

$S = k \cdot \log W$

AN INTRODUCTION TO UNSUPERVISED LEARNING AND BOLTZMANN MACHINES



Marco Baiesi

Department of Physics and Astronomy, Univ. Padova
& INFN

TNT seminar,
18.9.2020

Physics of Data, UNIPD

- Now starting 3rd year of master
- Growing number of students
 - Complex systems
 - high energy
 - astrophysics



Laboratory of Computational Physics B

- 24 hours theory & exercises
- This seminar is extracted from there
- Review:

Mehta et al, “A high-bias, low-variance introduction to Machine Learning for physicists”
Ph.Rep.810 (2019) 1–124

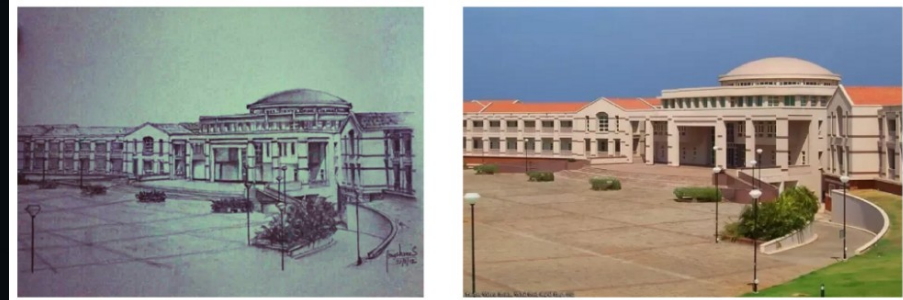
- 24 hours supervised projects in groups



Learning

- Supervised
 - Deep neural networks, regression,...
 - Labeled data
 - Optimization
 - Discriminative: obvious estimate of performances

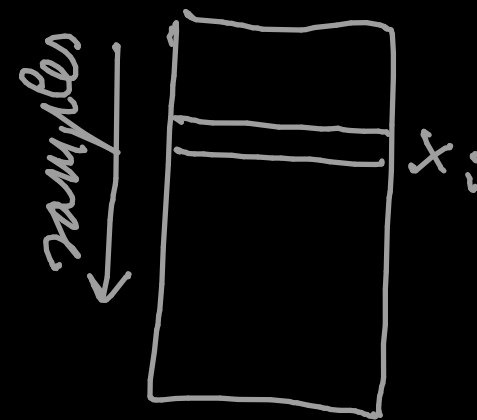
- Unsupervised
 - Boltzmann machines, variational autoencoders, generative adversarial nets
 - No labels
 - Generative



05-1

Unsupervised
Learning
(VL)

$\mathbf{x} = \{x_1, x_2, \dots\}$ data



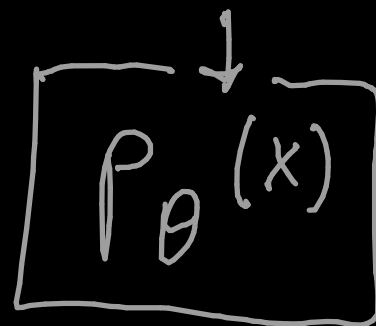
~~X~~

No labels

$p(\mathbf{x})$ probability distribution function (PDF)
of data, usually not known

$\mathbf{z} = \{z_1, z_2, \dots\}$ hidden, latent variables

$\theta = \{\theta_1, \theta_2, \dots\}$ parameters of the model



goal of UL:

represent "true" $p(x)$ of data
by approximate $p_\theta(x)$

generative models

- generating "fantasy" data x
- denoising
- filling missing data
- discrimination

Key quantities:

(information theory)

$$S_p = - \sum_i p(x_i) \log p(x_i)$$

Shannon
entropy

key quantities:

$$S_p = - \sum_i p(x_i) \log p(x_i)$$

Shannon
entropy

$$D_{KL}(p \parallel p') = \sum_i p(x_i) \log \frac{p(x_i)}{p'(x_i)}$$

Kullback
Leibler
divergence

key quantities:

$$S_p = - \sum_i p(x_i) \log p(x_i)$$

Shannon
entropy

$$D_{KL}(p \parallel p') = \sum_i p(x_i) \log \frac{p(x_i)}{p'(x_i)}$$

Kullback
Leibler
divergence

- $D_{KL} \geq 0$ (using $\log \frac{1}{\pi} \geq 1 - \pi$)
- $D_{KL}(p \parallel p') \neq D_{KL}(p' \parallel p)$

(relative
entropy)

Why physics in UL?

- similar problems

- variational free energy minimiz.
- Jaynes Max Ent
- Disordered systems, spin glasses, ...

Why physics in UL?

- similar problems

- useful setup

$$p(x) = \frac{1}{Z} e^{-\beta \bar{E}(x)}$$

$$\log p = -\beta \bar{E}(x) - Z$$

Why physics in UL?

- similar problems

- useful setup

$$p(x) = \frac{1}{Z} e^{-\beta \bar{E}(x)}$$

- training by physics methods, Monte Carlo (MC)
NOT via backpropagation, ~~Keras~~ (RBM)

however:

physics

$$\langle f \rangle_{\text{obs}}$$

conceptual average
with respect to model



UL

$$\langle f \rangle_{\text{data}}$$

empirical average
from data

however:

physics

$\langle f \rangle_{\text{obs}}$

conceptual average
with respect to model



UL

$\langle f \rangle_{\text{data}}$

empirical average
from data

- overfitting
- heterogeneity in precision

Log-likelihood maximization by varying params θ

$$\mathcal{L}(\theta) = \langle \log p_{\theta}(x) \rangle_{\text{data}}$$

review
(189)

$$= -\langle E_{\theta}(x) \rangle_{\text{data}} - \log Z_{\theta}$$

$$\boxed{\beta=1}$$

↑
No data in the
partition function!

$$Z_{\theta} = \sum_x p_{\theta}(x)$$

minus
log-likelihood ~~maximization~~
minimization

$$-\mathcal{L}(\theta) = -\langle \log p_{\theta}(x) \rangle_{\text{data}} \\ = +\langle E_{\theta}(x) \rangle_{\text{data}} + \log Z_{\theta}$$

energy is
minimized

↑
No data in the
partition function!

minus
log-likelihood maximization
minimization

$$-\mathcal{L}(\theta) = -\langle \log p_{\theta}(x) \rangle_{\text{data}} \quad (191)$$

$$= +\langle E_{\theta}(x) \rangle_{\text{data}} + \log Z_{\theta} + \underline{E_{\theta}^{\text{reg}}}$$

eventually, also
regularization term

$$E_{\theta}^{\text{reg}} = \lambda \sum_i |\theta_i| \quad (\text{LASSO})$$

$$E_{\theta}^{\text{reg}} = \lambda \sum_i |\theta_i|^2 \quad (\text{ridge})$$

Computing gradients

to minimize $-L(\theta)$ via e.g. stochastic gradient descent

define "operators"

$$O_j = \partial_{\theta_j} \bar{E}_{\theta}(x)$$

role of minus
force (193)

Computing gradients

to minimize $-L(\theta)$ via e.g. stochastic gradient descent

define "operators"

$$O_j = \partial_{\theta_j} \bar{E}_{\theta}(x)$$

role of minus
force (193)

$$\partial_{\theta_j}(-L(\theta)) = \langle \partial_{\theta_j} \bar{E}_{\theta}(x) \rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta}$$

Computing gradients

to minimize $-L(\theta)$ via e.g. stochastic gradient descent

define "operators"

$$O_j = \partial_{\theta_j} \bar{E}_{\theta}(x)$$

role of minus
force (193)

$$\partial_{\theta_j}(-L(\theta)) = \langle \partial_{\theta_j} \bar{E}_{\theta}(x) \rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta}$$

$$= \langle O_j(x) \rangle_{\text{data}} \longrightarrow \langle O_j(x) \rangle_{\text{model}} \quad (195)$$

$$Z_{\theta} = \sum_x p_{\theta}(x) = \sum_x e^{-\bar{E}_{\theta}(x)}$$

$$\partial_{\theta_j} \log Z_{\theta} = \frac{1}{Z_{\theta}} \sum_x \left(-\partial_{\theta_j} \bar{E}_{\theta}(x) \right) e^{-\bar{E}_{\theta}(x)}$$

$$= - \left\langle \partial_{\theta_j} \bar{E}_{\theta}(x) \right\rangle_{\text{model}}$$

$$= - \langle G_j \rangle_{\text{model}}$$

Computing gradient

$$\partial_{\theta_j}(-\mathcal{L}(\theta)) = \left\langle \partial_{\theta_j} E_{\theta}(x) \right\rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta}$$

$$= \left\langle O_j(x) \right\rangle_{\text{data}} \longrightarrow \left\langle O_j(x) \right\rangle_{\text{model}}$$

positive phase
of the gradient
(contains all info on data)

negative phase
...
(only model)

