

C) Restricted Boltzmann Machine

The Hamiltonian is given by

$$H[\vec{s}, \vec{z}] = - \sum_{i \in V} s_i w_{i\alpha} \tau_{\alpha} - \sum_{i=1}^{N_v} \theta_i s_i - \sum_{\alpha=1}^{N_h} \eta_{\alpha} \tau_{\alpha}$$

↓ visible nodes $\{0,1\}$ ↓ hidden nodes $\{0,1\}$
↓ vis bias (magn fields) ↓ hid bias (magn fields)

$$p(\vec{s}, \vec{\tau}) = \frac{1}{Z} e^{-H[\vec{s}, \vec{\tau}]} \quad \text{with} \quad Z = \sum_{\{s_i, \tau_{\alpha}\}} \exp(-H[\vec{s}, \vec{\tau}])$$

Why RBM can be better than Hopfield model?

- i) it does contain patterns: $w_{i\alpha}$
- ii) \rightarrow it can modelize higher-order correlations

Let's consider $\vec{\eta} = \vec{\theta} = 0$ (for simplicity)

$$p(\vec{s}) = \frac{1}{Z} \sum_{\{\tau_{\alpha}\}} \exp\left(\sum_{i \in V} s_i w_{i\alpha} \tau_{\alpha}\right) = \frac{1}{Z} \prod_{\alpha} \sum_{\tau_{\alpha} \in \{0,1\}} e^{\tau_{\alpha} (\sum_i w_{i\alpha} s_i)}$$

$$= \frac{1}{Z} \prod_{\alpha} (1 + e^{\sum_i w_{i\alpha} s_i}) = \frac{1}{Z} \exp\left[\sum_{\alpha} \log(1 + e^{\sum_i w_{i\alpha} s_i})\right]$$

$w_{i\alpha}$ small $\approx \frac{1}{Z} \exp\left[\sum_{\alpha} \log\left(1 + 1 + \sum_i w_{i\alpha} s_i + \frac{1}{2} \left(\sum_i w_{i\alpha} s_i\right)^2 + \frac{1}{3!} \left(\sum_i w_{i\alpha} s_i\right)^3 + \dots\right)\right]$

N_h \dots $\left[\dots \right]$ ↙ OR W

$$\begin{aligned}
&\approx \frac{2}{L} \sum_{i=1}^{N_h} \exp \left[\sum_{j=1}^{N_h} w_{ij} s_j \right] \\
&\quad + \frac{1}{4} \sum_{i,j} s_i \left(\sum_k w_{ik} v_{jk} \right) s_j - \frac{1}{8} (-) \quad O(w^4) \\
&\quad + \frac{1}{3!} \sum_{i,j,h} s_i s_j s_h \left(\sum_k w_{ik} w_{jk} w_{hk} \right) + \dots \quad O(w^3) \\
&\quad +
\end{aligned}$$

1) Sampling with RBD

Sampling here, is as simple as for the Hopfield model!

First we need to compute the cdt^d dist

$$p(s_i | \vec{\tau}) = \frac{e^{s_i (\sum_k w_{ik} \tau_k + \theta_i)}}{\sum_{s_i} e^{s_i (\sum_k w_{ik} \tau_k + \theta_i)}} = \frac{e^{s_i (\sum_k w_{ik} \tau_k + \theta_i)}}{1 + e^{\sum_k w_{ik} \tau_k + \theta_i}}$$

$$p(s_i = 1 | \vec{\tau}) = \frac{1}{1 + e^{-\sum_k w_{ik} \tau_k + \theta_i}} = \text{sig} \left(\sum_k w_{ik} \tau_k + \theta_i \right) = \langle s_i \rangle_{\vec{\tau}}$$

where: $\text{sig}(x) = \frac{1}{1 + e^{-x}}$ sigmoid fct.

