

Title: PSI 2017/2018 - Machine Learning for Many Body Physics - Lecture 10

Date: Apr 19, 2018 03:45 PM

URL: <http://pirsa.org/18040071>

Abstract:

## Probabilistic Graphical Models

$P_{\lambda}(\vec{x})$  → diagrammatic representation of a probability distribution.  
(parameterized by  $\lambda$  = weights, biases, etc.)

Generative modelling can occur via graphical models.



## Probabilistic Graphical Models

$P_{\lambda}(\vec{x})$  → diagrammatic representation of a probability distribution.  
(parameterized by  $\lambda$  = weights, biases, etc.)

Generative modelling can occur via graphical models.

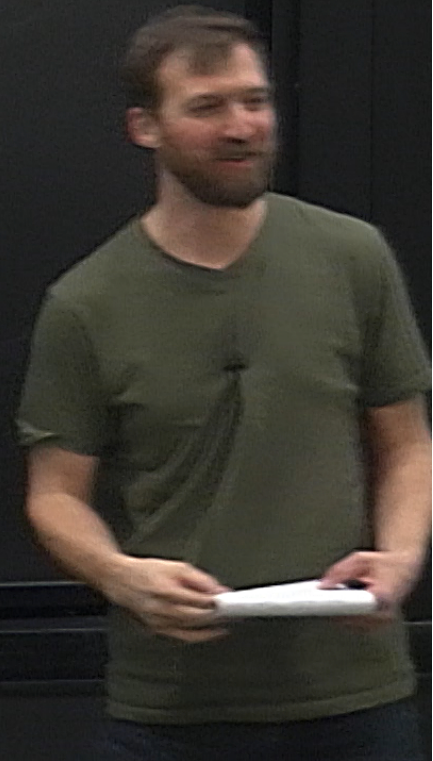
- directed connections "Bayesian" networks
- undirected "Markov random field models"



probability distribution,  
(e.g., biases, etc.)  
Bayesian models  
works  
"old models"



What is Bayesian statistics?



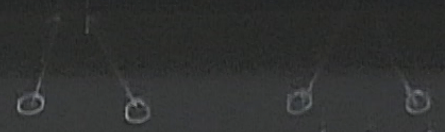


probability distribution.

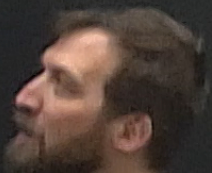
(, biases, etc.)

real models:

works  
"real models"

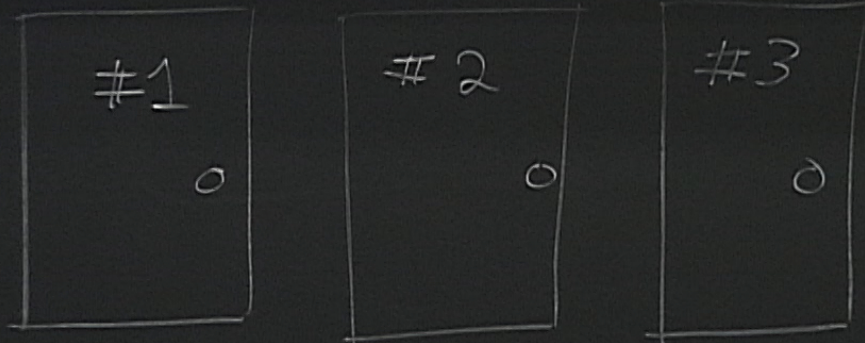


What is Bayesian statistics?  
Consider the Monte Hall problem:





- direction connections "Bayesian" networks
- undirected "Markov random field models"



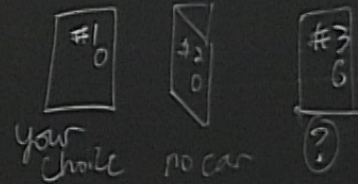
• you (contestant)

- one prize (a car)
- equally likely to be behind each door.



"Search" networks  
"random field models"

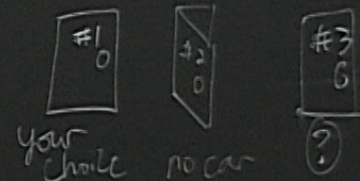
- you (contestant) chooses a door (say door #1)
- Me (the host) reveals a non-selected door that does not contain the car, (say door #2)





"Sean" networks  
"random field models"

- you (contestant) chooses a door (say door #1)
- Me (the host) reveals a non-selected door that does not contain the car, (say door #2)



Q: Do you stick with your original choice (#1)  
or switch to door #3 (?)



