# Deep Learning
# lecture 4
# Energy-Based Model

Yi Wu, IIS

Spring 2024

Mar-22

# Logistics

- Coding Project 2 due in 1 week
  - Use local compute for coding & Colab for testing
  - Cloud for long-term training
  - Any questions can be posted in Dingding channel

- Tips
  - Tricks: overfitting & regularization
  - First overfit!
    - **Learning rate decay, architecture, initialization, normalization and preprocess …**
  - Then regularize!
  - Be aware of your model size and computation (flops)!
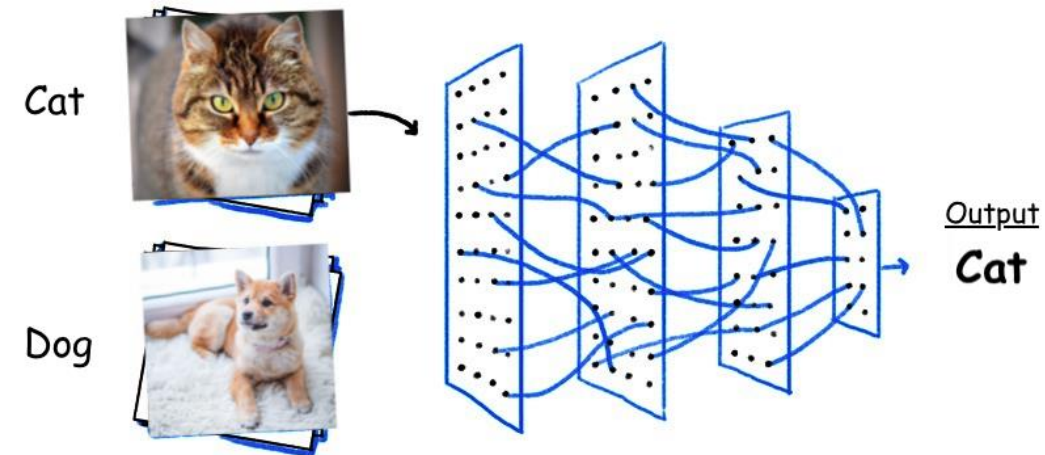
# Overview of Lecture 3

- Algorithm Design
  - Principle 1: adaptive and decoupled learning rate
  - Principle 2: momentum
  - Principle 3: second-order information is great! Let's approximate it

- Practical Algorithms
  - Mini-Batch SGD
  - Momentum SGD
  - AdaGrad, RMSProp, AdaDelta
  - Adam

# Overview of Lecture 3

- Regularizations
  - Goal: stabilize gradients and generalization!
  - Gradient Tricks:
    - Initialization, Gradient Clipping
  - Generalization Tricks:
    - Weight Decay, Dropout, Data Augmentation, Early Stopping
  - Normalization Layers:
    - BatchNorm, LayerNorm
  - Other:
    - Ensemble & practice makes perfect ☺
- Architecture
  - Residual Connection, Dense Connection
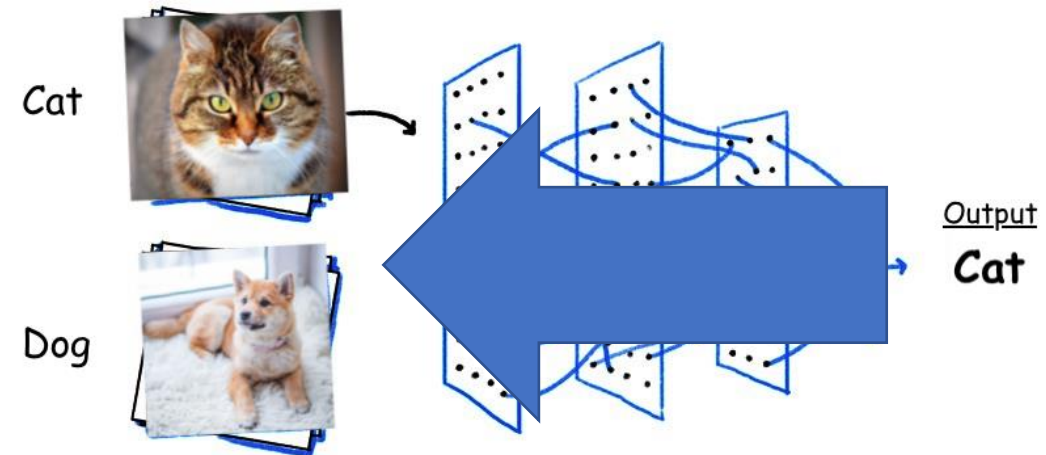  - Fully Convolutional Network

# Story So Far

- History
  - Lecture 1
    - first neural network (1943) to recent advances in deep learning
- Supervised Learning (Classification)
  - Lecture 2
    - MLP and basic components; Backpropagation
  - Lecture 3
    - Algorithms, Tricks and Architecture
- Discriminative Model
  - $P(y|X)$
  - Labeled data; $X \rightarrow y$

# Afterwards

- What if we want to generate $X$?

- Generative Model
  - $P(X, y) = P(y) * P(X|y)$
  - Or just $P(X)$

- Lecture 3~7
  - Deep Generative Models
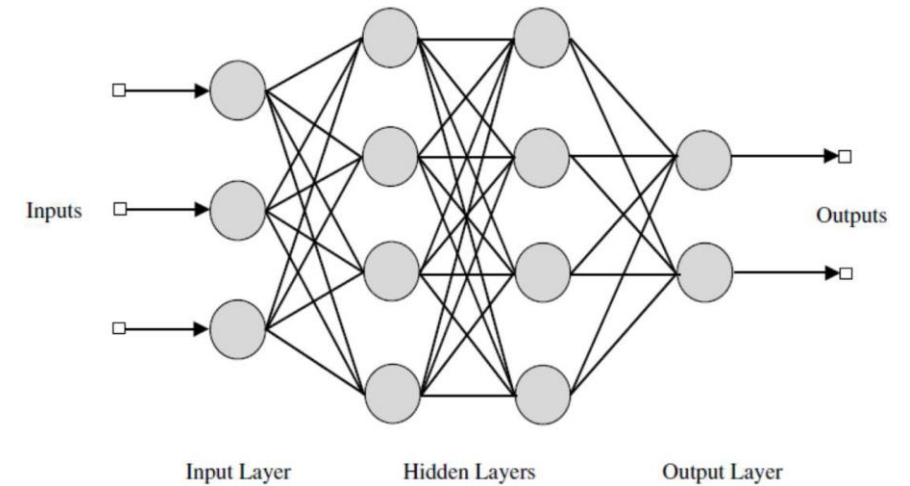  - Ask the neural network to generate a cat!

# Today's Lecture: Energy-Based Models

- A particularly flexible and general form of ***generative model***

- Part 1: Hopfield Network
  - The simplest model that can memorize and generate patterns

- Part 2: Boltzmann Machine
  - The first deep generative model

- Part 3: General Energy-Based Models

# Today's Lecture: Energy-Based Models

- A particularly flexible and general form of ***generative model***

- <span style="color:red">Part 1: Hopfield Network</span>
  - <span style="color:red">The simplest model that can memorize and generate patterns</span>

- Part 2: Boltzmann Machine
  - The first deep generative model
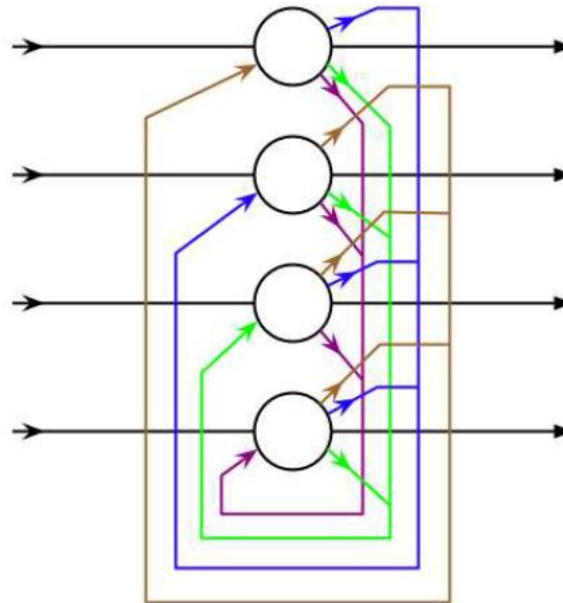
- Part 3: General Energy-Based Models

# Classification

- Recap: Classification
  - Layer-by-layer computation
  - Computation Graph: Directed Acyclic Graph
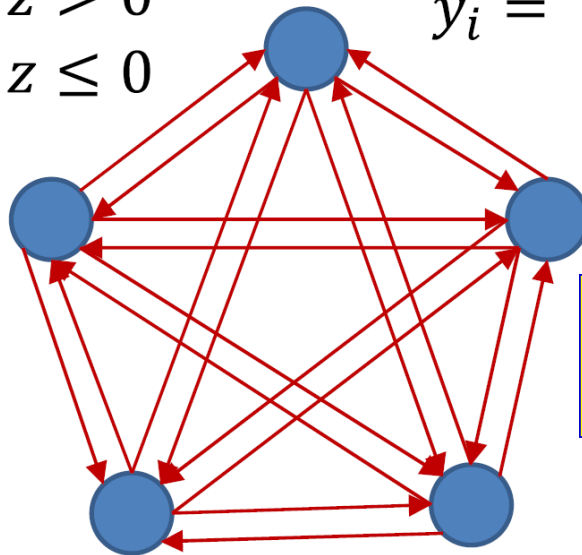  - Feedforward networks



- What about …
  - Loops!

# A Loopy Network

- A "fully-connected" network
  - Each neuron receives inputs from all the other neurons
  - $y_i = +1 \ or \ -1$ with hard thresholding

$$\Theta(z) = \begin{cases} +1 \ if \ z > 0 \\ -1 \ if \ z \leq 0 \end{cases}$$

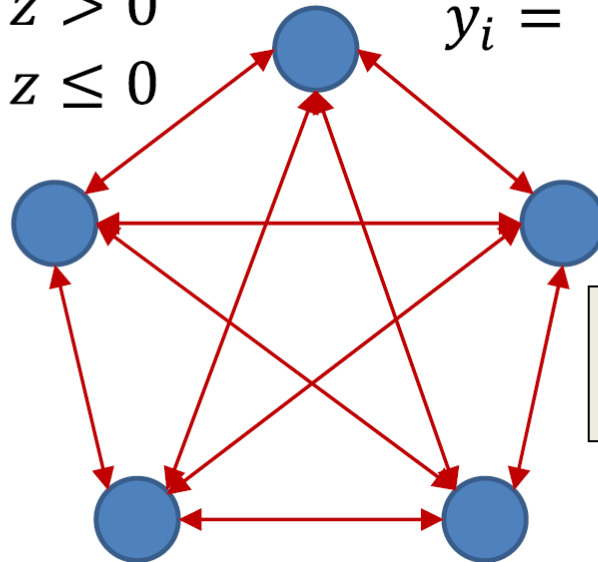$$y_i = \Theta\left(\sum_{j \neq i} w_{ji}y_j + b_i\right)$$

The output of a neuron affects the input to the neuron

# Hopfield Network

- A "fully-connected" network
  - Each neuron receives inputs from all the other neurons
  - $y_i = +1 \; or \; -1$ with hard thresholding
  - Symmetric weights

$$\Theta(z) = \begin{cases} +1 \; if \; z > 0 \\ -1 \; if \; z \leq 0 \end{cases} \qquad y_i = \Theta\left(\sum_{j \neq i} w_{ji} y_j + b_i\right)$$
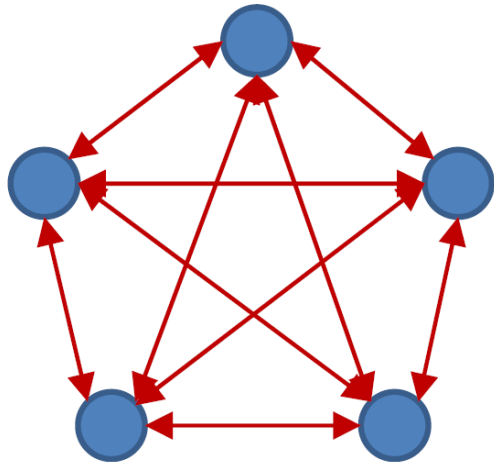
A symmetric network:
$$w_{ij} = w_{ji}$$

# Hopfield Network

- A Hopfield Network may not be stable!
  - At each time each neuron receives a "field" $\sum_{j \neq i} w_{ji} y_j + b_i$
  - If the sign of neuron matches the sign of the field, it flips

$$y_i \leftarrow -y_i \ \text{if} \ y_i \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) < 0$$

  - This can further cause other neurons to flip!