By symmetry, we have

$$p(\tau_a = 1 | \vec{s}) = \text{sig} \left(\sum_i w_{ia} s_i + \gamma_a \right) = \langle \tau_a \rangle_{\vec{s}}$$

Again we can sample configurations layerwise:

i) we start with a random config: $\vec{s}^0 \sim \text{bernoulli}(p=q_s)$

ii) $\vec{z}^0 \sim p(\tau | \vec{s}^0) \sim \text{bernoulli}(1-p = \text{sig}(\sum_i w_{i\tau} s_i^0 + \eta_\tau))$

$\vec{s}^1 \sim p(s_i | \vec{z}^0)$ $\vec{z}^1 \sim p(\tau | \vec{s}^1)$

$\vec{s}^2 \sim p(s_i | \vec{z}^1)$..

⋮

$\{\vec{s}^T, \vec{z}^T\}$

⇒ Learning of RBM: the procedure is as always:

we need to maximize the likelihood

Let's consider a dataset: $\{\vec{s}^{(m)}\}_{m=1, \dots, M}$

The log-likelihood is $L = \frac{1}{M} \sum_m \log [p(\vec{s}^{(m)})] = \frac{1}{M} \sum_m \log \left[\sum_{\vec{z}} p(\vec{s}^{(m)}, \vec{z}) \right]$

$L = \frac{1}{M} \sum_m \log \left[\sum_{\vec{z}} \exp \left(\sum_{i,\tau} s_i^{(m)} w_{i\tau} z_\tau + \sum_i \theta_i s_i^{(m)} + \sum_\tau \eta_\tau z_\tau \right) \right] - \log Z$

$\langle s_i \rangle_D =: \frac{1}{M} \sum_m \sum_i \theta_i s_i^{(m)} + \frac{1}{M} \sum_m \log \left(\prod_{\tau=0,1} e^{\tau (\sum_i w_{i\tau} s_i^{(m)} + \eta_\tau)} \right) - \log Z$

$$\langle f(\vec{s}, \vec{z}) \rangle := \frac{1}{\prod_m} \sum_{\tau \in \mathcal{I}} f(\vec{s}^{(m)}, \tau) \cdot p(\tau | \vec{s}^{(m)})$$

$$= \sum_i \theta_i \langle s_i \rangle_D - \frac{1}{\prod_m} \left[\sum_{\tau} \log \left(1 + e^{\left(\sum_i w_{i,c} s_i^{(m)} + \eta_c \right)} \right) \right] - \log Z$$

We compute the gradients.

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \langle s_i \rangle_D - \frac{1}{Z} \sum_{\tau \in \mathcal{I}} s_i \cdot e^{-H[\vec{s}, \vec{z}]} = \langle s_i \rangle_D - \underbrace{\langle s_i \rangle_{\text{RBN}}}_{= \langle s_i \rangle_{\text{RBN}}}$$

$$\frac{\partial \mathcal{L}}{\partial \eta_c} = \frac{1}{\prod_m} \sum \frac{e^{\sum_i w_{i,c} s_i^{(m)} + \eta_c}}{1 + e^{\sum_i w_{i,c} s_i^{(m)} + \eta_c}} - \langle \tau_c \rangle_{\text{RBN}}$$

$$= \frac{1}{\prod_m} \sum \underbrace{p(\tau_c = 1 | \vec{s}^{(m)})}_{\langle \tau_c \rangle_{\vec{s}^{(m)}}} - \langle \tau_c \rangle_{\text{RBN}}$$

$$= \langle \tau_c \rangle_D - \langle \tau_c \rangle_{\text{RBN}}$$

⚠ Non $\langle \tau_c \rangle_D$ depends on the model. ⚠

$$\frac{\partial \mathcal{L}}{\partial w_{i,c}} = \frac{1}{\prod_m} \sum s_i^{(m)} p(\tau_c = 1 | \vec{s}^{(m)}) - \langle s_i \tau_c \rangle_{\text{RBN}}$$

$$= \underbrace{\langle s_i \tau_c \rangle_D}_{\text{positive term}} - \underbrace{\langle s_i \tau_c \rangle_{\text{RBN}}}_{\text{neg. term}}$$

