

Deep Learning: An Overview

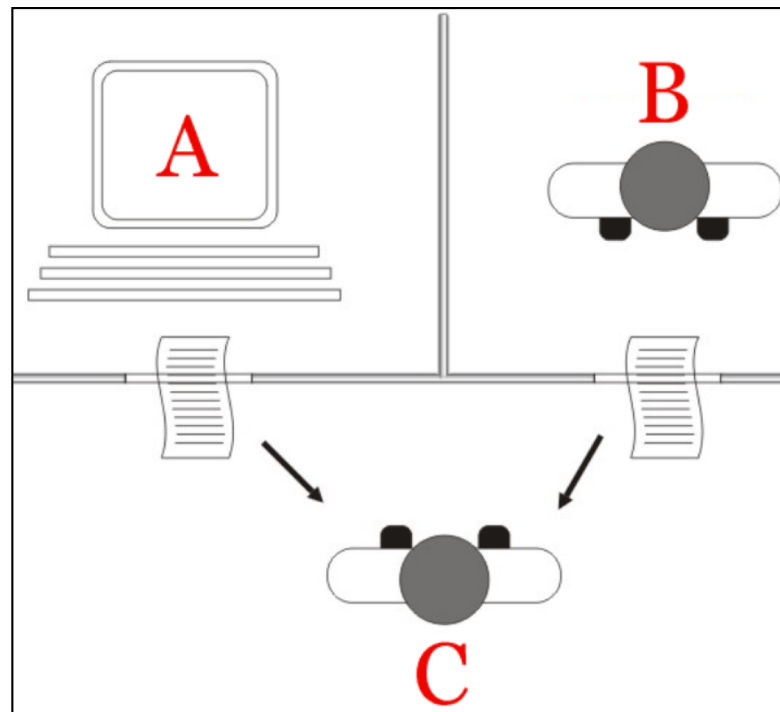
Sargur N. Srihari
srihari@buffalo.edu

Topics

- Artificial Intelligence
 - Definition, Challenges
- Comparison of approaches to AI
 - Knowledge-based, Simple ML, Deep Learning
- Deep Learning
 - Architectures
 - Supervised and Unsupervised Learning
 - Discriminative and Generative Models
- Teaching of AI
- Future of AI

Beginnings of AI

- Advent of *programmable* computers led people to wonder whether AI was possible
- Definition of AI: Turing Test 1950s



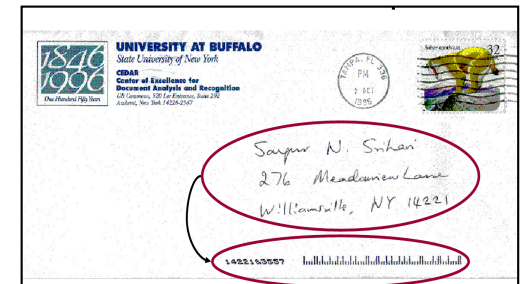
AI in the 1980s

- Moravec paradox
 - Hard problems are easy and easy problems hard
 - Early successes of AI
 - Solved problems intellectually difficult for humans
 - Chess- 64 squares, 32 pieces
 - True challenge of AI
 - Solve tasks easy that people do automatically
 - e,g recognize faces in images
- Artificial General Intelligence (Strong AI)
 - Perform any task that a human can
- Artificial Narrow Intelligence (Weak AI)
 - Specific problem-solving or reasoning tasks

Today AI is ubiquitous

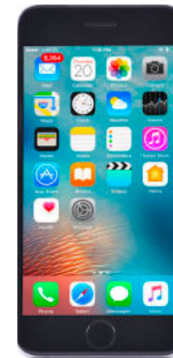
- Due to advances in chip design, processing power and big-data hosting, we rarely notice it
- Computers are already proficient in

- Picking stocks
- Translating speech
- Watson Diagnoses cancer
- Automate routine labor (USPS)



- We take it for granted when

- Siri schedules appointments, Alexa reads a book
- Facebook tags our photos (and subverts democracy)

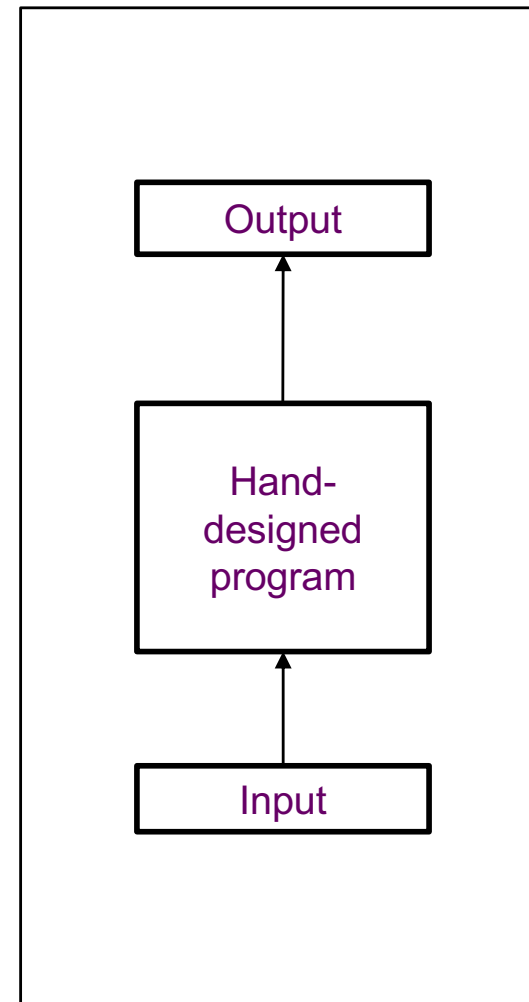


Today's AI research

- It is about solving more intuitive problems
- A person's everyday life requires immense amount of knowledge of the world
 - Much of this knowledge is intuitive and subjective
- Key challenge: how to get this informal knowledge into a computer

Knowledge-Based AI

- Hard-code knowledge in a formal language
- Computer can reason about statements using inference rules
- Example: Expert Systems
- Rule-based System



Disadvantage of Knowledge-Based AI

- Unwieldy process
 - Time of human experts
 - People struggle to formalize rules with enough complexity to describe the world

The Machine Learning approach

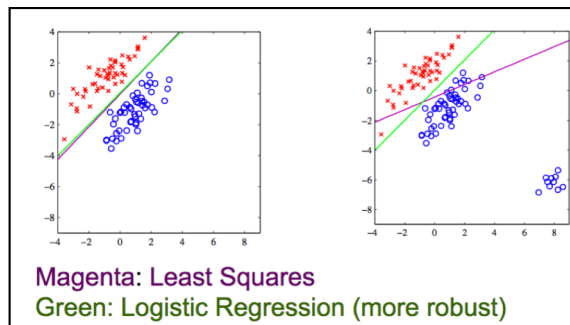
- Allow computers to learn from experience
 - Program said to learn from *experience* E
 - Wrt *task* T and *performance measure* P ,
 - if its performance improves with E
- The Task T
 - Tasks too difficult for human designed programs
 - ML leads to principles underlying intelligence
- Process of learning itself is not the task
 - Learning is a means to perform the task
 - E.g., making a walking robot
 - We could either program the robot to learn to walk or we could directly write a program manually for walking

The Machine Learning approach

- Simple ML algorithms for classification are:

