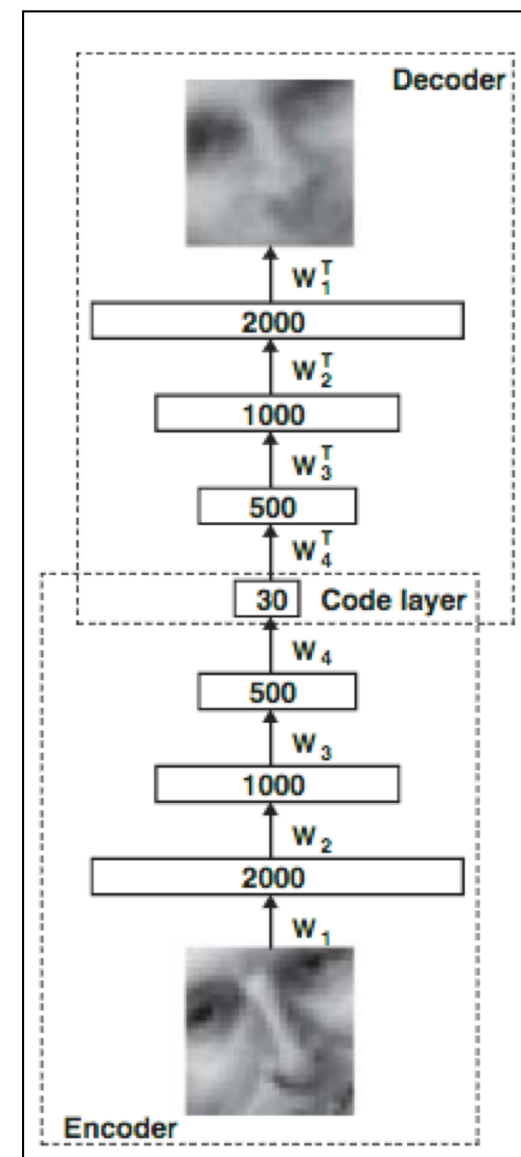


General structure of an autoencoder

- Maps input x to output r through code h



- The network has two parts
 - Encoder function $h=f(x)$
 - Decoder for reconstruction $r=g(h)$
- Model forced to prioritize which aspects of input are copied
 - Learns useful data properties



Autoencoder training using a loss function

- Encoder f and decoder g

$$f : X \rightarrow h$$

$$g : h \rightarrow X$$

$$\arg \min_{f,g} \|X - (f \circ g)X\|^2$$

$$h = \sigma_1(Wx + b)$$

$$x' = \sigma_2(W'h + b')$$

Trained to minimize reconstruction error

$$L(x, x') = \|x - x'\|^2 = \|x - \sigma_2(W^t(\sigma_1(Wx + b)) + b')\|^2$$

Provides a compressed representation of the input x

Regularized autoencoders properties:

Sparsity of representation

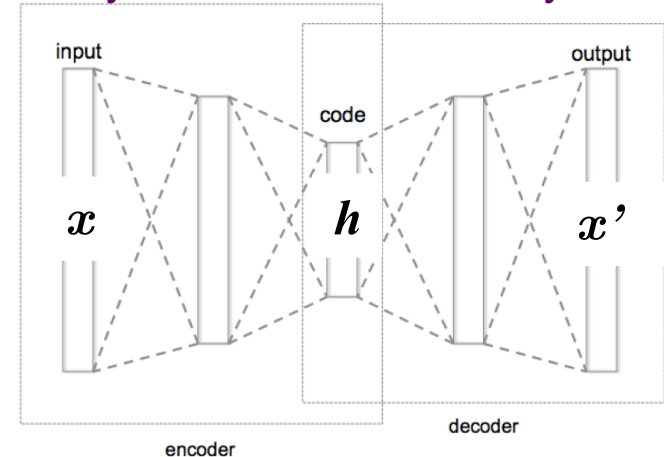
Smallness of the derivative of the representation

Robustness to noise

Robustness to missing inputs

Autoencoder:

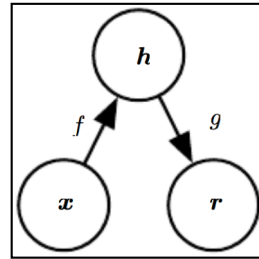
3 fully connected hidden layers



Regularization: Sparse Autoencoders

- Loss includes a sparsity penalty $\Omega(\mathbf{h})$ on code \mathbf{h}

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$



- Penalty limits the capacity of the encoder
- Together they force the the hidden representation to capture information about the data generating distribution
- Typically used to learn features for another task such as classification
 - Sparsity penalty can yield a model that has learned useful features as a byproduct

Autoencoders learn manifolds

- Autoencoders concentrate data around a low-dimensional manifold
 - Unusual behavior if input is off of the manifold
- MNIST with 784 pixels per digit