This is an example of Bayesian probabilities:



This is an example of Bayesian probabilities:

Notation:  $P(A|B)$  "probability of A given B"



This is an example of Bayesian probabilities:

Notation:  $P(A|B)$  "probability of A given B"

$P(A) = \sum_B p(A,B)$  "marginal probability" (all states of B are traced out)

$p(A,B)$  "joint probability"

$P(A,B) = P(B|A)p(A)$  "product rule"

From these, plus

$$p(A,B) = p(B,A)$$

$$p(B|A) = p(A|B)$$



given B "

probability" (all states of B are traced out)

From these, plus

$$P(A, B) = P(B, A)$$

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

Bayes' Theorem



$$P(A, B) = P(B|A)P(A) \quad \text{"product rule"}$$

$$P(A, B) = P(B, A) \quad P$$

Bayes' Theorem

Monte Hall: Assumed that you pick #1 and I reveal door #2

A: the event that the car is behind door #1

B: the event that I show an empty door #2

$$P(A|B)$$



product rule

$$P(A, B) = P(B, A) \quad P(B|A) = \frac{P(A, B)}{P(A)}$$

Bayes' Theorem

$c \neq 1$  and I reveal door  $\neq a$

behind door  $\neq 1$

empty door  $\neq 2$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} =$$



product rule

$$P(A, B) = P(B, A) \quad P(B|A) = \frac{P(A, B)}{P(A)}$$

Bayes' Theorem

$c \neq 1$  and I reveal door #  $a$   
behind door  
empty door  $\neq$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{P(B)}$$



$$P(A, B) = P(B|A)P(A) \quad \text{"product rule"}$$

$$P(A, B) = P(B, A)$$

Bayes' Theo

Monte Hall: Assumed that you pick #1 and I reveal door #2

A: the event that the car is behind door #1

B: the event that I show an empty door #2

$P(A|B)$

$$\text{so: } P(B) = \sum_A P(B|A)P(A)$$

$\frac{1}{3}$



# Bayes' Theorem

#1 and I reveal door #2

nd door #1

ty door #2

$$P(A/B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\frac{1}{2}}{P(B)}$$

- If car  $\rightarrow$  #1, door #2 opened 50% of the time.
- If car  $\rightarrow$  #2, door #2 never opened.



# Bayes' Theorem

and I reveal door #2

door #1

door #2

$$P(A/B) = \frac{P(B|A)P(A)}{P(B)} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{P(B)}$$

• If car  $\rightarrow$  #1, door #2 opened 50% of the time.

• If car  $\rightarrow$  #2, door #2 never opened.

• If car  $\rightarrow$  #3, door #2 opened 100% of the time



$$P(B) = \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 = \frac{1}{2} \Rightarrow P(A|B) = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}$$

there is a  $1 - \frac{1}{3} = \frac{2}{3}$  chance the car is behind door #3. So yes: change your pick to door #3.



$$P(B) = \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 = \frac{1}{2} \Rightarrow P(A|B) = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}$$

here  $B$  is a  $1 - \frac{1}{3} = \frac{2}{3}$  chance the car is behind door #3. So yes: change your pick door #3.

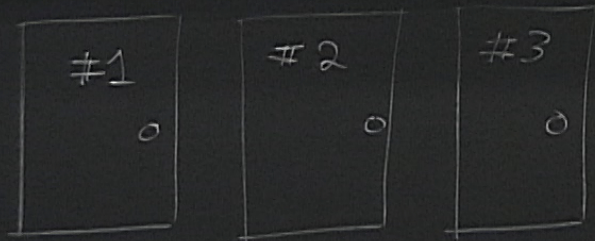
Goal of generative modelling  $\rightarrow$  produce a "good" estimate for  $p_i(x) \sim q(x)$  ("physical" distribution underlying your data set  $\{X\}$ )



$$P(B) = \frac{1}{3} \cdot \frac{1}{2} + \frac{1}{3} \cdot 0 + \frac{1}{3} \cdot 1 = \frac{1}{2} \Rightarrow P(A|B) = \frac{\frac{1}{3} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}$$

there is a  $1 - \frac{1}{3} = \frac{2}{3}$  chance the car is behind door #3. So yes: change

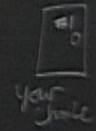
Goal of generative modelling  $\rightarrow$  produce a "good" estimate for  $p_{\theta}(\vec{X}) \sim q(\vec{X})$  "physical" distribution underlying your data  
 $\rightarrow$  generate new instances of  $\vec{X}$  (possibly not contained in the original data set.)