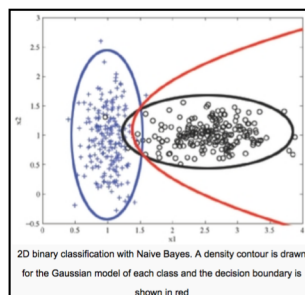
1. Logistic Regression

- Caesarian delivery using uterine scar feature

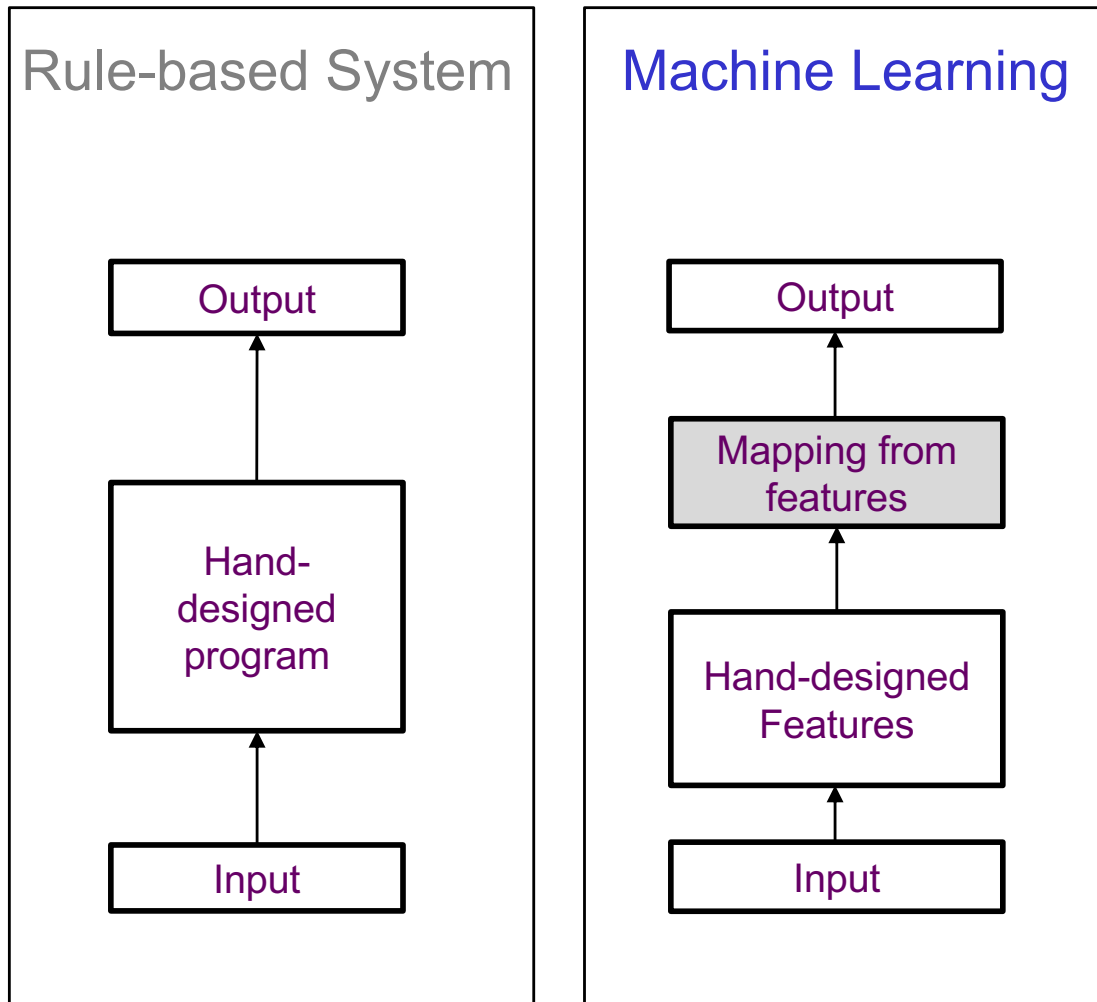


2. Naive Bayes

- Decide whether email is spam using keyword features



Simple Machine Learning AI Model



■ Shaded boxes indicate components that can learn from data

Designing right set of features

- Simple ML algorithms depend heavily on *representation* of given data
- Difficult to know what set of features are good
- For detecting a car in photographs
 - Simple shape but difficult in terms of pixel values
 - Shadows, glare from metal parts, fender or object in front obscures

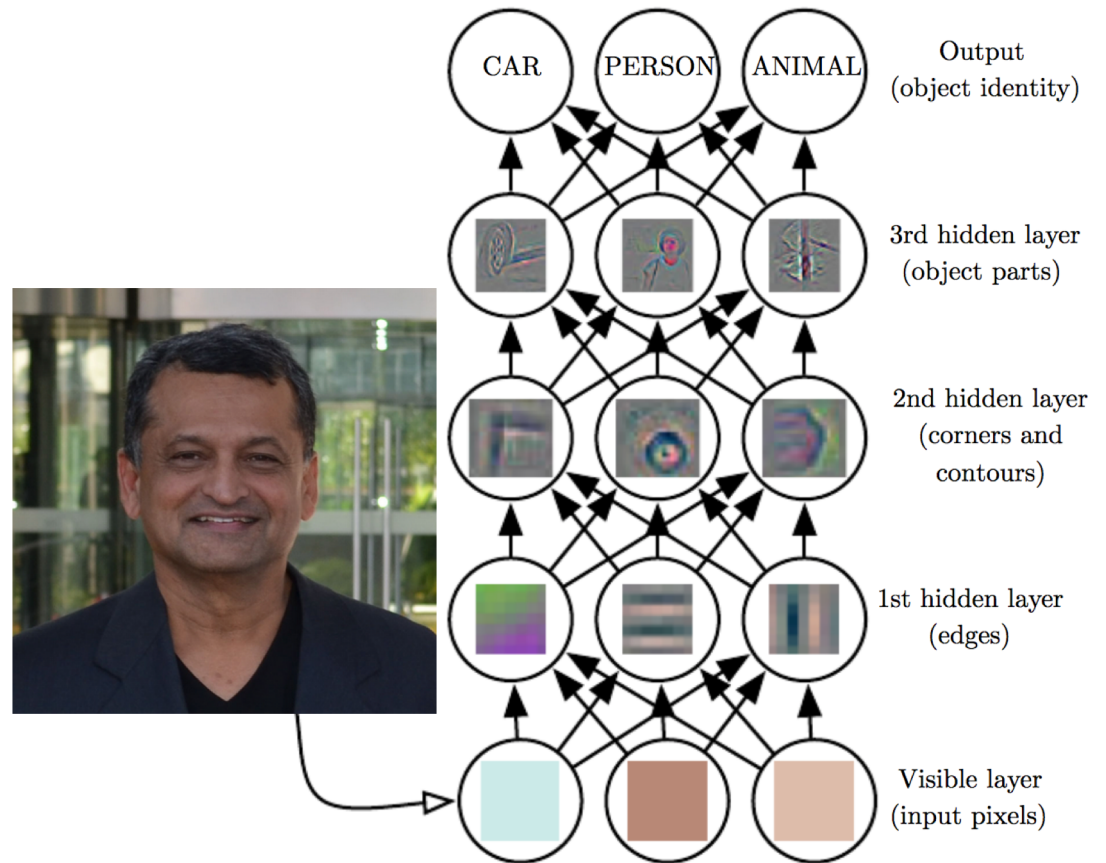


Representation Learning

- Solution: use ML to not only learn mapping from representation to output but representation itself
 - Better results than hand-coded representations
- Allows AI systems to rapidly adapt to new tasks
 - Designing features can take great human effort
 - Can take decades for a community of researchers

Feature Learning for Classification

- Function to map pixels to object identity is complicated
- Series of hidden layers extract increasingly abstract features
- Final decision made by a simple classifier

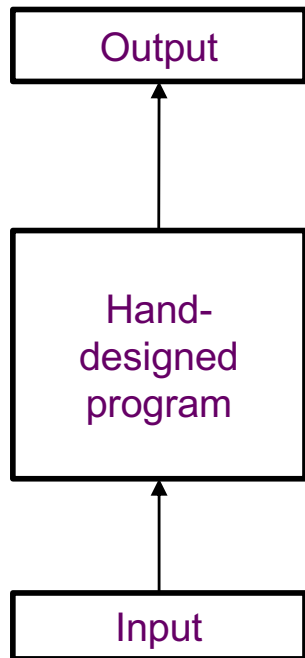


Depth of a Model

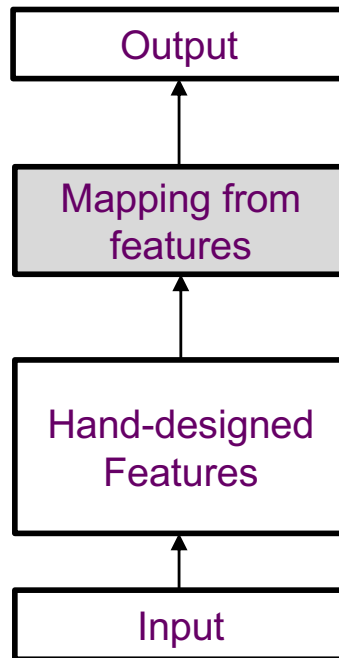
- Deep learning differs from traditional ML in that models involve
 - (a) composition of learned functions
 - (b) learned concepts
- Network describes concept relationships
 - Often concept with more complex underlying graph
 - Eye more complex than face (which consists of two eyes)
- No single correct value for depth of architecture just as no single correct length of a program

Summary of AI Models

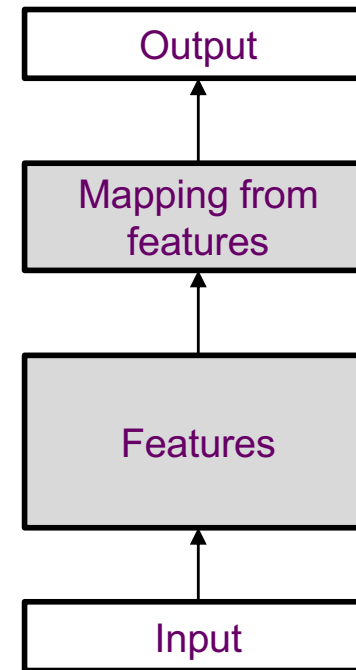
Rule-based System



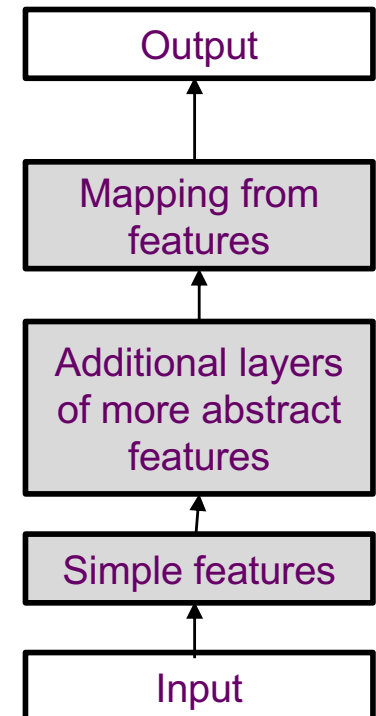
Classic Machine learning



Representation Learning



Deep Learning



■ Shaded boxes indicate components that can learn from data

Challenge of Unstructured Modeling

- The classification task is a limited goal of ML
- Goal of deep learning is to scale ML to the kind of challenges needed to solve AI
 - This means being able to understand high-dimensional data with rich structure
 - E.g., understand natural images,
 - Audio waveforms representing speech
 - Documents containing multiple words and punctuations



