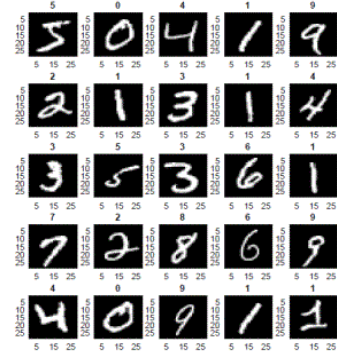


Goal: construct a network that can recognize handwritten numbers $\in \{0, 1, \dots, 9\}$ from MNIST (Modified National Institute of Standards and Technology) data set.



- contains 60000 training images, labeled by 'image ID' $n = 1, \dots, N$ and 10000 testing images

- 28x28 pixels, labeled by 'pixel ID' $l = 1, \dots, 784 \equiv L$

- each pixel contains grey-scale value $x_n^l \in (0, 1) \equiv I \subset \mathbb{R}$
 white black unit interval

- image n is represented by 'image vector' $\vec{x}_n = (x_n^1, \dots, x_n^L) \in I^L$

- each image has been assigned a 'target name' $\vec{t}_n \in \{\vec{e}_0, \dots, \vec{e}_9\}$,

where $\vec{e}_j = (0, 0, \dots, 1, \dots, 0)$, a basis vector in \mathbb{N}^{10} , represents the number $j \in \{0, \dots, 9\}$
 position j

Goal: find 'decision function' \vec{f} that maps image vector to 'predicted name',

$$\vec{f}: I^L \rightarrow \mathbb{N}^{10}, \quad \vec{x}_n \mapsto \vec{f}(\vec{x}_n) \equiv \vec{f}_n$$

'predicted name'

while minimizing the cost function $C = \sum_{n=1}^N (\vec{f}_n - \vec{t}_n)^2$

target name

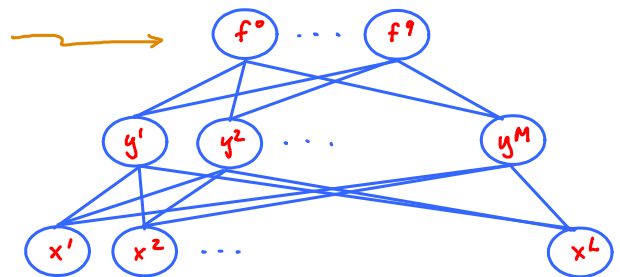
[Alternatively, choose $\vec{f} \in I^{10}$, $|\vec{f}| = 1$ then f_j = probability that image is the number j]

1. Neural network

'output layer': $\vec{f} = (f^0, \dots, f^9) \in \mathbb{N}^{10}$

'hidden layer': $\vec{y} = (y^1, \dots, y^M) \in I^L$

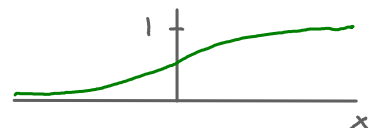
'input layer': $\vec{x} = (x^1, \dots, x^L) \in I^L$



Non-linear transformation: $y^k = \sigma(b^k + w_{\ell}^k x^{\ell})$
 bias weight input

with $\sigma(x) = \frac{1}{1 + e^{-x}}$

'sigmoid function'



mimics neuron: 'fires' when input is above threshold

'soft-max layer':
$$f^j = \frac{e^{(a^j + u^j_l y^l)}}{\sum_{i=0}^9 e^{(a^i + u^i_l y^l)}}$$

use of exponentials emphasizes largest output at expense of others

$\vec{v} = (b, w, a, u)$ are variational parameters, used to minimize C (e.g. by gradient descent)
 \Rightarrow 'train the network' = 'supervised learning'

Multilayer networks (many layers = 'deep learning')

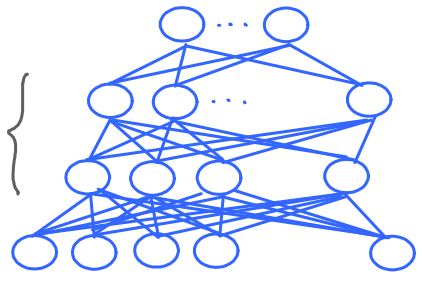
All of the above is just one possible Ansatz.
 Many others can and have been tried.

E.g.: multilayer networks:
 hope is: will capture hierarchical structure better

output layer:

hidden layers: {

input layer:



As before, sigmoid functions can be used to map input to output from one layer to the next.