• one prize (a car)

• 1/3 chance to be behind each door

- you (contestant) chooses a door (say door #1)
- Me (the host) reveals a non-selected door contain the car, (say door #2)

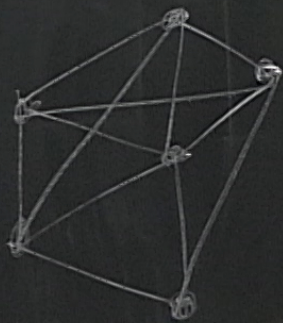


Q: Do you stick with your original choice or switch to door #3?



Goal of generative modelling  $\rightarrow$  produce a "good" estimate for  $P_{\lambda}(X)$   
 $\rightarrow$  generate new instances of  $\vec{X}$  (possibly  
 $\rightarrow$  derive from generative rules

### Hopfield network (1982)



$x_i = 0, 1$



→ produce a good estimate for  $P_n(\vec{x}) \sim q(\vec{x})$  (physical underlying)

→ generate new instances of  $\vec{x}$  (possibly not contained in the original data)

known random field models

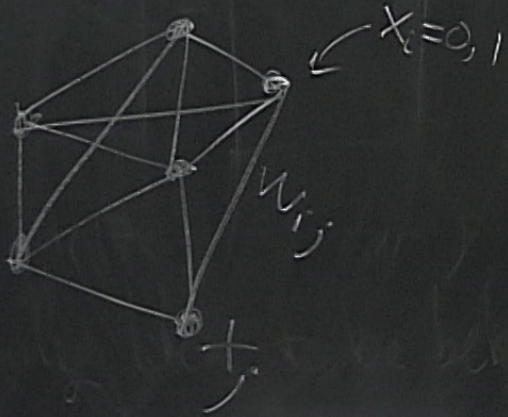
Like an Ising model with "Energy"

$$E(\vec{x}) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n W_{ij} x_i x_j + \sum_{i=1}^n h_i x_i$$



Goal of generative modelling  $\rightarrow$  produce a good estimate for  $P_n(x) \sim q$   
 $\rightarrow$  generate new instances of  $\vec{x}$  (possibly not in  $\mathcal{X}$ )  
 - often "Markov random field" models

Hopfield network (1982)



Like an Ising model with "Energy"

$$E_\lambda(\vec{x}) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n W_{ij} x_i x_j + \sum_{i=1}^n b_i x_i$$

where  $\lambda = \{W, \vec{b}\}$

and  $P_\lambda = \frac{1}{Z_\lambda} e^{-E_\lambda}$

"Machine Learning" IS

take  $S = \{\vec{x}_1, \vec{x}_2, \dots\}$



new instances of  $\vec{x}$  (possibly not contained in the original data set.)  
 underlying your data set  $\mathcal{X}$

Using model with "Energy"

$$-\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n W_{ij} x_i x_j + \sum_{i=1}^n b_i x_i$$

$$\lambda = \{W, \vec{b}\}$$

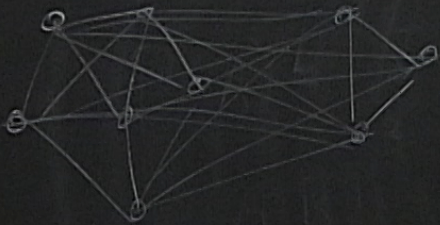
"Machine Learning" is tuning  $\lambda$  so that  $p_\lambda \approx q$

take  $S = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_{10,000}\}$

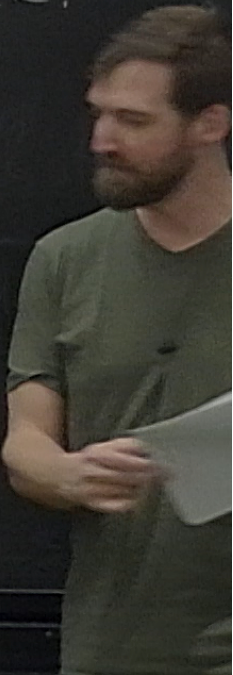


Boltzmann Machine

Like an extended Hopfield network.