$$y_i = \Theta \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)$$

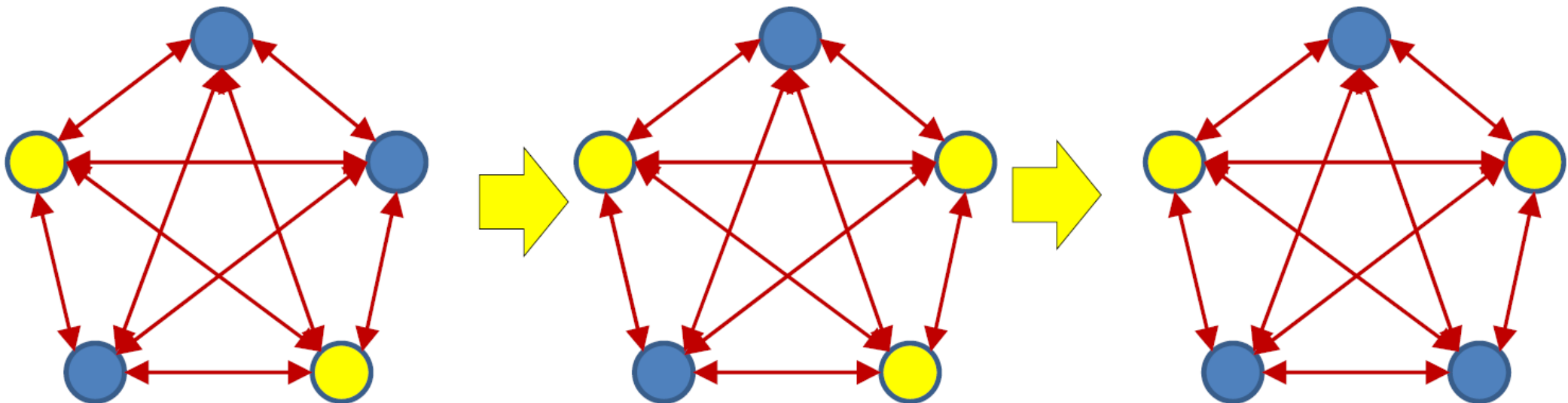$$\Theta(z) = \begin{cases} +1 \ if \ z > 0 \\ -1 \ if \ z \leq 0 \end{cases}$$

# Hopfield Network

- Neurons flip if its sign does not match its local "field"
  - $y_i \leftarrow -y_i$ if $y_i\left(\sum_{j \neq i} w_{ji} y_j + b_i\right) < 0$ for all neurons
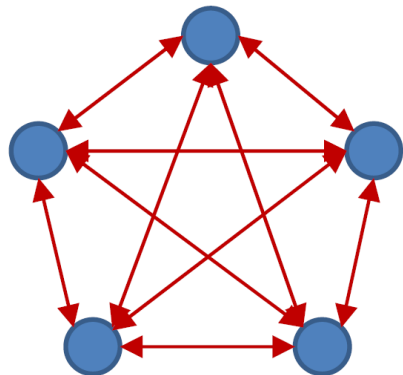  - Repeat until no neuron flips
  - Will this process converge?

$$\Theta(z) = \begin{cases} +1 \ if \ z > 0 \\ -1 \ if \ z \leq 0 \end{cases}$$

$$y_i = \Theta\left(\sum_{j \neq i} w_{ji} y_j + b_i\right)$$

# Hopfield Network

- Let $y_i^-$ denote the value of $y_i$ before a "flip"

- Let $y_i^+$ denote the value of $y_i$ after a "flip"

- If $y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) \geq 0$, nothing happen

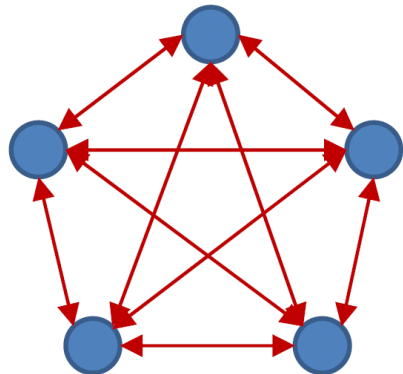$$y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) - y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) = 0$$

$$y_i = \Theta \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)$$

$$\Theta(z) = \begin{cases} +1 \ if \ z > 0 \\ -1 \ if \ z \leq 0 \end{cases}$$

# Hopfield Network

- Let $y_i^-$ denote the value of $y_i$ before a "flip"

- Let $y_i^+$ denote the value of $y_i$ after a "flip"

- If $y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) \geq 0$, nothing happen

- If $y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) < 0$, $y_i^+ = -y_i^-$

$$y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) - y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) = 2 y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)$$

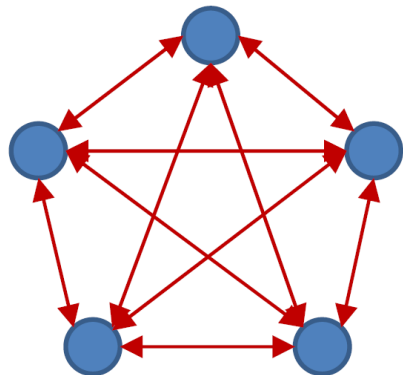$$y_i = \Theta \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)$$

$$\Theta(z) = \begin{cases} +1 \; if \; z > 0 \\ -1 \; if \; z \leq 0 \end{cases}$$

# Hopfield Network

- Let $y_i^-$ denote the value of $y_i$ before a "flip"
- Let $y_i^+$ denote the value of $y_i$ after a "flip"
- If $y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) \geq 0$, nothing happen
- If $y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) < 0$, $y_i^+ = -y_i^-$

$$y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) - y_i^- \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) = \textcolor{red}{2 y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)} \quad \textbf{Positive!}$$

*Every flip increases*
$$\textcolor{red}{2 y_i \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)}$$

$$y_i = \Theta \left( \sum_{j \neq i} w_{ji} y_j + b_i \right)$$

$$\Theta(z) = \begin{cases} +1 \ if \ z > 0 \\ -1 \ if \ z \leq 0 \end{cases}$$
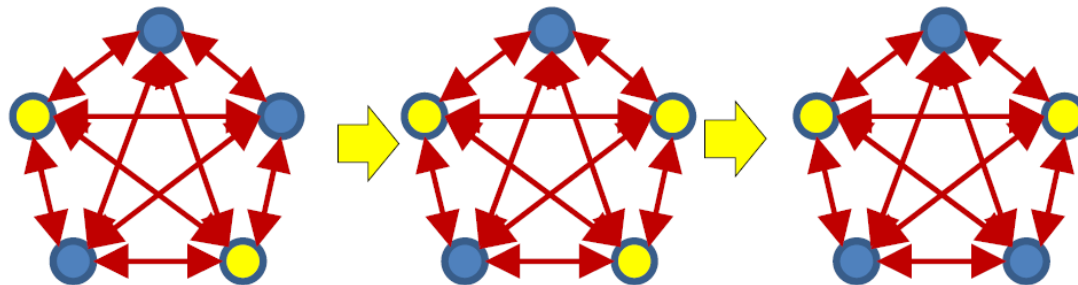
# Hopfield Network

- Consider the sum over every pair of neurons (assume $w_{ii} = 0$)

$$D(y_1, \ldots, y_N) = \sum_{i<j} y_i w_{ij} y_j + y_i b_i$$

- Any flip that changes $y_i^-$ to $y_i^+$ increases $D(y_1, \ldots, y_N)$

$$\Delta D = D(\ldots, y_i^+, \ldots) - D(\ldots, y_i^-, \ldots) = 2y_i^+ \left( \sum_{j \neq i} w_{ji} y_j + b_i \right) > 0$$
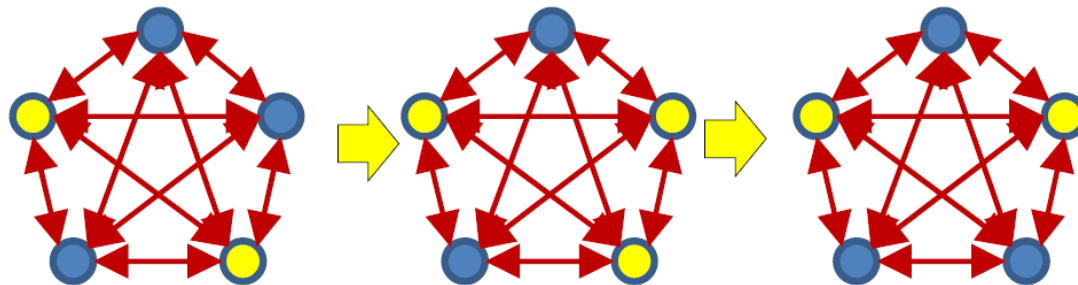
- Convergence?

# Hopfield Network

- $D$ is upper-bounded (we only change $y_i$)

$$D(y_1, \ldots, y_N) = \sum_{i<j} w_{ij} y_i y_j + \sum_i b_i y_i \leq \sum_{i<j} |w_{ij}| + \sum_i |b_i|$$

- $\Delta D$ is lower-bounded

$$\Delta D_{\min} = \min_{i, \{y_j\}} 2 \left| \sum_j w_{ij} y_j + b_i \right| > 0$$

- $\{y_i\}$ converges with a finite number of iterations!
  - $\{y_i\}$: *state*
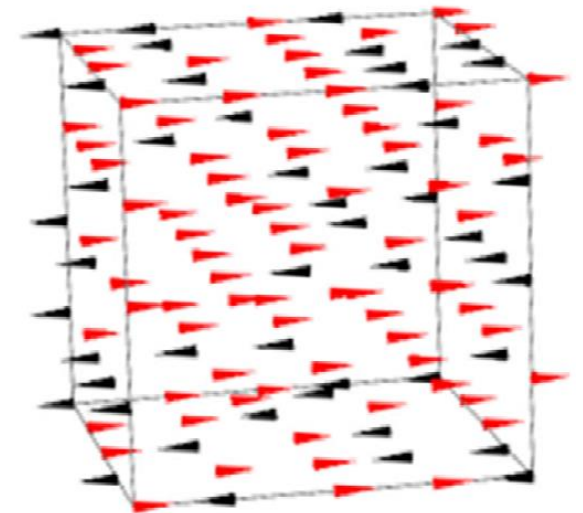
# Hopfield Network

- The **Energy** of Hopfield Network

$$E = -D = -\sum_{i<j} w_{ij} y_i y_j - \sum_i b_i y_i$$

  - The evolution of Hopfield network always decreases its energy!

- The concept of Energy
  - Magnetic dipoles in a disordered magnetic material
  - Each dipole tries to align itself to the local field
  - Field at a particular dipole $f(p_i)$, $p_i$ is the position of $x_i$

$$f(p_i) = \sum_{j \neq i} J_j x_j + b_i$$

  - **Ising model** of magnetic materials (Ising and Lenz, 1924)

# Hopfield Network

- Ising model for magnetic materials
  - Total field for a dipole

$$f(p_i) = \sum_{j \neq i} J_j x_j + b_i$$

  - Response of a dipole
    - $x_i \leftarrow -x_i$ if $\text{sign}(x_i f(p_i)) \neq 1$
  - Hamiltonian (total energy) of the system

$$E = -\frac{1}{2} \sum_i x_i f(p_i) = -\sum_{i<j} J_{ij} x_i x_j - \sum_i b_i x_i$$
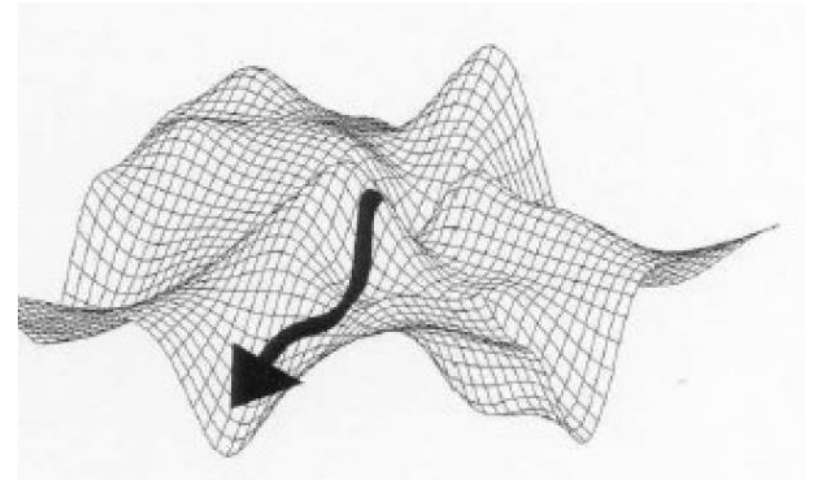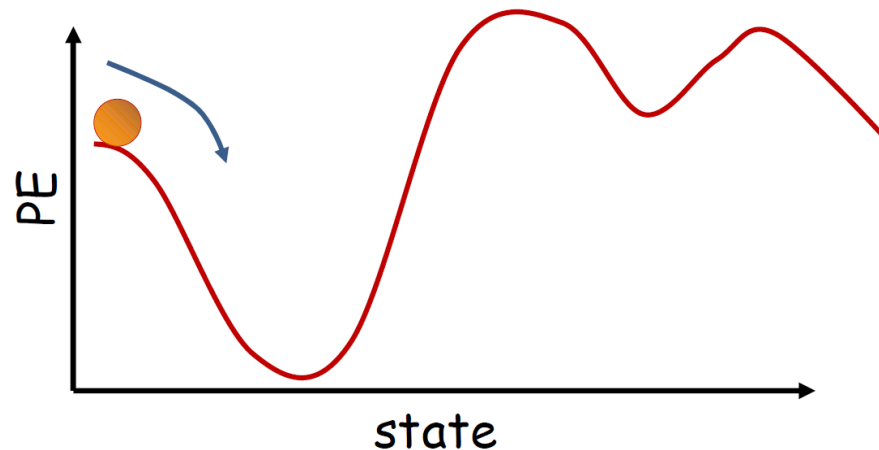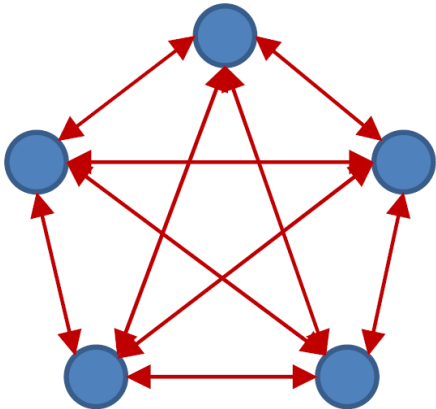
  - The system *evolves* to minimize the energy
    - A dipole stop filling if that flip increases the energy → a local minimum