A 1-D manifold in 784-D space

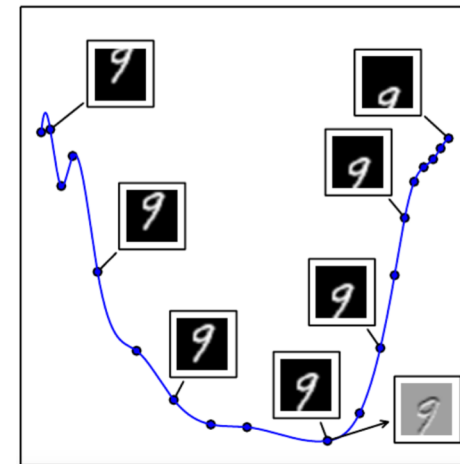
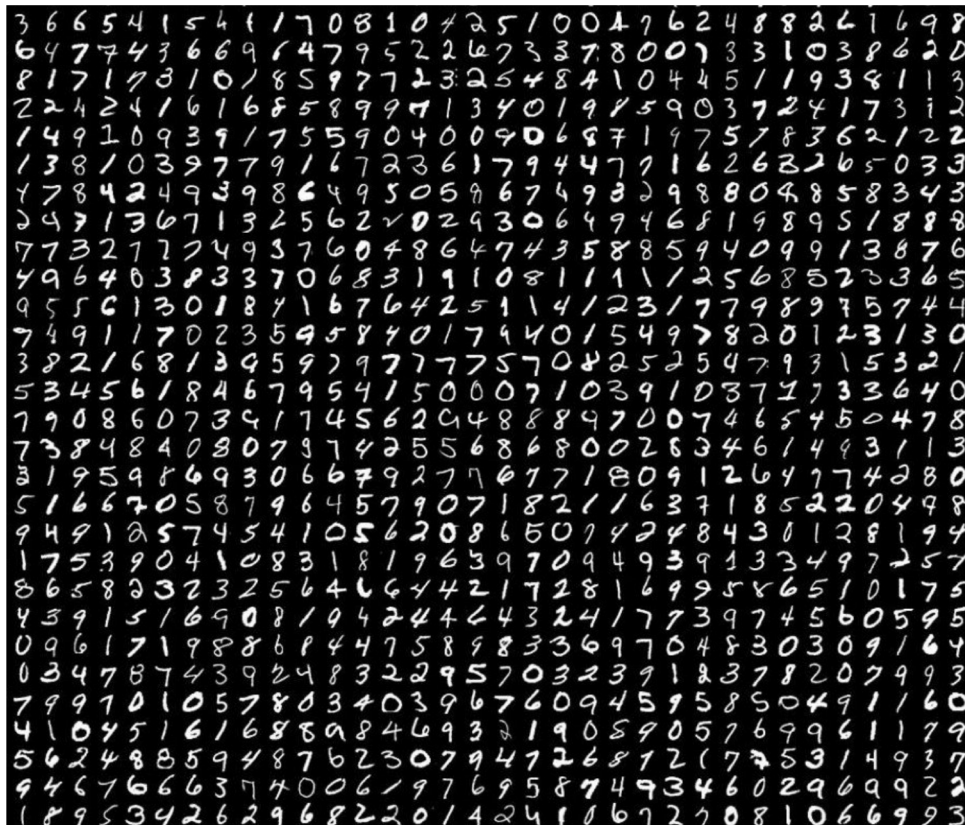
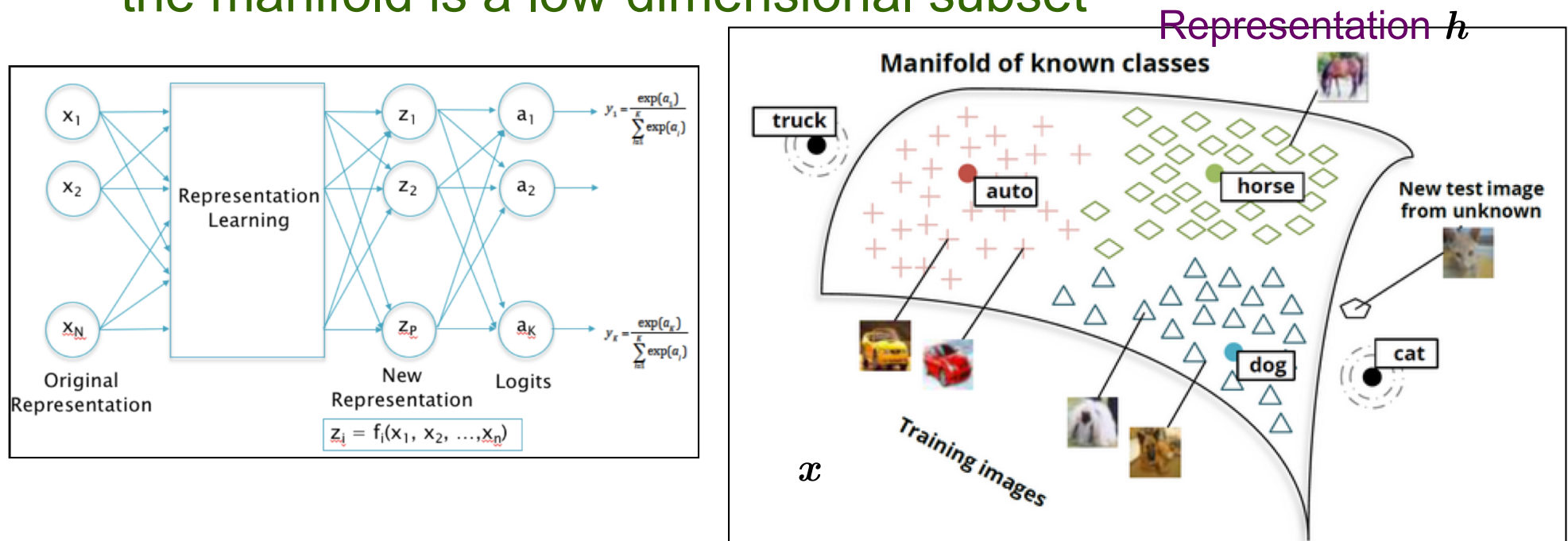


Figure is projection into 2-D space using PCA

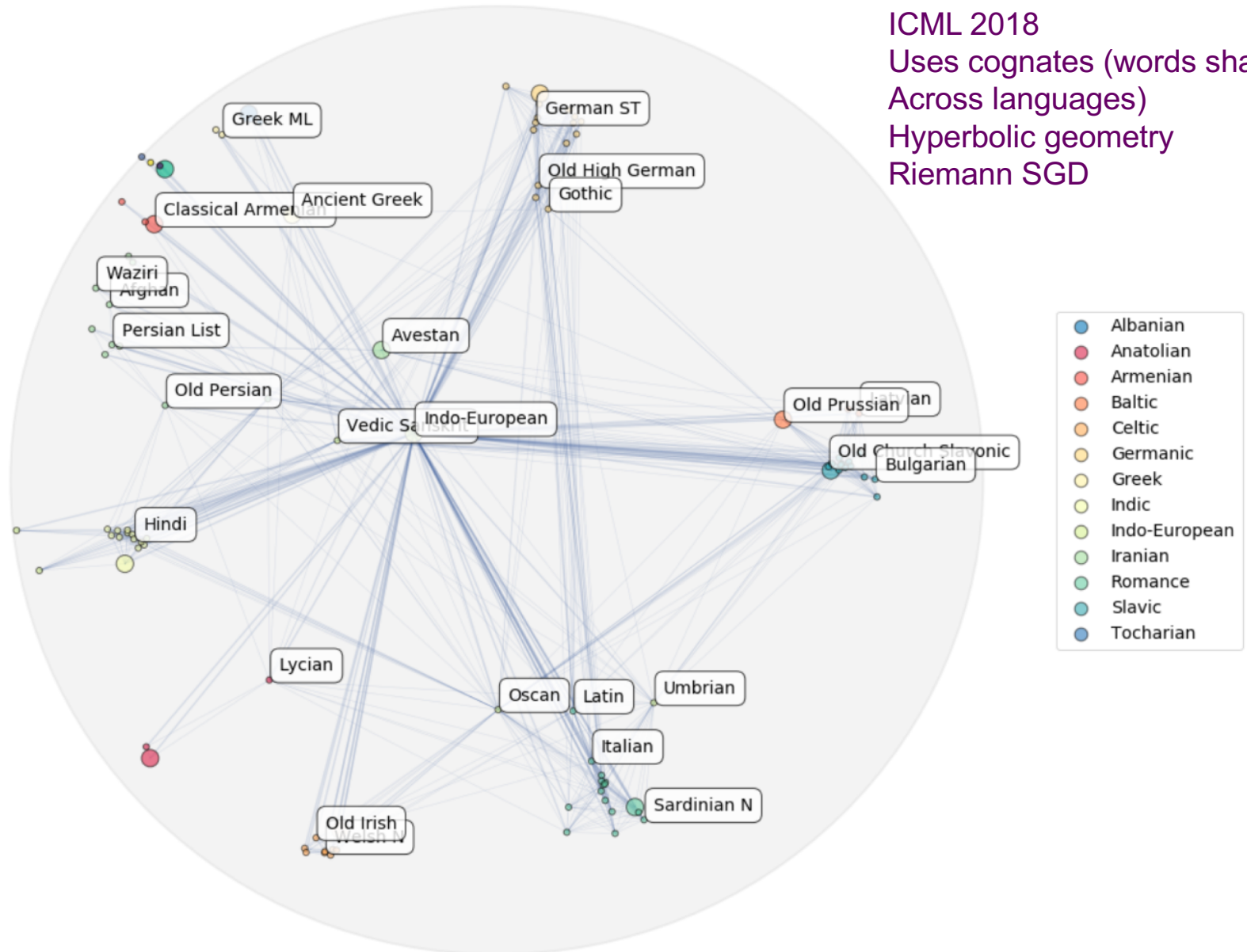
Embedding is a point on a manifold

- An embedding is a low-dimensional vector
 - With fewer dimensions than the ambient space of which the manifold is a low-dimensional subset