Computing gradients

$$\begin{aligned}\partial_{\theta_j}(-\mathcal{L}(\theta)) &= \left\langle \partial_{\theta_j} E_{\theta}(x) \right\rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta} \\ &= \left\langle O_j(x) \right\rangle_{\text{data}} \longrightarrow \left\langle O_j(x) \right\rangle_{\text{model}}\end{aligned}$$

Nice physical interpretation:

optimum when zero "force", i.e.
when expectation from model
equals " data

Computing gradients

$$\begin{aligned}\partial_{\theta_j}(-\mathcal{L}(\theta)) &= \left\langle \partial_{\theta_j} E_{\theta}(x) \right\rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta} \\ &= \left\langle O_j(x) \right\rangle_{\text{data}} \longrightarrow \left\langle O_j(x) \right\rangle_{\text{model}}\end{aligned}$$

- only in some Gaussian cases we have analytic solutions
- in general, intractable likelihood

to evaluate

$$\langle f(x) \rangle_{\text{model}} = \sum_x p_{\theta}(x) f(x)$$

to evaluate

$$\langle f(x) \rangle_{\text{model}} = \sum_x p_\theta(x) f(x) \approx \frac{\sum_{x'_i} f(x'_i)}{\sum_{x'_i} 1}$$

draw samples x'_i
from the model
according to p_θ ,
following Monte Carlo procedure

to evaluate

$$\langle f(x) \rangle_{\text{model}} = \sum_x p_\theta(x) f(x) \approx \frac{\sum_{x'_i} f(x'_i)}{\sum_{x'_i} 1}$$

draw samples x'_i
from the model
according to p_θ ,
following Monte Carlo procedure

x'_i , "fantasy particle"

to evaluate

$$\langle f(x) \rangle_{\text{model}} = \sum_x p_\theta(x) f(x) \approx \frac{\sum_{x'_i} f(x'_i)}{\sum_{x'_i} 1}$$

draw samples x'_i
from the model
according to p_θ ,
following Monte Carlo procedure

Normalization

x'_i , "fantasy particle"


log-derivative trick

to compute gradient of any $f(x)$

$$\partial_{\theta_j} \langle f(x) \rangle_{\text{model}} = \sum_i \partial_{\theta_j} p_{\theta}(x_i) f(x_i)$$

log-derivative trick

to compute gradient of any $f(x)$

$$\begin{aligned}\partial_{\theta_j} \langle f(x) \rangle_{\text{model}} &= \sum_i \partial_{\theta_j} p_{\theta}(x_i) f(x_i) \\ &= \langle \partial_{\theta_j} \log p_{\theta}(x) f(x) \rangle_{\text{model}}\end{aligned}$$


log-derivative trick

to compute gradient of any $f(x)$

$$\begin{aligned}\partial_{\theta_j} \langle f(x) \rangle_{\text{model}} &= \sum_i \partial_{\theta_j} p_{\theta}(x_i) f(x_i) \\ &= \langle \partial_{\theta_j} \log p_{\theta}(x) f(x) \rangle_{\text{model}} \quad \triangle! \\ &= \langle O_j(x) f(x) \rangle_{\text{model}}\end{aligned}$$

log-derivative trick

to compute gradient of any $f(x)$

$$\begin{aligned}\partial_{\theta_j} \langle f(x) \rangle_{\text{model}} &= \sum_i \partial_{\theta_j} p_{\theta}(x_i) f(x_i) \\ &= \langle \partial_{\theta_j} \log p_{\theta}(x) f(x) \rangle_{\text{model}} \quad \triangle! \\ &= \langle O_j(x) f(x) \rangle_{\text{model}} \\ &\approx \frac{\sum_{x'_i} O_j(x'_i) f(x'_i)}{\sum_{x'_i} 1} \quad (198)\end{aligned}$$

Summary of training procedure

goal: fit $\{\theta\}$ of model $p_{\theta}(x) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$

- train:
- 1) read minibatch B of data, $\{x\}_B$
 - 2) generate fantasy particles $\{x'\}_B \sim p_{\theta}(x)$ with MC
 - 3) compute gradients (195) (for B)
 - 4) update θ with gradient descent

05-2

Latent variables

and

Restricted Boltzmann Machines

Latent variables

enhance expressive power of generative models
by encoding complex correlations between data

Latent variables

enhance expressive power of generative models
by encoding complex correlations between data

$z \rightarrow h$ for "hidden" in this case

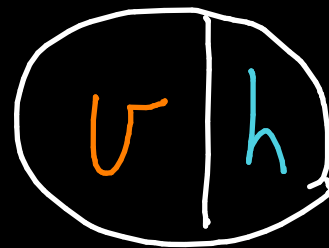
$x \rightarrow v$ for "visible"

Latent variables

enhance expressive power of generative models
by encoding complex correlations between data

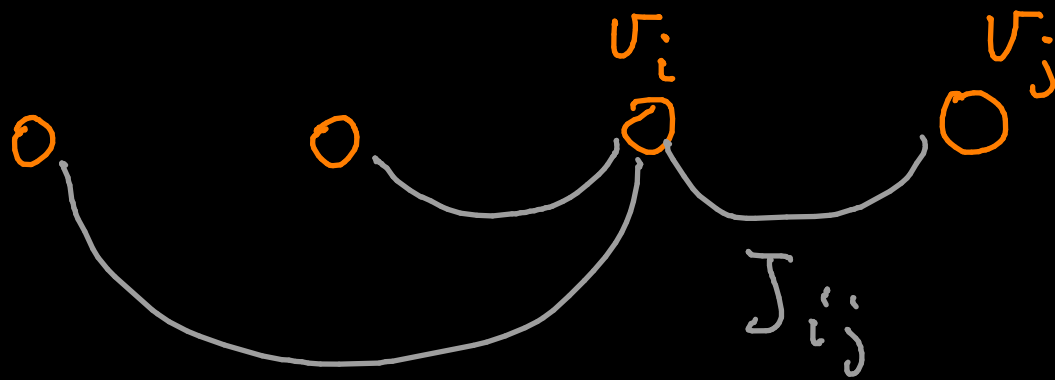
$z \rightarrow \underline{h}$ for "hidden" in this case

$x \rightarrow \underline{v}$ for "visible"



$v \cup h$ system

- spin systems (physics again relevant for VL...)



$i = \text{index of the spin (} i = 1, 2, \dots \text{)}$

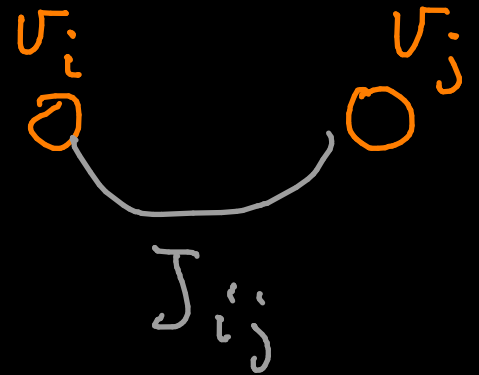
mean field: all couplings $J_{ij} \neq 0$

energy
$$E(v) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij} J_{ij} v_i v_j$$

