

Latent variables

and

Restricted Boltzmann Machines

# Latent variables

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

K-means

$$\pi_{mm} \in 0, 1 \quad \leftarrow z$$

$$\mu_m \quad \leftarrow \theta$$

# Latent variables

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

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

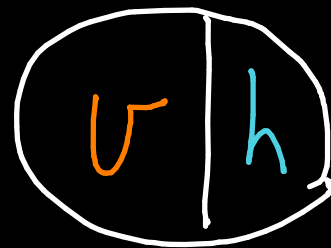
$x \rightarrow v$  for "visible"

# Latent variables

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

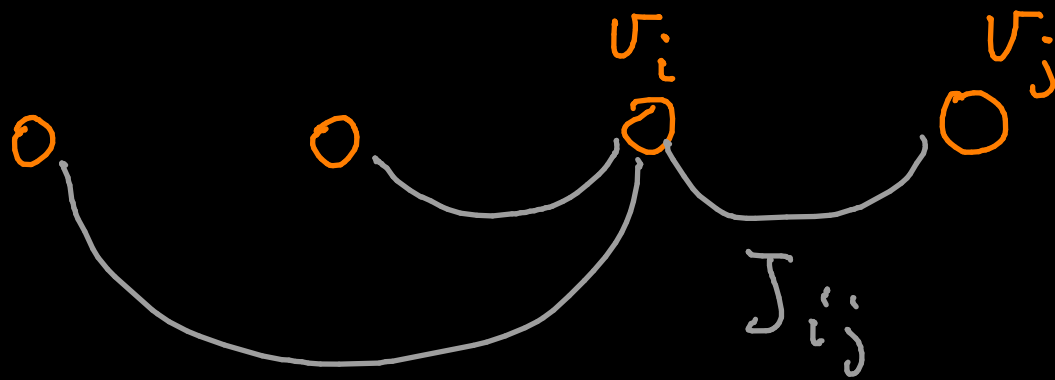
$z \rightarrow \underline{h}$  for "hidden" in this case

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



$v \cup h$  system

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



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

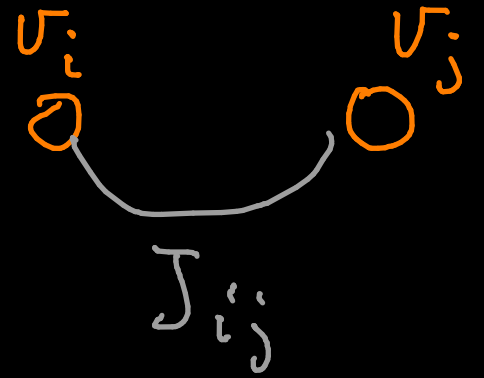
mean field: all couplings  $J_{ij} \neq 0$