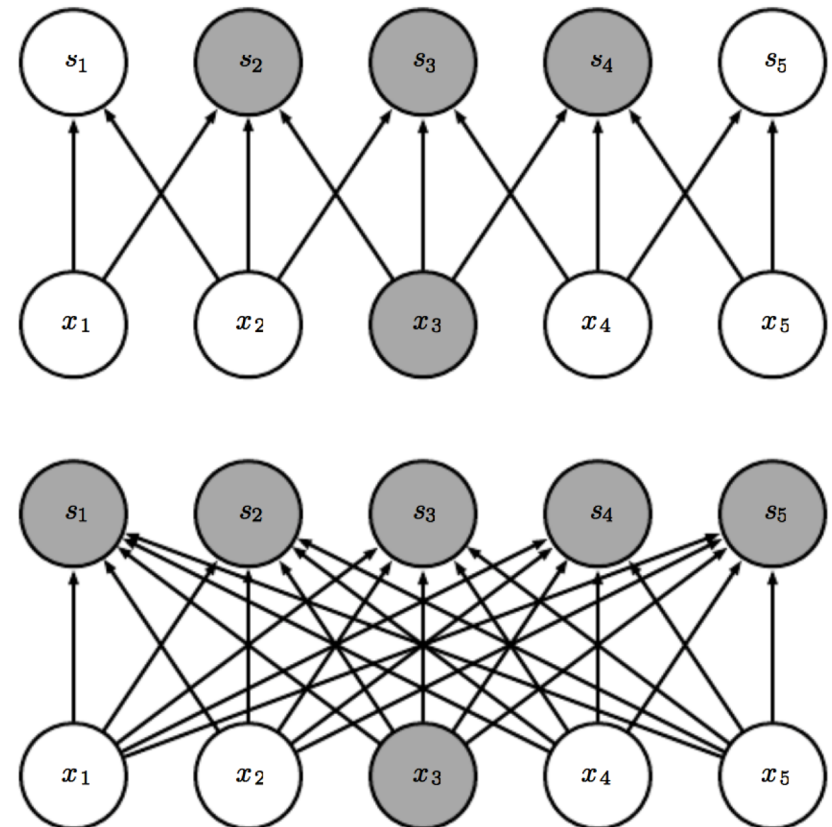
- Autoencoders learn a representation function that maps any point in ambient space to its embedding

Embedding Historical Linguistics Data



Sparse Connectivity, viewed from below

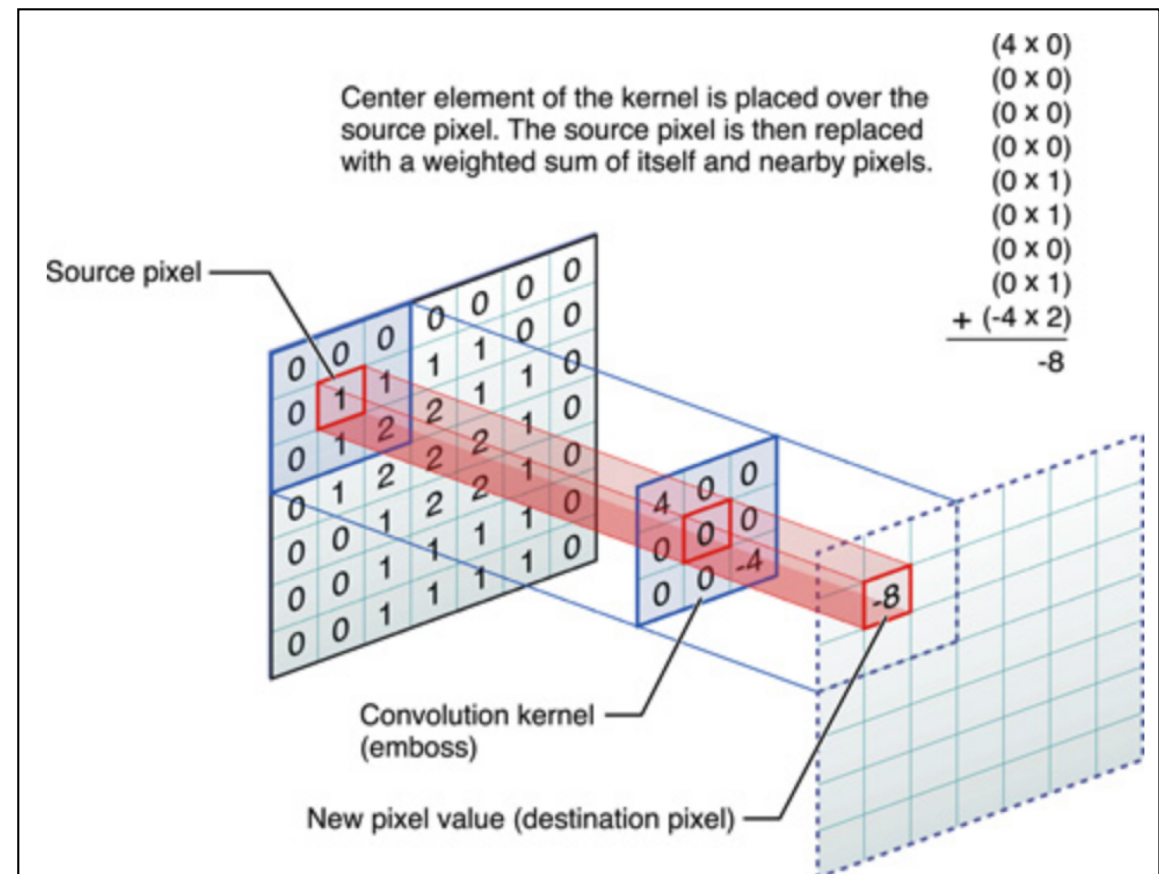
- Highlight one input x_3 and output units s affected by it
- Top: when s is formed by convolution with a kernel of width 3, only three outputs are affected by x_3
- Bottom: when s is formed by matrix multiplication



Convolutional Neural Network

- Input image may have *millions* of pixels,
- But we can detect edges with kernels of *hundreds* of pixels
- If we limit no of connections for each input to k
 - we need $k \times n$ parameters and $O(k \times n)$ runtime
- It is possible to get good performance with $k \ll n$