# Hopfield Network

- The Hopfield network (simplified)

$$E = -\sum_{i<j} w_{ij} y_i y_j$$

  - Network evolution arrives at a local optimum in the energy contour
    - Every change in the network state decreases the energy
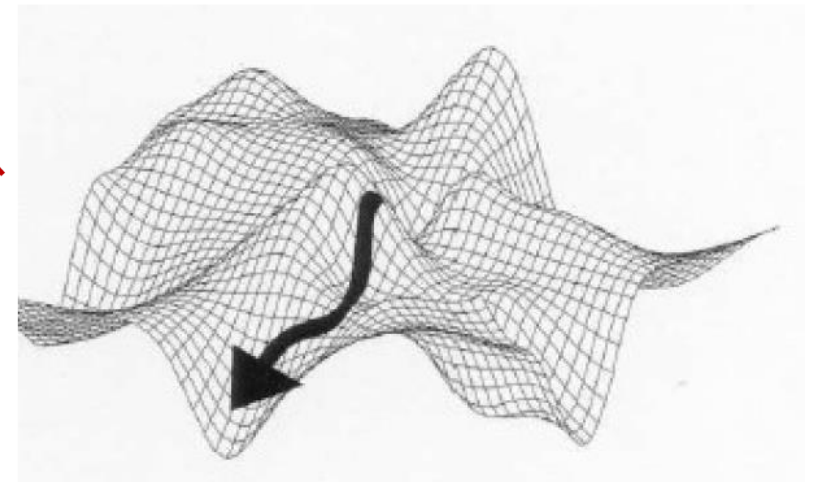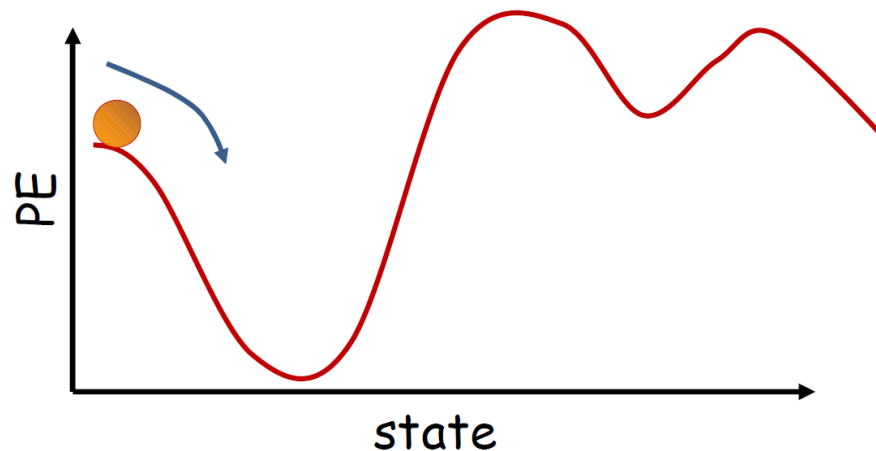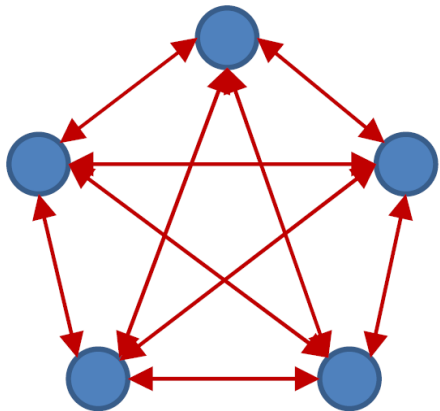  - Any small jitter from this stable state returns it to the stable state

# Hopfield Network

- The Hopfield network (simplified)

$$E = -\sum_{i<j} w_{ij} y_i y_j$$

  - Each local optimum state is a "stored" pattern
    - If the network is initialized close to a stored pattern, it evolves to the pattern
  - *Associated Memory (content addressable memory)*

# Hopfield Network

- Image Reconstruction by Hopfield Network (1982)



Hopfield network reconstructing degraded images from noisy (top) or partial (bottom) cues.

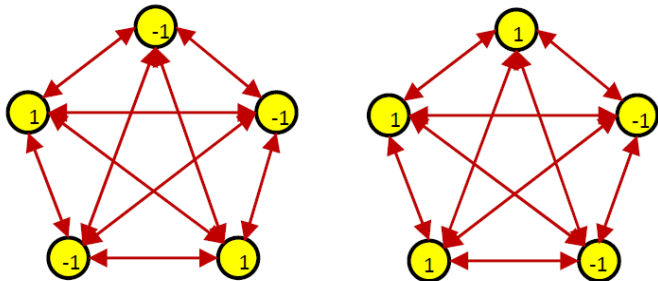- *How can we store the desired patterns?*

# Hopfield Network: Training

- Let's teach the network to store this image
  - $N$ pixels $\rightarrow$ $N$ neurons
  - Symmetric weights $\rightarrow \frac{1}{2} N(N-1)$ parameters to learn

- Design $\{w_{ij}\}$ such that the energy is at a local minimum for a desired pattern $y$
  - Redundancy! $y$ & $-y$ will be both stored

$$E = -\sum_i \sum_{j<i} w_{ji} y_j y_i$$
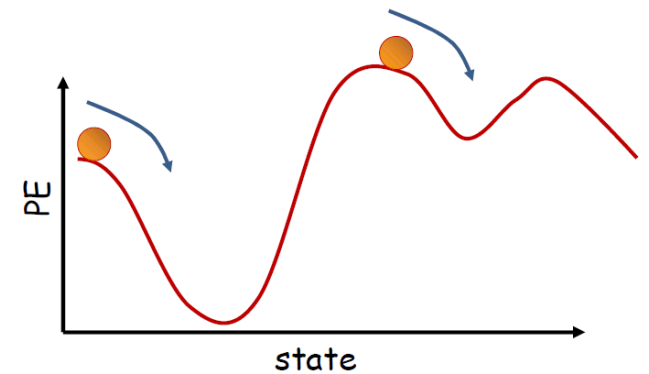
# Hopfield Network: Training



- Let's teach the network to store this image
  - $N$ pixels $\rightarrow$ $N$ neurons
  - Symmetric weights $\rightarrow \frac{1}{2}N(N-1)$ parameters to learn

- Design $\{w_{ij}\}$ such that the energy is at a local minimum for a desired pattern $y$
  - Hebbian Learning Rule $w_{ij} \leftarrow y_i y_j$ (1949)
  - $E = -\sum_{i<j} w_{ij} y_i y_j = -\frac{1}{2}N(N-1)$ $\rightarrow$ <span style="color:red">lowest possible energy!</span>

# Hopfield Network: Training

- What if we want to store **multiple** patterns?
    - $P = \{y^p\}$ $N_p$ patterns
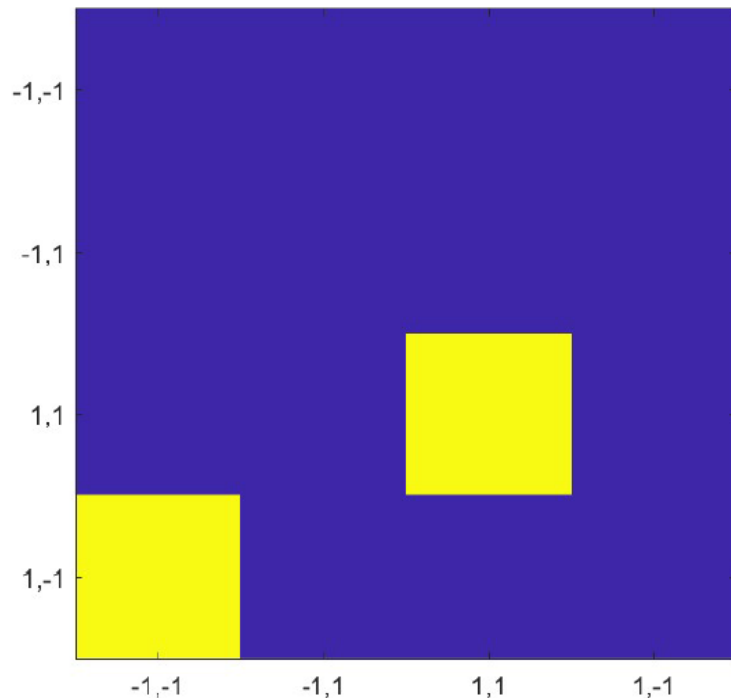    - Hebbian Learning Rule

$$w_{ij} = \frac{1}{N_p} \sum_p y_i^p y_j^p$$

- The issue of Hebbian Learning
    - Spurious local optima

# Hopfield Network: Training

- Example: 4-dimensional Hopfield Network with Hebbian Learning
  - Two orthogonal patterns to store
    - *Let's assume the value of each neuron is 1 or -1*

**Left: desired patterns**



**Right: stored patterns**

# Hopfield Network: Training

- Example: 4-dimensional Hopfield Network with Hebbian Learning
  - Three patterns
  - *How many patterns can a Hopfield network store?*



*Left: desired patterns*

*Right: stored patterns*

# Hopfield Network: Training

- How many patterns can a Hopfield network store?

- A fact: you can store all the $2^N$ patterns!

- Solution: find any $N$ orthogonal patterns
  - Prove this fact in your homework ☺

*A "flatten" landscape does not help evolve desired patterns!*

# Hopfield Network: Training

- We want to construct a network with desired ***stable*** local optimum
  - A pattern can be recovered after 1-bit change
- For a specific set of $K$ patterns, we can always build a network for which all patterns are stable provided $K \leq N$
  - Mostafa and St. Jacques (1985)
  - For large $N$, the upper bound on $K$ is actually $\frac{N}{4}\log N$
    - McElice et. al. (1987)
  - Still possible with undesired local minimum
- **How can we find the weights?**
  - $K$ patterns remembered
  - Avoid undesired local minimum as much as we can

# Hopfield Network: Optimization

- Problem Formulation
  - Desired patterns $P = \{y^p\}$
  - Energy function $E(y) = -\frac{1}{2} y^T W y$  (we omit bias for simplicity)

- Objective for $W$
  - Minimize $E$ for all $y^p$
  - It should also maximize $E$ for all non-desired patterns!

$$W = \arg\min_{W} \sum_{y \in P} E(y) - \sum_{y' \notin P} E(y')$$

  - Gradient Descent

$$W \leftarrow W - \eta \left( \sum_{y \in P} y y^T - \sum_{y' \notin P} y' y'^T \right)$$

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} y y^T - \sum_{y' \notin P} y' y'^T \right)$$



Energy

state

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} y y^T - \sum_{y' \notin P} y' y'^T \right)$$

  - The first term is minimizing the energy of desired patterns!

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \boxed{\sum_{y' \notin P} y'y'^T} \right)$$

  - The second term essentially raises all the patterns in the space
    - **Issue??**

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \mathbf{Valley}} y'y'^T \right)$$

- Let's just focus on the valleys!

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \boldsymbol{Valley}} y'y'^T \right)$$

  - Let's just focus on the valleys!
  - **But how can we find the valleys?**

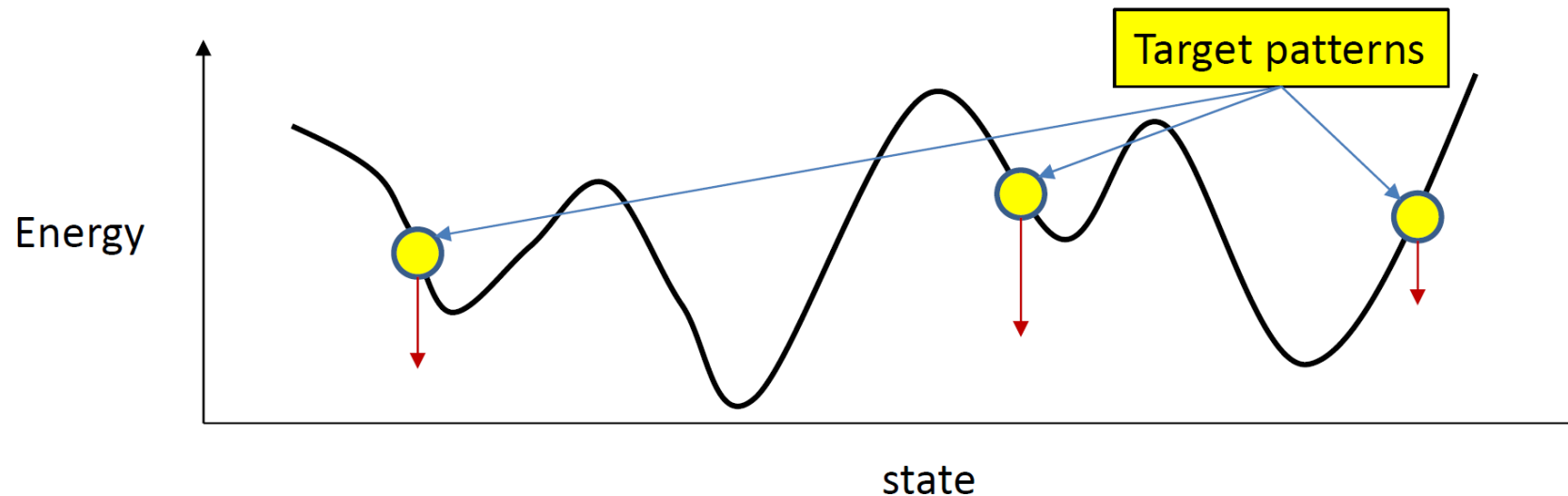# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \textbf{Valley}} y'y'^T \right)$$