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

$$J_{ij} = \sum_{\mu=1}^M W_{i\mu} W_{\mu j}$$

o

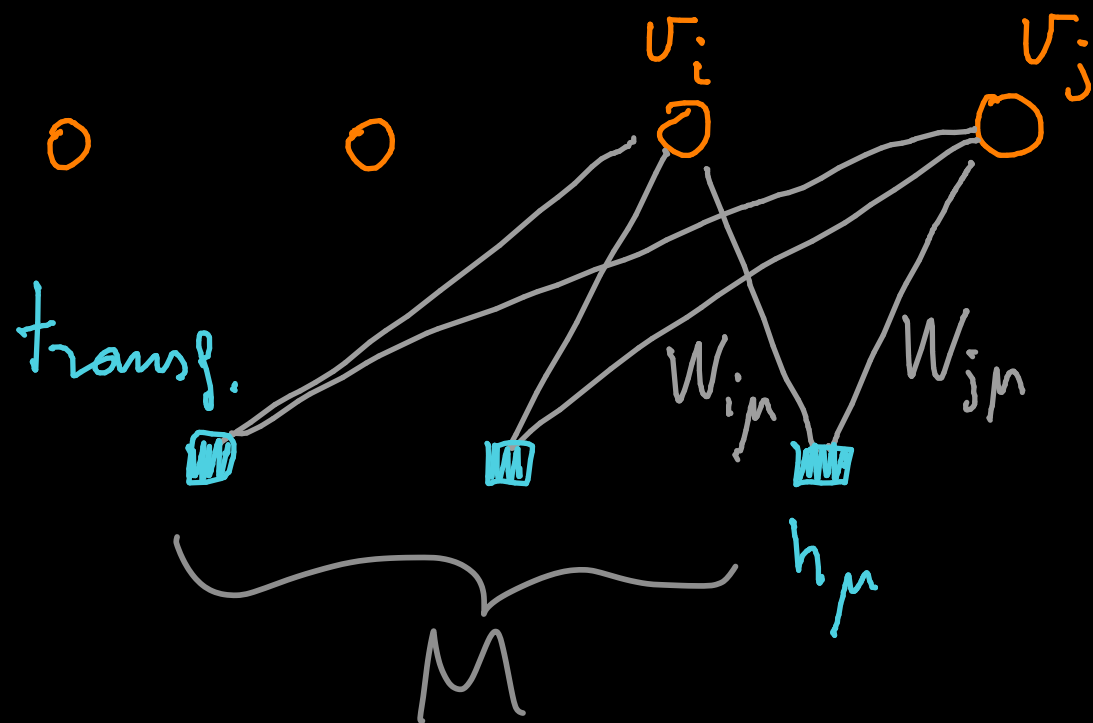
o



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

$$J_{ij} = \sum_{\mu=1}^M W_{i\mu} W_{j\mu}$$

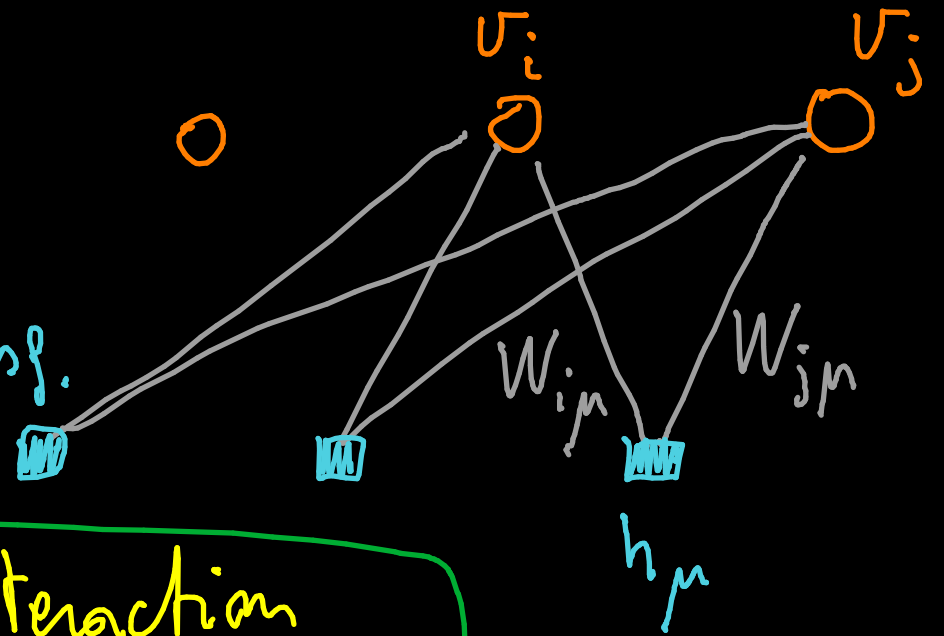
Hubbard-Stratonovich transg.  
( $h_\mu$ 's with Gaussian stat.)



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

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

Hubbard-Stratonovich transf.



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

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

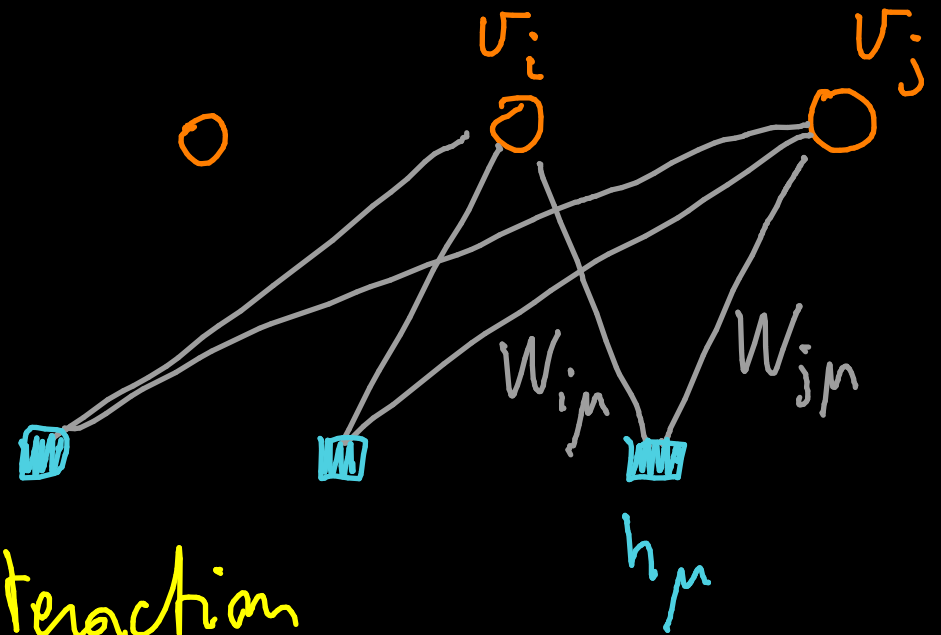
$\Downarrow$

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

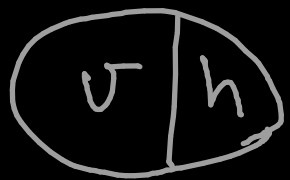


Visible Layer

hidden layer



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

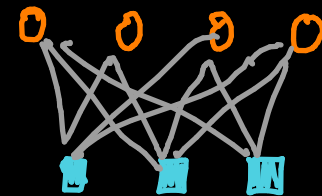


bipartite system

"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

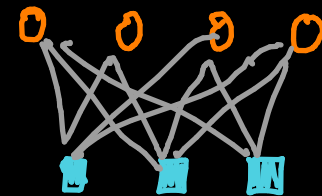
$$E(v, h) = - \sum_i a_i(v_i) - \sum_{\mu} b_{\mu}(h_{\mu}) - \sum_{i, \mu} v_i W_{i, \mu} h_{\mu}$$

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

$v_i \in \mathbb{R}$

$a_i(v_i)$

$a_i v_i$

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

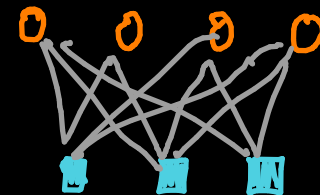
$b_{\mu}(h_{\mu})$

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

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

(also other versions,  
see Mnason et al.)

# Restricted Boltzmann Machines



inspired by previous considerations,

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

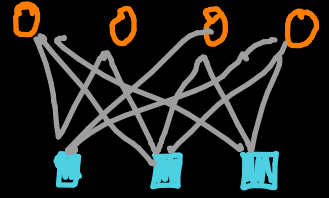
functions

$a_i(\cdot)$

$b_{\mu}(\cdot)$

	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

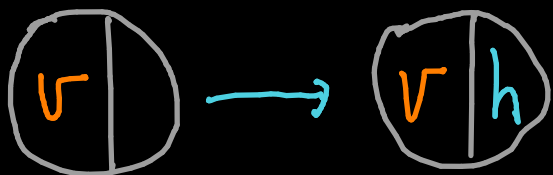
maximize log-likelihood

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

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

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

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

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

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

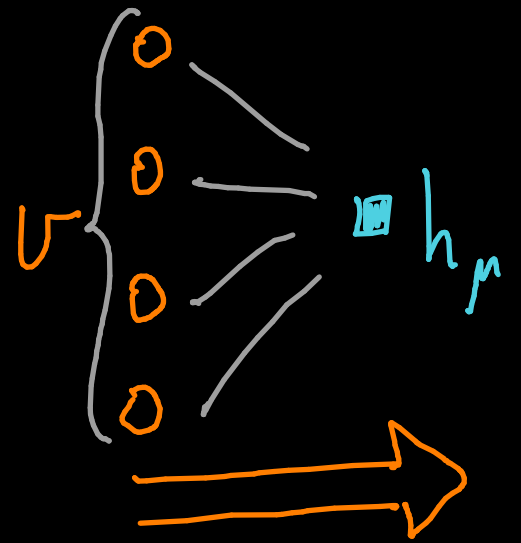
$$\begin{aligned} v &= \{v_i\} \\ h &= \{h_{\mu}\} \end{aligned}$$

