# Energy-based models

# Energy-based models for distributions

o Distribution as: $p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \dfrac{1}{\int_x g_{\boldsymbol{\theta}}(x)dx} g_{\boldsymbol{\theta}}(\boldsymbol{x})$

o $p_{\boldsymbol{\theta}}$ as known probability distributions (Gaussian, exp.) can be restrictive
  ◦ Maybe I want to encode domain knowledge of how variables interact

o We can also define an energy function and divide by its volume

$$g_{\boldsymbol{\theta}}(\boldsymbol{x}) = \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x})) \Rightarrow p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x}))$$

# Energy-based models for distributions

$$g_{\boldsymbol{\theta}}(\boldsymbol{x}) = \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x})) \Rightarrow p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x}))$$

o  $-f_{\boldsymbol{\theta}}(\boldsymbol{x})$ is the energy function

o  Partition function is the hard bit

$$Z(\boldsymbol{\theta}) = \int_{\boldsymbol{x}} \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x})) \, d\boldsymbol{x}$$

◦ Note the multi-dimensional integral due to $\boldsymbol{x}$

# Why exponential?

- Why $g_\theta(x) = \exp(f_\theta(x))$ and not $g_\theta(x) = f_\theta^2(x)$?


- Couples well with maximum likelihood and natural logarithms

- Many existing distributions are exponential-based

- They arise often in statistical physics $\rightarrow$ Good inspiration

# Advantages & disadvantages

$$p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x}))$$

o  Very flexible in defining our energy function

o  Sampling from $p_{\boldsymbol{\theta}}(\boldsymbol{x})$ can be very hard
   ◦ The CDF introduces another integral

o  Evaluating and optimizing likelihood can be hard $\Rightarrow$ Learning is hard
   ◦ Must be able to compute the partition function

o  In vanilla case no latent variables $\Rightarrow$ no representation learning
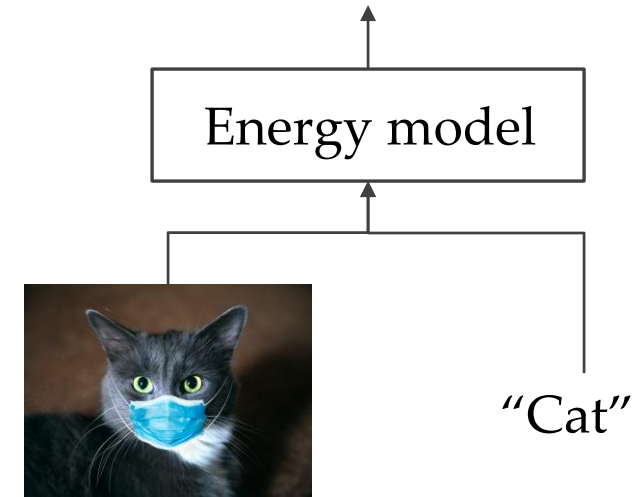   ◦ Latent variables can be added though

# Ratio of likelihoods
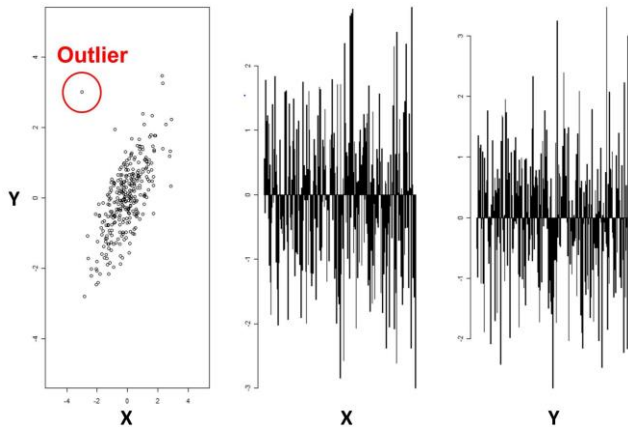
o The partition function is often very hard to compute analytically

o But if we have pairs of inputs

$$\frac{p_{\boldsymbol{\theta}}(\boldsymbol{x_a})}{p_{\boldsymbol{\theta}}(\boldsymbol{x_b})} = \exp(f_{\boldsymbol{\theta}}(\boldsymbol{x}_a) - f_{\boldsymbol{\theta}}(\boldsymbol{x}_b))$$

o No partition function anymore

# Applications

o Given trained model
  ◦ Anomaly detection
  ◦ Denoising & restoration
  ◦ Classification

UNIVERSITY OF AMSTERDAM

VISLab

# Examples of energy models

- Ising model

$$p_{\boldsymbol{\theta}}(\boldsymbol{y}, \boldsymbol{x}) = \frac{1}{Z} \exp\left(\sum_i \psi_i(x_i, y_i) + \sum_{i,j \in E} \psi_{ij}(y_i, y_j)\right)$$



- Product of experts (similar to AND)

$$p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{Z(\boldsymbol{\theta}, \boldsymbol{\varphi}, \boldsymbol{\omega})} q_\theta(\boldsymbol{x}) r_\varphi(\boldsymbol{x}) s_\omega(\boldsymbol{x})$$

- Hopfield networks

- Boltzmann machines

- Deep belief networks