$$J_{ij} = \sum_{\mu} W_{i\mu} W_{\mu j}$$

o

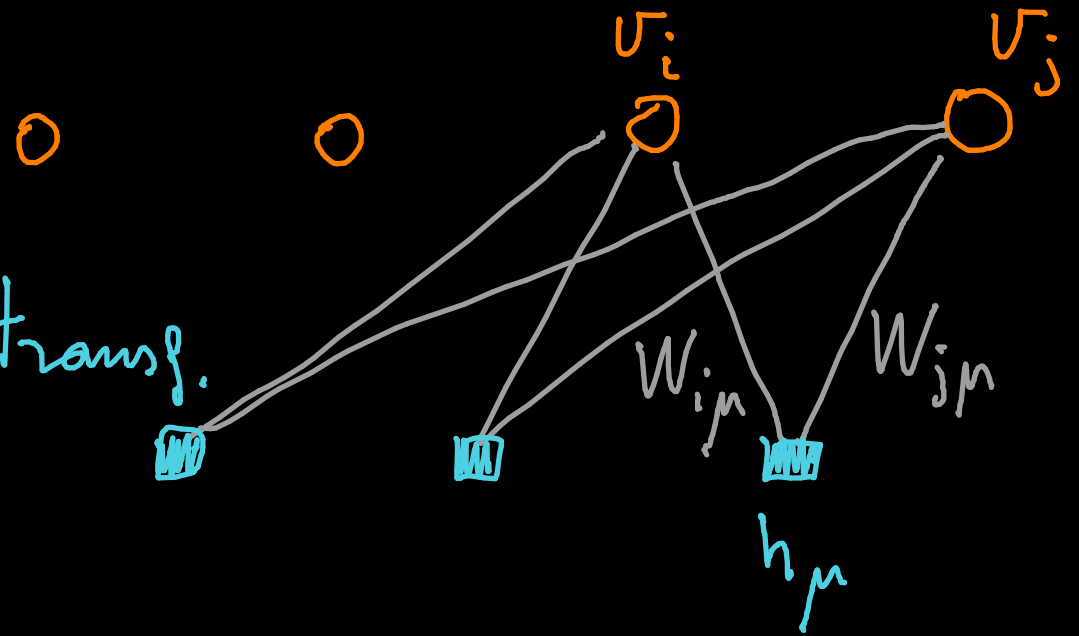
o



$$E(v) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij} J_{ij} v_i v_j$$

$$J_{ij} = \sum_{\mu} W_{i\mu} W_{j\mu}$$

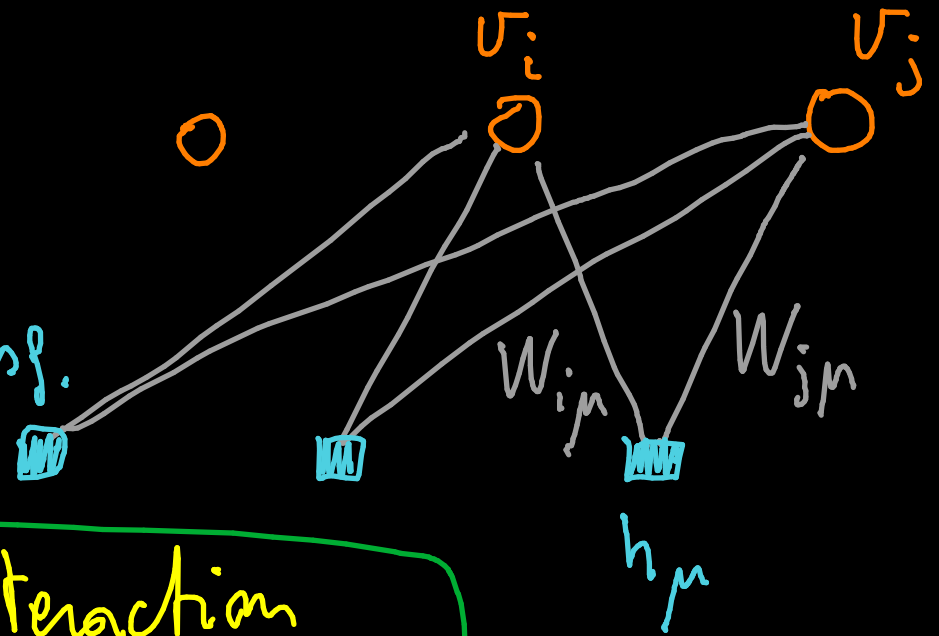
Hubbard-Stratonovich transf.
(h_{μ} 's with Gaussian stat.)



$$E(v) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij} J_{ij} v_i v_j$$

$$J_{ij} = \sum_{\mu} W_{i\mu} W_{j\mu}$$

Hubbard-Stratonovich transf.



J_{ij} removed: no direct interaction between "spins" v_i & v_j

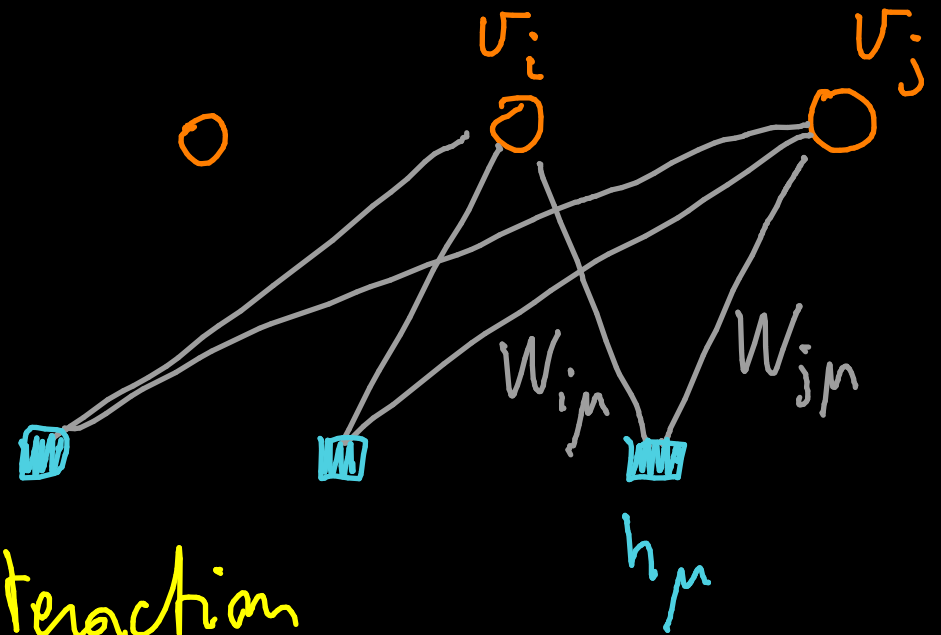
$$E(v) = - \sum_i a_i v_i - \frac{1}{2} \sum_{ij} J_{ij} v_i v_j$$

\Downarrow

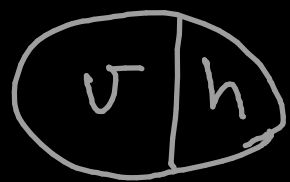
$$E(v, h) = - \sum_i a_i v_i + \frac{1}{2} \sum_{\mu} h_{\mu}^2 - \sum_{i\mu} v_i W_{i\mu} h_{\mu}$$

Visible Layer

hidden layer



\tilde{J}_{ij} removed: no direct interaction between "spins" v_i & v_j

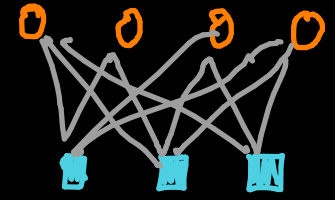


bipartite system

"Restricted"
(also NO $h_\mu h_\nu$ interaction)

$$E(v, h) = - \sum_i a_i v_i + \frac{1}{2} \sum_\mu h_\mu^2 - \sum_{i\mu} v_i W_{i\mu} h_\mu$$

Restricted Boltzmann Machines



inspired by previous considerations,

energy

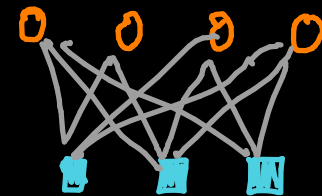
$$E(v, h) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i, \mu} v_i W_{i, \mu} h_{\mu}$$

functions

$a_i(\cdot)$

$b_{\mu}(\cdot)$

Restricted Boltzmann Machines



inspired by previous considerations,

energy $E(v, h) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i, \mu} v_i W_{i, \mu} h_{\mu}$

functions

$a_i(\cdot)$

$b_{\mu}(\cdot)$

Bernoulli layers

binary

$v_i \in \{0, 1\}$

Gaussian

$v_i \in \mathbb{R}$

$a_i(v_i)$

$a_i v_i$

$\frac{v_i^2}{2 \sigma_i^2}$

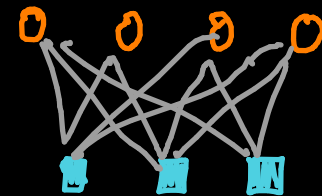
$b_{\mu}(h_{\mu})$

$b_{\mu} h_{\mu}$

$\frac{h_{\mu}^2}{2 \sigma_{\mu}^2}$

(also other versions,
see Mnason et al.)