Unsupervised Representation Learning

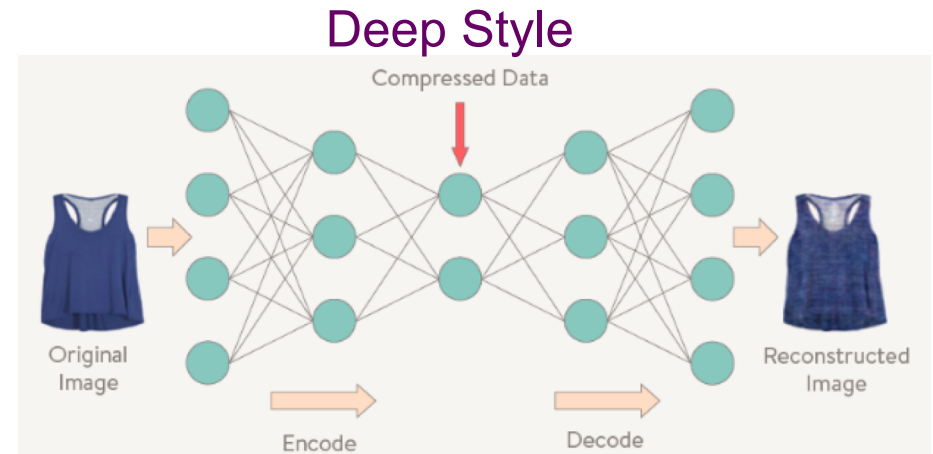
- Autoencoder

- Encoder:

- Converts input into a representation

- Decoder:

- Converts representation back to input










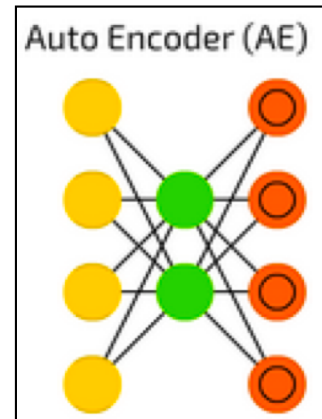
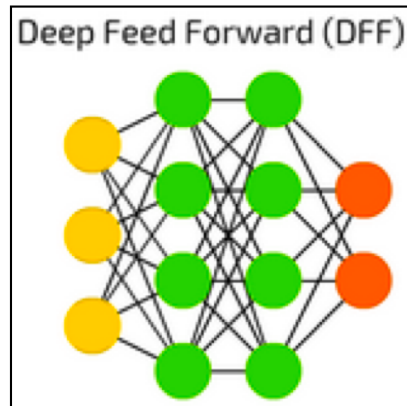
New designs from representation



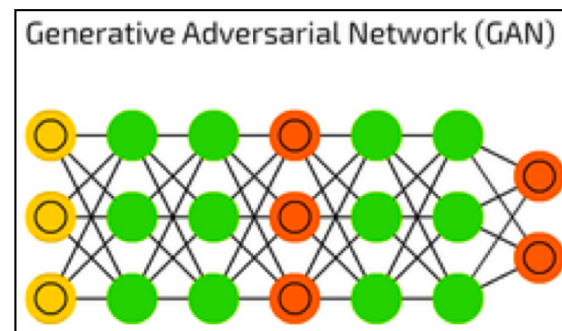
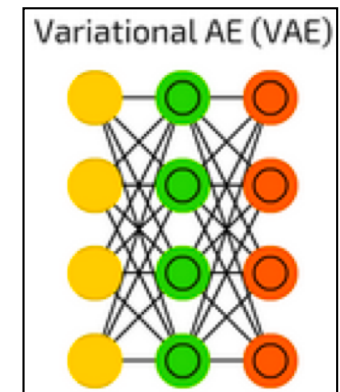
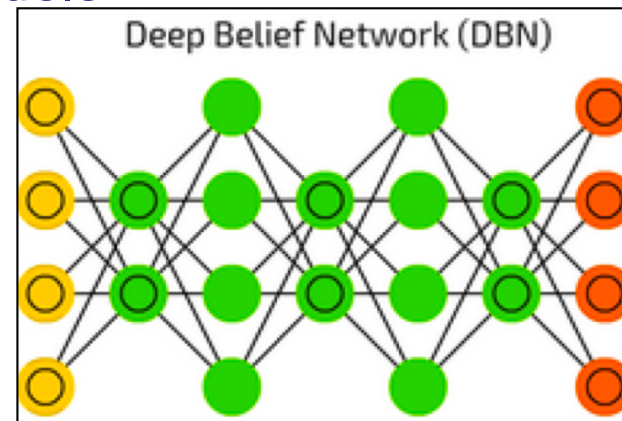
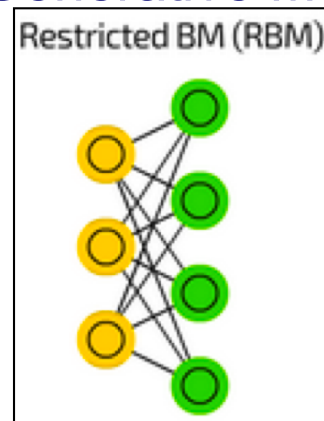
Deep Learning Architecture Examples

Discriminative Models

-  Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Different Memory Cell
-  Kernel
-  Convolution or Pool



Generative Models

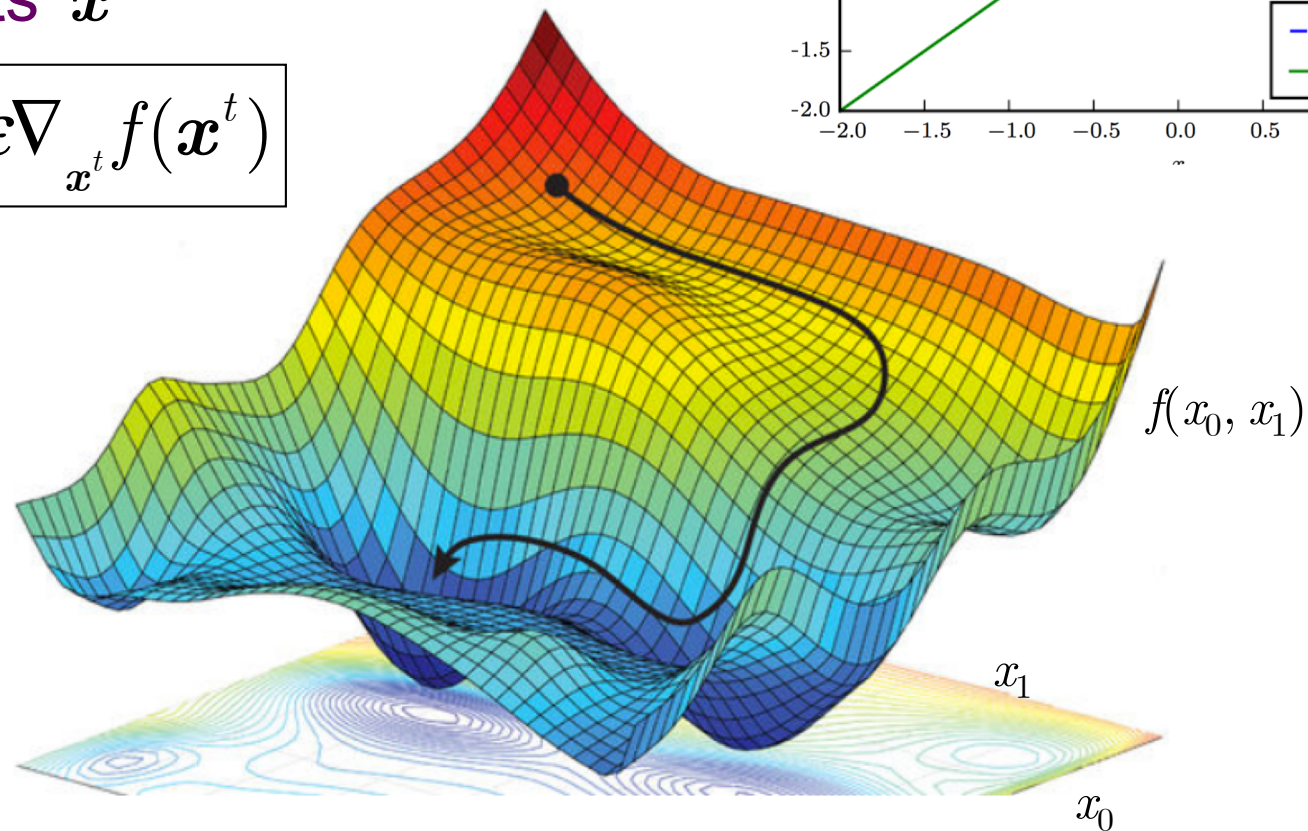
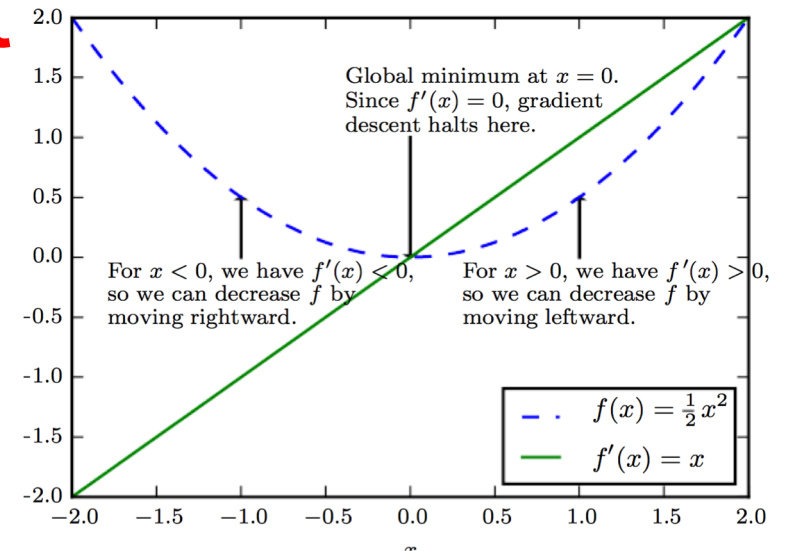