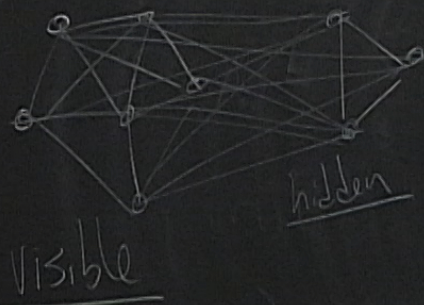


Visible





## Boltzmann Machine



Like an extended Hopfield network.

a generalization that has some hidden nodes  
"more powerful" than a Hopfield.



Like an extended Hopfield network.

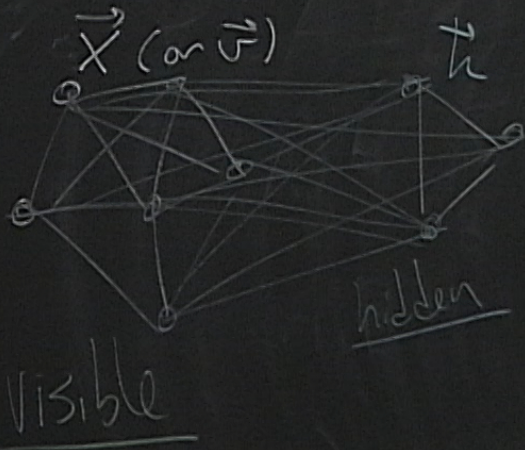
a generalization that has some hidden nodes.

"more powerful" than a Hopfield.

Imagine "learning" the  $\lambda = \{W, b\}$  so that  $P_{\lambda}(X) \approx q_{\lambda}(X)$