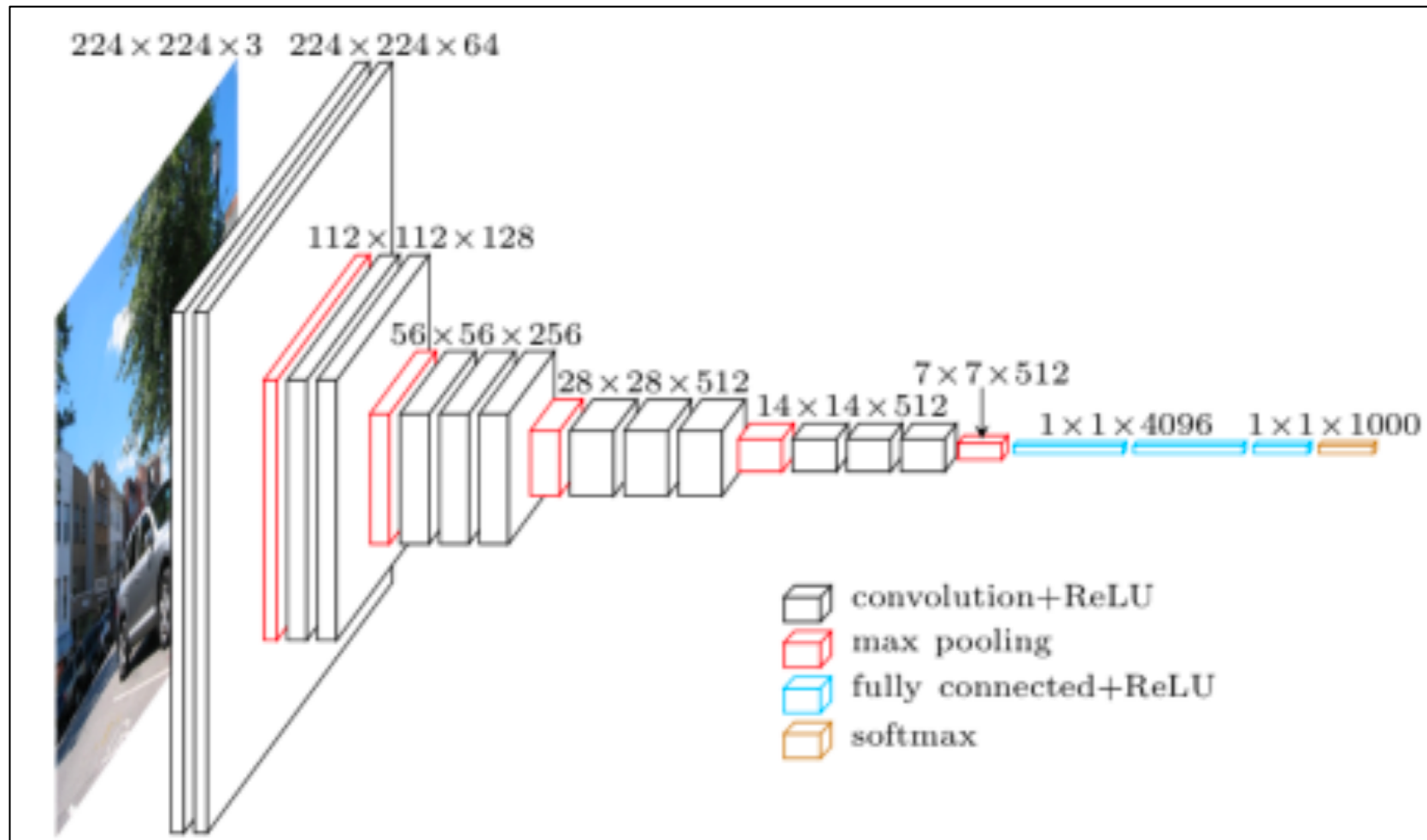
- Convolutional networks have sparse interactions
 - Accomplished by making the kernel smaller than the input
- Next slide shows graphical depiction



Types of Pooling functions

- A pooling function replaces the output of the net at a certain location with a summary statistic of the nearby inputs
- Popular pooling functions:
 1. *max pooling* operation reports the maximum output within a rectangular neighborhood
 2. Average of a rectangular neighborhood
 3. L^2 norm of a rectangular neighborhood
 4. Weighted average based on the distance from the central pixel

VGG 16



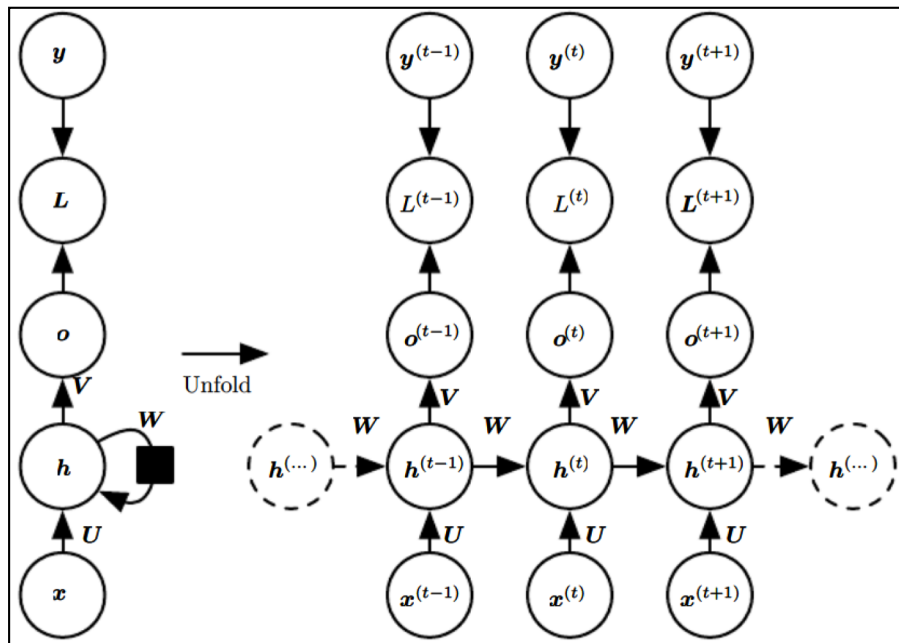
VGG is a convolutional neural network model

The model achieves 93% top-5 test accuracy in ImageNet

which is a dataset of over 14 million images belonging to 1000 classes.

Computational Graphs for RNNs

- Computational graphs include cycles
 - Cycles represent the influence of the present value of a variable on its own value at a future time step
 - RNN to map input sequence of x values to output sequence of o values
 - Loss L measures how far each output o is from the training target y



- Forward propagation is given as follows:

For each time step t , $t=1$ to $t=\tau$

Apply the following equations

$$o^{(t)} = c + Vh^{(t)}$$

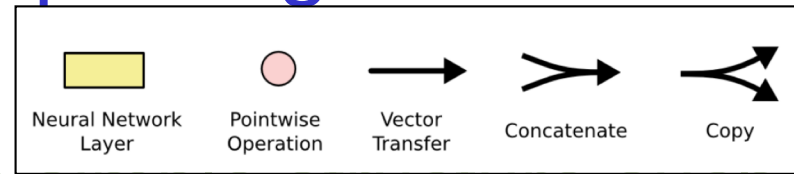
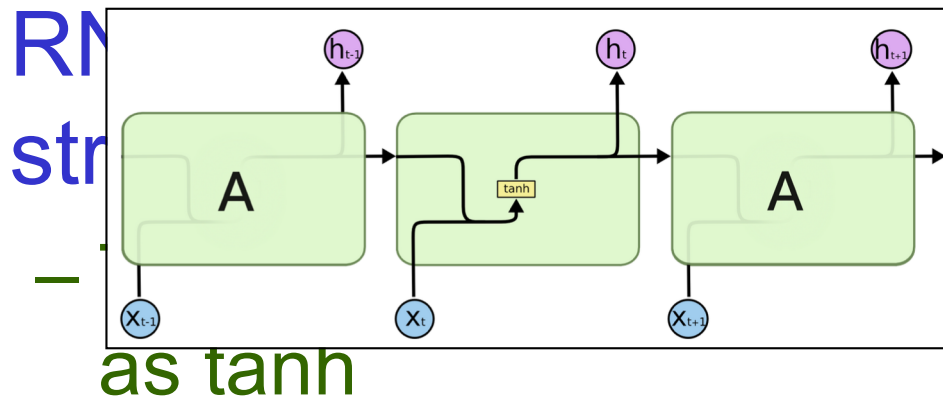
$$h^{(t)} = \tanh(a^{(t)})$$