Learning algo: For a mini-batch

1) compute $p(\tau_c | \vec{s}^{(m)}) \rightarrow s_i^{(m)} p(\tau_c = 1 | \vec{s}^{(m)})$

2

i) compute $p(\tau_i | \vec{s}^{(m)}) \rightarrow s_i^{(m)} p(\tau_i = 1 | \vec{s}^{(m)})$
also $\langle s_i \rangle_D$

ii) Sample N_s conf. for T timesteps
starting from rdun init. cdk.

$$\rightarrow \left\{ \vec{s}^{T, (m)}, \tau^{T, (m)} \right\}_{m=1, \dots, N_s}$$

$$\left\langle \tau_i \right\rangle_{RBN} \approx \frac{1}{N_s} \sum_m p(\tau_i = 1 | \vec{s}^{T, (m)})$$

$$\left\langle s_i \right\rangle_{RBN} \approx \frac{1}{N_s} \sum_m s_i^{T, (m)}$$

$$\left\langle s_i \tau_i \right\rangle_{RBN} \approx \frac{1}{N_s} \sum_m s_i^{T, (m)} \tau_i^{T, (m)}$$

iii) Update your parameters:

$$\theta_i \leftarrow \theta_i + \alpha (\langle s_i \rangle_D - \langle s_i \rangle_{RBN})$$

$$\eta_i \leftarrow \eta_i + \alpha (\langle \tau_i \rangle_D - \langle \tau_i \rangle_{RBN})$$

$$w_{ic} \leftarrow w_{ic} + \alpha (\langle s_i \tau_i \rangle_D - \langle s_i \tau_i \rangle_{RBN})$$

α is
the learning rate

iii) Mean-Field Phase diagram

can we understand what happens in the small couplings limit?

assume w_{ic} is small : $s_i \approx \langle s_i \rangle + \delta_i$ $\langle s_i \rangle = \langle \tau_c \rangle = 1/2$
 $\tau_c \approx \langle \tau_c \rangle + \delta_c$

$$\langle \tau_c \rangle_{1/2} + \delta_c = \text{sig} \left(\sum_i w_{ic} (\langle s_i \rangle + \delta_i) + \eta_c \right)$$

$$\left(\right) \delta_c \quad \eta_c = - \sum_i w_{ic} \langle s_i \rangle$$

$$= \frac{1}{1 + e^{-\sum_i w_{ic} \delta_i}} \approx \frac{1}{1 + (1 - \sum_i w_{ic} \delta_i)}$$

$$\approx \frac{1}{2(1 - \sum_i \frac{w_{ic} \delta_i}{2})} \approx \frac{1}{2} \left[1 + \frac{1}{2} \sum_i w_{ic} \delta_i \right]$$

$$\frac{1}{2} + \delta_c \approx \frac{1}{2} + \frac{1}{4} \sum_i w_{ic} \delta_i$$

$$\delta_c \approx \frac{1}{4} \sum_i w_{ic} \delta_i \quad \rightarrow \quad \vec{\delta}^{hid} = \frac{1}{4} W^T \cdot \vec{\delta}^{vis}$$

$$\vec{\delta}^{vis} = \frac{1}{4} W \cdot \vec{\delta}^{hid}$$


These perturbations are amplified if the strongest eigenvalue of W :