Supervised Learning by Gradient Descent

Classification Task:

Loss function $f(\mathbf{x})$,
e.g., sum of squared errors,
given weights \mathbf{x}

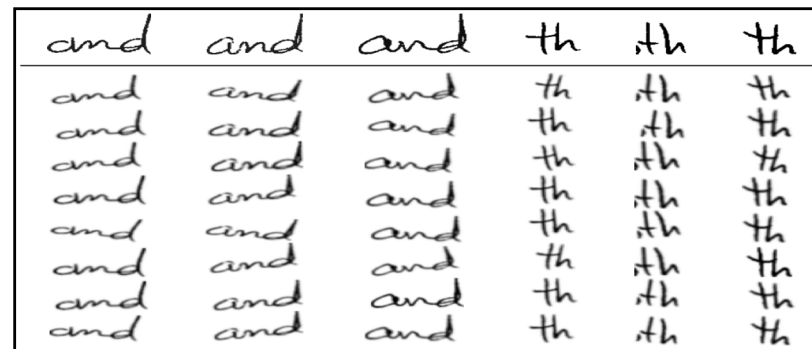
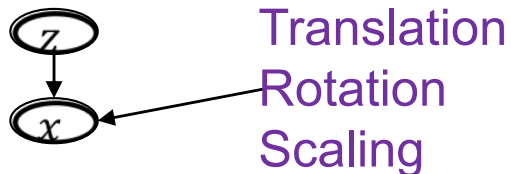
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \epsilon \nabla_{\mathbf{x}^t} f(\mathbf{x}^t)$$



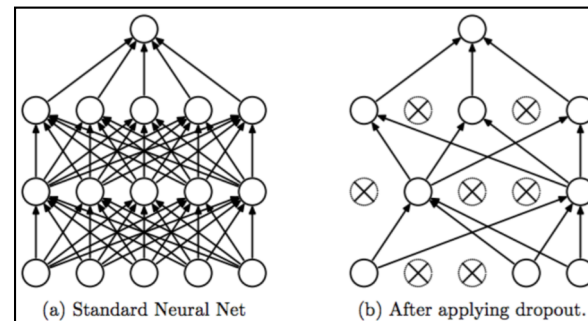
Regularization

- ML algorithms need to perform well not just on training data but on new inputs as well

- Parameter Norm Penalties (L^2 - and L^1 - regularization)
- Data Set Augmentation



- Early Stopping
- Dropout



Generative Model: RBM

- Joint configuration (\mathbf{v}, \mathbf{h})

- visible and hidden units has an energy

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij}$$

- Network assigns a probability to every pair of hidden and visible vectors

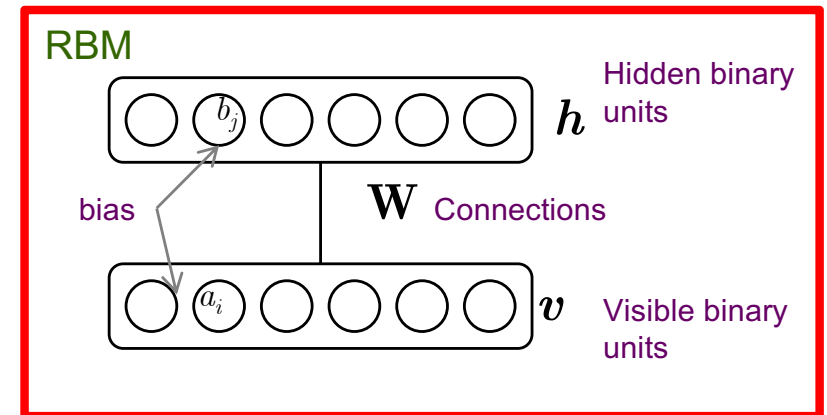
$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$

- where partition function Z is a sum over all possible pairs of visible/hidden vectors

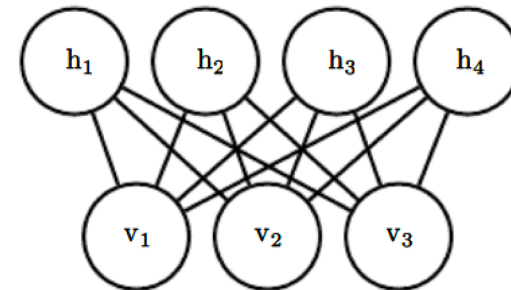
$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

- Probability that network assigns to a visible vector \mathbf{v} is

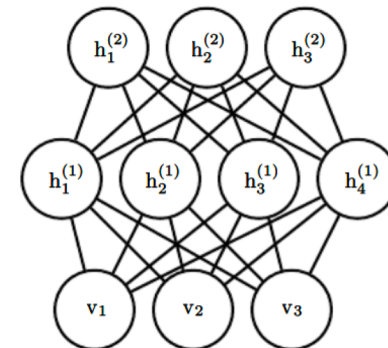
$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$



Stochastic binary pixels \mathbf{v} connected to stochastic binary feature detectors \mathbf{h} using symmetrically weighted connections



Deep Boltzmann Machine



Summary of RBM training

Probability Distribution of Undirected model (Gibbs)

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{x}, \boldsymbol{\theta})$$

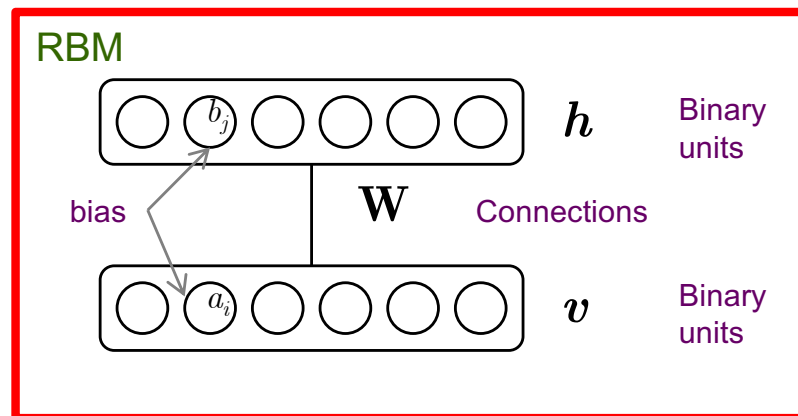
$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}, \boldsymbol{\theta})$$

Intractable
Partition function

For an RBM: $\mathbf{x} = \{\mathbf{v}, \mathbf{h}\}$ $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$

Determine parameters $\boldsymbol{\theta}$ that maximize log-likelihood (negative loss)
 $\max_{\boldsymbol{\theta}} L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}; \boldsymbol{\theta}) = \sum_m \log p(\mathbf{x}^{(m)}; \boldsymbol{\theta})$

$$L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}; \boldsymbol{\theta}) = \sum_m \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \sum_m \log Z(\boldsymbol{\theta})$$



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} = -\sum_i \sum_j W_{i,j} v_i h_j - \sum_i a_i v_i - \sum_j b_j h_j$$

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$Z = \sum_{\mathbf{v}, \mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$$

For stochastic gradient ascent, take derivatives:

$$g_m = \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}^{(m)}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) \quad \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \varepsilon g$$

$$\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$$

Derivative of positive phase:

$$\frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta})$$

Summation is over samples from the training set

Since it is summed m times $1/m$ has no effect

Derivative of negative phase:

$$\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})$$

An identity

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta})$$

Summation is over samples from the RBM

Contrastive Divergence Algorithm

- Using gradient ascent for optimization

Set ϵ , the step size, to a small positive number.

Set k , the number of Gibbs steps, high enough to allow a Markov chain sampling from $p(\mathbf{x}; \boldsymbol{\theta})$ to mix when initialized from p_{data} . Perhaps 1-20 to train an RBM on a small image patch.

while not converged **do**

Sample a minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from the training set.

$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$.

Positive phase

for $i = 1$ to m **do**

$\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}$.

end for

for $i = 1$ to k **do**

for $j = 1$ to m **do**

$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs_update}(\tilde{\mathbf{x}}^{(j)})$.

end for

end for

$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$.

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$.

end while

Negative Phase:

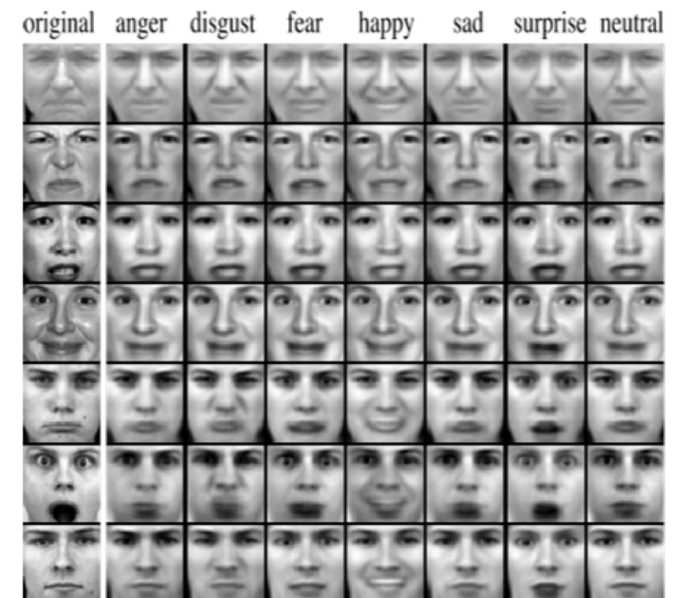
Gibbs samples from the model initialized with training data are used to determine gradient update

What makes one representation better than another?