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

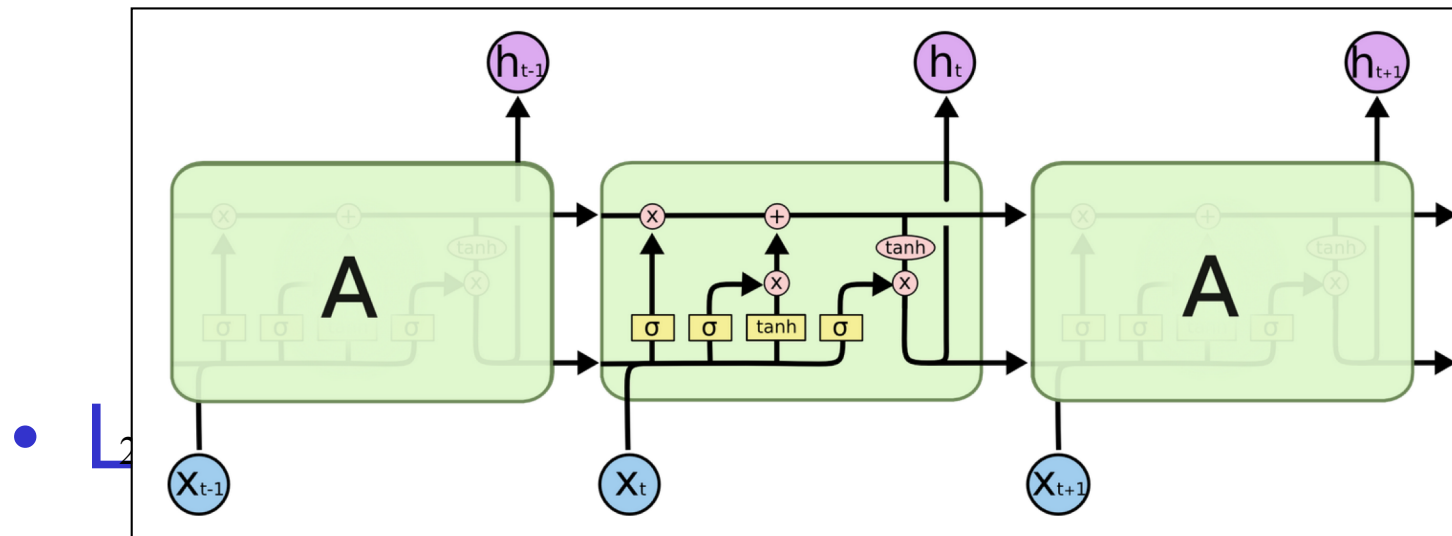
Long Short Term Memory

- Explicitly designed to avoid the long-term dependency problem

- RNN as a repeating chain



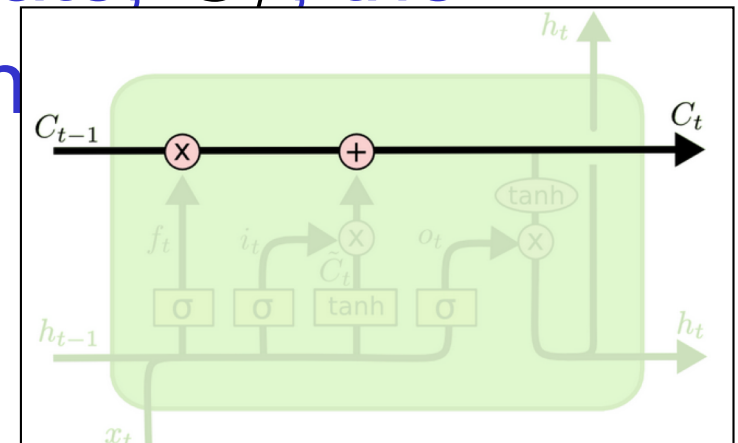
- as a simple structure such as tanh



- but the repeating module has a different structure

Core idea behind LSTM

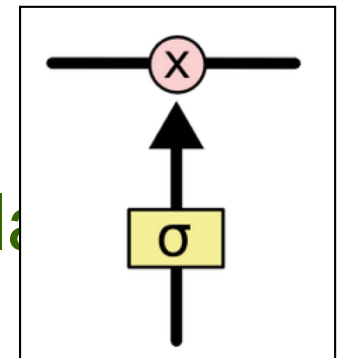
- The key to LSTM is the cell state, C_t , the horizontal line running through diagram



- Like a conveyor belt

- Runs through entire chain with minor interactions

- LSTM does have the ability to remove/add information to cell state regularly through structures called gates



- Gates are an optional way to let information through

Transfer Learning

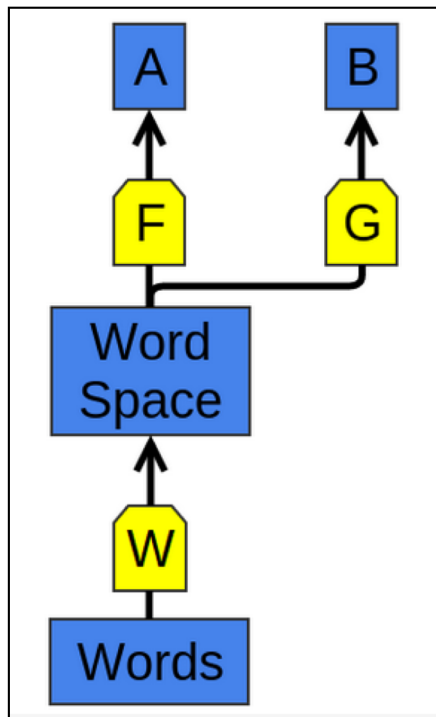
Learner must perform two or more different tasks

But we assume that the factors that explain the variations in A are relevant to the variations to be captured for learning B

Visual classification

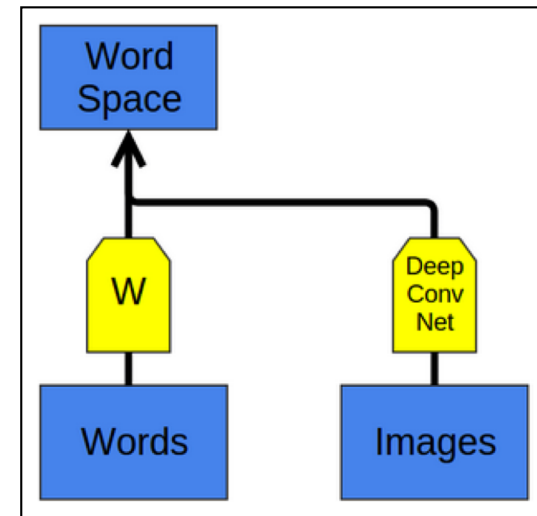
Significantly more data in distribution sampled from cats and dogs

Then learn to quickly generalize to ants and wasps



Shared Representation:

W and F are used to learn to perform task A
Later, G can learn to perform task B based on the representation of W



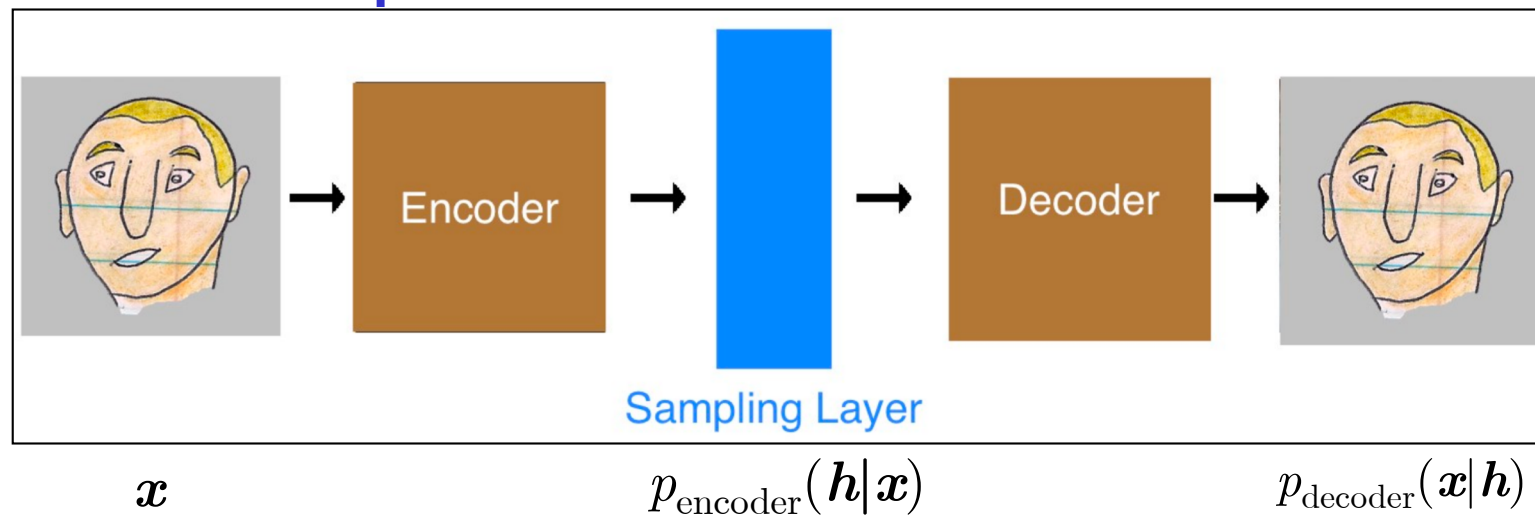
Embedding words and images in a single representation