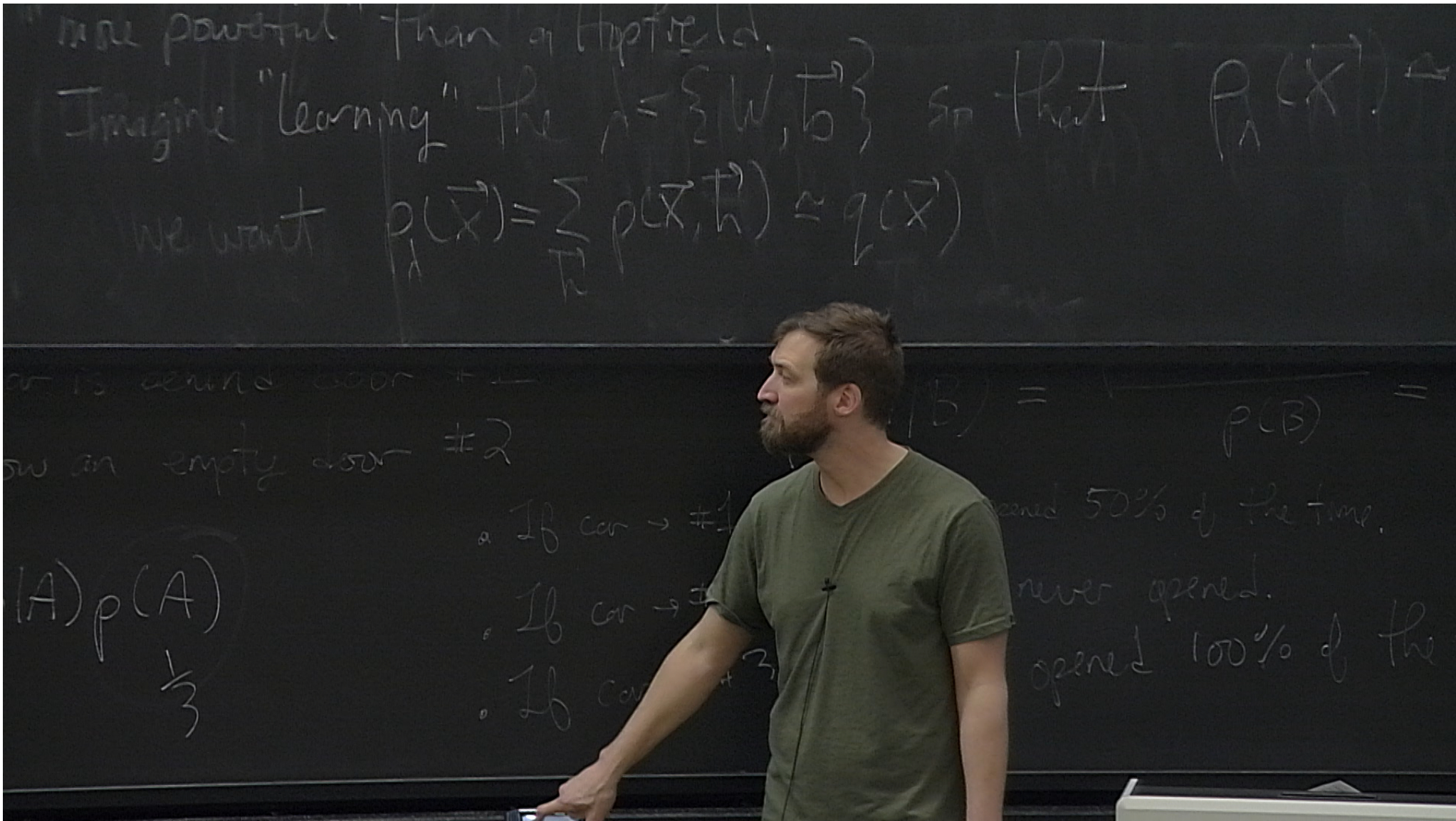


# Boltzmann Machine



Like an extended Hopfield net  
a generalization that has so  
"more powerful" than a Hopf  
Imagine "learning" the  $n$ -







an extended Hopfield network.

a generalization that has some hidden nodes.

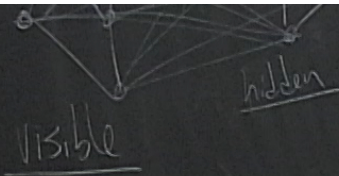
"more powerful" than a Hopfield.

Imagine "learning" the  $\lambda = \{W, \vec{b}\}$  so that  $P_{\lambda}(\vec{x}) \approx$

$$\text{we want } P_{\lambda}(\vec{x}) = \sum_{\vec{h}} p(\vec{x}, \vec{h}) \approx q(\vec{x})$$

u pick #1 and I reveal door #2



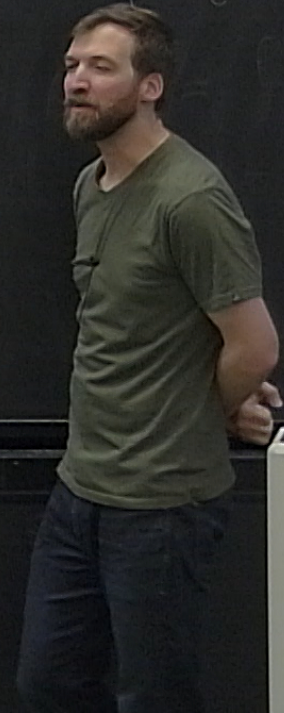


more powerful than a Hopfield.

Imagine "learning" the  $\lambda = \{W, b\}$  so that  $p_{\lambda}(\vec{x}) \approx q_{\lambda}$

We want 
$$p_{\lambda}(\vec{x}) = \sum_{\vec{h}} p(\vec{x}, \vec{h}) \approx q(\vec{x})$$

Let's restrict the architecture: Restricted Boltzmann Machine (Smolensky '86)





Visible

we want  $p_{\lambda}(\vec{x}) = \sum_{\vec{h}} p(\vec{x}, \vec{h}) \approx q(\vec{x})$

Let's restrict the architecture: Restricted Boltzmann Machine (RBM)

Let's consider binary units  $v_i = 0, 1$  (visible),  $h_i = 0, 1$  (for hidden)