- Ideal representation h is one where features correspond to the underlying causes of the observed x
 - With features h_i correspond to different causes
 - Representation disentangles causes from one another
- This motivates approaches in which we seek a good representation for $p(x)$
 - Which may also be good for representing $p(y|x)$ if y is among the most salient causes of x

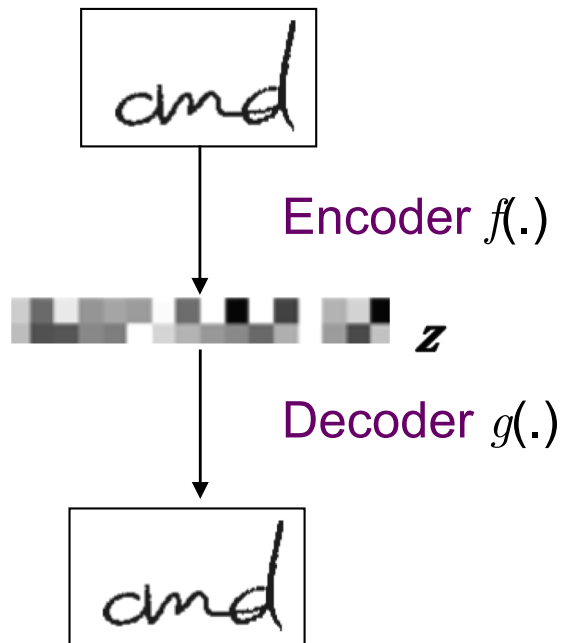
Factors of Variation (FoVs)

- When designing features or algorithms for learning features,
 - we seek invariance to FoVs in the observed data
- Factors refer to separate sources of influence
 - Often not directly observed
- Examples of FoVs
 - Face:
 - pose, expression, morphology
 - Speech:
 - age, sex, accent, words spoken
 - Car image:
 - position, color, angle of sun



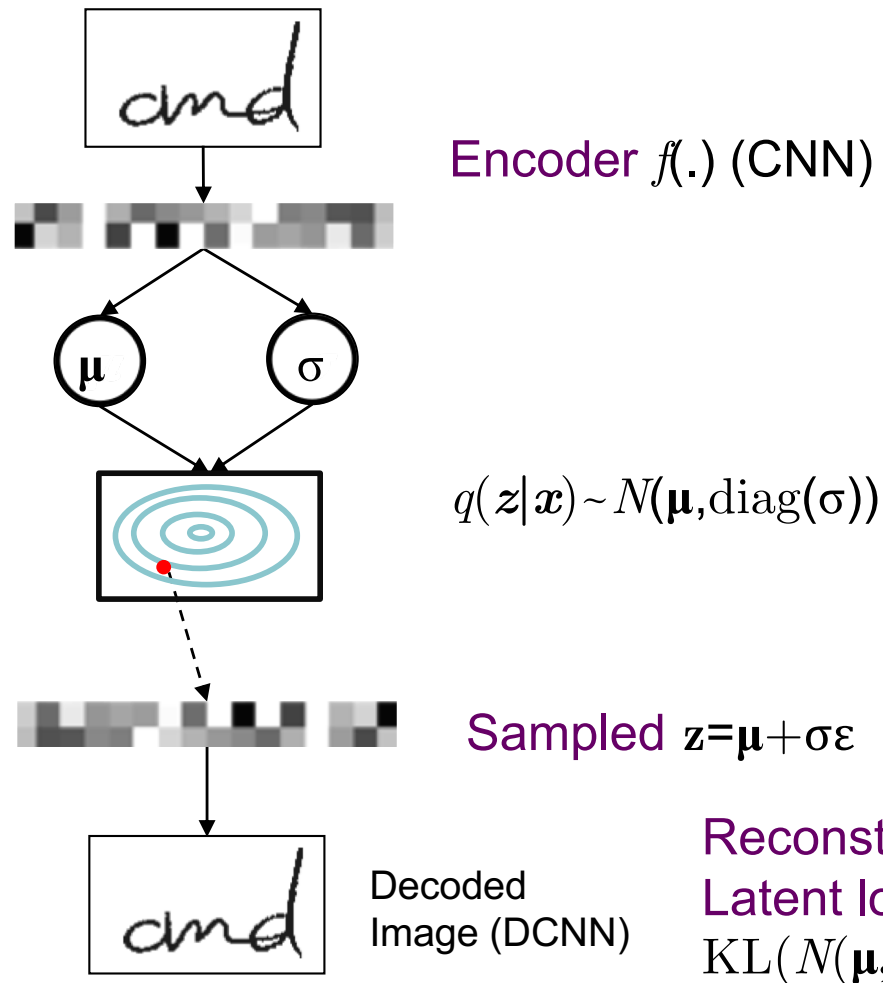
Unsupervised Representation Learning Using Generative Models

Autoencoder



Loss: Reconstruction loss

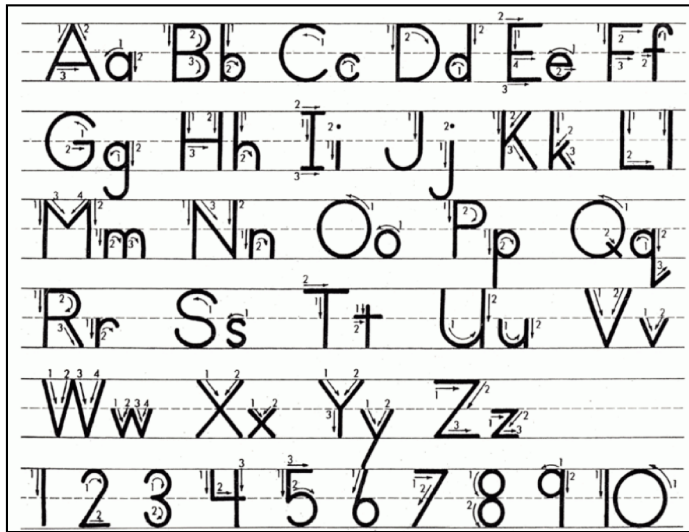
Variational Autoencoder



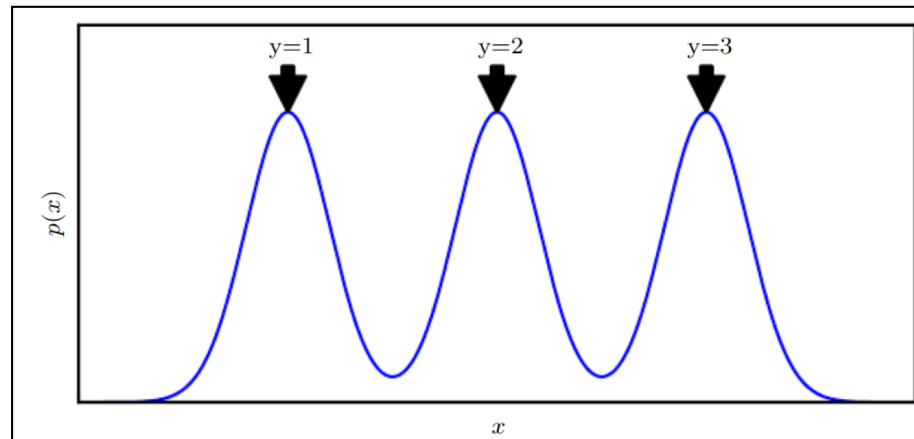
Reconstruction loss +
Latent loss:
 $\text{KL}(N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), N(\mathbf{0}, \mathbf{I}))$

How semi-supervised can succeed

- Ex: density over \mathbf{x} is a mixture over three components, one per value of \mathbf{y}
- If components well-separated:
 - modeling $p(\mathbf{x})$ reveals where each component is
 - A single labeled example per class enough to learn $p(\mathbf{y}|\mathbf{x})$