- Let's just focus on the valleys!
- But how can we find the valleys?
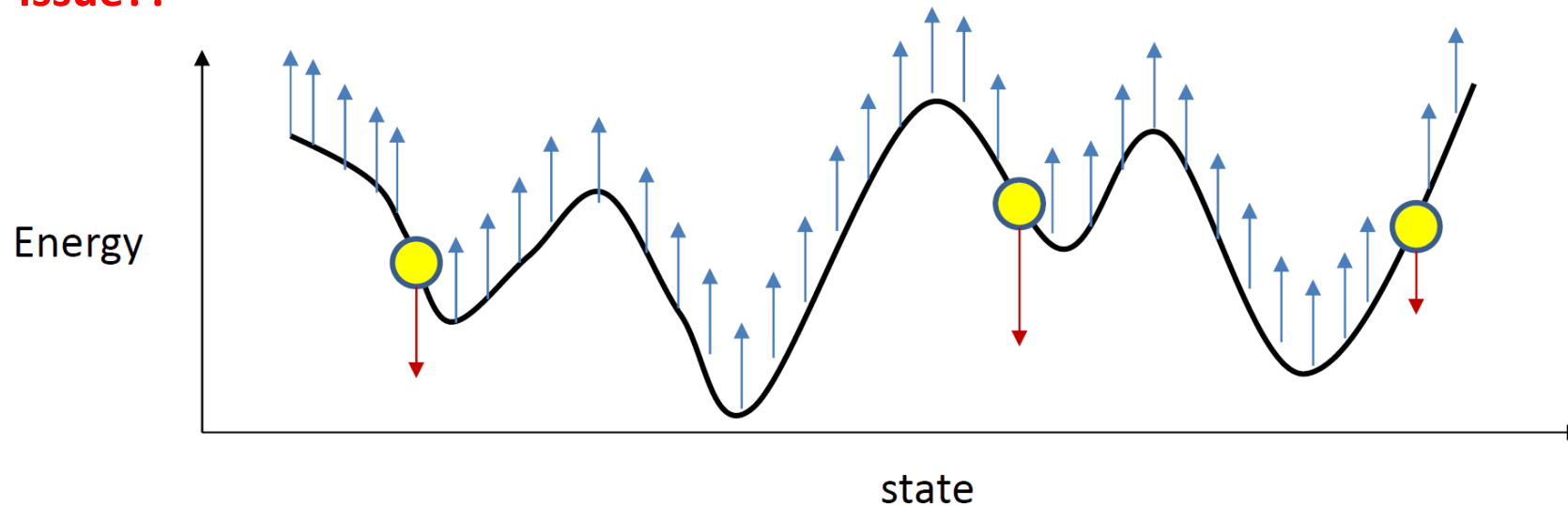- **Evolution of Hopfield Network will converge to a valley**

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \boldsymbol{Valley}} y'y'^T \right)$$

- Compute outer-products of desired patterns $y$
- Randomly initialize $y'$ for multiple times
  - Run evolution for random $y'$ until convergence
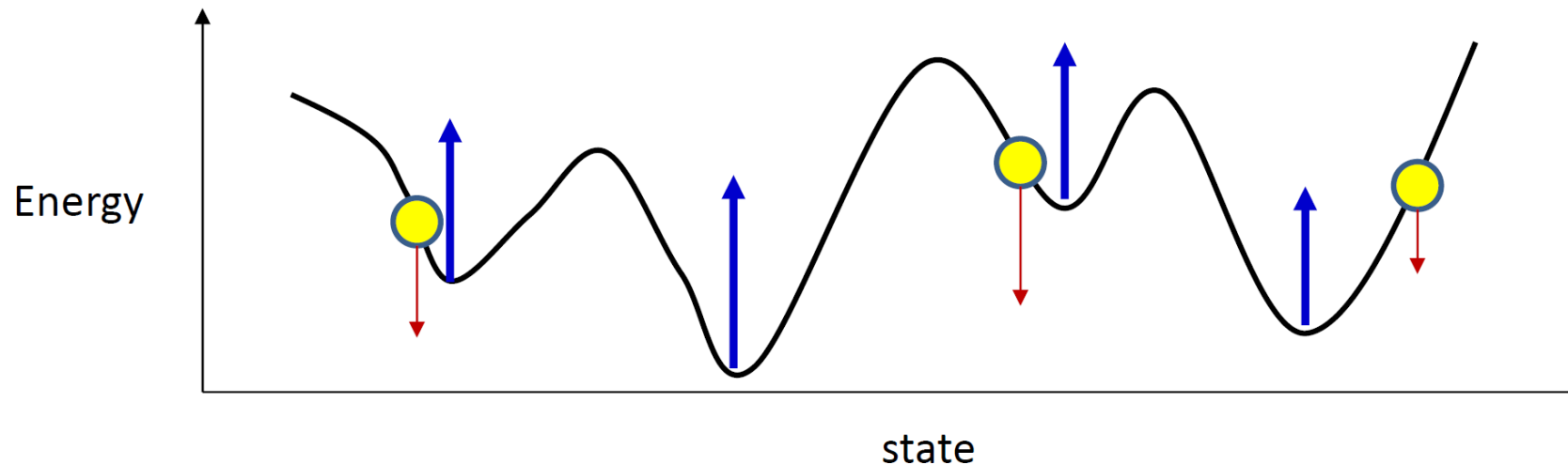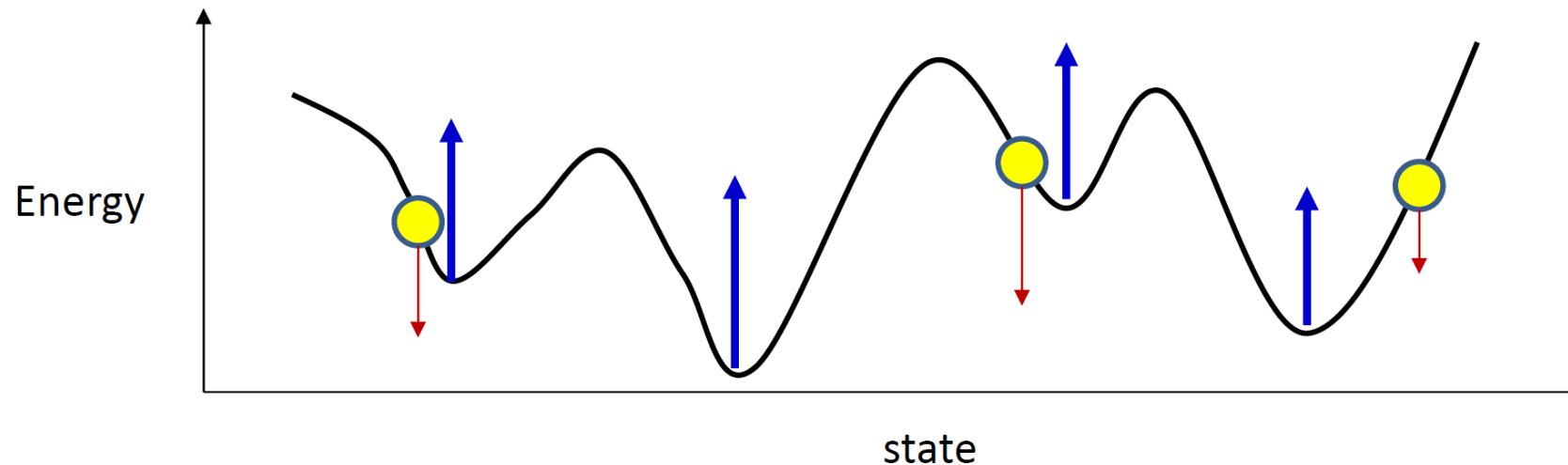  - Calculate outer-product of $y'$
- Compute gradient and update $W$

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \textbf{\textit{Valley}}} y'y'^T \right)$$

- Compute outer-products of desired patterns $y$
- Randomly initialize $y'$ for multiple times
  - Run evolution for random $y'$ until convergence
  - Calculate outer-product of $y'$
- Compute gradient and update $W$

- **Valleys are NOT equivalently important…**

# Hopfield Network: Optimization

- Which valleys are important?


- Primary object: ensure desired pattens stable
    - We want to ensure desired patterns are in broad valleys

# Hopfield Network: Optimization

- Which valleys are important?

- Primary object: ensure desired pattens stable
  - We want to ensure desired patterns are in broad valleys
  - **Spurious valleys around desired patterns are more important to eliminate**

# Hopfield Network: Optimization

- Which valleys are important?

- Primary object: ensure desired pattens stable
    - We want to ensure desired patterns are in broad valleys
    - Spurious valleys around desired patterns are more important to eliminate
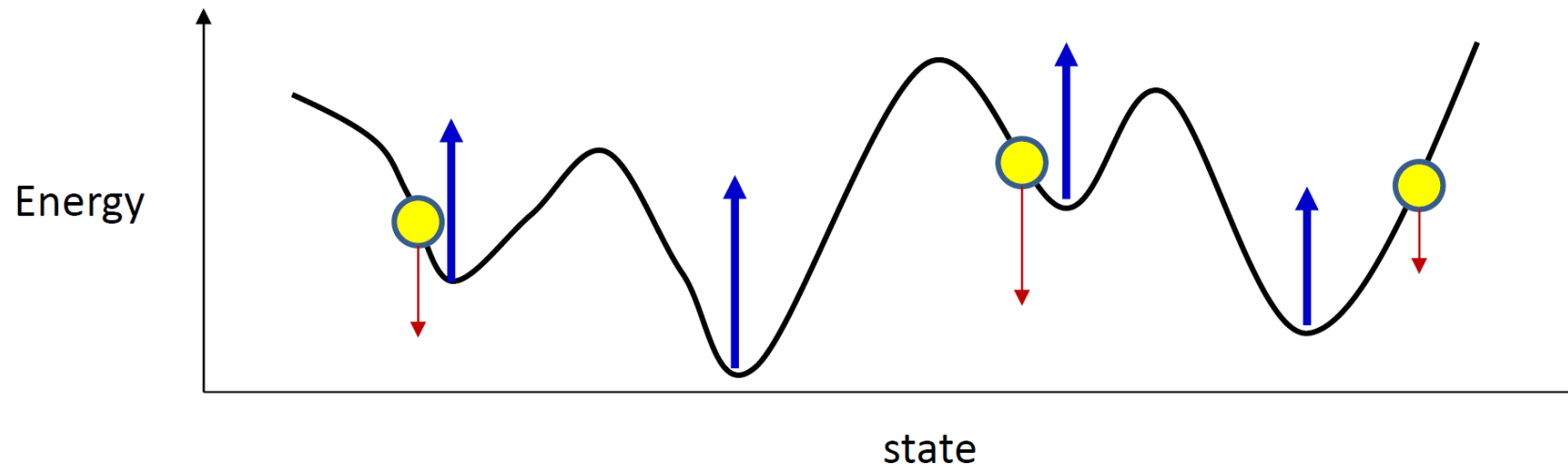    - **Evolution from desired patterns**

# Hopfield Network: Optimization

- Update rule for $W$

$$W \leftarrow W - \eta \left( \sum_{y \in P} yy^T - \sum_{y' \notin P \ \& \ y' \in \textbf{Valley}} y'y'^T \right)$$

  - Compute outer-products of desired patterns $y$
  - Initialize $y'$ by all the desired patterns
    - Run evolution for random $y'$ until convergence
    - Calculate outer-product of $y'$
  - Compute gradient and update $W$

  - **Still issues?**

# Hopfield Network: Optimization

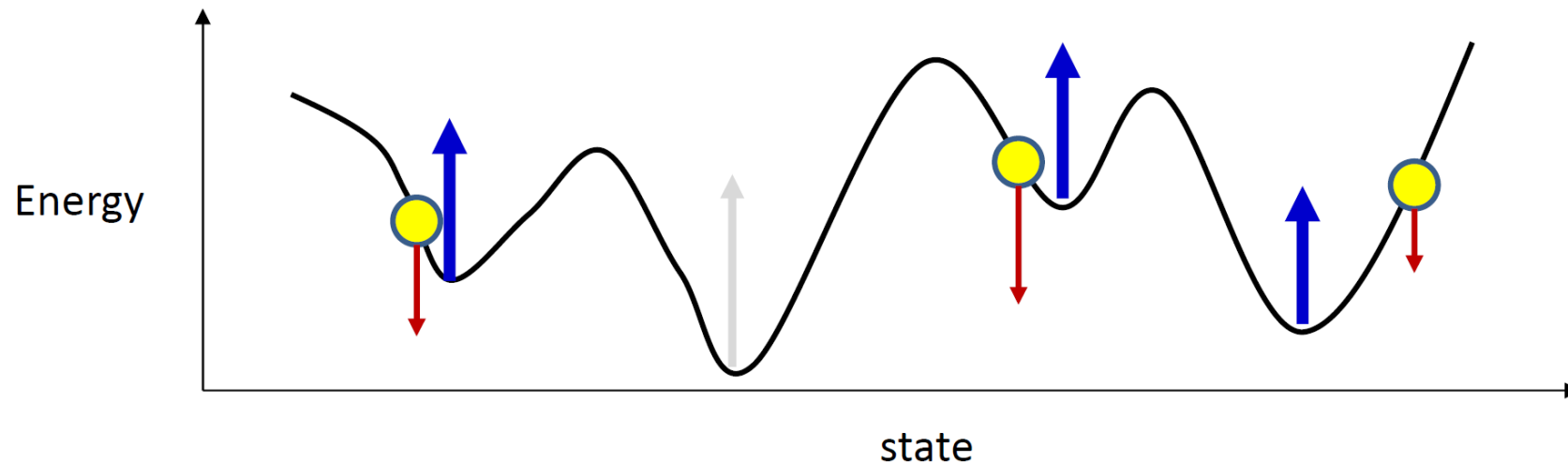• Recap: we raise the valleys next to the desired patterns

# Hopfield Network: Optimization

- Recap: we raise the valleys next to the desired patterns

- What if a pattern is close to the valley?
    - Naively forcing a valley to raise may hurt the learned representation
    - Particularly challenging when $y$ are continuously valued (e.g., tanh activation)

# Hopfield Network: Optimization

- New idea: we only raise the "neighborhood" of desired patterns!
    - It is sufficient to make each desired pattern a valley
    - Note: we want to raise the "decent" neighborhood

# Hopfield Network: Optimization

- New idea: we only raise the "neighborhood" of desired patterns!
  - It is sufficient to make each desired pattern a valley
  - Note: we want to raise the "decent" neighborhood

- Implementation
  - We initialize $y'$ by the desired patterns
  - **Only perform evolution for a few steps!**

# Hopfield Network: SGD Optimization

- SGD update rule for $W$
$$W \leftarrow W - \eta\left(\mathrm{E}_{y \in P}[yy^T] - \mathrm{E}_{y'}[y'y'^T]\right)$$

  - Compute outer-products of random desired pattern $y$
  - Initialize $y'$ by a random desired pattern
    - Run evolution for random $y'$ for a few timesteps (2~4)
    - Calculate outer-product of $y'$
  - Compute gradient and update $W$

- In theory, $O(N)$ patterns can be stored in the network (with undesired valleys)
  - How to store more patterns?