(2.2)

probabilities are  
factorized

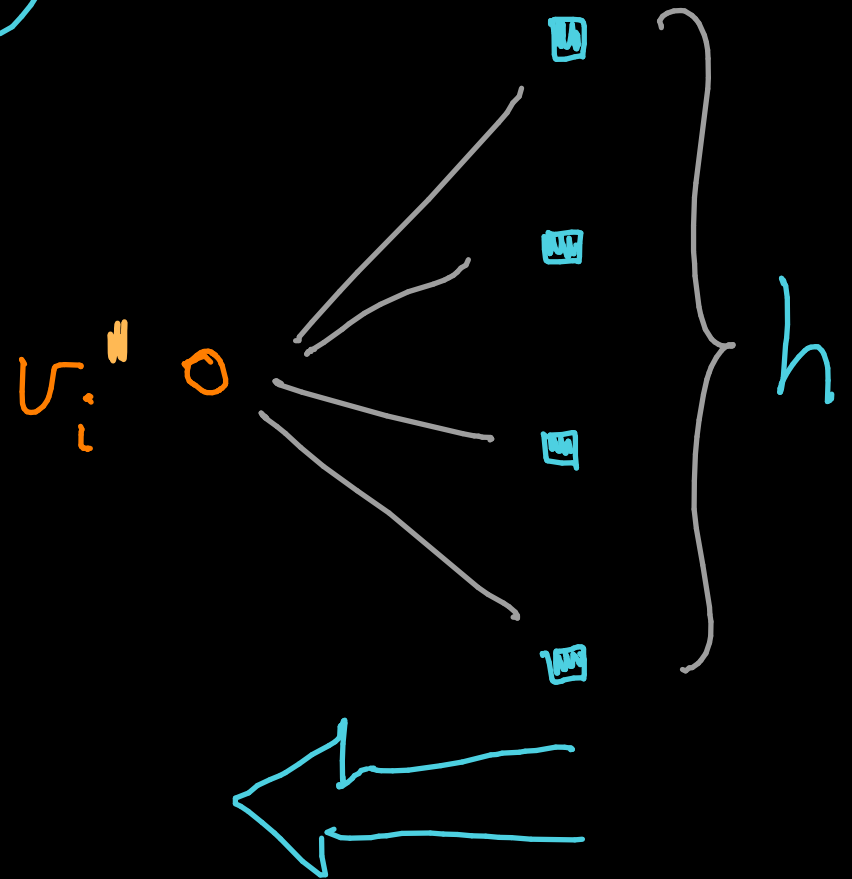
we can draw each  $h_\mu$  independently  
from the others ("restricted"!)  
according to its  
 $p(h_\mu | v)$

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



probabilities are  
factorized

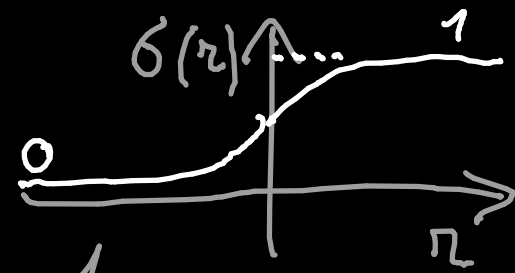
we can draw each  $U_i$  independently  
from the others ("restricted"!)  
according to its  
 $p(U_i | h)$



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)

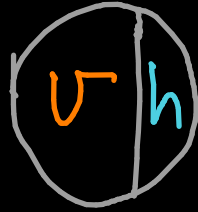
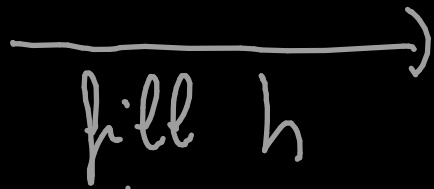
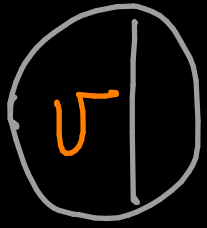


5

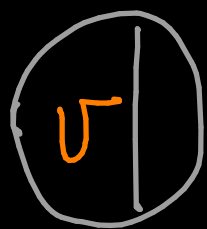
→

fill h

with  
probabilities  
from (213)



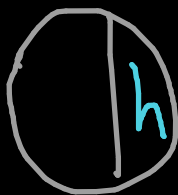
fill h  
with  
probabilities  
from (213)

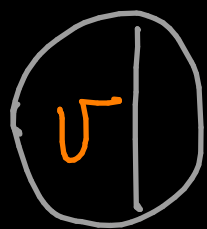


fill h  
with  
probabilities  
from (213)



erase v

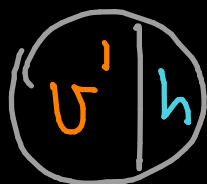
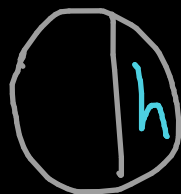