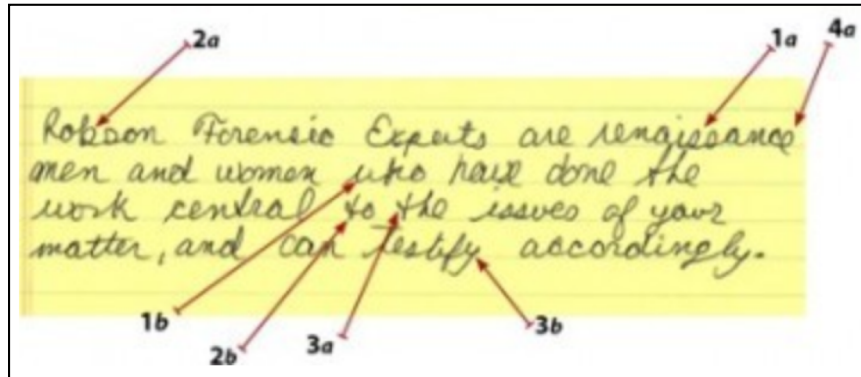
\mathbf{x} = no. of black pixels



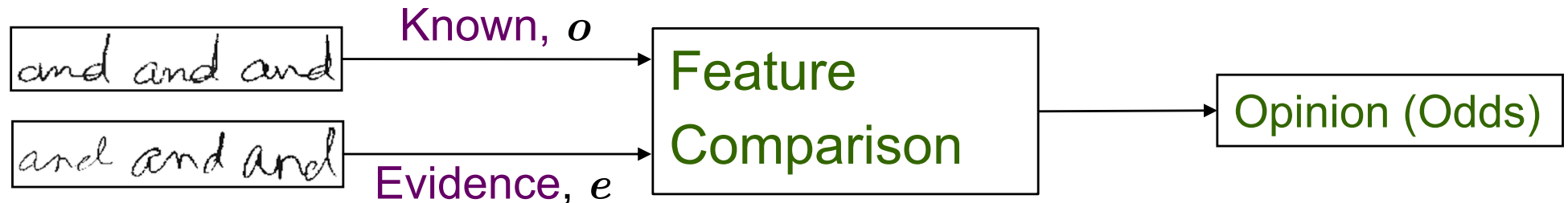
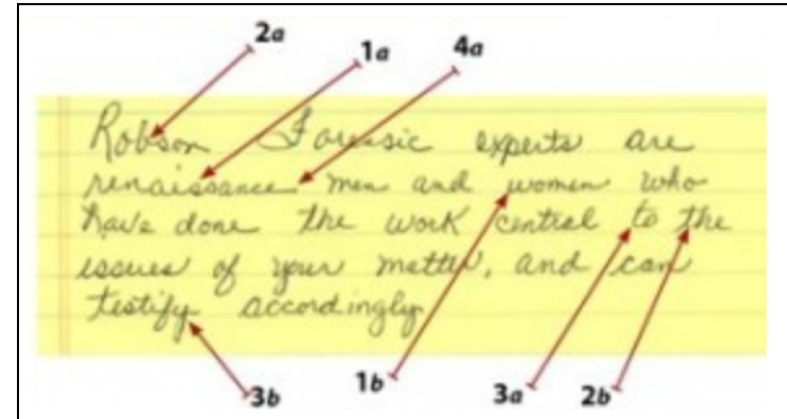
In this case $p(\mathbf{y}|\mathbf{x})$ is a univariate Gaussian for $\mathbf{y}=1,2,3$

A Computer Vision Task

Known Writer Sample, o



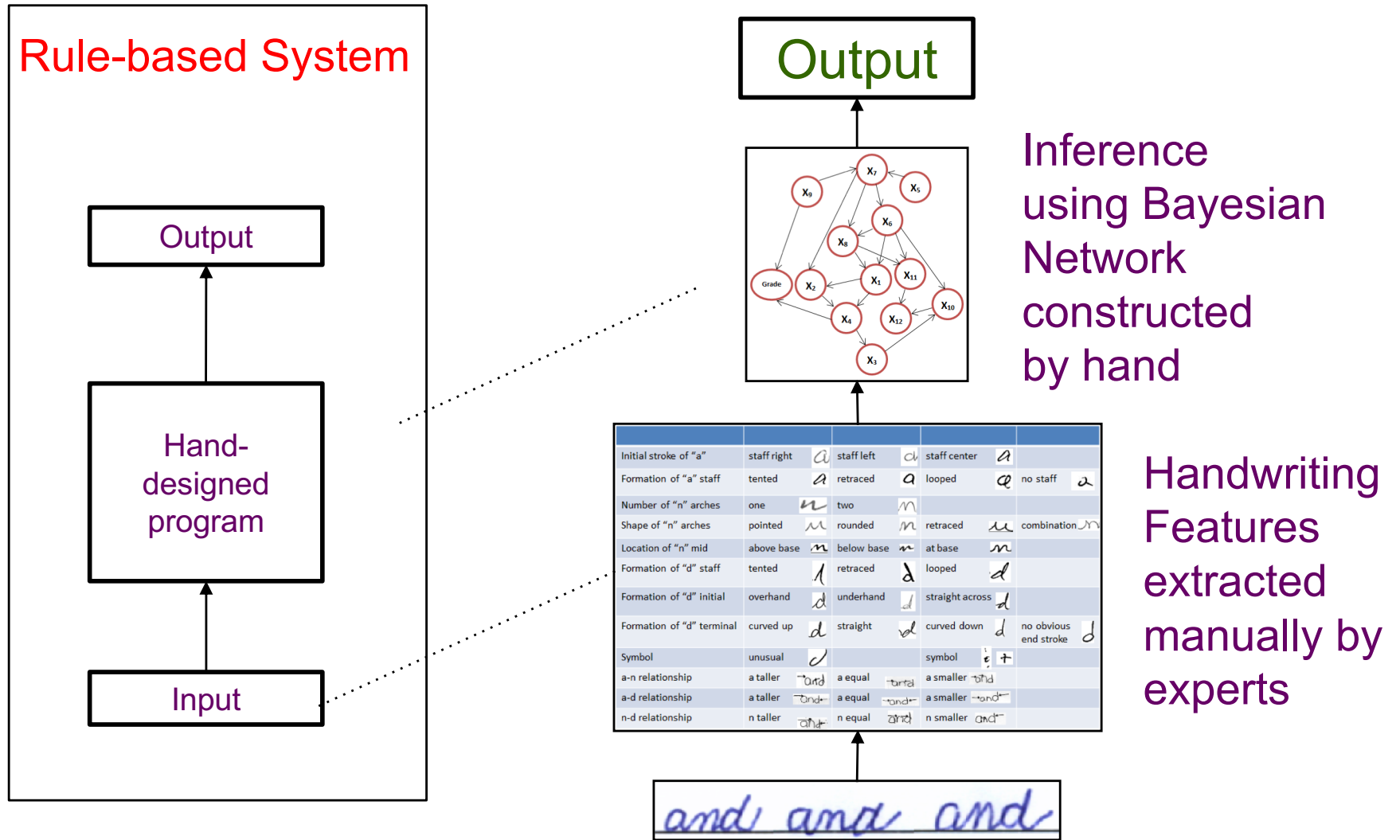
Evidence, e



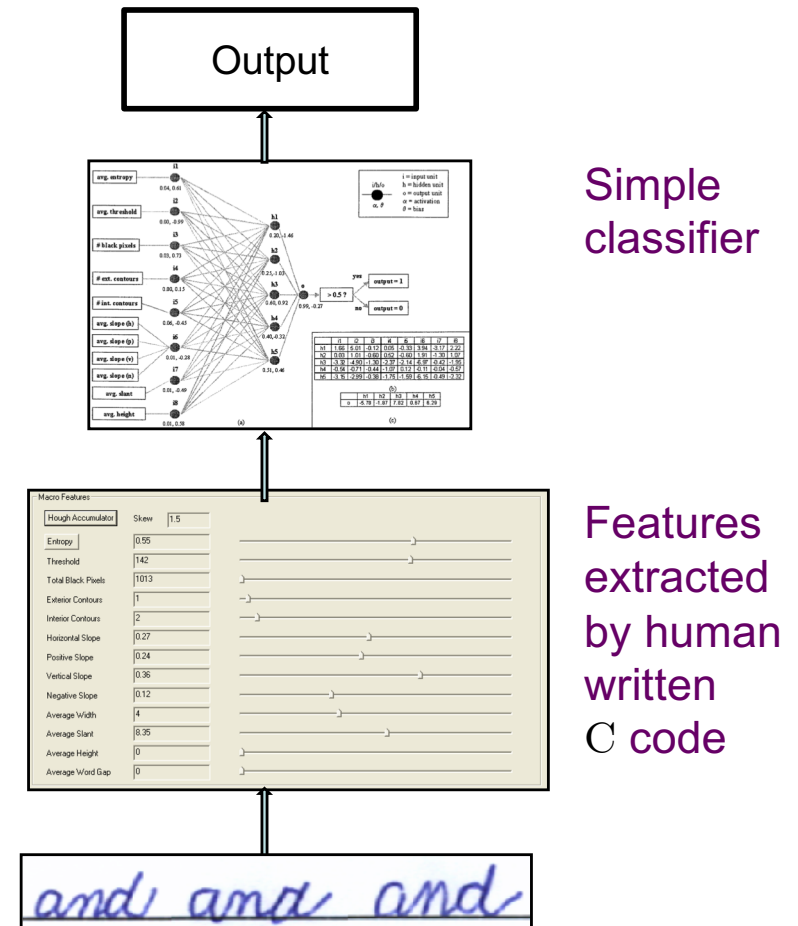
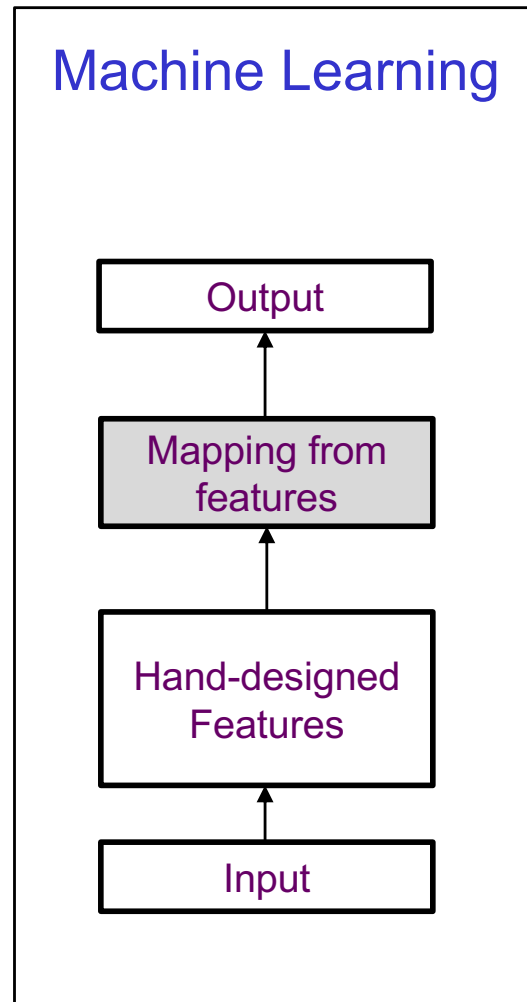
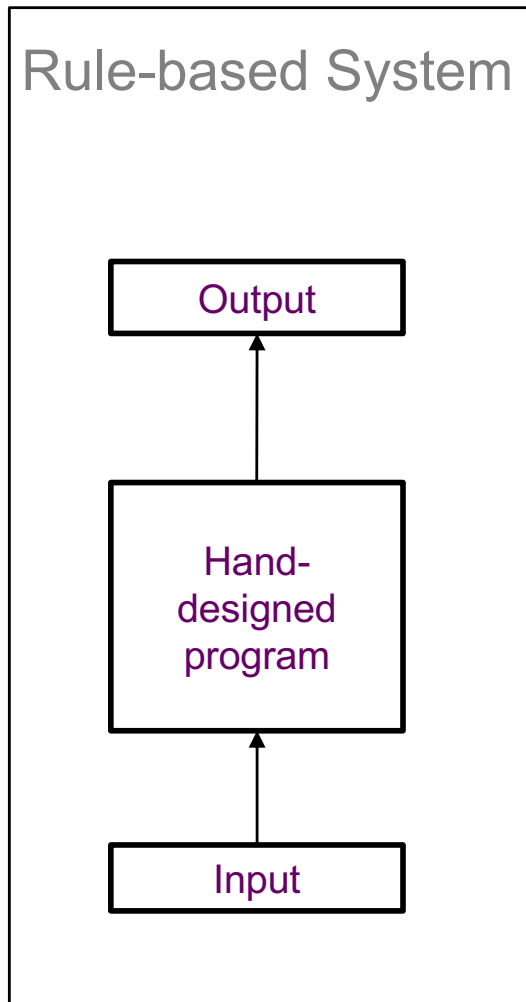
- h^0 : Evidence e was generated by known o
- h^1 : Evidence e was not generated by o