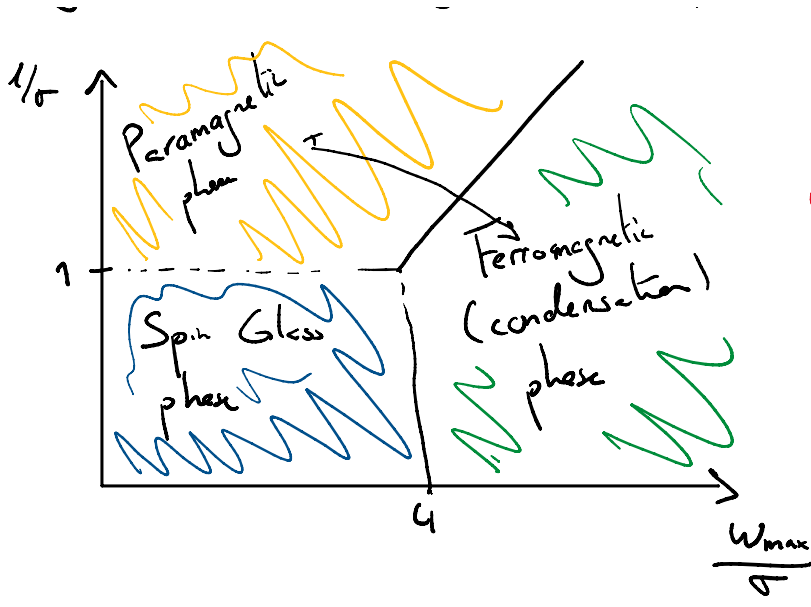
$$\boxed{w_{max} > 4}$$

otherwise they are killed

If we assume the following decomposition : $w_{ic} = \sum_{\alpha=1}^k u_i^\alpha w_\alpha v_c^\alpha + r_{ic}$

1/2 ↑ 

 ↓
rdm, gauss.



a contribution of $K=O(1)$
isolated eigenvalues

rdm Gauss. noise with variance σ

U & V are left and right eigenvectors

Singular Value Decomposition

Ferromagnetic phase: the system condenses over the strongest mode of w_{ic}

$$\begin{cases} \bar{m}_x = \frac{1}{\sqrt{N_s}} \sum_i \langle s_i \rangle u_i^x \\ \bar{m}_x = \frac{1}{\sqrt{N_h}} \sum_a \langle \tau_a \rangle v_a^x \end{cases} \quad \text{becomes } O(1)$$

IV) Learning behavior

⚠ for this section I will consider $\begin{cases} s_i = \pm 1 \\ \tau_a = \pm 1 \end{cases}; \begin{cases} \theta_i = 0 \\ \eta_a = 0 \end{cases}$

We have seen that the SVD decomposition of w_i plays an important role.

$w_{ic} = \sum_{\alpha=1}^{\min(N_s, N_h)} u_i^{\alpha} w_{\alpha} v_j^{\alpha}$

→ project the gradient over \vec{u}^{α} and \vec{v}^{α}

$$\left(\frac{\partial \mathcal{L}}{\partial w_{ic}} \right) = \sum_{\alpha=1}^{\min(N_s, N_h)} u_i^{\alpha} \left[\langle \tau_a \rangle - \langle s_i \tau_a \rangle \right] v_j^{\alpha}$$

$$\left(\frac{\partial L}{\partial \underline{w}}\right)_{\alpha p} := \sum_{i, r} u_i^{\alpha} \frac{\partial L}{\partial w_{i r}} v_r^{\beta} = \sum_{i, r} u_i^{\alpha} \left[\langle s_i \tau_r \rangle_0 - \langle s_i \tau_r \rangle_{\text{RNN}} \right] v_r^{\beta}$$

$$\begin{bmatrix} s_{\alpha} = \sum_i s_i u_i^{\alpha} \\ \tau_{\alpha} = \sum_r \tau_r v_r^{\alpha} \end{bmatrix} = \langle s \tau \rangle_0 - \langle s \tau \rangle_{\text{RNN}} \quad (**)$$

Gradient upd

$$\underline{w}(t + \Delta t) \approx \underline{w}(t) + \Delta t \left[\frac{d\underline{w}}{dt} \right] = \frac{\partial L}{\partial \underline{w}}$$

$$\left(\frac{d\underline{w}}{dt}\right)_{\alpha p} = \sum_{i, r} u_i^{\alpha} \left(\frac{d}{dt} w_{i r}\right) v_r^{\beta} = \sum_{i, r} u_i^{\alpha} \left(\frac{d}{dt} \sum_{\delta} u_i^{\delta} w_{\delta r} v_r^{\beta}\right) v_r^{\beta}$$