# The Expanded Network

- Idea: introduce redundant neurons to increase network capacity
- Original $N$ neurons for patterns: visible neurons
- Additional $K$ neurons: hidden neurons

# The Expanded Network

- Idea: introduce redundant neurons to increase network capacity
- Original $N$ neurons for patterns: visible neurons
- Additional $K$ neurons: hidden neurons



Visible bits          Hidden bits

$N$

$N + K$

# The Expanded Network

- $N$ dimensional pattern $\rightarrow$ $N + K$ dimension
    - How can we store the patterns with $K$ additional units?
    - Possible solution: random filling of $K$ units (not great but okay...)

# The Expanded Network

- How to retrieve the stored pattern?
  - Still evolution?
  - Evolution is performed on the entire network but we only care about $N$ units

A mechanism that can decouple $N$ visible units when needed!

Visible bits          Hidden bits



$N$          $N + K$

# The Expanded Network

- How to retrieve the stored pattern?
  - Idea: Probabilistic Framework $P_w(v, h)$ **Let's borrow some ideas from physics!**
    - Consider desired patterns by computing the marginal $P_w(v) = \sum_h P(v, h)$
    - How to convert Hopfield network to a distribution?

Visible bits                    Hidden bits



$N$                    $N + K$

# The Helmholtz Free Energy of a System

- Recap: A thermodynamic （热力学） system
  - We previously discussed a discrete-time version
  - In fact, it is a probabilistic system
- A thermodynamic system at temperature $T$
  - $P_T(S)$ the probability of the system at state $S$
  - $E_T(S)$ the potential energy at state $S$
  - $U_T$ the internal energy, the capability to do work
  - $H_T$ the entropy, internal disorder of the system
  - $k$ Boltzmann constant
  - Free energy $F_T = U_T - kTH_T$

# The Helmholtz Free Energy of a System

- A thermodynamic system at temperature $T$
  - Internal energy $U_T = \sum_S P_T(S)E_T(S)$
  - Entropy $H_T = -\sum_S P_T(S)\log P_T(S)$
  - Free energy $F_T = \sum_S P_T(S)E_T(S) - kT\sum_S P_T(S)\log P_T(S)$
- Minimum Free-Energy Principle
  - A system held at temperature T anneals by varying the rate at which it visits the various states until a minimum free-energy state is achieved
- Boltzmann Distribution
  - The probability distribution of states at equilibrium

# The Helmholtz Free Energy of a System

- Free energy

$$F_T = \sum_S P_T(S) E_T(S) + kT \sum_S P_T(S) \log P_T(S)$$

- Boltzmann distribution (minimize $F_T$ w.r.t. $P_T(S)$)

$$P_T(S) = \frac{1}{Z} \exp\left(-\frac{E_T(S)}{kT}\right)$$

  - It is also known as Gibbs distribution
  - $Z$ normalizing constant

Given an energy function $E_T(S)$, if we follow a proper physical evolution process, the system will converge to the Boltzmann distribution

# Stochastic Hopfield Network

- Let's model our Hopfield network as a thermodynamic system
  - $T = k = 1$ for simplicity
  - Energy

$$E(y) = -\sum_{i<j} w_{ij} y_i y_j - b_i y_i$$

  - Boltzmann Probability

$$P(y) = \frac{1}{Z} \exp\left(\sum_{i<j} w_{ij} y_i y_j + b_i y_i\right)$$

- Stochastic Hopfield Network
  - Models the stationary probability distribution of states
  - Generative model: generate state from $P(y)$

# Stochastic Hopfield Network



- Let's consider the "flip" operation
  - Deterministic $\rightarrow$ probabilistic
  - Goal: change $y_i$ to 1 with probability $P(y_i = 1 | y_{j \neq i})$

- Assume $y$ and $y'$ only differ at position $i$ and $y_i' = -1$
  - $\log P(y) = -E(y) + C$
  - $E(y) = -\sum_{i<j} w_{ij} y_i y_j - b_i y_i$
  - $\log P(y) - \log P(y') = E(y') - E(y) = -\sum_j w_{ij} y_j - 2b_i$

$$\log \frac{P(y)}{P(y')} = \log \frac{P(y_i = 1 | y_{j \neq i}) P(y_{j \neq i})}{P(y_i' = -1 | y_{j \neq i}') P(y_{j \neq i}')} = \log \frac{P(y_i = 1 | y_{j \neq i})}{1 - P(y_i = 1 | y_{j \neq i})} = -\sum_j w_{ij} y_j - 2b_i$$

# Stochastic Hopfield Network



- Let's consider the "flip" operation
  - Deterministic $\rightarrow$ probabilistic
  - Goal: change $y_i$ to 1 with probability $P(y_i = 1|y_{j \neq i})$

- Assume $y$ and $y'$ only differ at position $i$ and $y_i' = -1$
  - $\log P(y) = -E(y) + C$
  - $E(y) = -\sum_{i<j} w_{ij} y_i y_j - b_i y_i$
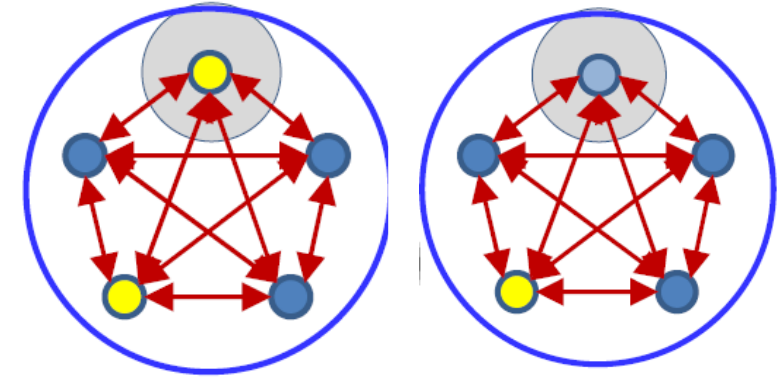  - $\log P(y) - \log P(y') = E(y') - E(y) = -\sum_j w_{ij} y_j - 2b_i$

$$\log \frac{P(y)}{P(y')} = \log \frac{P(y_i = 1|y_{j \neq i})P(y_{j \neq i})}{P(y_i' = -1|y_{j \neq i}')P(y_{j \neq i}')} = \log \frac{P(y_i = 1|y_{j \neq i})}{1 - P(y_i = 1|y_{j \neq i})} = -\sum_j w_{ij} y_j - 2b_i$$
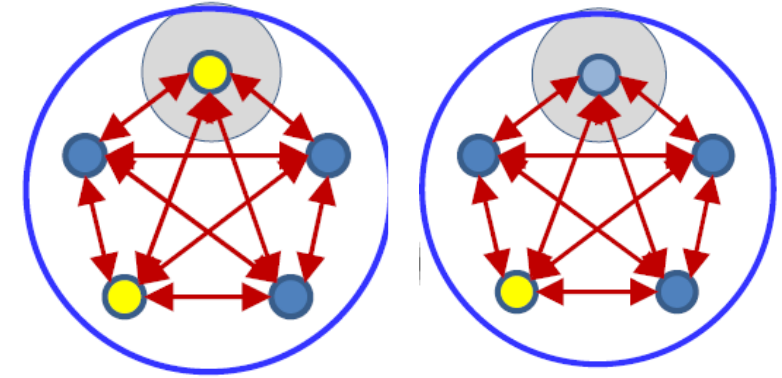
- A sigmoid conditional: $P(y_i = 1|y_{j \neq i}) = \dfrac{1}{1 + \exp\left(-\sum_j w_{ij} y_j - 2b_i\right)}$

*This is also called Gibbs sampling (remember the name for now ☺)*

# Stochastic Hopfield Network



- The whole update rule
  - Field at $y_i$: $z_i = \sum_j w_{ij} y_j + 2b_i$
  - $P\left(y_i = 1 \middle| y_{j \neq i}\right) = \frac{1}{1 + \exp(-z_i)} = \sigma(z_i)$

  <span style="color:red">Field quantifies the delta energy of flip</span>

- Running the network
  - Randomly initialize $y$
  - Cycle over $y_i$, fixed other variables fixed and sample $y_i$ according to the conditional probability
  - After "convergence", we can get samples of $y$ according to $P(y)$
  - *This sampling procedure is called Gibbs sampling*
  - **<span style="color:red">How can we retrieve the stored pattern???</span>**
    - **<span style="color:red">This is a stochastic process!</span>**

# Stochastic Hopfield Network

- Network evolution
  - initialize $y_0$
  - For $1 \le i \le N, y_i(t+1) \sim Bernoulli(\sigma(z_i(t)))$
  - Until convergence

- Retrieve a stored pattern $y$
  - Given sequence of samples $y_0, \dots, y_L$
  - Simply take the average of final $M$ samples

$$y_i = I\left[\frac{1}{M}\sum_{t=L-M+1}^{L} y_i(t) > 0\right]$$

  - If you want a probability instead of a single vector, you can use the frequency derived from $\{y_{L-M+1}, \dots, y_L\}$ to approximate the stationary distribution
  - **In many applications, we simply take $M = 1$ (output $y_L$)**

# Stochastic Hopfield Network: Annealing

- Find the state with lowest energy?

- Network evolution
  - initialize $y_0, T \leftarrow T_{\max}$
  - Repeat
    - Repeat a few cycles
      - For $1 \leq i \leq N, y_i(T) \sim Bernoulli\left(\sigma\left(\frac{1}{T}z_i(T)\right)\right)$
    - $y_i(\alpha T) \leftarrow y_i(T); T \leftarrow \alpha\, T$
  - Until convergence

- Final state as the retrieved pattern
  - With temperature annealing, the system will converge to the most likely state
  - Possibly local minimum in practice

# Boltzmann Machine

- A generative Model
  - $E(y) = -\frac{1}{2} y^T W y$
  - $P(y) = \frac{1}{Z} \exp\left(-\frac{E(y)}{T}\right)$
  - Parameter $W$

**How to learn $W$ for desired patterns?**

- It has a probability for producing any binary pattern $y$
  - We assume $y_i = 0$ or 1 (or $\pm 1$)

$$z_i = \frac{1}{T} \sum_j w_{ji} s_j$$

$$P(s_i = 1 | s_{j \neq i}) = \frac{1}{1 + e^{-z_i}}$$

# Boltzmann Machine: Training

- Goal
  - Remember a set of desired patterns $P = \{y^p\}$
  - Now we have a probability distribution

- Objective: maximum likelihood learning (assume $T = 1$)
  - Probability of a particular pattern

$$P(y) = \frac{\exp\left(\frac{1}{2}y^T W y\right)}{\sum_{y'} \exp\left(\frac{1}{2}y'^T W y'\right)}$$

  - Maximize log-likelihood

$$L(W) = \frac{1}{N_P} \sum_{y \in P} \frac{1}{2} y^T W y - \log \sum_{y'} \exp\left(\frac{1}{2}y'^T W y'\right)$$

# Boltzmann Machine: Training

- Maximize log-likelihood

$$L(W) = \frac{1}{N_P} \sum_{y \in P} \frac{1}{2} y^T W y - \log \sum_{y'} \exp\left(\frac{1}{2} y'^T W y'\right)$$

- Gradient Ascent $\nabla_{w_{ij}} L$

# Boltzmann Machine: Training

- Maximize log-likelihood

$$L(W) = \frac{1}{N_P} \sum_{y \in P} \frac{1}{2} y^T W y - \log \sum_{y'} \exp \left( \frac{1}{2} y'^T W y' \right)$$

- Gradient Ascent $\nabla_{w_{ij}} L$

  - $\nabla_{w_{ij}} L = \frac{1}{N_P} \sum_{y \in P} y_i y_j$

# Boltzmann Machine: Training

- Maximize log-likelihood

$$L(W) = \frac{1}{N_P} \sum_{y \in P} \frac{1}{2} y^T W y - \boxed{\log \sum_{y'} \exp\left(\frac{1}{2} y'^T W y'\right)}$$

- Gradient Ascent $\nabla_{w_{ij}} L$

  - $\nabla_{w_{ij}} L = \frac{1}{N_P} \sum_{y \in P} y_i y_j - \sum_{y'} \frac{\exp\left(\frac{1}{2} y'^T W y'\right)}{Z} \cdot y_i' y_j'$    Exponentially many terms!