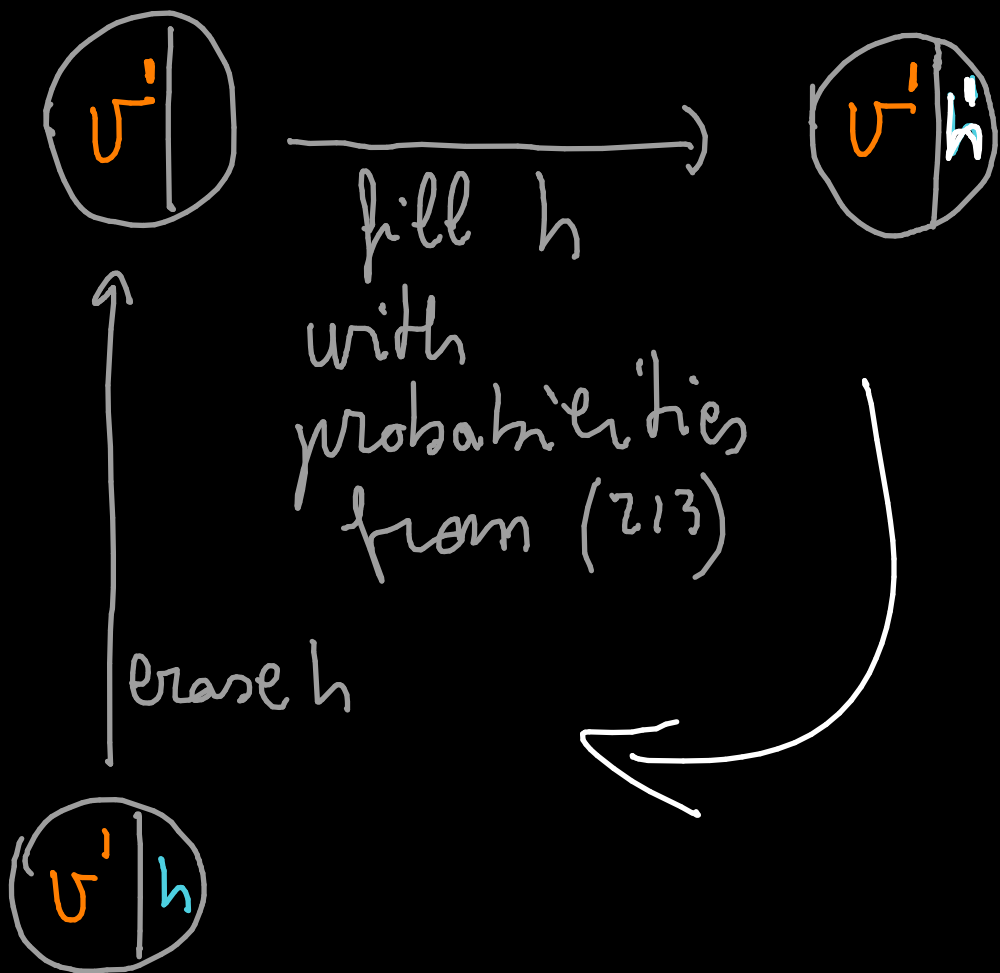
fill h  
with  
probabilities  
from (213)

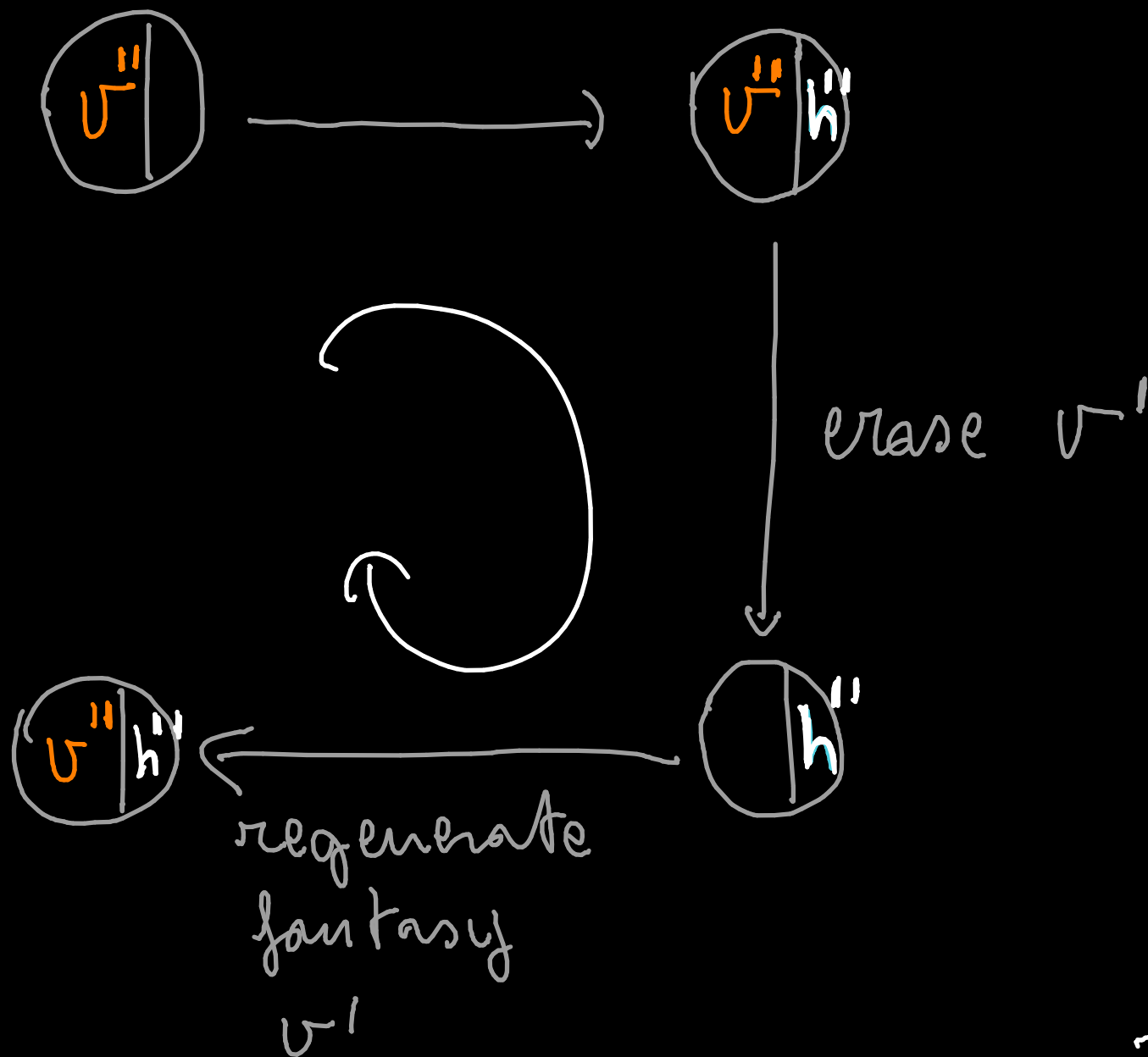


erase v



regenerate  
fantasy  
v'