Restricted Boltzmann Machines



inspired by previous considerations,

energy $E(v, h) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i, \mu} v_i W_{i, \mu} h_{\mu}$

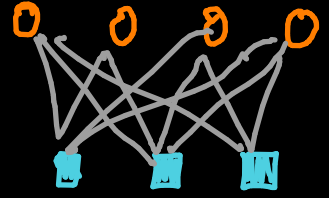
functions

$a_i(\cdot)$

$b_{\mu}(\cdot)$

	Bernoulli layers	Gaussian
	binary $v_i \in \{0, 1\}$	$v_i \in \mathbb{R}$
$a_i(v_i)$	$a_i v_i$	$\frac{v_i^2}{2 \sigma_i^2}$
$b_{\mu}(h_{\mu})$	$b_{\mu} h_{\mu}$	$\frac{h_{\mu}^2}{2 \sigma_{\mu}^2}$

Restricted Boltzmann Machines



energy $E(v, h) = - \sum_i a_i v_i - \sum_{\mu} b_{\mu} h_{\mu} - \sum_i v_i W_{i, \mu} h_{\mu}$

$$v_i, h_{\mu} = \begin{matrix} -1, 1 \\ \text{OR} \\ 0, 1 \end{matrix}$$

correlations induced by latent variables \rightarrow see the review

training

parameters $\theta = \{W_{i\mu}, a_i, b_\mu\}$

$$O_j = \partial_{\theta_j} E_\theta(v, h)$$

$$O_j(x) = O_j(v, h)$$

$$\partial_{\theta_j} (-\mathcal{L}(\theta)) = \langle O_j \rangle_{\text{data}} - \langle O_j \rangle_{\text{model}} \quad (195)$$

for example $\partial_{W_{i\mu}} E = -v_i h_\mu$

thanks to the
simple linear
appearance of
term $v_i W_{i\mu} h_\mu$

hence training via (195) follows these gradient components of $-L(\theta)$ to minimize it:

$$- \partial_{w_{i\mu}} L = \langle -v_i h_\mu \rangle_{\text{data}} - \langle -v_i h_\mu \rangle_{\text{model}}$$

$$- \partial_{a_i} L = \langle -v_i \rangle_{\text{data}} - \langle -v_i \rangle_{\text{model}}$$

$$- \partial_{b_\mu} L = \langle -h_\mu \rangle_{\text{data}} - \langle -h_\mu \rangle_{\text{model}}$$

hence training via (195) follows these
gradient components of $L(\theta)$ to maximize

$$\partial_{w_{i\mu}} L = \langle v_i h_\mu \rangle_{\text{data}} - \langle v_i h_\mu \rangle_{\text{model}}$$

$$\partial_{a_i} L = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}}$$

$$\partial_{b_\mu} L = \langle h_\mu \rangle_{\text{data}} - \langle h_\mu \rangle_{\text{model}}$$

maximize log-likelihood

$$\partial_{w_{i\mu}} \mathcal{L} = \langle v_i h_\mu \rangle_{\text{data}} - \langle v_i h_\mu \rangle_{\text{model}}$$

$$\partial_{a_i} \mathcal{L} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}}$$

$$\partial_{b_\mu} \mathcal{L} = \langle h_\mu \rangle_{\text{data}} - \langle h_\mu \rangle_{\text{model}}$$

same interpretation: optimum
where predictions of model match
the averages from data

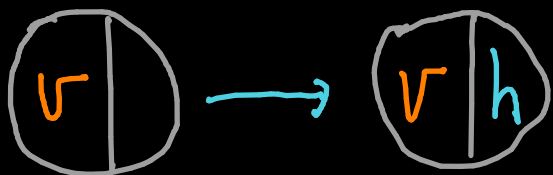
maximize log-likelihood

$$\partial_{w_{i\mu}} \mathcal{L} = \langle v_i h_\mu \rangle_{\text{data}} - \langle v_i h_\mu \rangle_{\text{model}}$$

$$\partial_{a_i} \mathcal{L} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}}$$

$$\partial_{b_\mu} \mathcal{L} = \langle h_\mu \rangle_{\text{data}} - \langle h_\mu \rangle_{\text{model}}$$

from "data"
 $v \cup h$



maximize log-likelihood

$$\partial_{w_{i\mu}} \mathcal{L} = \langle v_i h_\mu \rangle_{\text{data}} - \langle v_i h_\mu \rangle_{\text{model}}$$

$$\partial_{a_i} \mathcal{L} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}}$$

$$\partial_{b_\mu} \mathcal{L} = \langle h_\mu \rangle_{\text{data}} - \langle h_\mu \rangle_{\text{model}}$$

run MC to
generate v' & h'

Gibbs sampling

much simplified by bipartite structure of
restricted B.M. (no interaction between
 v 's and between h 's)

\Rightarrow conditionally independent variables

Gibbs sampling

much simplified by bipartite structure of
restricted B.M. (no interaction between
 v 's and between h 's)

\Rightarrow conditionally independent variables

$$p(v|h) = \prod_i p(v_i|h)$$

$$p(h|v) = \prod_{\mu} p(h_{\mu}|v)$$

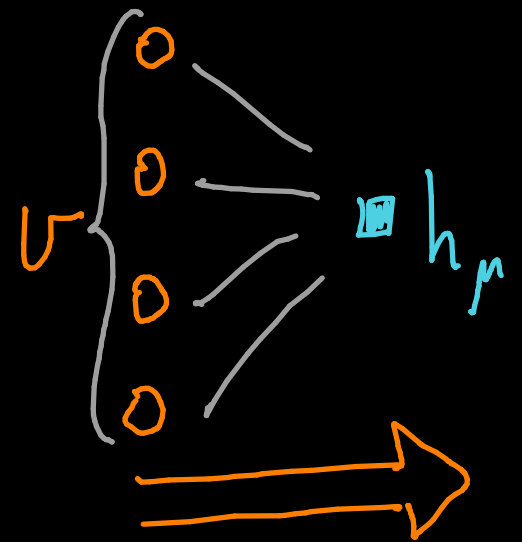
$$\begin{aligned} v &= \{v_i\} \\ h &= \{h_{\mu}\} \end{aligned}$$

(2.2)

probabilities are
factorized

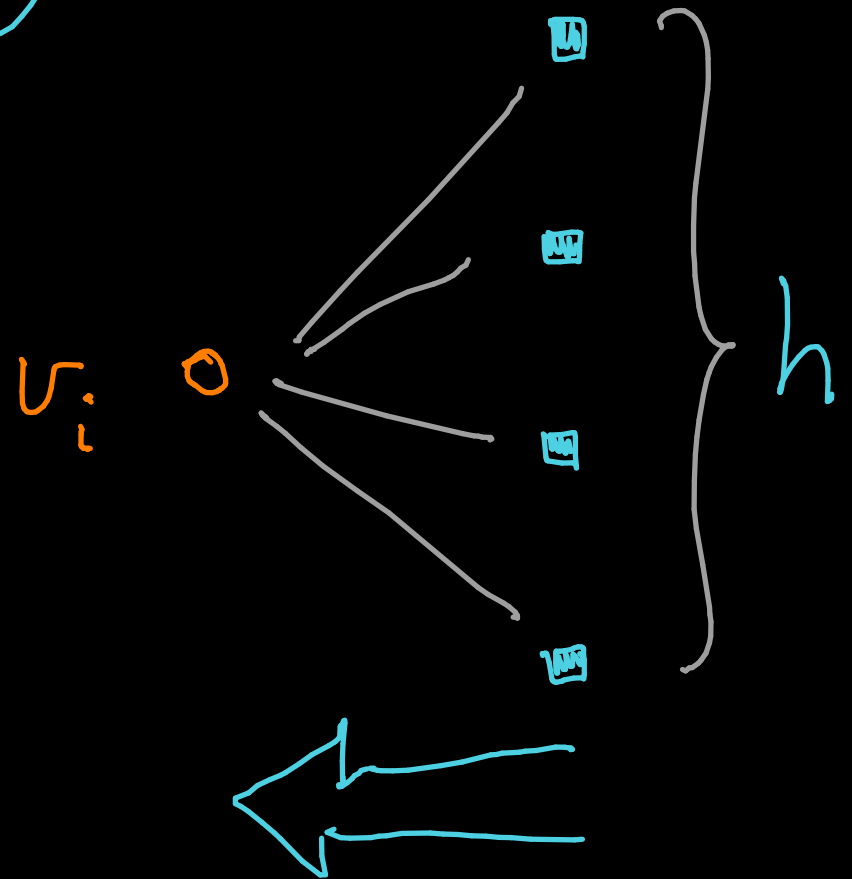
we can draw each h_μ independently
from the others ("restricted"!)
according to its
 $p(h_\mu | v)$