Success of Transfer Learning

- Unsupervised deep learning for transfer learning has found success in ML competitions
 - Each participant is given data from distribution P_1 illustrating some set of categories
 - Participants learn a feature space
 - Mapping raw input to a representation space
 - This transformation is applied to samples from P_2
 - A linear classifier is trained from very few samples
- As deeper representations used (learned purely unsupervised from P_1) performance improves
 - For deeper representations fewer samples needed

Generative Models as Autoencoders

- Generative models with latent variables can be viewed as a particular form of autoencoder



- Generative modeling approaches which emphasize connection with autoencoder:
 - Variational autoencoder
 - Generative stochastic networks

What is a good representation for ML?

- Answer:
 - Representation makes subsequent learning easier
- Choice usually depends on subsequent learning task
- A feed-forward network $f(\mathbf{x}) = f^{(3)} [f^{(2)} [f^{(1)}(\mathbf{x})]]$ trained by supervised learning performs representation learning
 - Every hidden layer makes the classification easier
 - Ex: classes not linearly separable in input features may become linearly separable in the last layer

What is an ideal representation?

- One answer: ideal representation is one that disentangles the underlying *causal factors of variation* that generated the data, especially those that are relevant to our application
- Most strategies for representation learning are based on introducing clues that help the learning to find these underlying factors of variation

Qualities of Representations

1. They should be as informative as the data
2. They should be invariant, or unaffected by nuisance factors, or “noise”
3. They should be disentangled
4. Be as simple as possible and easy to work with
 - In other words, for a representation to be optimized to the task at hand, it should find the right tradeoff between accuracy and complexity.

Untangling

- Make (something complicated or confusing) easier to understand or deal with
 - Tangled inadvertently



Disentangle
→

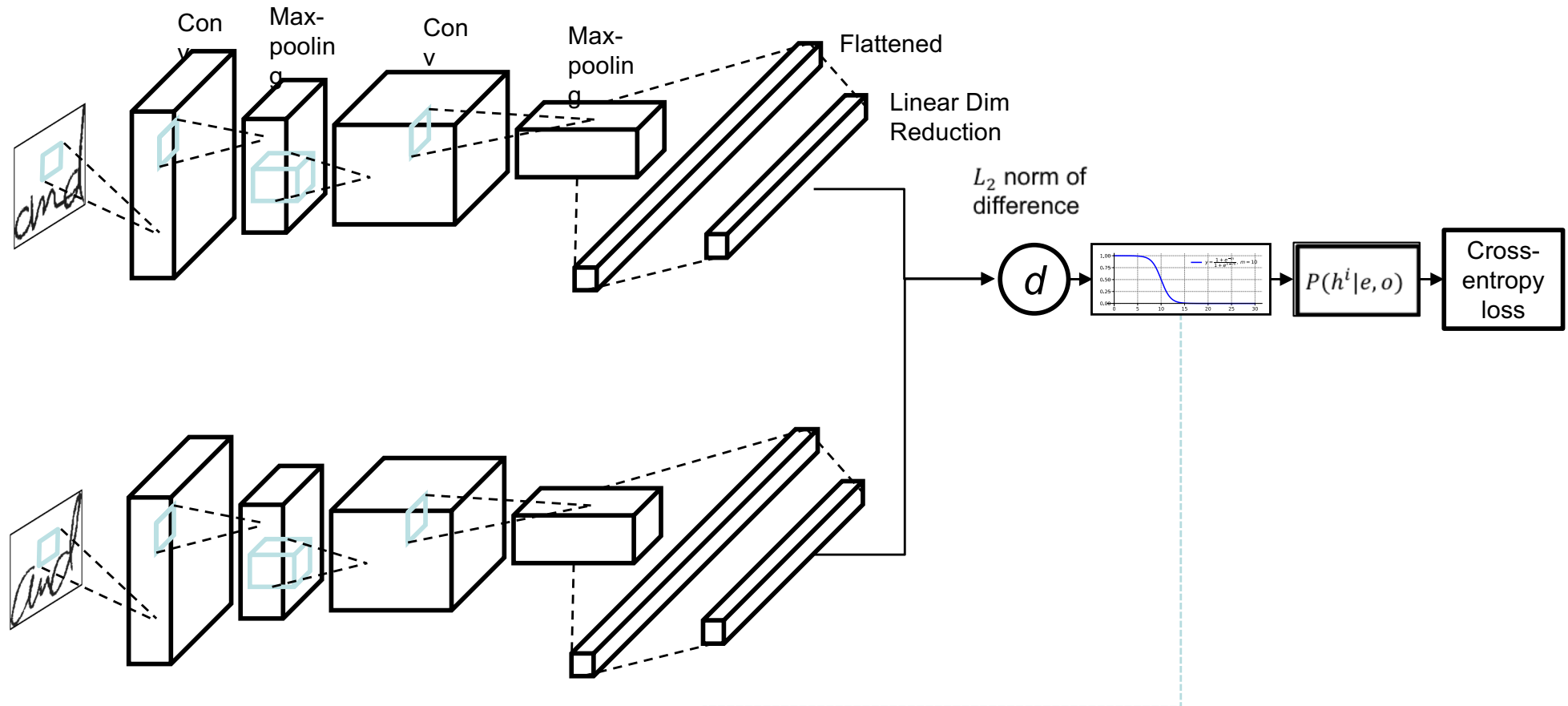


- Tangled purposefully

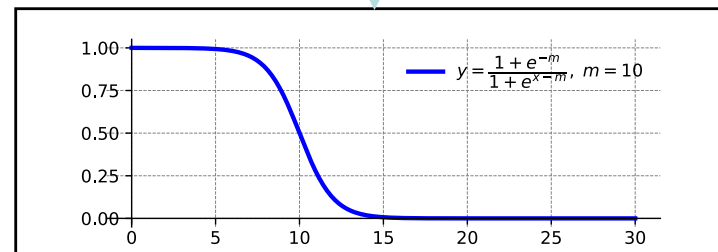
- “O what a tangled web we weave when we aim to deceive”: Walter Scott

- “He knew he'd have to untangle a financial mess and somehow restore investor confidence”