repeat the  
process to  
sample the  
model averages

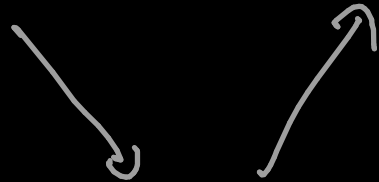
Alternating

Gibbs

sampling

$U(0)$

$U(1)$

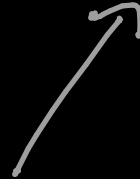


$h(0)$

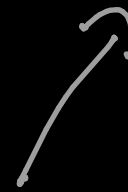


$h(1)$

$U(2)$

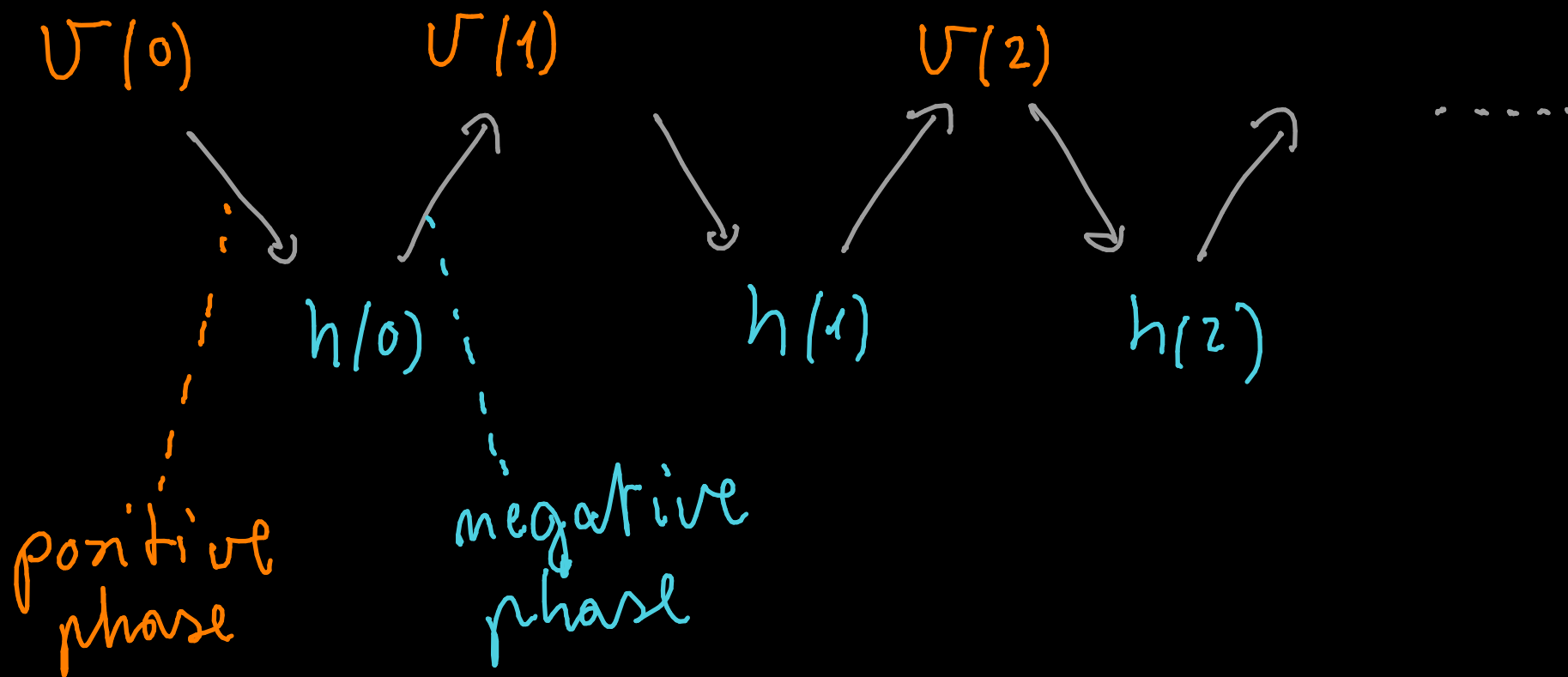


$h(2)$



...