Joint dist:  $p_{\lambda}(\vec{v}, \vec{h}) = \frac{1}{Z_{\lambda}} e^{-E_{\lambda}(\vec{v}, \vec{h})}$ ,  $E_{\lambda} = - \sum_{ij} W_{ij} v_i h_j$



visible

hidden

Imagine learning the  $\lambda = \{W, b\}$  so that

we want  $p_{\lambda}(\vec{x}) = \sum_{\vec{h}} p(\vec{x}, \vec{h}) \approx q(\vec{x})$

Let's restrict the architecture: Restricted Boltzmann Machine (Smolensky)

Let's consider binary units  $v_i = 0, 1$  (visible),  $h_i = 0, 1$  (for hidden)

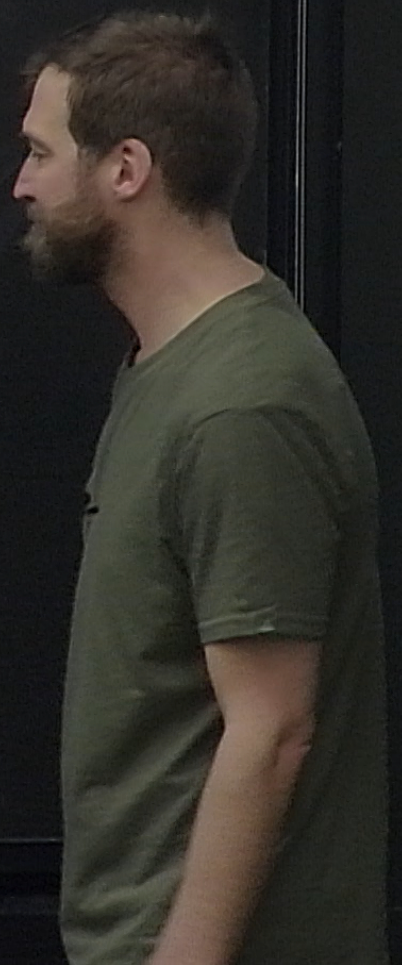
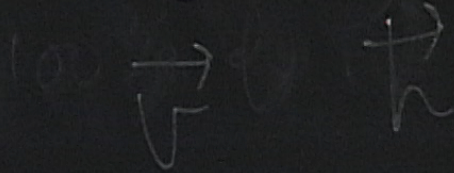
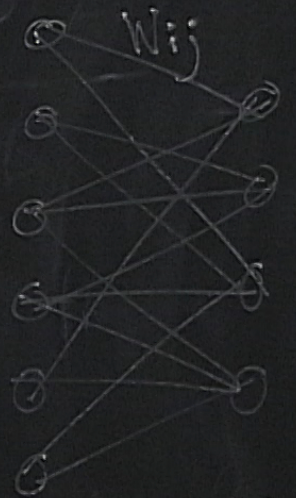
Joint dist:  $P_{\lambda}(\vec{v}, \vec{h}) = \frac{1}{Z_{\lambda}} e^{-E_{\lambda}(\vec{v}, \vec{h})}$ ,  $E_{\lambda} = - \sum_{ij} W_{ij} v_i h_j - \sum_i b_i v_i - \sum_j c_j h_j$