# Boltzmann Machine: Training

- Maximize log-likelihood

$$L(W) = \frac{1}{N_P} \sum_{y \in P} \frac{1}{2} y^T W y - \boxed{\log \sum_{y'} \exp\left(\frac{1}{2} y'^T W y'\right)}$$

- Gradient Ascent $\nabla_{w_{ij}} L$

  - $\nabla_{w_{ij}} L = \frac{1}{N_P} \sum_{y \in P} y_i y_j - \sum_{y'} \frac{\exp\left(\frac{1}{2} y'^T W y'\right)}{Z} \cdot y_i' y_j'$

  - $\nabla_{w_{ij}} L = \frac{1}{N_P} \sum_{y \in P} y_i y_j - E_{y'}\left[y_i' y_j'\right]$  **Monte-Carlo Approximation**

  - Draw a set of samples $S$ for $y'$ according to the probability,

  - $\nabla_{w_{ij}} L = \frac{1}{N_P} \sum_{y \in P} y_i y_j - \frac{1}{|S|} \sum_{y' \in S} y_i' y_j'$

# Boltzmann Machine: Training

- Maximize log-likelihood with $M$ Monte-Carlo samples

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P} \sum_{y \in P} y_i y_j - \frac{1}{M} \sum_{y' \in S} y_i' y_j'$$

- How to draw samples from $P(y)$?
  - Running the stochastic network (Gibbs sampling)
  - Randomly initialize $y(0)$
  - Cycle over $y_i(t)$, sampling according to $P(y_i(t)|y_{j \neq i}(t))$
  - After convergence, we get a sequence of samples $\{y(0), \dots, y(L)\}$
  - Get the final $M$ states as samples $S = \{y(L - M + 1), \dots, y(L)\}$

# Boltzmann Machine: Training
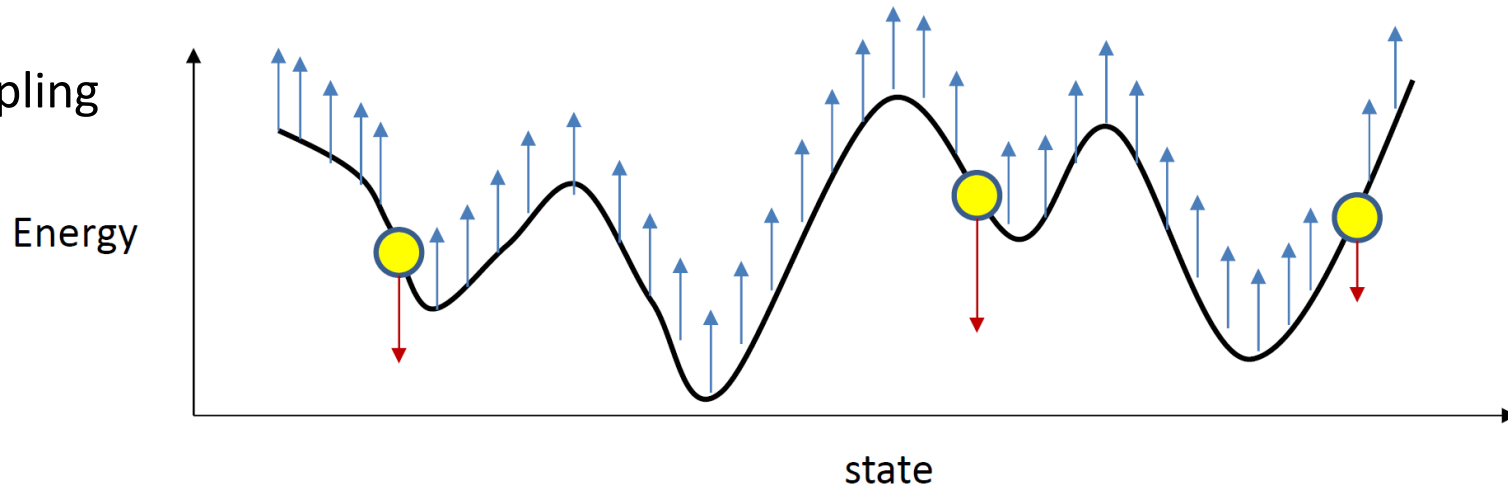
- Overall Training
  - Initialize $W$
  - Maximize log-likelihood with $M$ Monte-Carlo samples

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P} \sum_{y \in P} y_i y_j - \frac{1}{M} \sum_{y\prime \in S} y_i' y_j'$$

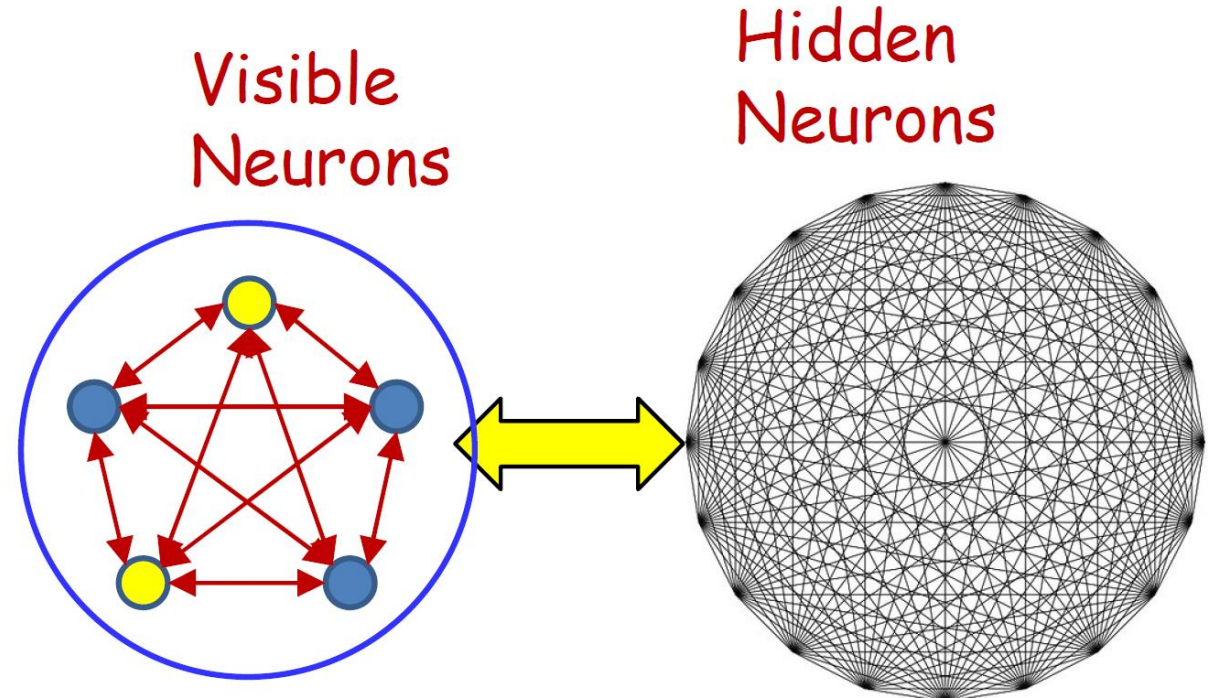  - $w_{ij} \leftarrow w_{ij} + \eta \nabla_{w_{ij}} L(W)$ (*we are maximizing likelihood*)
  - Remark
    - We can also run Gibbs sampling from states in $P$ (will discuss later...)

Energy

state

# Boltzmann Machine with Hidden Neurons

- Let's get back to hidden neurons!
  - $v$ visible neurons (pattern), $h$ hidden neurons
  - $y = (v, h)$
- A joint probability distribution
  - $P(y) = P(v, h)$
  - $P(v) = \sum_h P(v, h)$
    - **The marginal distribution!**
  - $h$: latent representation
- New objective
  - Maximize the marginal probability

Visible
Neurons

Hidden
Neurons

# Boltzmann Machine with Hidden Neurons

- Maximum log-likelihood learning

$$P(v) = \sum_h P(v,h) = \sum_h \frac{\exp(y^T W y)}{\sum_{y'} \exp(y'^T W y')}$$

$$L(W) = \frac{1}{|P|} \sum_{v \in P} \log\left(\sum_h \exp(y^T W y)\right) - \log\left(\sum_{y'} \exp(y'^T W y')\right)$$

- Gradient $\nabla L(W)$?

# Boltzmann Machine with Hidden Neurons

- Maximum log-likelihood learning

$$P(v) = \sum_h P(v,h) = \sum_h \frac{\exp(y^T W y)}{\sum_{y'} \exp(y'^T W y')}$$

$$L(W) = \frac{1}{|P|} \sum_{v \in P} \log\left(\sum_h \exp(y^T W y)\right) - \color{red}{\log\left(\sum_{y'} \exp(y'^T W y')\right)}$$

- Gradient $\nabla L(W)$?

**Monte-Carlo Estimate!**

# Boltzmann Machine with Hidden Neurons

- Maximum log-likelihood learning

$$P(v) = \sum_h P(v,h) = \sum_h \frac{\exp(y^T W y)}{\sum_{y'} \exp(y'^T W y')}$$

$$L(W) = \frac{1}{|P|} \sum_{v \in P} \log\left(\sum_h \exp(y^T W y)\right) - \log\left(\sum_{y'} \exp(y'^T W y')\right)$$

- Gradient $\nabla L(W)$?
  - The first term is also in the form of log-sum
  - Monte Carlo Estimate for each $v \in P$!

# Boltzmann Machine with Hidden Neurons

- Maximum log-likelihood learning

$$\nabla_{w_{ij}} L(W) = \frac{1}{|P|} \sum_{v \in P} E_h[y_i y_j] - E_{y'}[y_i' y_j']$$

- Second term
  - Freely generate samples w.r.t. $p(y)$
  - Random initialization, cyclic Gibbs sampling

- First term
  - Generate samples w.r.t. $p(y)$ conditioned on a fixed $v$
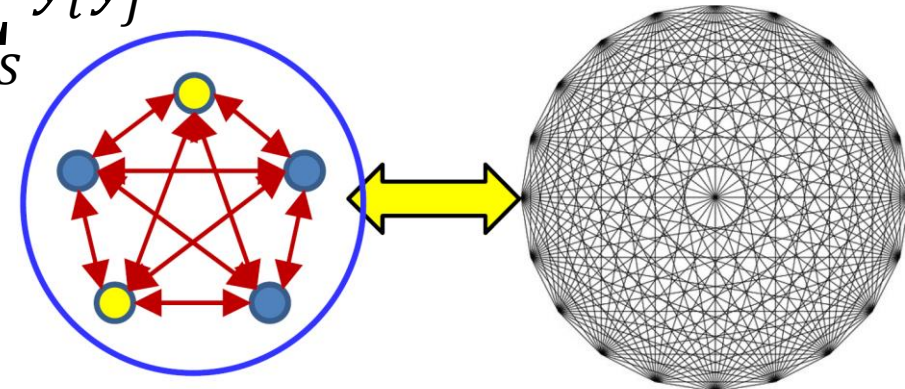  - Randomly initialize $h$, run Gibbs sampling over $h$

# Boltzmann Machine with Hidden Neurons

- Overall Training
  - Initialize $W$
  - For $v \in P$, fixed the visible neurons, run Gibbs sampling to get $K$ samples
    - Collect all conditioned samples as $S_c$
  - Randomly initialize all neurons, run Gibbs sampling to get $M$ samples
    - Collect free samples as $S$
  - Maximize log-likelihood with $M$ Monte-Carlo samples

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{y \in S_c} y_i y_j - \frac{1}{M} \sum_{y' \in S} y_i' y_j'$$

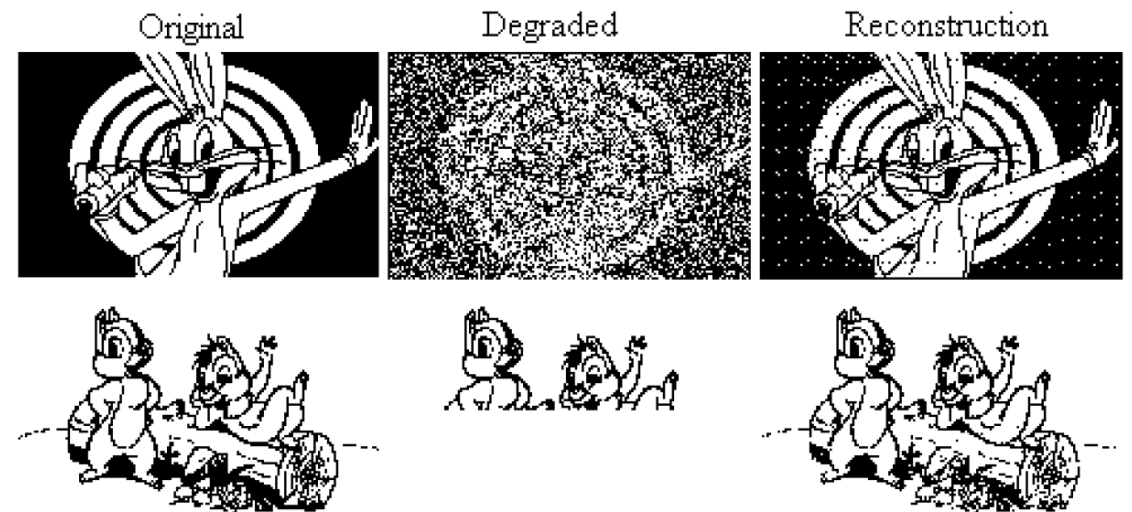  - $w_{ij} \leftarrow w_{ij} + \eta \nabla_{w_{ij}} L(W)$

# Boltzmann Machine

- Summary
  - A stochastic version of Hopfield Network
  - Nice mathematical properties
  - Large capacity for storing patterns (with hidden neurons)
  - Pattern generation
    - Gibbs sampling
  - Pattern completion
    - Conditioned Gibbs sampling
- ***Classification??***
  - *$y = (v, h, c)$, c is label*
  - *c as a one-hot vector (0-1 variables)*
  - Posterior $P(c|v)$
  - Even conditional generation: $P(v|c)$!