Deep Learning for HW Comparison

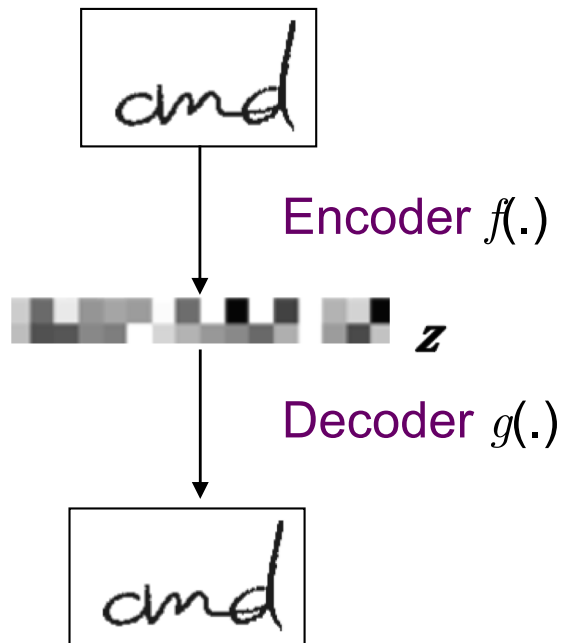


$$P(h^0 | o, e) = \frac{1 + e^{-m}}{1 + e^{d-m}}$$



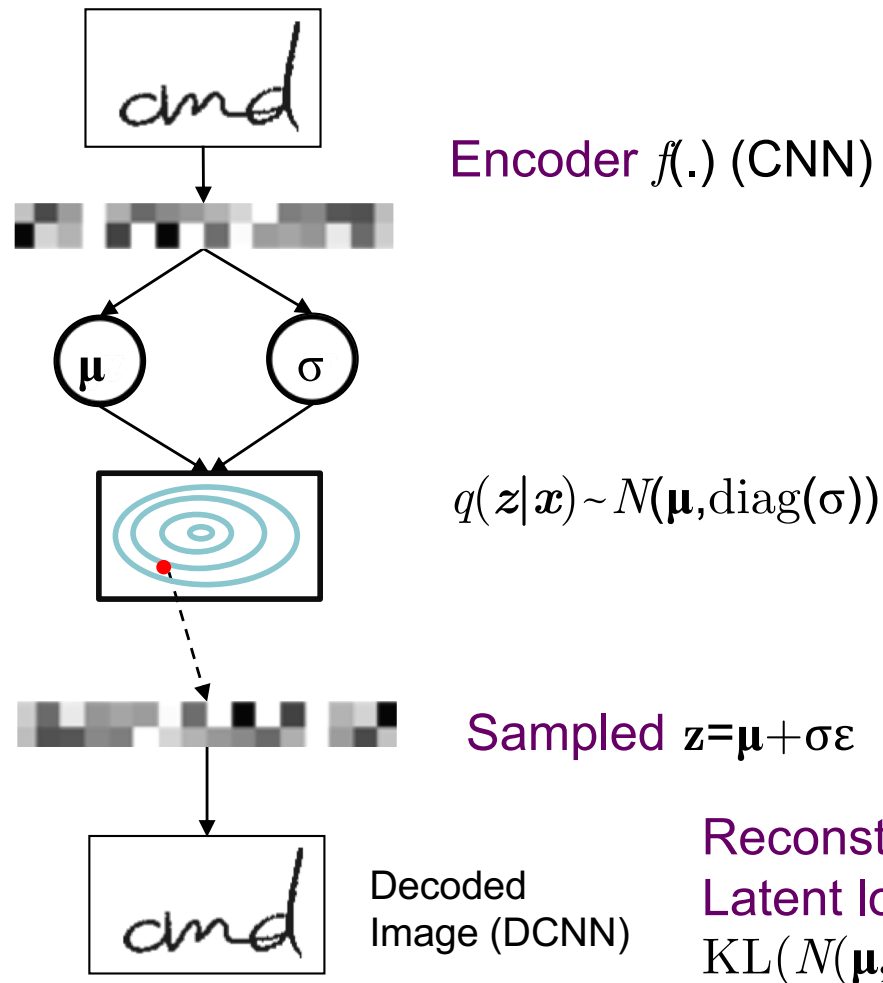
Unsupervised Representation Learning Using Generative Models

Autoencoder



Loss: Reconstruction loss

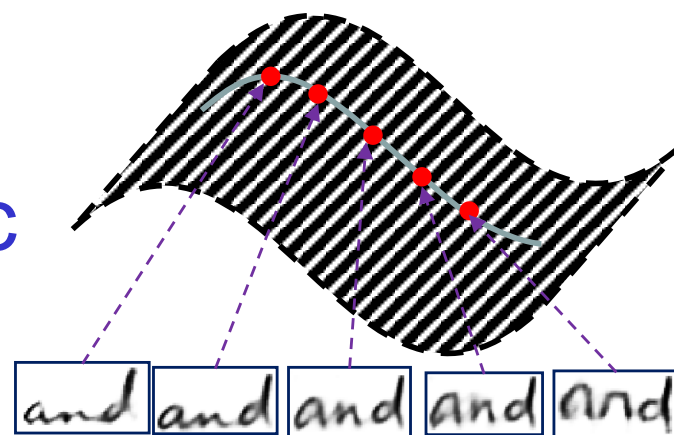
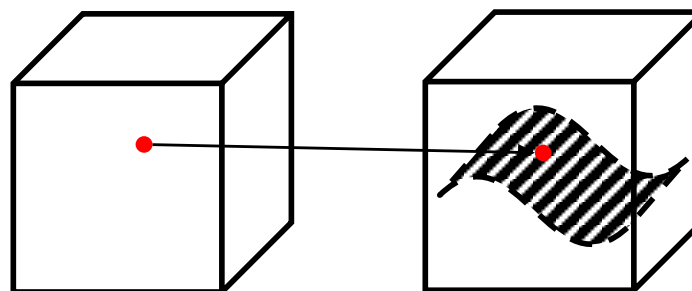
Variational Autoencoder



Reconstruction loss +
Latent loss:
 $\text{KL}(N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), N(\mathbf{0}, \mathbf{I}))$

Manifold Learning

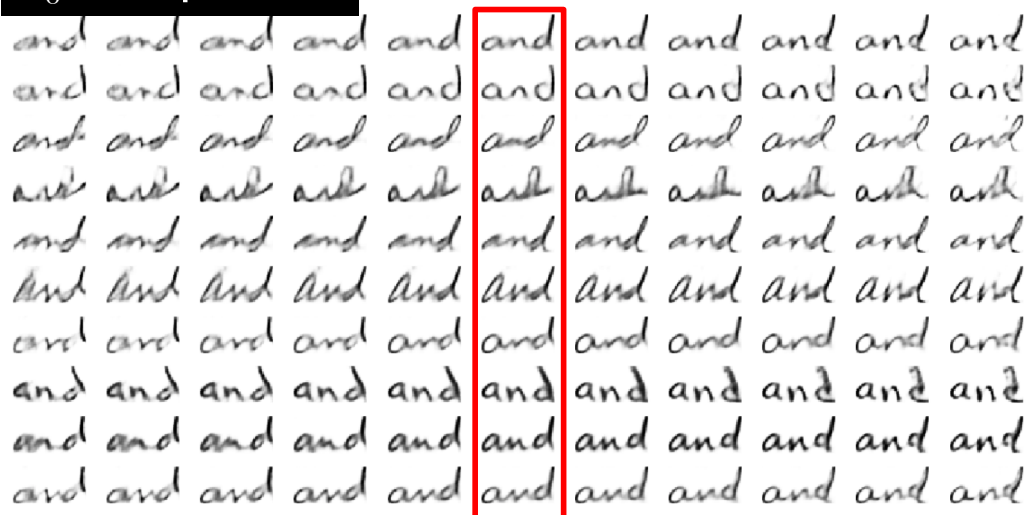
- Semantic manifold
 - Low-dim sub-space
- Interpretability
- Transferability
- Invariance
- Contribute to forensic comparison