$$\text{Odds Ratio} = \frac{P(o, e|h^0)}{P(o, e|h^1)}$$

Knowledge-Based AI: Handwriting

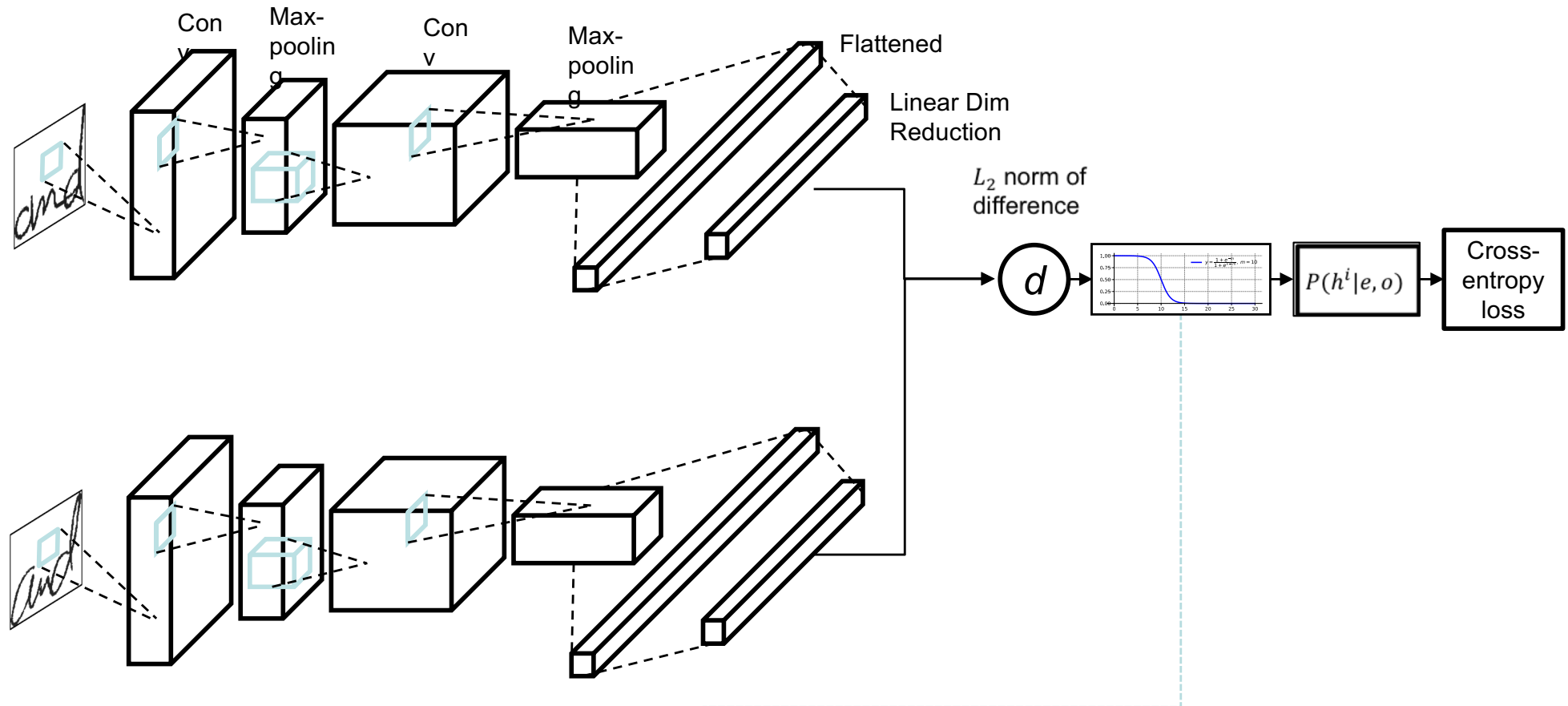


Simple ML for Handwriting

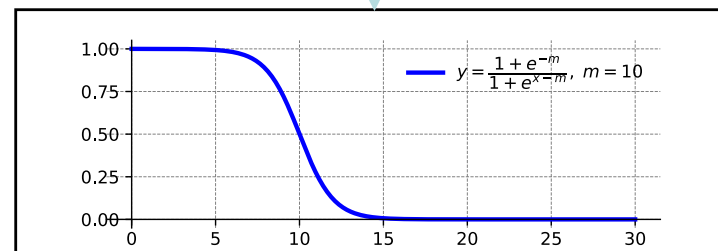


Shaded boxes indicate components that can learn from data

Deep Learning for HW Comparison

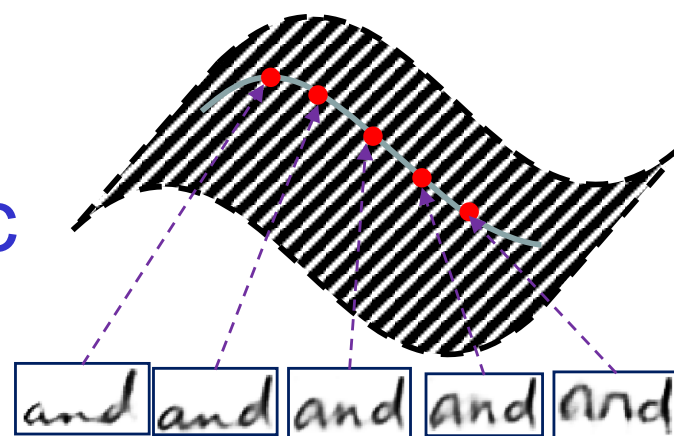
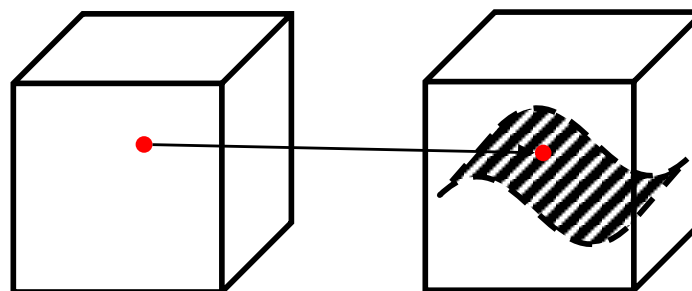