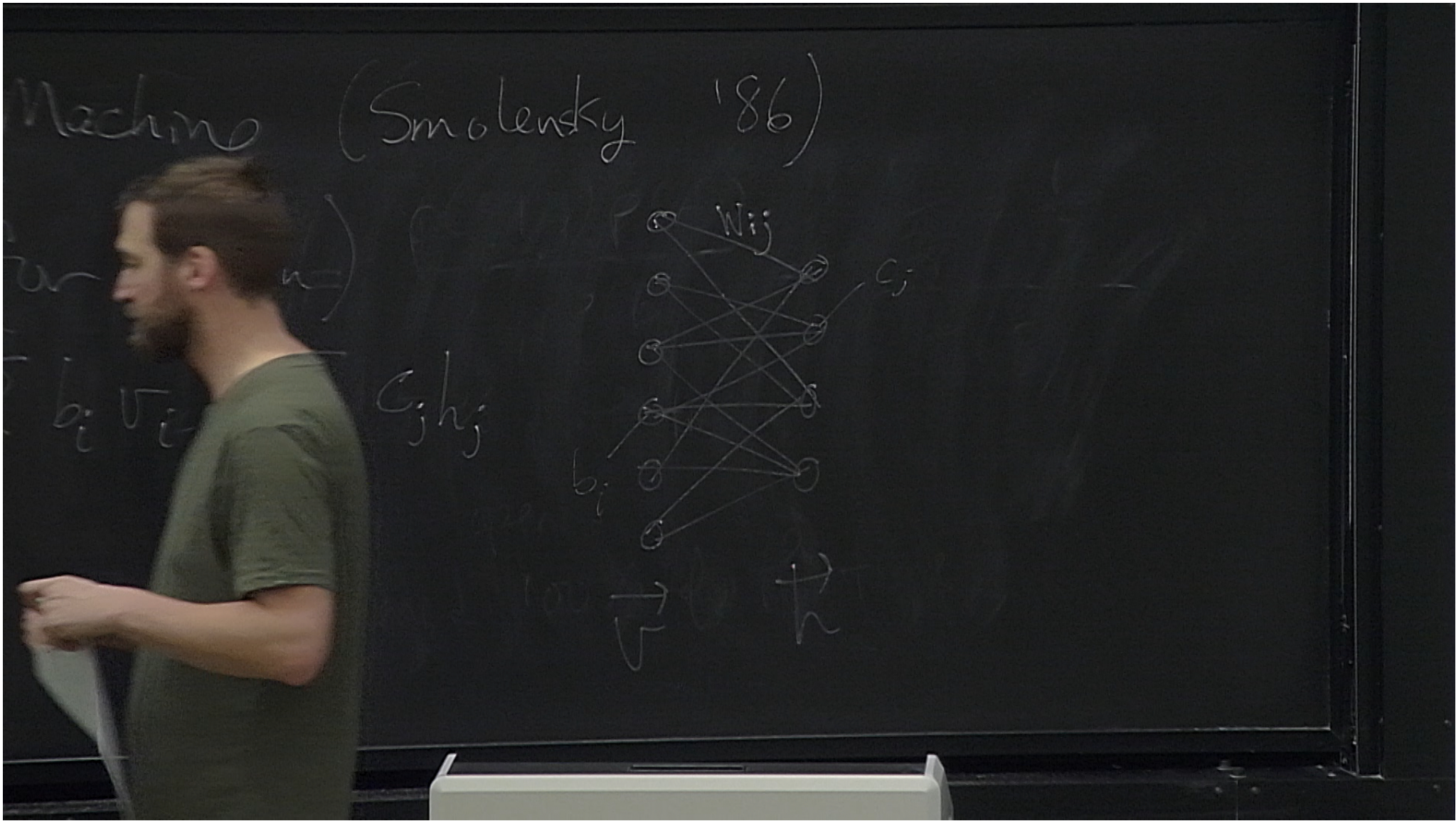
Machine (Smolensky '86)

for hidden

$$b_i v_i = \sum_j c_j h_j$$









We want  $p_{\lambda}(x) = \sum_{h} p(x, h) = q(x)$

restrict the architecture? Restricted Boltzmann Machine (Smolensky)

consider binary units  $v_i = 0, 1$  (visible),  $h_i = 0, 1$  (for hidden)

$$p_{\lambda}(\vec{v}, \vec{h}) = \frac{1}{Z_{\lambda}} e^{-E_{\lambda}(\vec{v}, \vec{h})}, \quad E_{\lambda} = -\sum_{ij} W_{ij} v_i h_j - \sum_i b_i v_i - \sum_j c_j h_j$$

$$Z_{\lambda} = \sum_{\vec{v}, \vec{h}} e^{-E(\vec{v}, \vec{h})}$$

$$\lambda = (W, \vec{b}, \vec{c})$$



## Industry uses of RBMs

pre 2012 were important for unsupervised pre-training for deep learning.

provided to that were used for supervised learning (fed in a classifier, SVM, FFNN)

Collaborative filtering (Salakhutdinov, Mnih, Hinton 2007)



## Industry uses of RBMs

- pre 2012 were important for unsupervised pre-training for deep learning.
- provided  $h^0$  that were used for supervised learning (fed in a classifier, SVM, FFNN)
- Collaborative filtering (Salakhutdinov, Mnih, Hinton 2007)  
"Netflix competition" → given users ranking → predict the ranking of an unseen movie.