# Alternating Gibbs sampling

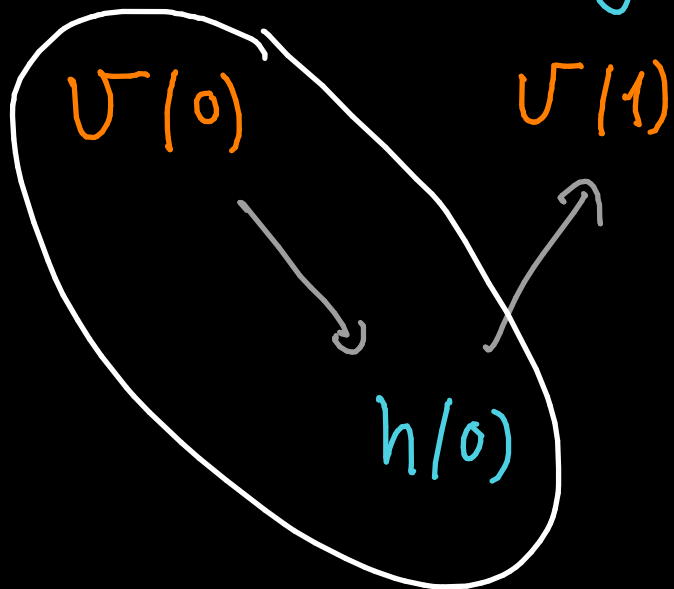




Alternating

Gibbs

sampling

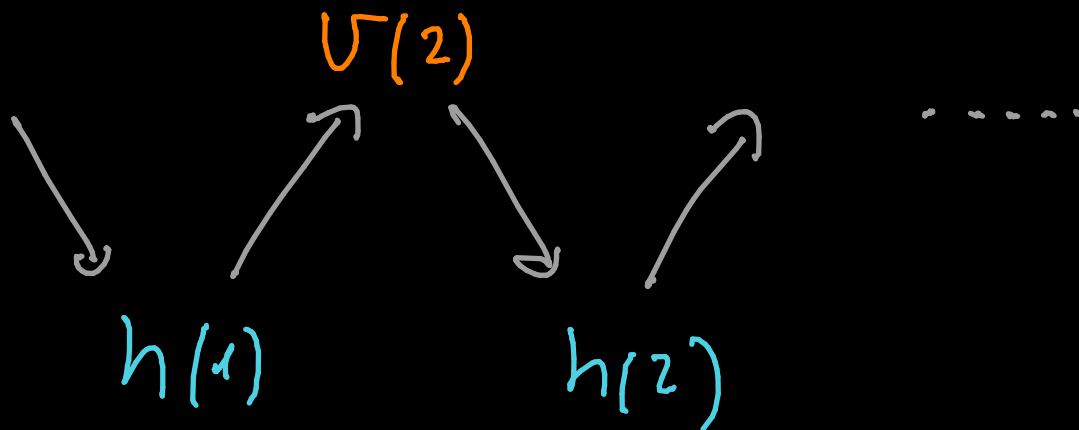


"data"

at  $t=0$



$\langle \dots \rangle_{data}$



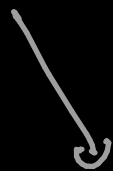
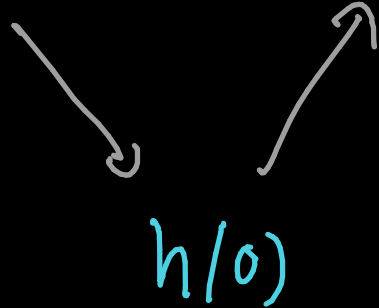
Alternating

Gibbs

sampling

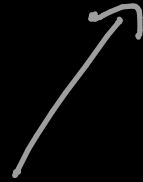
$U(0)$

$U(1)$



$h(1)$

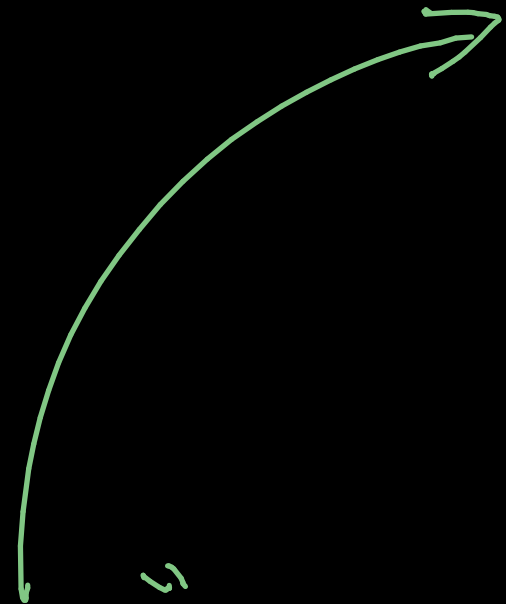
$U(2)$



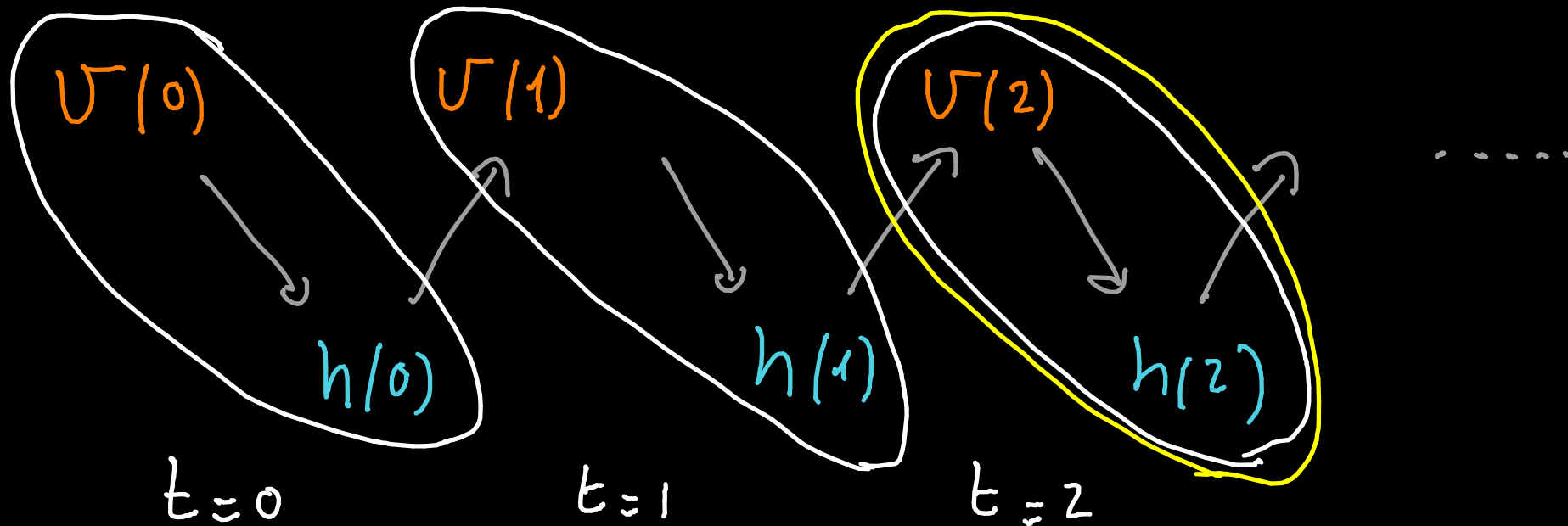
$h(2)$

...