$$p(h | v) = \prod_{\mu} p(h_{\mu} | v)$$



probabilities are
factorized

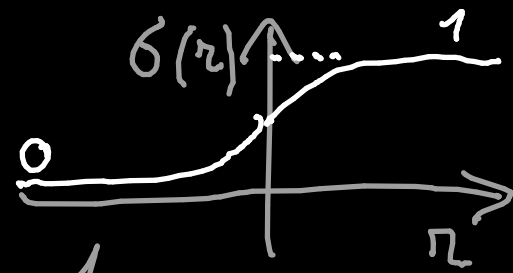
we can draw each U_i independently
from the others ("restricted"!)
according to its
 $p(U_i | h)$



for Bernoulli layers

(re)defining sigmoid

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



if $v_i = 0, 1$

$$p(v_i = 1 \mid h) = \sigma\left(a_i + \sum_{\mu} W_{i\mu} h_{\mu}\right)$$

$$p(h_{\mu} = 1 \mid v) = \sigma\left(b_{\mu} + \sum_i W_{i\mu} v_i\right)$$

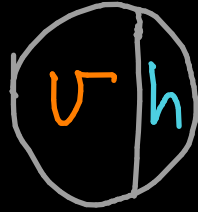
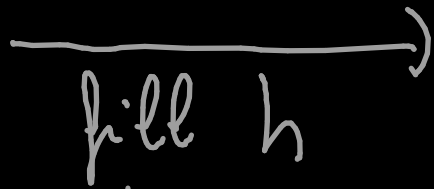
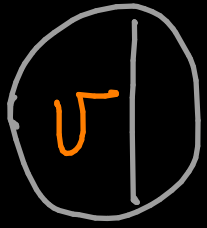
(213)

5

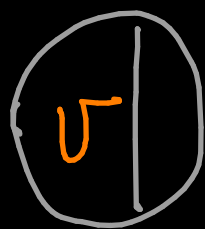
→

fill h

with
probabilities
from (213)



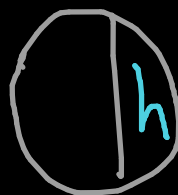
with
probabilities
from (213)

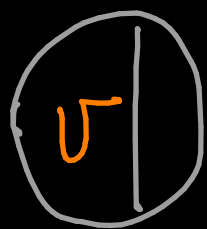


fill h
with
probabilities
from (213)

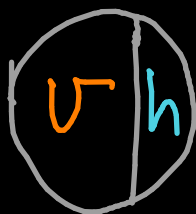


erase v

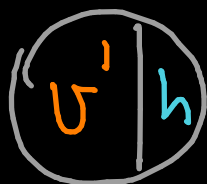
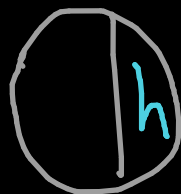




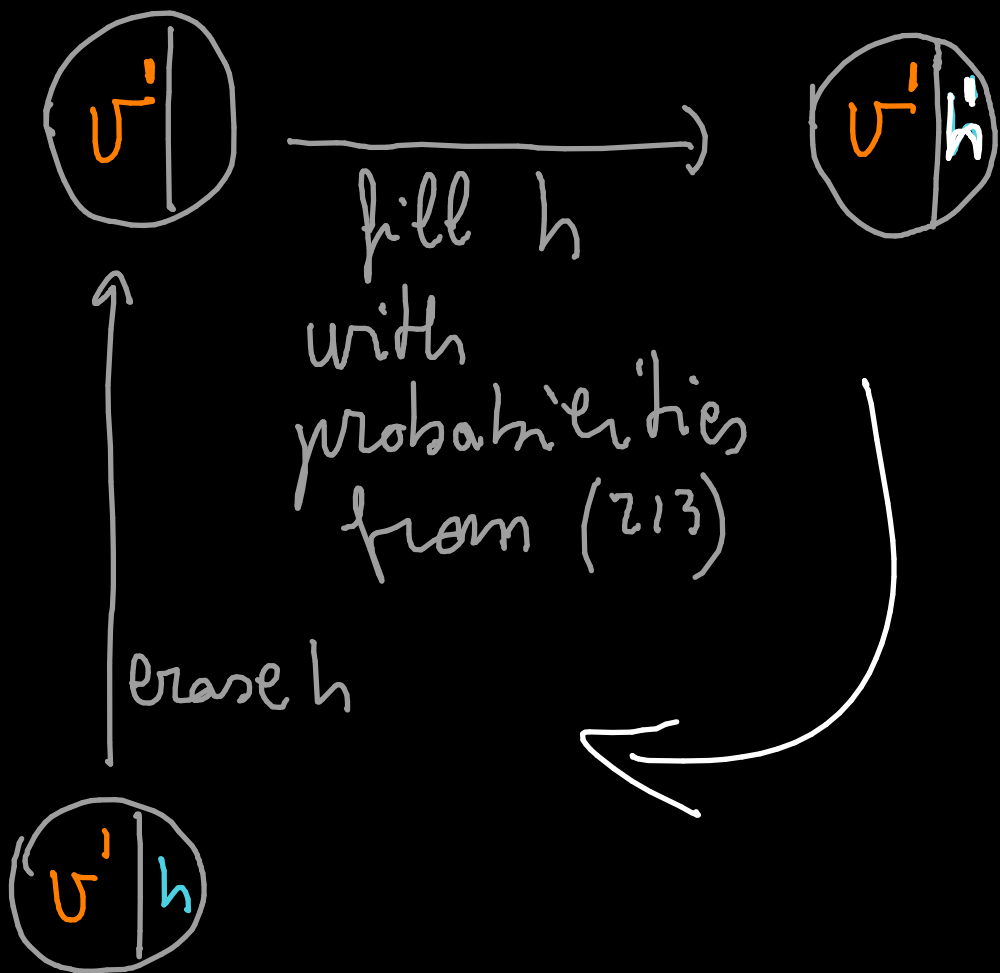
fill h
with
probabilities
from (213)