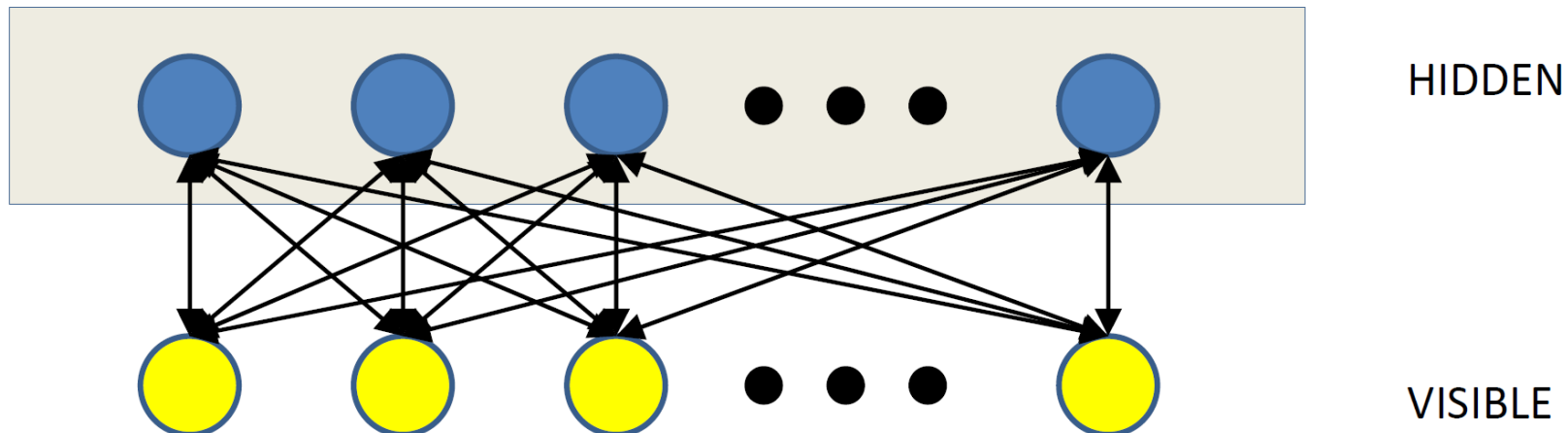

Original    Degraded    Reconstruction

Hopfield network reconstructing degraded images
from noisy (top) or partial (bottom) cues.

# Boltzmann Machine

- The issue
  - Training is hard!
  - Gibbs sampling may take a very long time to converge
    - also called ***mixing-time***
  - Not really applicable for large problems

- Can we design a better structure for faster Gibbs sampling mixing?

# Restricted Boltzmann Machine

- A particularly structured Boltzmann Machine
    - A partitioned structure
    - Hidden neurons are only connected to visible neurons
    - No intra-layer connections
    - *Invented under the name Harmonium by Paul Smolensky in 1986*
    - *Became promise after Hinton invented fast learning algorithms in mid-2000*



HIDDEN

VISIBLE

# Restricted Boltzmann Machine

- Computation Rules: same as Boltzmann machine
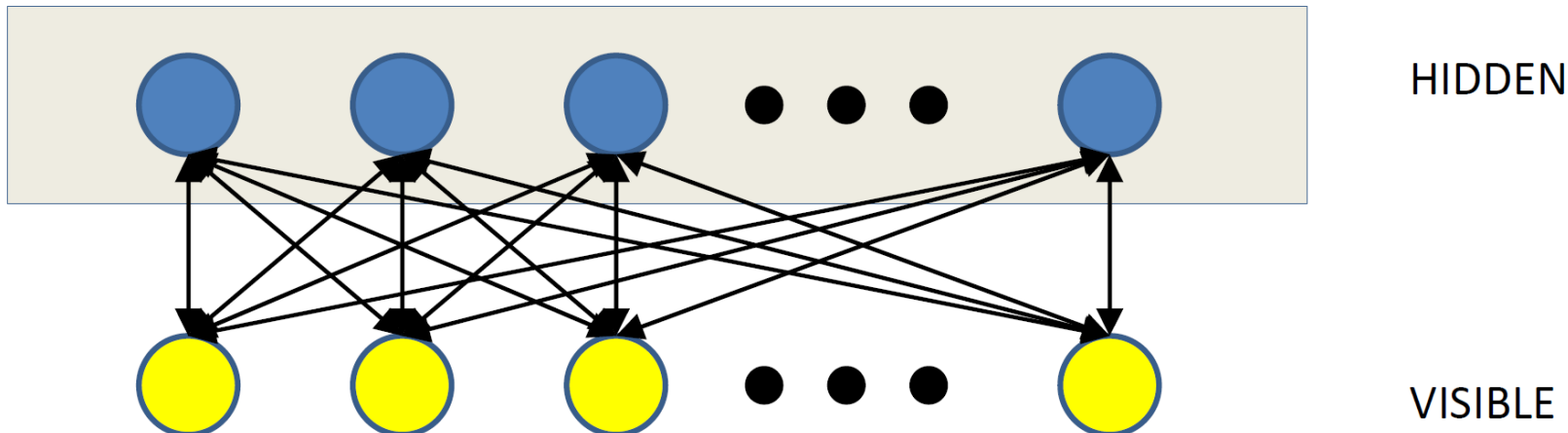  - Hidden neurons $h_i$

  $$z_i = \sum_j w_{ij} v_j, \qquad P(h_i = 1 | v_j) = \frac{1}{1 + \exp(-z_i)}$$

  - Visible neurons $v_j$

  $$z_j = \sum_i w_{ij} h_i, \qquad P(v_j = 1 | h_i) = \frac{1}{1 + \exp(-z_j)}$$
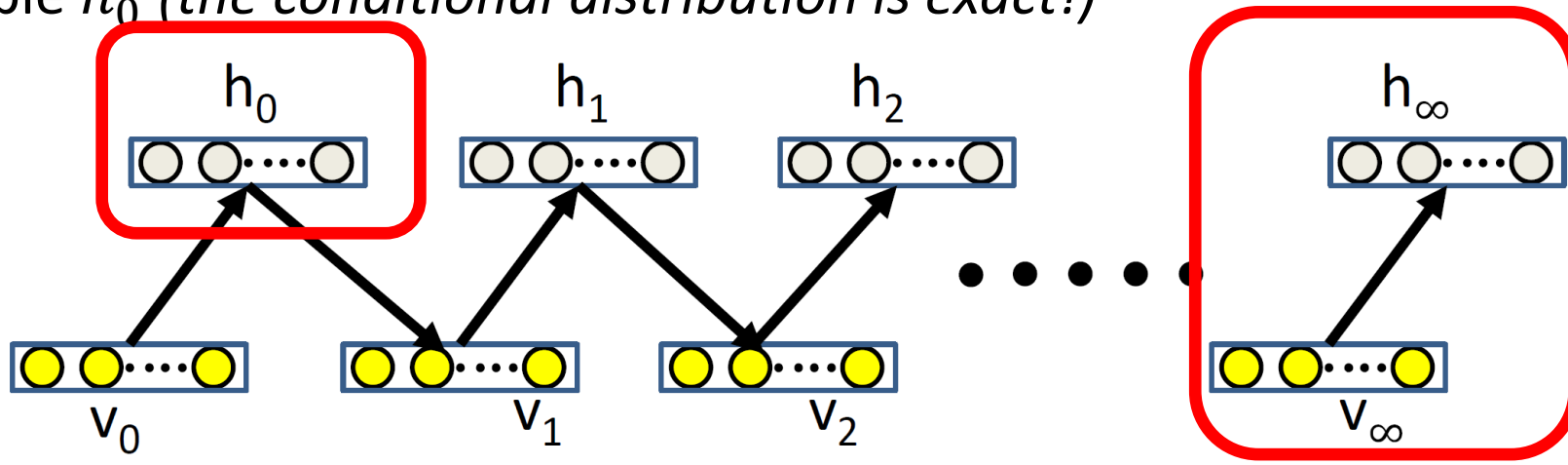
Iterative Sampling!



HIDDEN

VISIBLE

# Restricted Boltzmann Machine

- Sampling
  - Randomly initialize visible neurons $v_0$
  - Iterative between hidden neurons and visible neurons
  - Get final sample $(v_\infty, h_\infty)$

- Conditioned sampling?
  - Initialize $v_0$ as the desired pattern
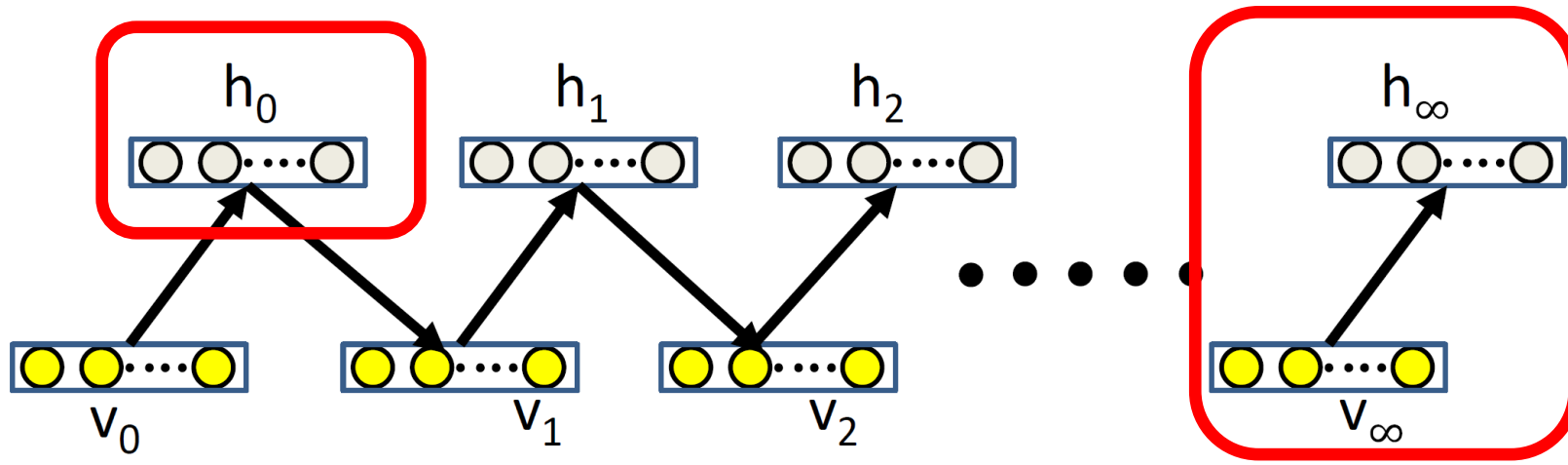  - Sample $h_0$ *(the conditional distribution is exact!)*

# Restricted Boltzmann Machine

- Maximum Likelihood Estimate

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{v \in P} v_{0_i} h_{0_j} - \frac{1}{M} \sum v_{\infty_i} h_{\infty_j}$$

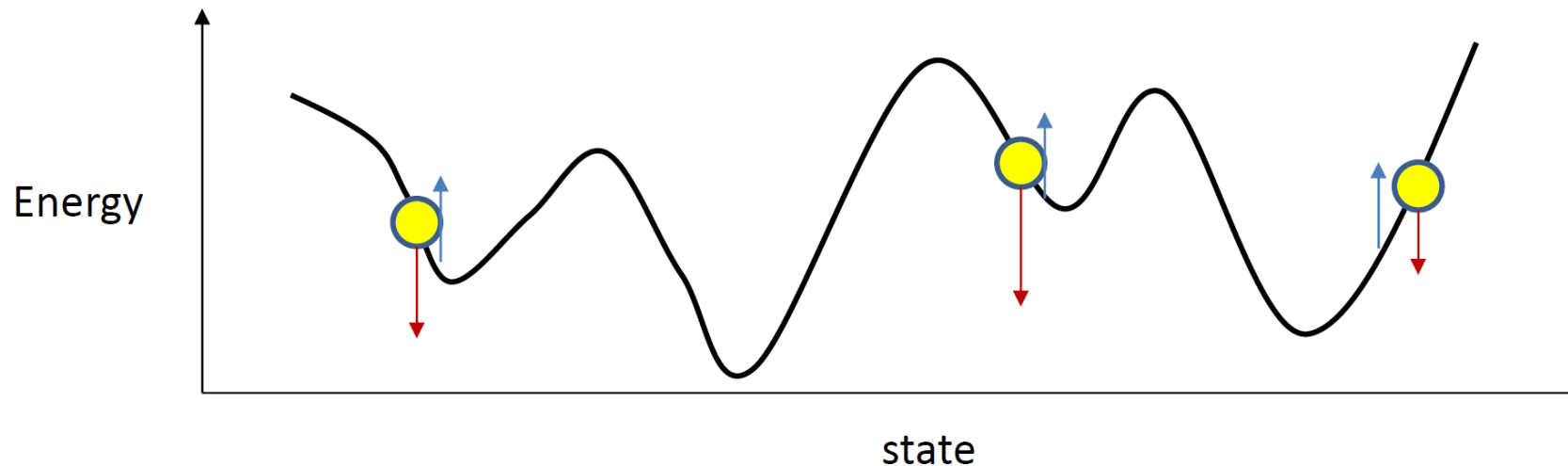  - No need to lift up the entire energy landscape! (recap)