"model"  
at  $t \rightarrow \infty$



# Contrastive Divergence (CD-m)

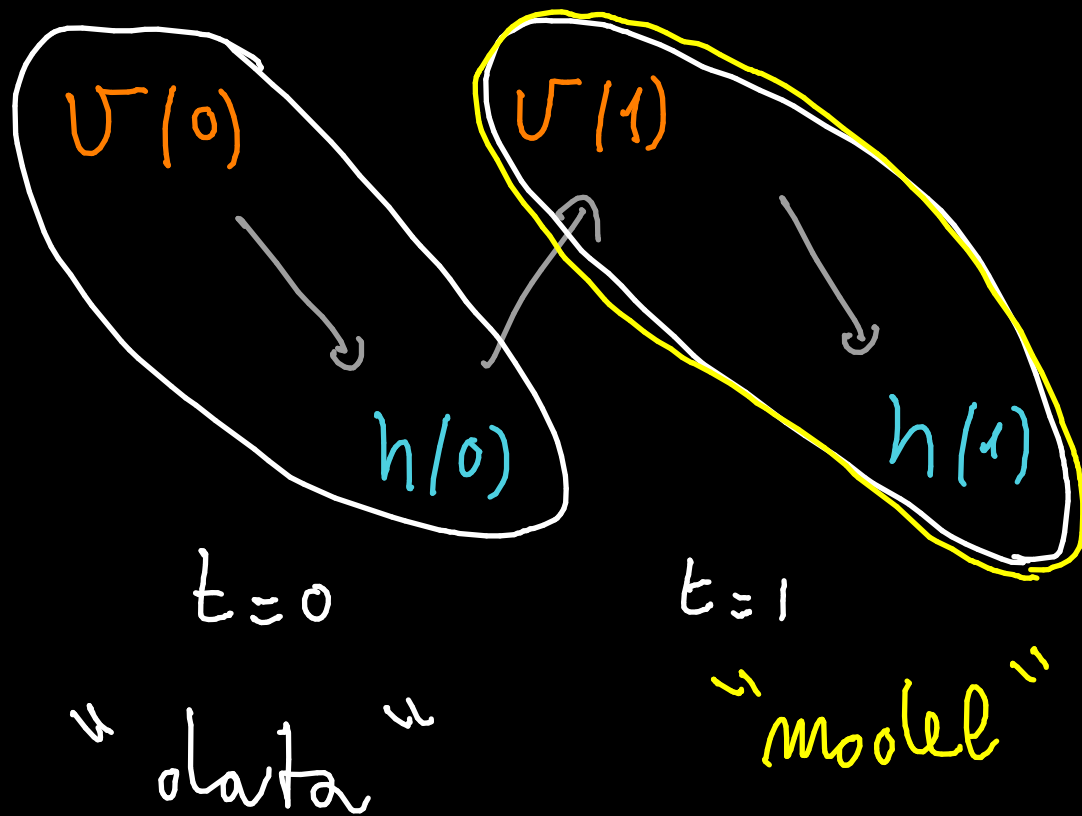


"data"

$m=2$  (for example)

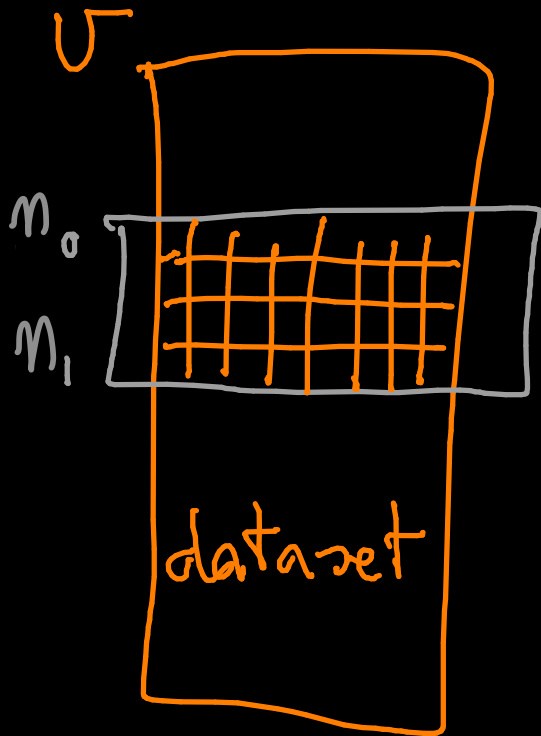
↓  
"model" evaluated  
at  $t=m$   
rather than  $t \rightarrow \infty$

# Contrastive Divergence (CD-1)



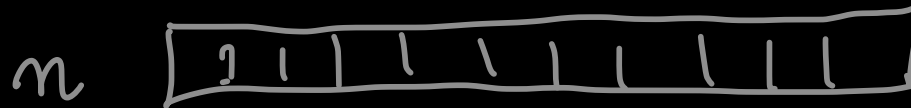
- most extreme example of CD
- fastest
- it works...