Optimize cost function using gradient descent:

$C = C(\vec{v})$
 ↑ parameters of network

Gradient: $-\vec{\nabla} C = -\left(\frac{\partial C}{\partial v^1}, \frac{\partial C}{\partial v^2}, \dots\right)$

points in direction of steepest descent:

New variables: $\vec{v}' = \vec{v} - \eta \vec{\nabla} C$

↑ 'learning rate' (should be neither too small, nor too large)

2. Supervised learning with tensor networks

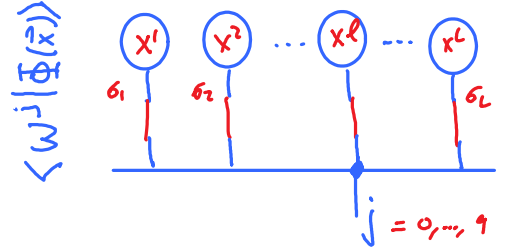
[Novikov2016], [Stoudenmire2017] with Schwab; [Maier2017] Bachelor thesis of David Maier

Goal: construct decision function \vec{f} using a tensor network (here MPS);
 train network using optimization techniques familiar from DMRG

Ansatz: $\vec{f} : \mathbb{I}^L \mapsto \mathbb{I}^{10}$, (1)

$\vec{x} \mapsto \vec{f}(\vec{x}) \equiv \langle \vec{w} | \Phi(\vec{x}) \rangle$ (2)

image vector predicted name



where right-hand side involves two separate maps:

'feature map' $\Phi : \vec{x} \mapsto |\Phi(\vec{x})\rangle$: encodes greyscale input data into L -leg MPS, $|\Phi(\vec{x})\rangle$ (3)

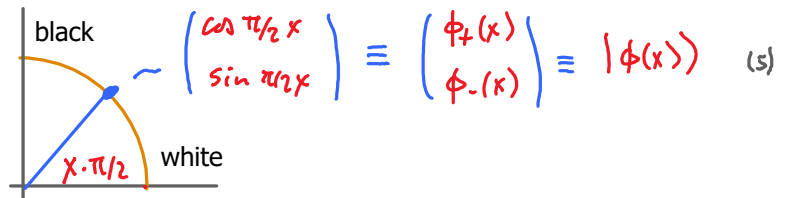
'weight vector' $\vec{w} : |\Phi(\vec{x})\rangle \mapsto f^j(\vec{x}) \equiv \langle w^j | \Phi(\vec{x}) \rangle$, $j = 0, \dots, 9$ (4)

converts feature map into predicted name via inner product with an L -leg MPS,

'predicted name': that label j for which f^j is maximal.

Feature map: encoding input data

map color range
 (0,1) = (white, black)
 to quarter-unit-circle,



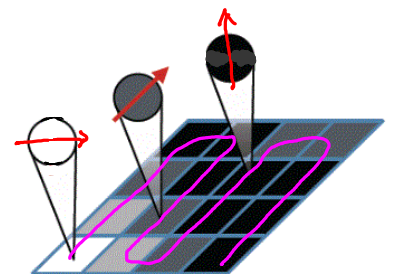
so that $\langle \phi(x') | \phi(x) \rangle = \sum_{\sigma=\pm} \phi_{\sigma}(x') \phi_{\sigma}(x) \approx \begin{cases} 1 & \text{if } x \approx x' \\ 0 & \text{if } x \approx \text{white}, x' \approx \text{black} \end{cases}$ (6)

Choose 'snake-ordering' of pixels,
 and encode image in a product state MPS:

$|\Phi(\vec{x})\rangle = |\phi(x^1)\rangle \otimes |\phi(x^2)\rangle \otimes \dots \otimes |\phi(x^L)\rangle$ (7)

$= \begin{matrix} \textcircled{x^1} & \textcircled{x^2} & \dots & \textcircled{x^L} & \dots & \textcircled{x^L} \\ \downarrow & \downarrow & & \downarrow & & \downarrow \\ \sigma_1 & \sigma_2 & & \sigma_L & & \sigma_L \end{matrix}$ (8)

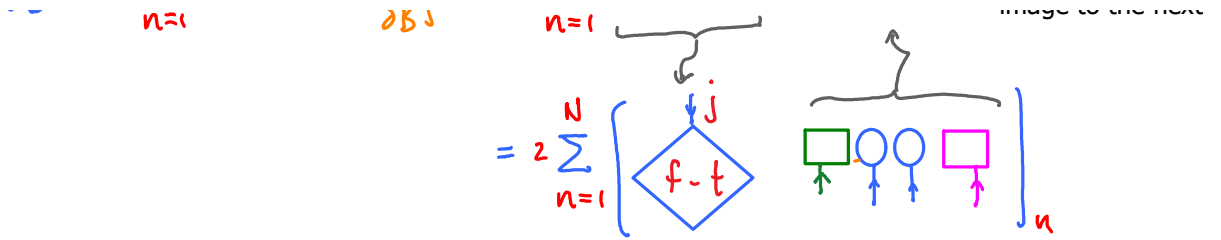
(d=2)



This construction for $|\Phi(\vec{x})\rangle$ is not unique. Other constructions are possible, provided that

$\langle \Phi(\vec{x}') | \Phi(\vec{x}) \rangle$ is a smooth and slowly varying function of \vec{x} and \vec{x}'

which induces a 'distance matrix' in feature space which tends to cluster similar images together.