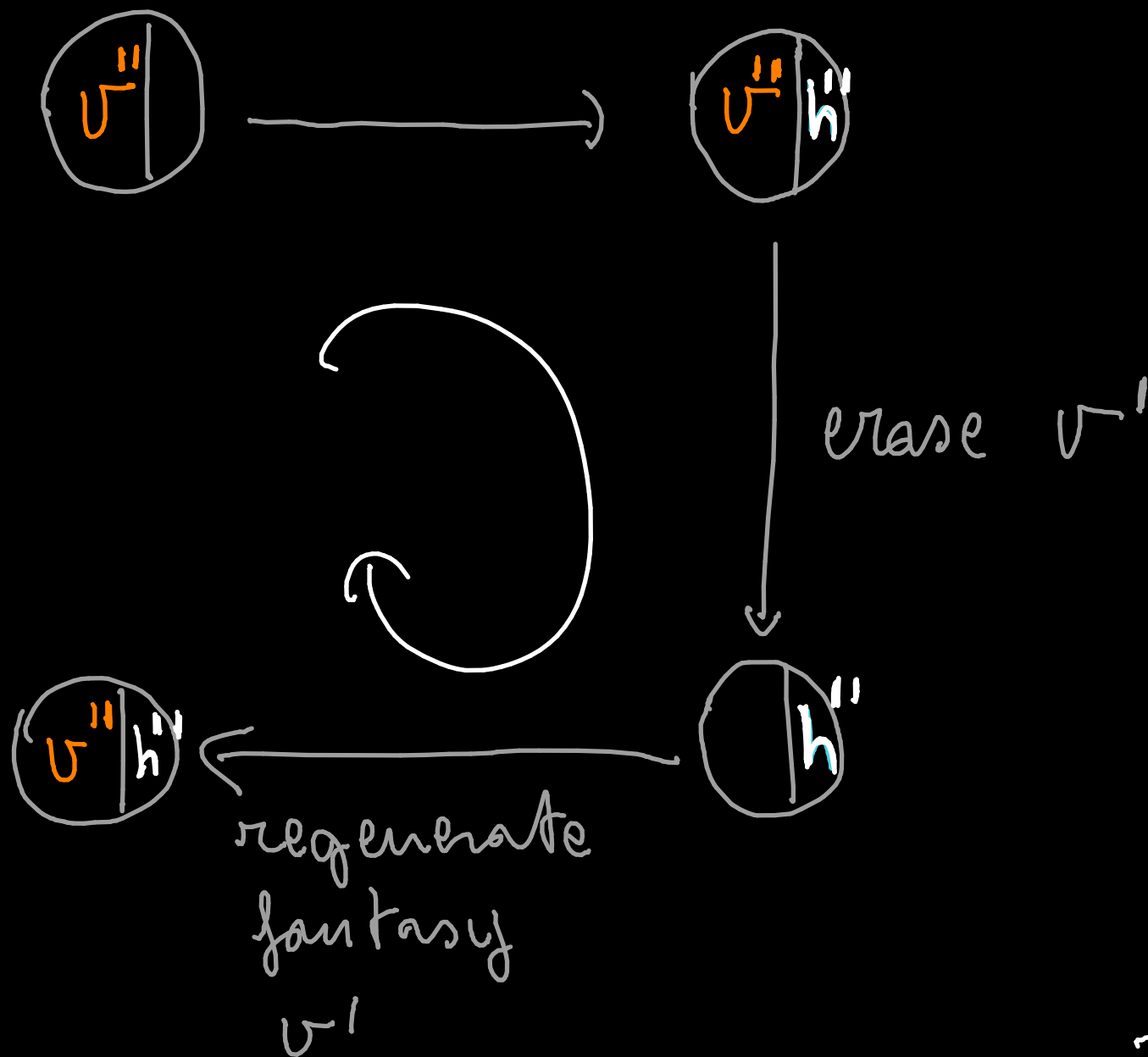


erase v



regenerate
fantasy
v'





repeat the process to sample the model averages

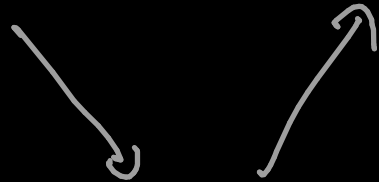
Alternating

Gibbs

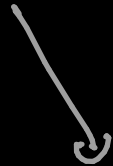
sampling

$U(0)$

$U(1)$

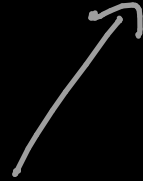


$h(0)$

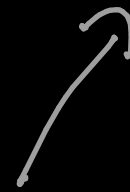


$h(1)$

$U(2)$

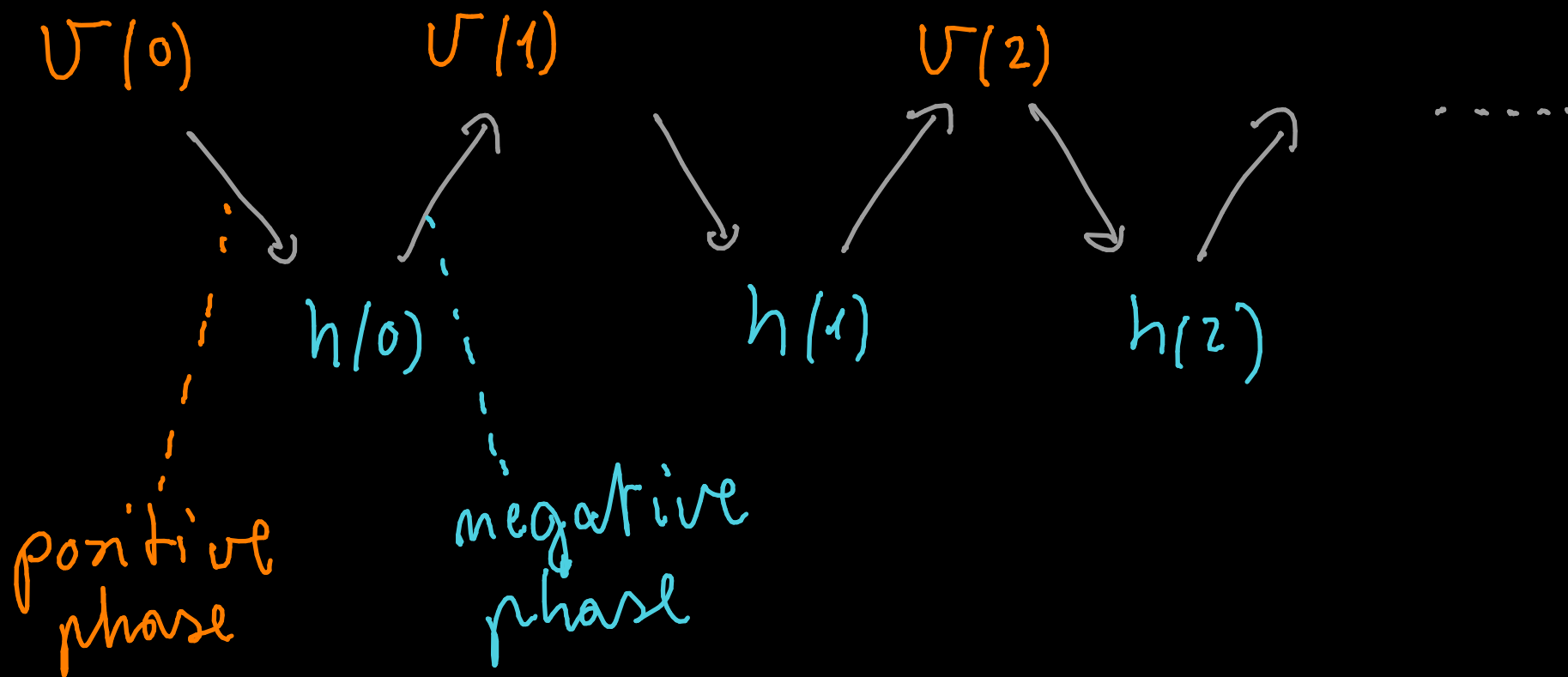


$h(2)$



...