# Mini batches

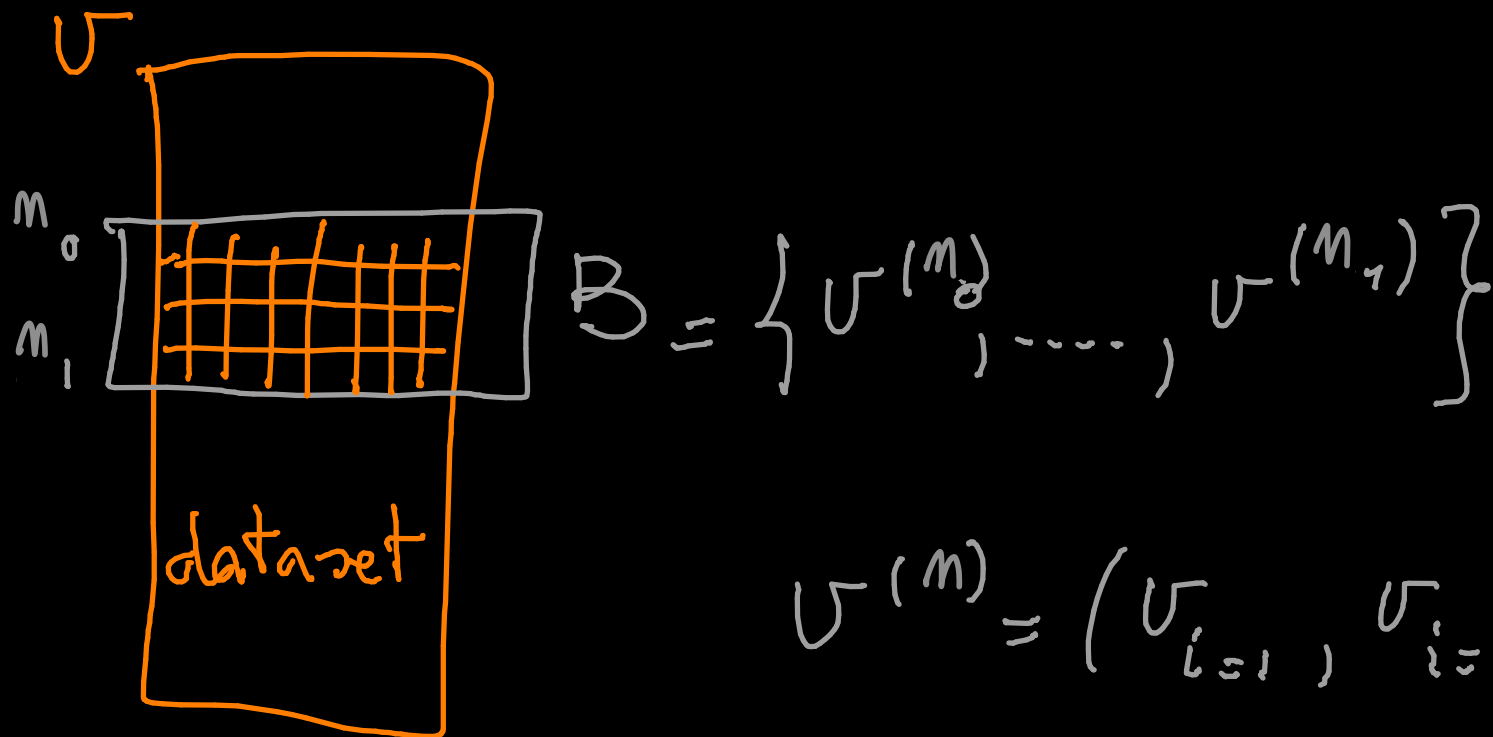


$$\mathcal{B} = \{U^{(n_0)}, \dots, U^{(n_1)}\}$$

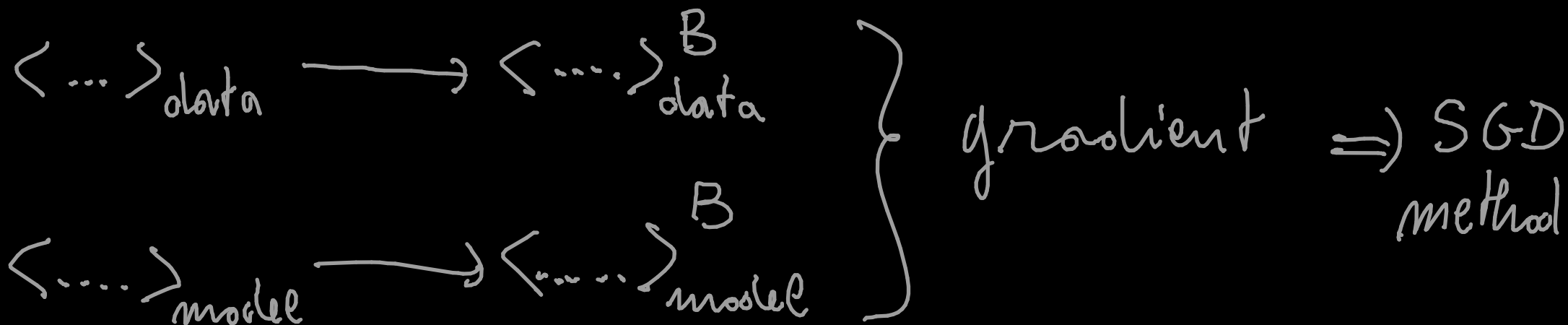
$$U^{(n)} = (U_{i=1}, U_{i=2}, \dots, U_{i=L})^{(n)}$$



# Mini batches

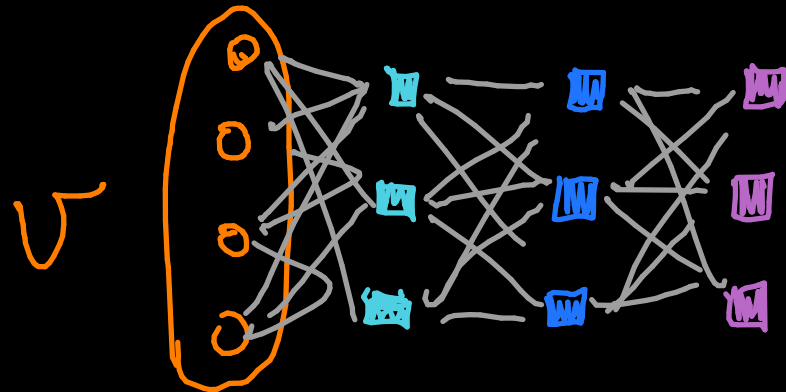


$$v^{(m)} = (v_{i=1}, v_{i=2}, \dots, v_{i=L})^{(m)}$$



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
- ...
- "reading"  $W$ 's  $\Rightarrow$  understand data