$$= \sum_{i, r, \delta} \left[\underbrace{u_i^{\alpha} u_i^{\delta}}_{\delta_{\alpha \delta}} \frac{d w_{\delta r}}{dt} v_r^{\beta} + \underbrace{u_i^{\alpha} \frac{d u_i^{\delta}}{dt}}_{\vec{u}^{\alpha} \frac{d \vec{u}^{\delta}}{dt}} w_{\delta r} v_r^{\beta} + \underbrace{u_i^{\alpha} u_i^{\delta} w_{\delta r}}_{(1 - \delta_{\alpha p}) w_p} \frac{d v_r^{\beta}}{dt} v_r^{\beta} \right]$$

$\delta_{\alpha p} \frac{d w_{\alpha r}}{dt} \quad \vec{u}^{\alpha} \frac{d \vec{u}^{\delta}}{dt} w_{\alpha} (1 - \delta_{\alpha p}) \quad (1 - \delta_{\alpha p}) w_p \frac{d \vec{v}^{\beta}}{dt} \cdot \vec{v}^{\beta}$

$$= \delta_{\alpha p} \frac{d w_{\alpha r}}{dt} + (1 - \delta_{\alpha p}) \left[\vec{u}^{\alpha} \cdot \frac{d \vec{u}^{\beta}}{dt} w_{\alpha} + w_p v_r^{\beta} \cdot \frac{d \vec{v}^{\alpha}}{dt} \right] \quad (***)$$

using (*) & (**)

$$\frac{d w_{\alpha}}{dt} = \langle s \tau \rangle_0 - \langle s \tau \rangle_{\text{RNN}}$$

$$\underbrace{\vec{u}^{\alpha} \frac{d \vec{u}^{\beta}}{dt}}_{\text{infinitesimal}} w_{\alpha} = \int \left(\langle s \tau \rangle_{\delta p} - \langle s \tau \rangle_{\text{RNN}} \right) \rightarrow \text{rotation of the orthog. matrix } u$$

→ infinitesimal rotation operator over the vec. u

→ the same for v .

The orthog. matrix. u

→ when w_{ic} is small

$$\frac{\partial \mathcal{L}}{\partial w_{ic}} = \frac{1}{n} \sum_m s_i^{(m)} \text{th} \left(\sum_j w_{jc} s_j^{(m)} \right) - \langle s_i z_c \rangle_{\text{train}}$$

$$\approx \frac{1}{n} \sum_m s_i^{(m)} \sum_j w_{jc} s_j^{(m)} - w_{ic}$$

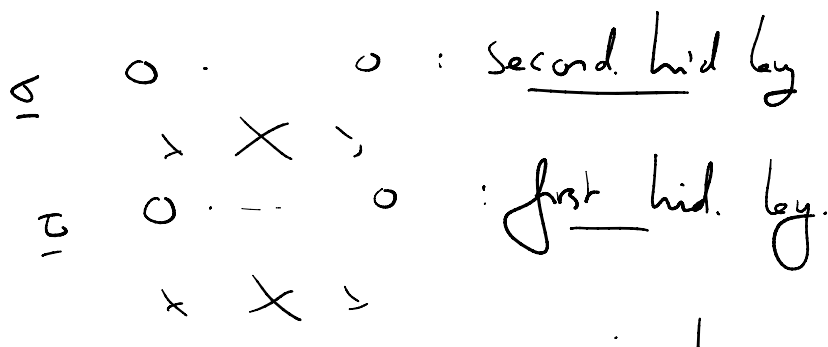
$$\approx \sum_j C_{ij} \cdot w_{jc} - w_{ic} \quad \text{when } C_{ij} = \frac{1}{n} \sum_m s_i^{(m)} s_j^{(m)}$$

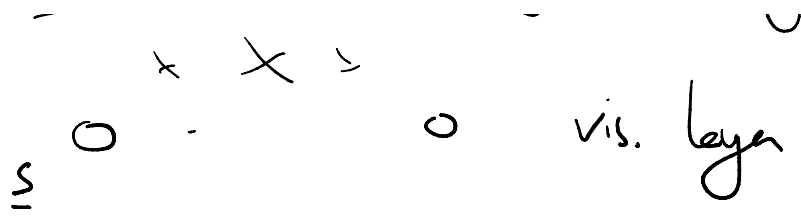
$$\frac{dw_{ic}}{dt} \approx w_{ic} \left[\langle s_i^2 \rangle - 1 \right]$$

→ transition behavior very similar to the one we get in GNN

C_{ij} projected on the U matrix

✓ Deep RBN?