$$P(h^0 | o, e) = \frac{1 + e^{-m}}{1 + e^{d-m}}$$



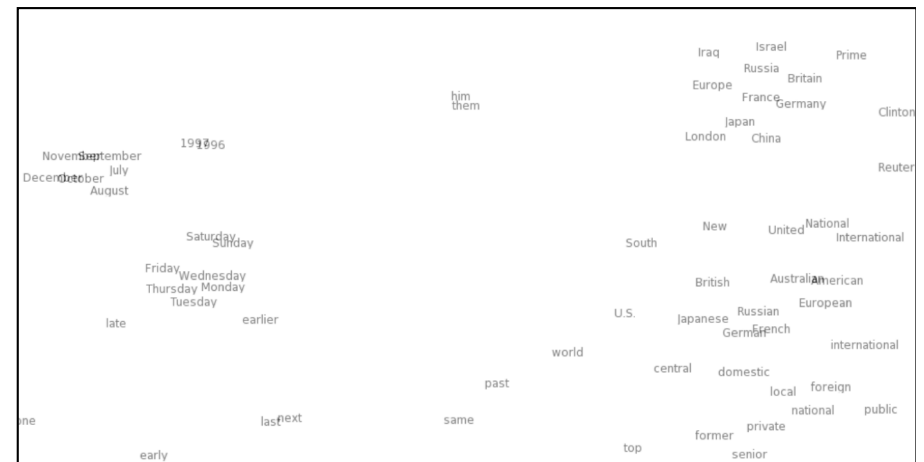
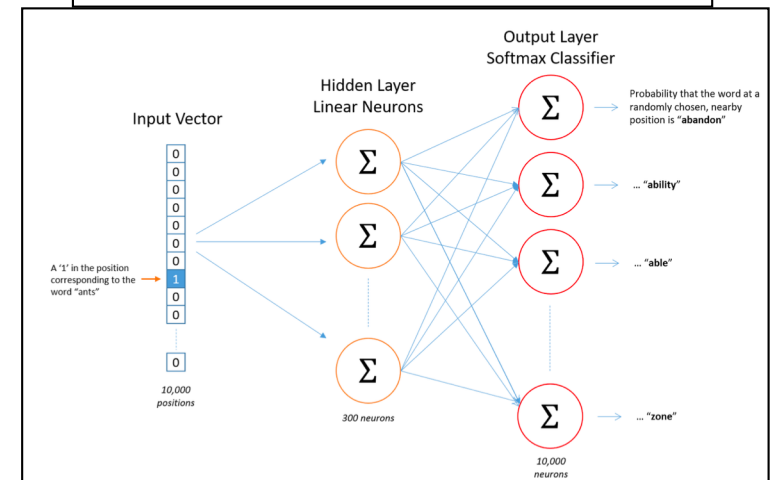
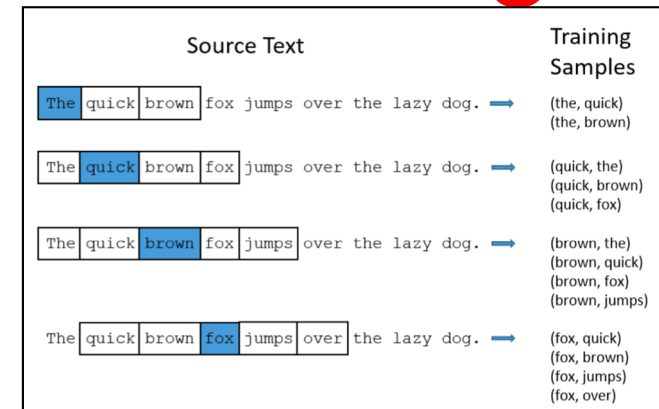
Manifold Learning

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



Natural Language Processing

- Training Data
- Word-to-vec
 - One-hot vector mapped to vector of 300
- Word embedding
 - Similar words are close together



Word-to-vec:

Represent noun by co-occurrences with 25 verbs*

Semantic feature values:

“celery”

0.8368, eat
0.3461, taste
0.3153, fill
0.2430, see
0.1145, clean
0.0600, open
0.0586, smell
0.0286, touch
...
...
0.0000, drive
0.0000, wear
0.0000, lift
0.0000, break
0.0000, ride

Semantic feature values:

“airplane”

0.8673, ride
0.2891, see
0.2851, say
0.1689, near
0.1228, open
0.0883, hear
0.0771, run
0.0749, lift
...
...
0.0049, smell
0.0010, wear
0.0000, taste
0.0000, rub
0.0000, manipulate

* in a trillion word text collection ³⁶

Teaching of AI

– Mathematics

– Probability Theory, Calculus, Linear Algebra, Graph Theory

– Algorithms

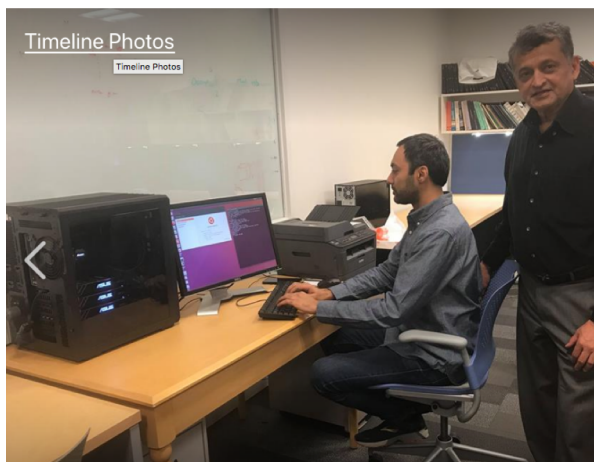
– Neural networks, Sampling, Variational Approximation

– Programming Languages

– Python, Tensorflow

– Hardware

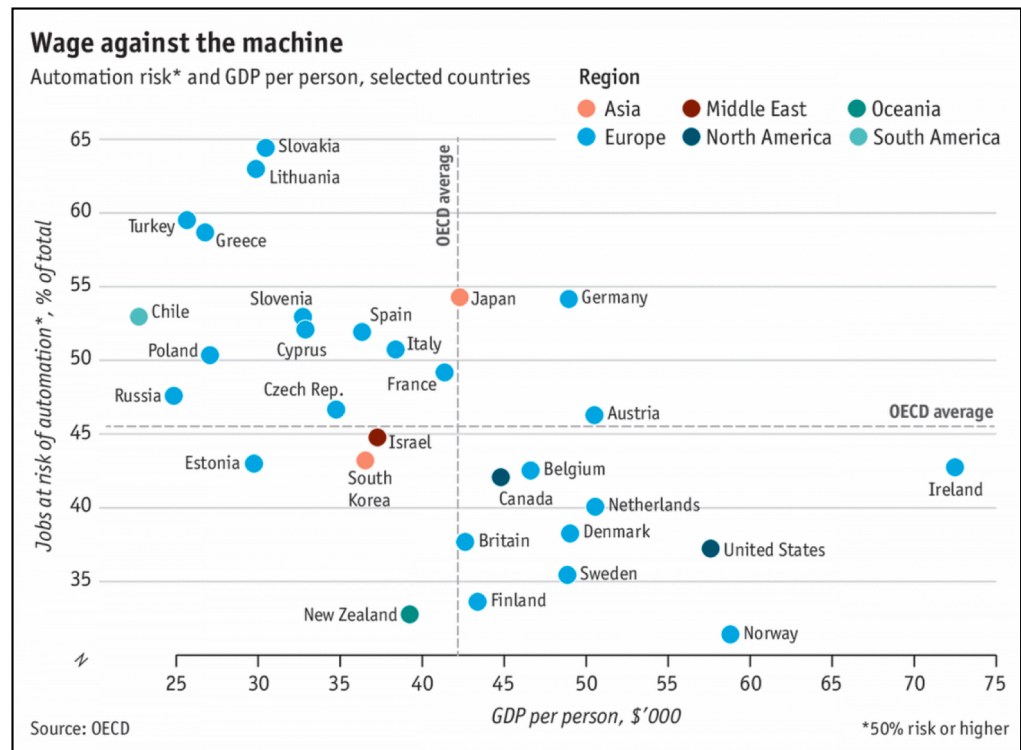
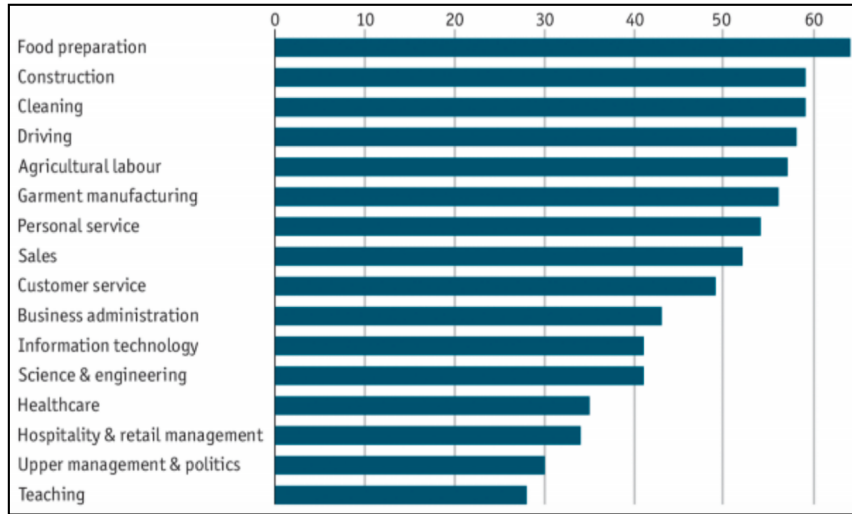
– GPUs



Jobs at risk due to automation

42% of US jobs at risk over next two decades.

International (Rich and mid-income):
 14% of jobs in 32 countries have 70% risk
 For another 32% of jobs, that probability is 60%
 Thus 210m jobs at risk in those countries alone.
 Rich-countries less at risk than mid-income states



Future prospects of AI

- Survey of AI researchers

- AI will outperform humans in:

- Translating languages 2024
 - Writing high-school essays 2026
 - Driving a truck 2027
 - Working in retail 2031
 - Writing a best-selling book 2049
 - Working as a surgeon 2053
 - Outperform humans in all tasks: 50% chance in 45 years
 - Automating all human jobs 120 years

- Survey population: 2015 NIPS/ICML authors

- Questions on AI capabilities (e.g. folding laundry, language translation), superiority at specific occupations (e.g. truck driver, surgeon), superiority over humans at all tasks.