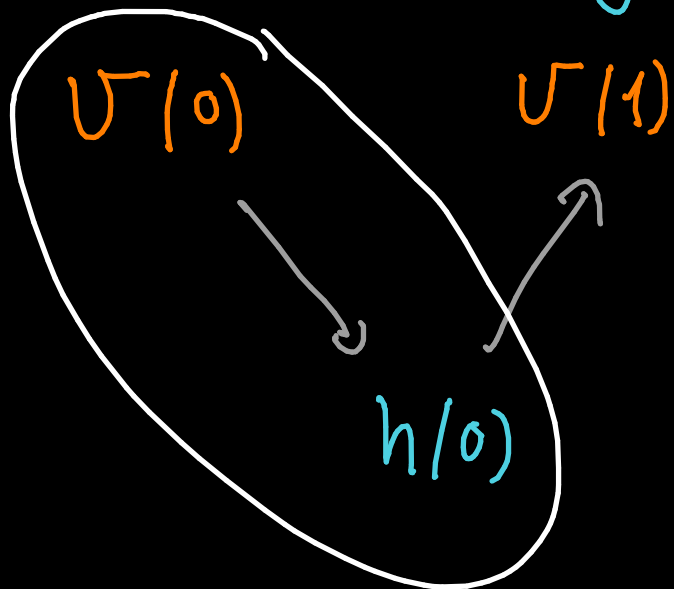
Alternating Gibbs sampling



Alternating

Gibbs

sampling

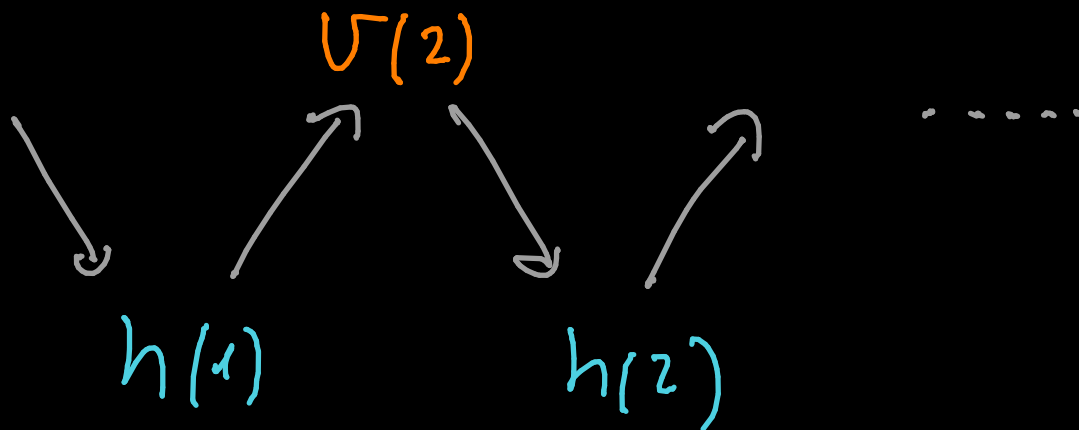


"data"

at $t=0$



$\langle \dots \rangle_{data}$



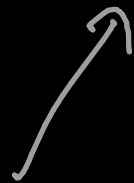
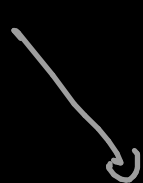
Alternating

Gibbs

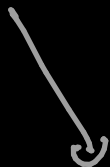
sampling

$U(0)$

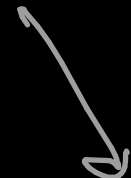
$U(1)$



$h(0)$



$h(1)$

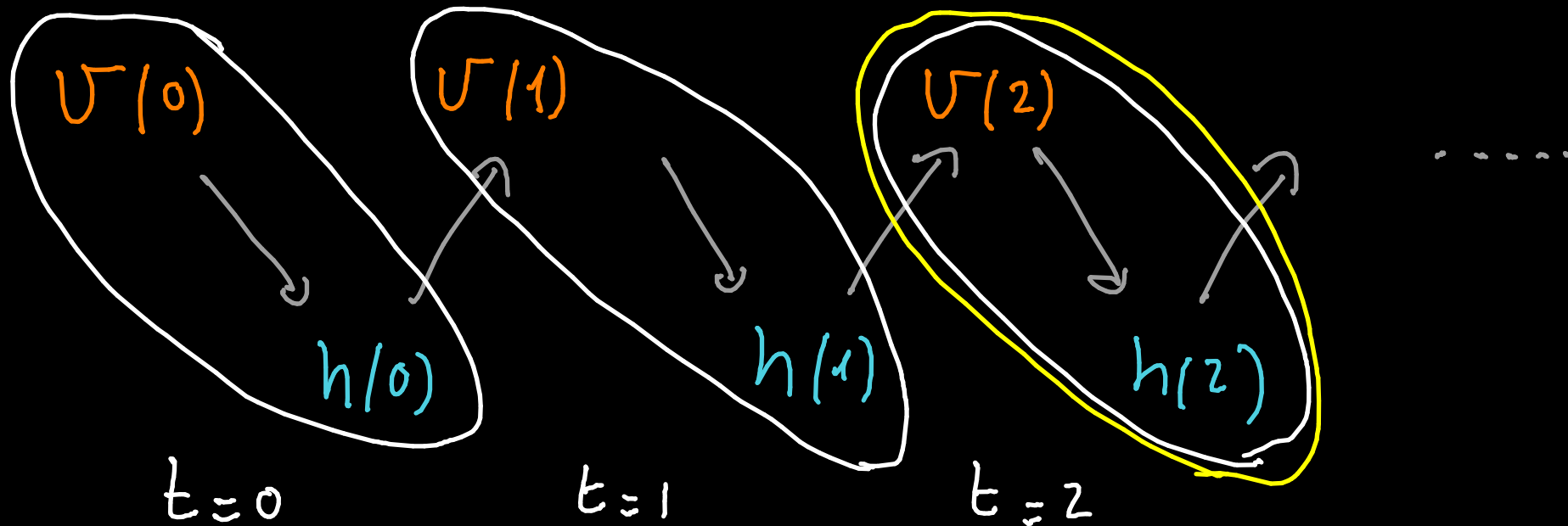


$h(2)$

...

"model"
at $t \rightarrow \infty$

Contrastive Divergence (CD-m)

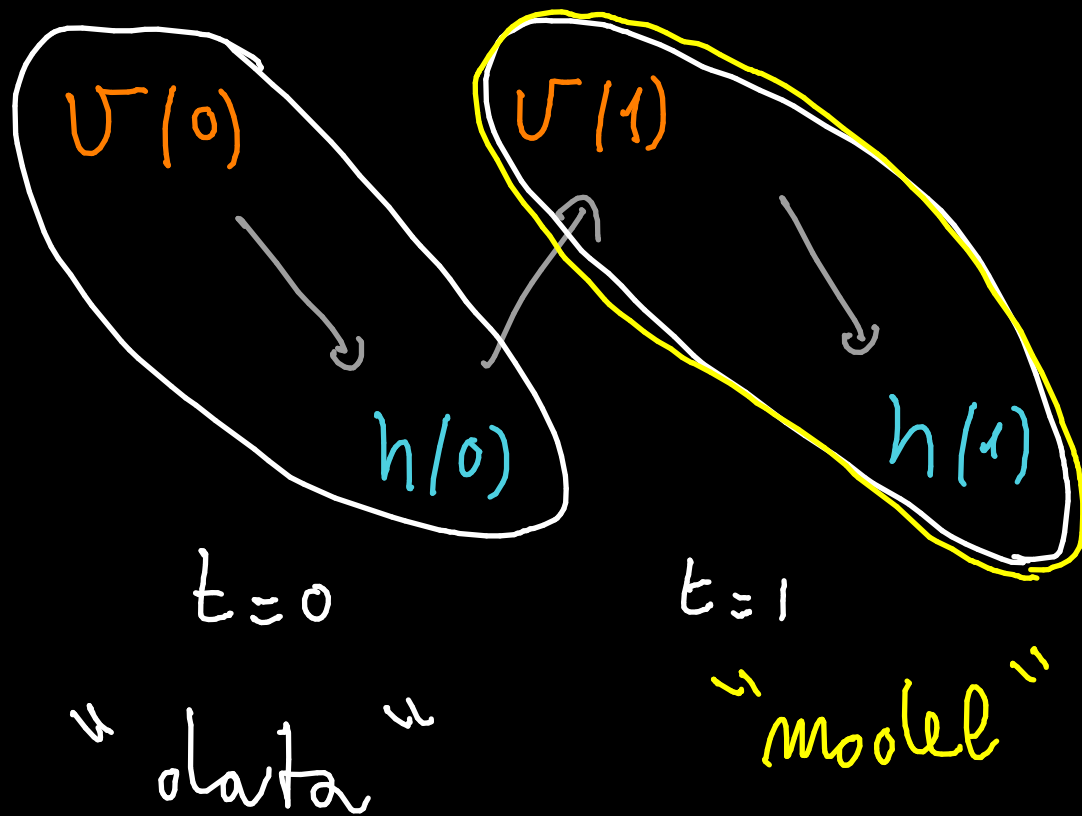


"data"

$m=2$ (for example)