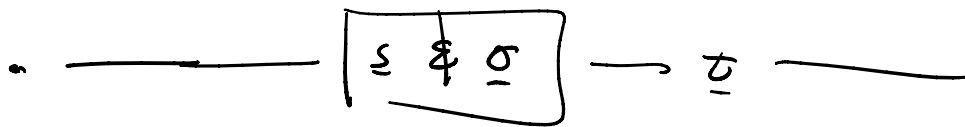




$$H_z = - \sum_{i < j} s_i w_{ij} \tau_j - \sum_i \theta_i s_i - \sum_j \eta_j \tau_j - \sum_{i < j} \tau_i w_{ij} \sigma_j - \sum_f \delta_f \sigma_f$$

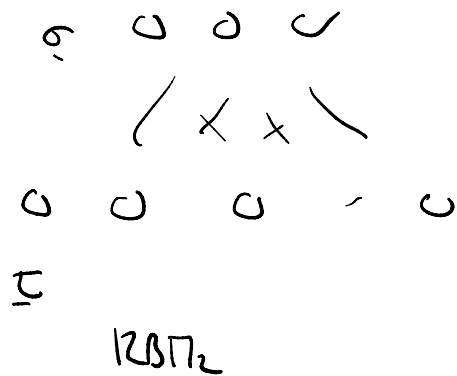
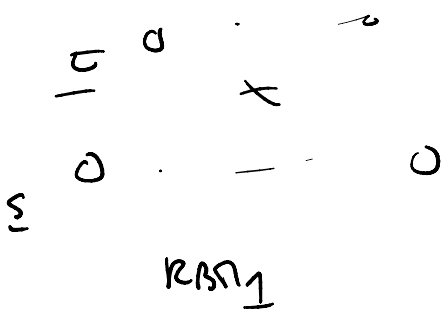
→ you still have conditional dependence:

- if you fix $\underline{\tau} \rightarrow \{ \underline{s}, \sigma \}$ becomes indep



→ Simplification: Deep. Belief Network

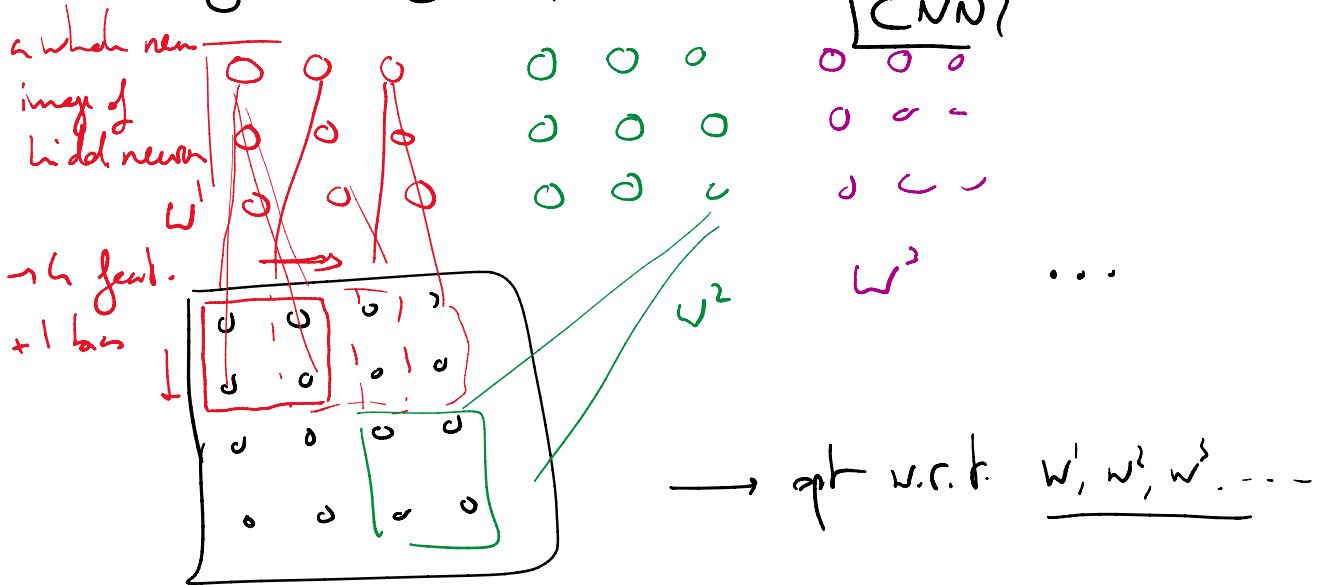
They consist in stacking RBM.