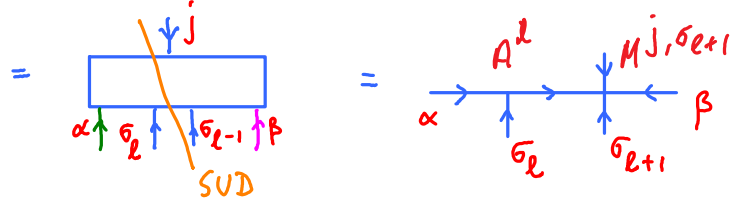


Then update the MPS:

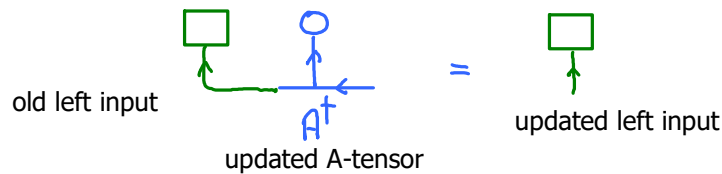
$$|B^j\rangle = |B^j\rangle - \eta (\nabla B^j) = \text{[Diagram of } B^j \text{]} - \eta \sum_{n=1}^N \left[\text{[Diagram of } f-t \text{]} \right]$$

learning rate η (must be chosen very carefully!)

Advance to next site:



Update training input:



Sweep back and forth until A-tensors no longer change -- then 'training of network' is complete.

Comments

Costs: $\mathcal{O}(d^3 D^3 N \cdot L \cdot 10)$

d : physical bond dimension (here: 2) N : number of training images

D : MPS bond dimension (free parameter) L : number of pixels per image

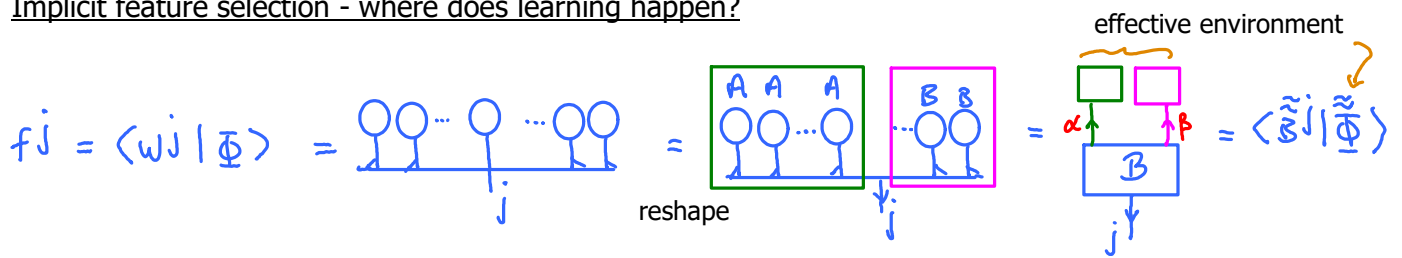
Once network has been trained, prediction of a new image \vec{x} proceeds simply via

$$f^j(\vec{x}) = \langle W^j | \Phi(\vec{x}) \rangle, \text{ predicted name is the } j \text{ yielding maximal } f^j$$

MNIST test:

- 28 x 28 was coarse-grained to 14 x 14 (to save resources)
- at most 5 sweeps were needed before training converges
- bond dimension $D = 10 \Rightarrow$ 5% error rate
- $20 \rightarrow$ 2% error rate
- $120 \rightarrow$ 0.97% error rate

Implicit feature selection - where does learning happen?



- $|\tilde{\Phi}\rangle$ is projection of $|\Phi\rangle$ onto space spanned by orthonormal basis, encoded in $\langle \tilde{B}_j |$
 has just D^2 components

- So, training an MPS model uncovers relatively small set of features, and simultaneously trains decision function using only those features.
- 'Feature selection' occurs when computing SVD: basis elements which do not contribute optimally to bond tensors are discarded

Future prospects

- try tensor networks that are designed for 2D (PEPS, TRG, MERA,)
- *exploit symmetries*
- try other sampling schemes
- ~~try other sampling schemes~~
- 'unsupervised learning' with tensor networks
- ...