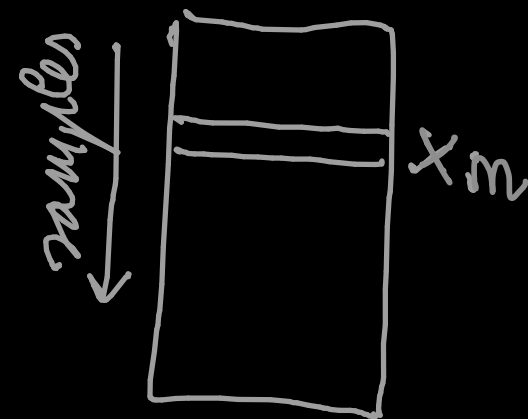


Unsupervised  
Learning  
(VL)

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



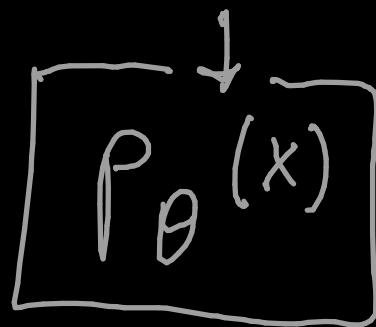
~~$\mathbf{y}$~~

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

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

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



# Why physics in UL?

- similar problems

- useful setup

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

- 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

$\langle f \rangle_{\text{model}}$  will be compared to  $\langle f \rangle_{\text{data}}$

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 $\theta$

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

review  
(189)

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

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

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

↑  
No data in the  
partition function!  
(only parameters  $\theta$  and  
possible  $x$ 's)

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)

$$\begin{aligned} \partial_{\theta_j}(-L(\theta)) &= \langle \partial_{\theta_j} \bar{E}_{\theta}(x) \rangle_{\text{data}} + \partial_{\theta_j} \log Z_{\theta} \\ &= \langle O_j(x) \rangle_{\text{data}} \longrightarrow \langle O_j(x) \rangle_{\text{model}} \end{aligned} \quad \begin{matrix} (*) \\ (195) \end{matrix}$$

GIVEN

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

(\*)

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

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

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

## 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  $p_\theta$ ,  
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  $p_\theta$ ,  
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  $p_\theta$ ,  
following Monte Carlo procedure

Normalization

$x'_i$ , "fantasy particle"


## log-derivative trick

to compute gradient of any  $f(x)$

$$\partial_{\theta_j} \langle f(x) \rangle_{\text{model}} = \sum_i \partial_{\theta_j} p_{\theta}(x_i) f(x_i) \quad \left( \partial_{\theta_j} = \frac{\partial}{\partial \theta_j} \right)$$

# 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  $f(x)$

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

# 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  $\theta$  with gradient descent

# 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



NOT

ANALYTICALLY

3) compute gradients (195) (for  $B$ )

4) update  $\theta$  with gradient descent

(as in back propagation)