D] Few words on Conv. neural Network

& Generative adversarial Network

1) One of the very important layers: convolutional NN



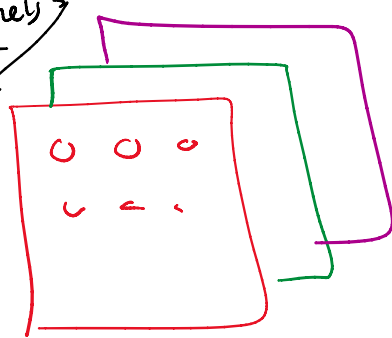
4x4 images

Invariant by translation

3 channels

after this:

this the output of your conv layer



we take this as a new input.

→ we make = perception layer



Im → conv → perception → softmax

Im \rightarrow conv \rightarrow perceptron \rightarrow softmax
lay (Dense)

This to CNN \rightarrow impressive results on image classification

1) GAN: generative model without sampling

Good Fellow

Ya design 2 neural networks

1 is the generator: from a set of $\sim N_I$ random input (e.g. Gaussian iid.)

it transforms it using a net: $G(\vec{z})$

$$\vec{z} \in \mathbb{R}^{M_1}$$

$$G(\vec{z}) \in \mathbb{R}^{M_x \times M_y}$$

2 The discriminator D it takes an image $N_x \times N_y$
 \rightarrow it outputs the prob that this image
belongs to the dataset or not

Training: 2 pass

1) you fix the weights of D

you opt G such that it "forces" D to answer that the samples from G are in the dataset

2) you fix G

you opt D to answer correctly

EJ Approximation to compute the negative terms

o] Normal Gibbs sampling

$\vec{s}^{(0)} \sim \text{unif}\{0,1\}$ \rightarrow MC sampling starting from $\vec{s}^{(0)}$
up to $\{\vec{s}^{(M)}\}$

$$\langle s_i; \tau_i \rangle_{\text{RSD}} \approx \frac{1}{N_s} \sum_{m=1}^{N_s} s_i^{(m)} p(\tau_i = 1 | \vec{s}^{(m)})$$

I] Hinton : Constructive Divergence (CD-h)

$\{\vec{s}^0\}$ = data used
in the minibatch

$$\langle s_i, t_i \rangle_{\text{RSP}} \approx \frac{1}{N} \sum_m \vec{s}^{(h), (m)} \cdot p(t_i = 1 | \vec{s}^{(h), (m)})$$

$\vec{s}^{(h)}$: vis after h NC steps

II] Persistent CD-h

for the first update : $\vec{s}^0 \sim \text{unif}\{0,1\}$

→ make h NC steps

$$\langle s_i, t_i \rangle_{\text{RSP}} \approx \frac{1}{N_s} \sum_m s_i^{(m)} p(t_i = 1 | \vec{s}^{(h), (m)})$$

Iter: for the successive updates

→ you take the previous final state
of the NC chain as initial ckt

for $\vec{s}^{(0)}$