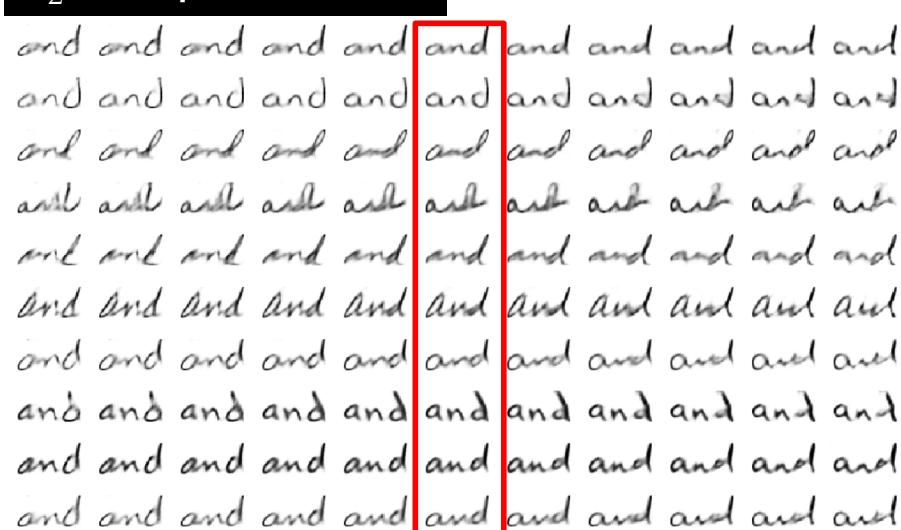
Images generated from manifold

Traversing along a dimension in latent space (keeping other dimensions fixed)

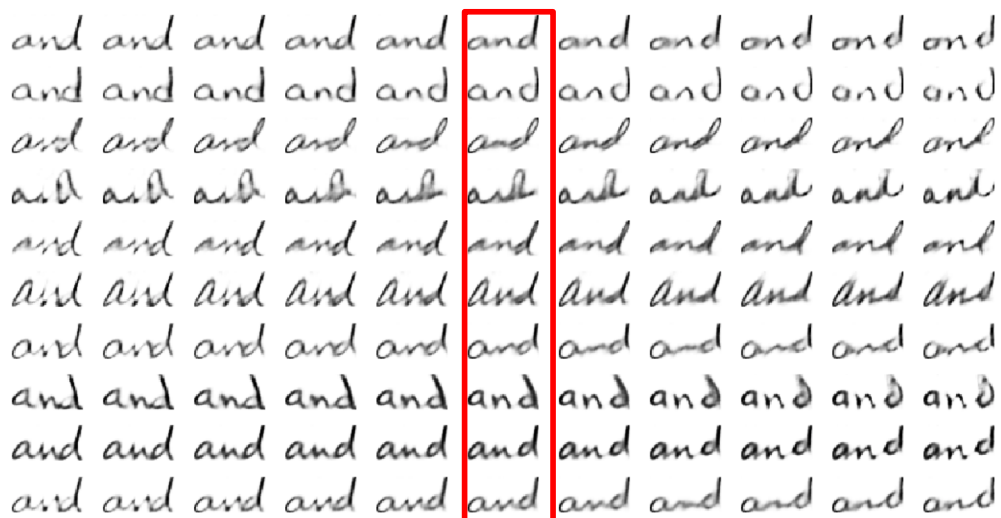
z_0 : Shape of d



z_2 : Loop size of d



z_1 : Vertical position of d



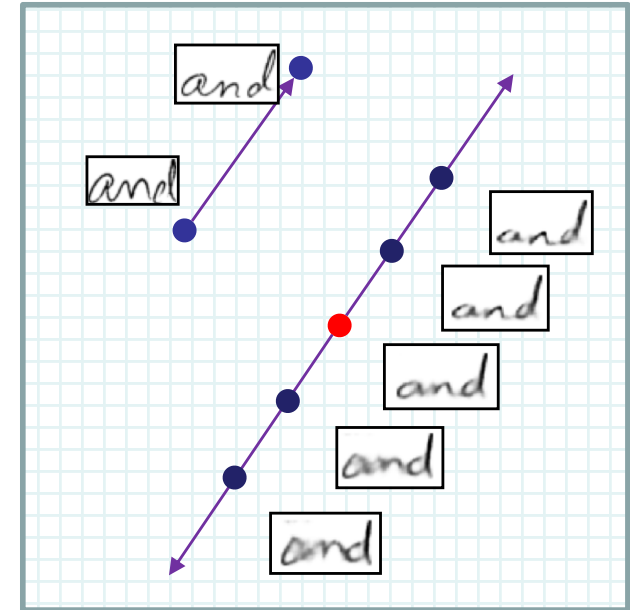
Red boxes contain original images.
Not necessarily from same or different people.

Images in each row generated by increasing/decreasing z_i

Attribute vector from data

and → and

and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and
and and and and and and and and and and and and



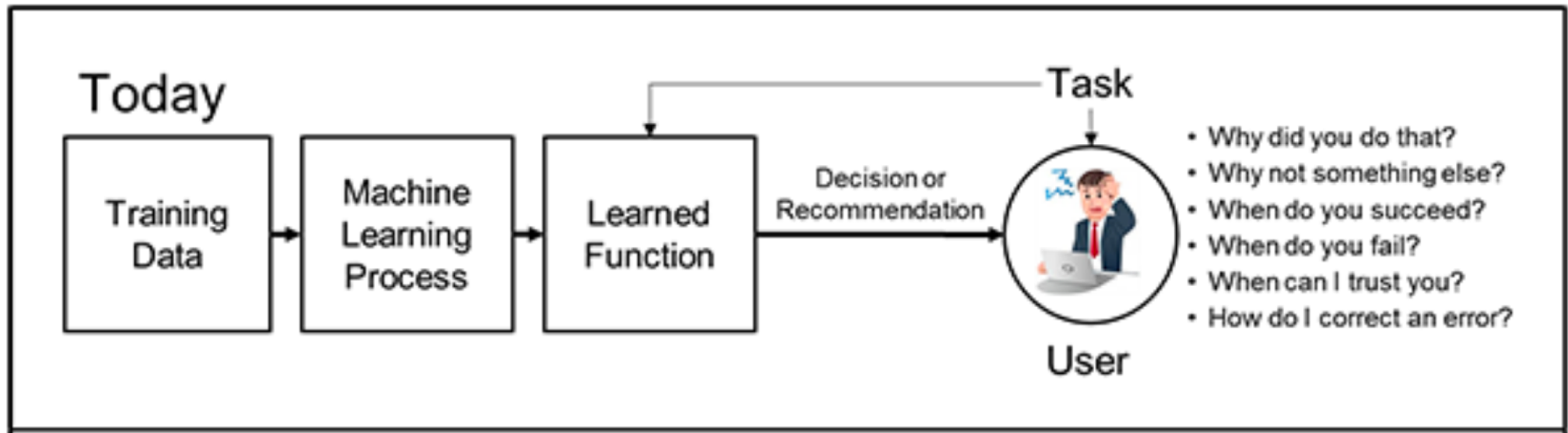
In previous slide traverse is along one dimension in latent space.

Here the traverse is along a so-called *attribute vector*, which points from the first "and" image to the second "and" image from below to the upper one

The "attribute vector" points in direction that makes the writing shorter.

If we apply this "attribute vector" on other "and" images, we see the same effect.

Explainable Artificial Intelligence



Explainable Artificial Intelligence