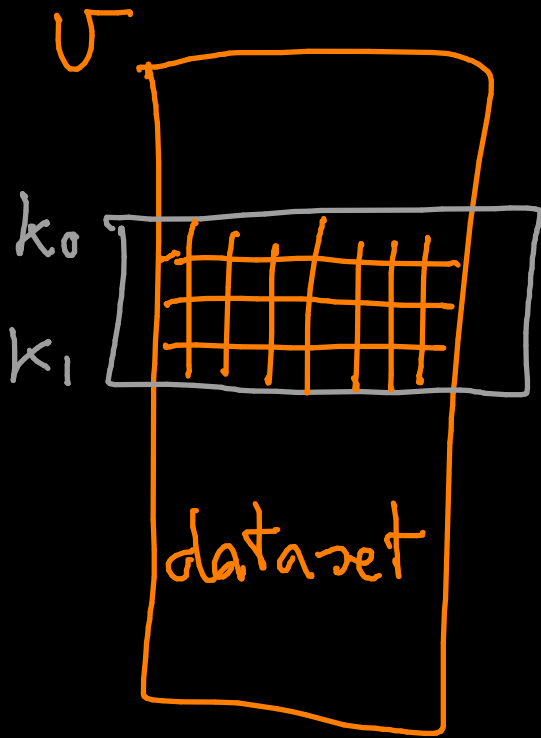
↓
"model" evaluated
at $t=m$
rather than $t \rightarrow \infty$

Contrastive Divergence (CD-1)



- most extreme example of CD
- fastest
- it works...

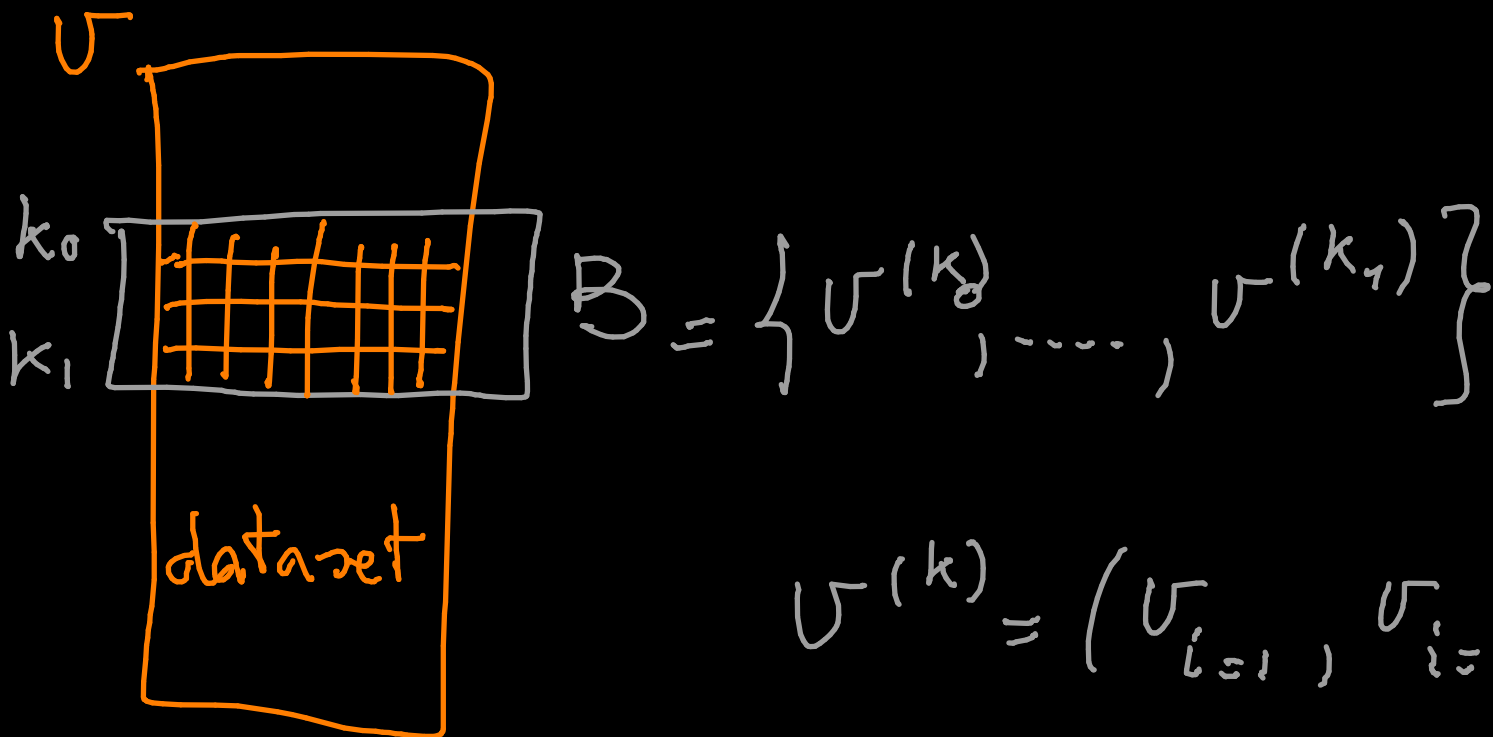
Mini batches



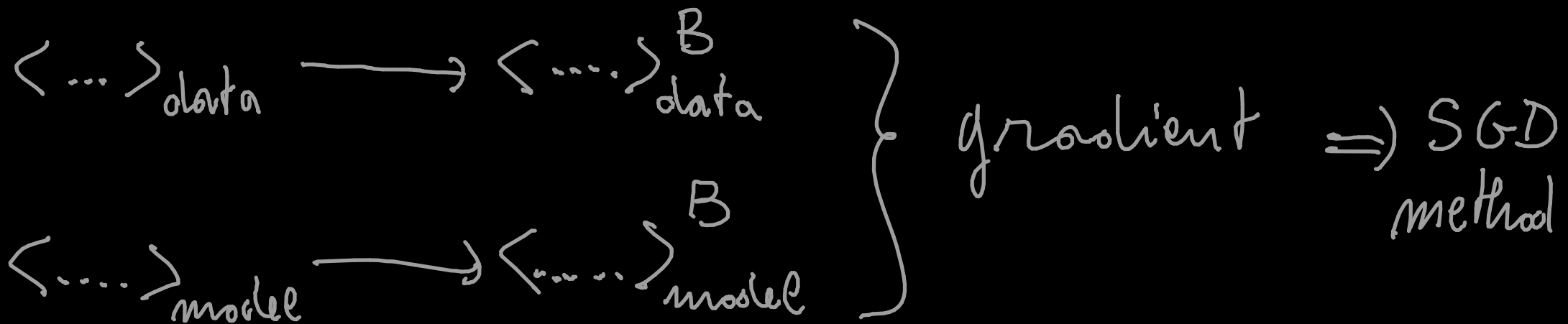
$$B = \{U^{(k_0)}, \dots, U^{(k_1)}\}$$

$$U^{(k)} = (U_{i=1}, U_{i=2}, \dots, U_{i=L})^{(k)}$$

Mini batches

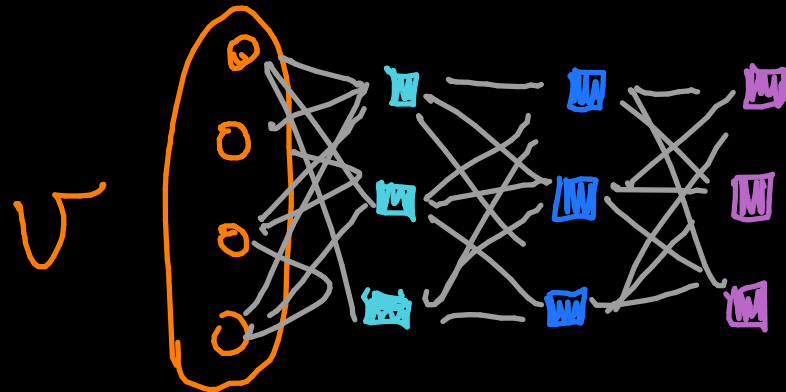


$$U^{(k)} = (U_{i=1}, U_{i=2}, \dots, U_{i=L})^{(k)}$$



More reading in the review:

- initialization
- regularization
- learning rates
- persistent contrastive divergence
- deep Boltzmann machines
(many hidden layers)



Summary: after training

- RBM has hidden layer that responds to data and can send back "fantasy data" with similar features
- generative
- denoising
- ...



Final example (Tubiana, Cocco, Monasson)