# Restricted Boltzmann Machine

- Maximum Likelihood Estimate

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{v \in P} v_{0_i} h_{0_j} - \frac{1}{M} \sum v_{\infty_i} h_{\infty_j}$$

- No need to lift up the entire energy landscape! (recap)
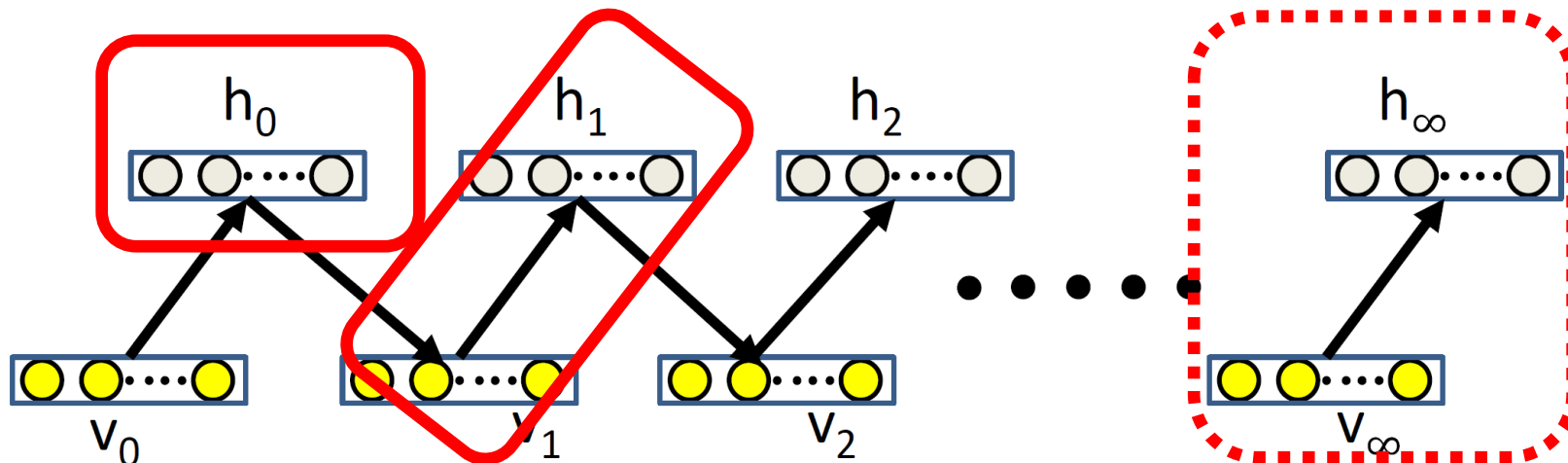  - Raising the neighborhood of desired patterns will be sufficient

# Restricted Boltzmann Machine

- Maximum Likelihood Estimate

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{v \in P} v_{0_i} h_{0_j} - \frac{1}{M} \sum v_{\infty_i} h_{\infty_j}$$

  - No need to lift up the entire energy landscape!
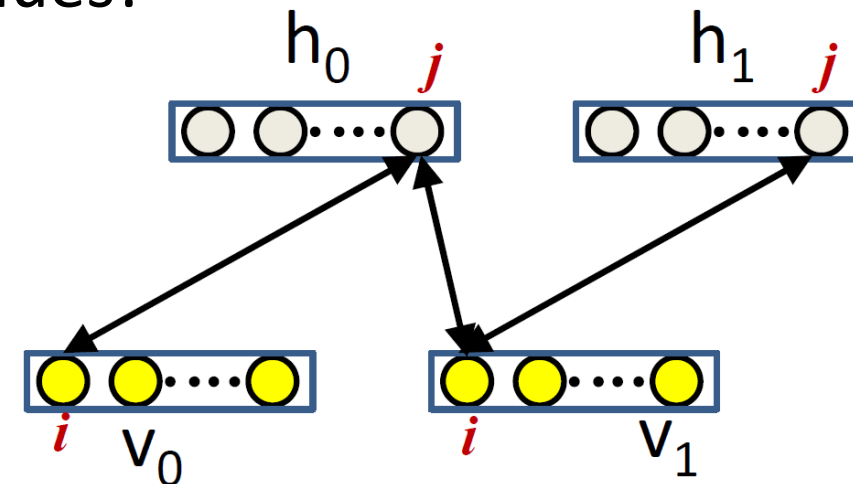    - One Gibbs sampling will be sufficient

# Restricted Boltzmann Machine
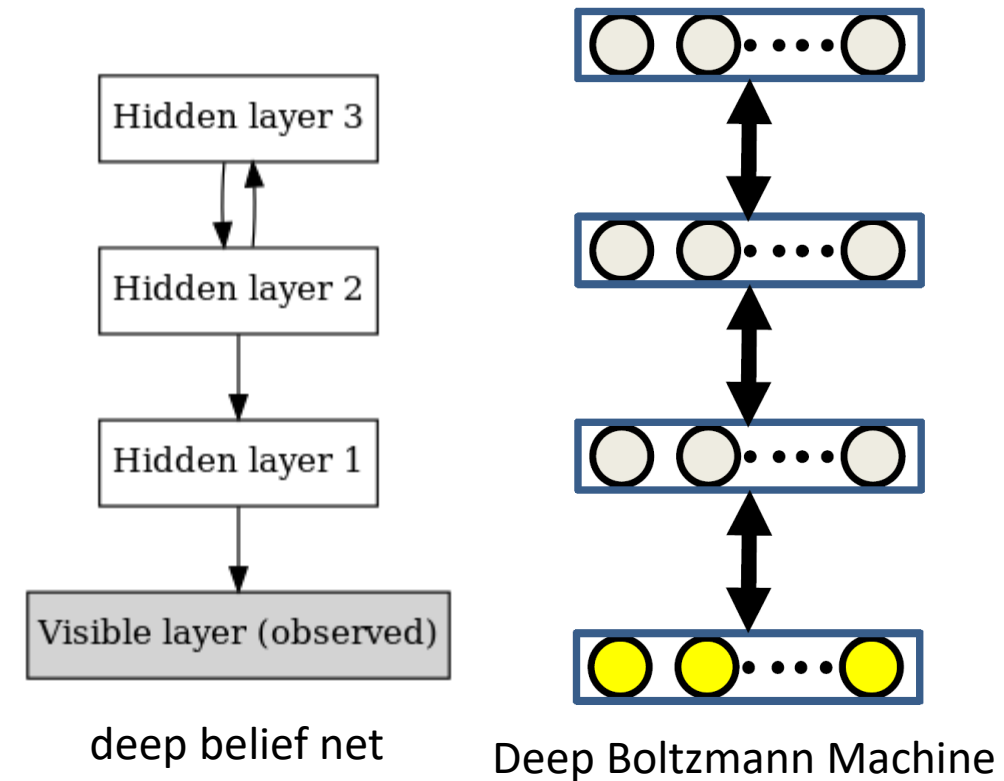
- Maximum Likelihood Estimate

$$\nabla_{w_{ij}} L(W) = \frac{1}{N_P} \sum_{v \in P} v_{0_i} h_{0_j} - v_{1_i} h_{1_j}$$

  - Only 3 Gibbs sampling steps are needed!

- We can also extend (R)BMs to to continuous values!
  - If we can explicitly sample from $P(y_i | y_{j \neq i})$
  - Exponential family! (FYI ☺)
    - "Exponential Family Harmoniums with an Application to Information Retrieval", Welling et al., 2004
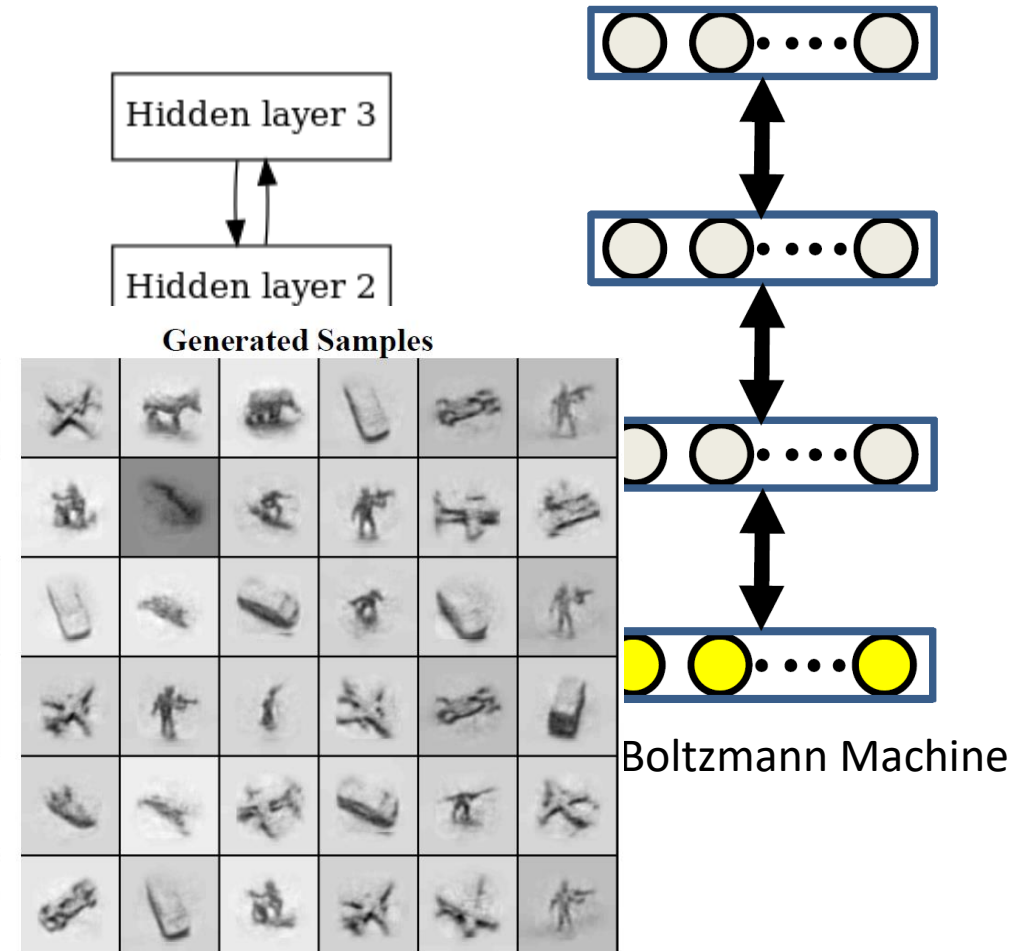
# Deep Boltzmann Machine
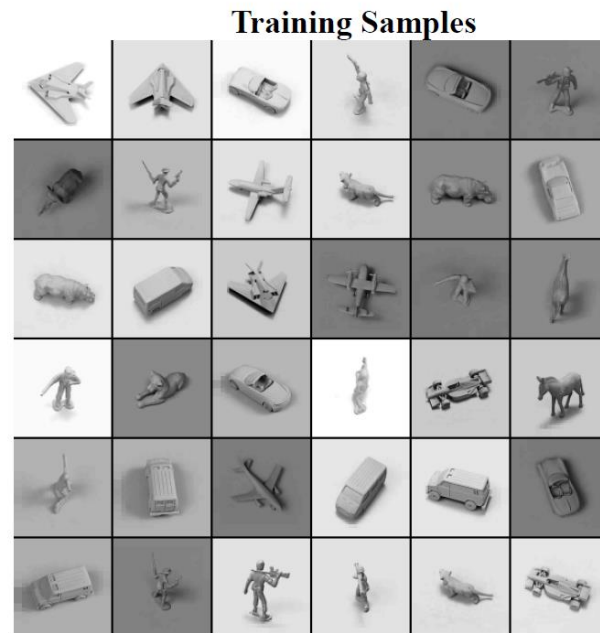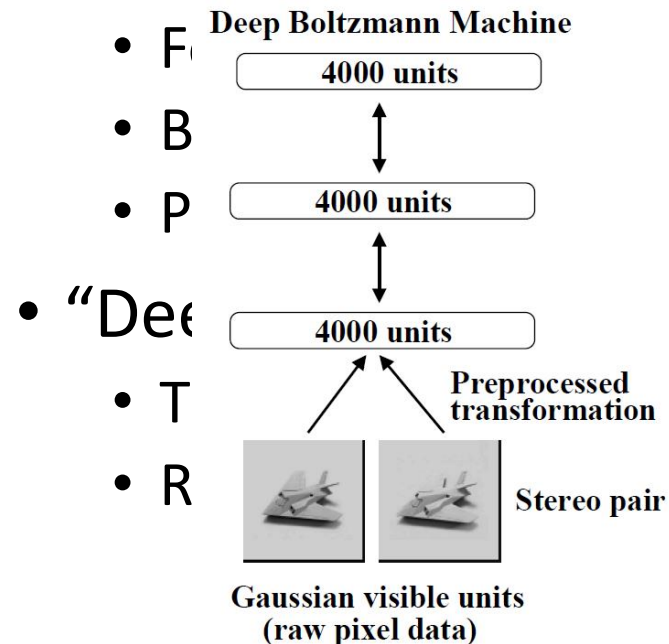
- Can we have a **deep** version of RBM?
  - Deep Belief Net (2006)
  - Deep Boltzmann Machine (2009)

- Sampling?
  - Forward pass: bottom-up
  - Backward pass: top-down
  - Practical Trick: Layer-by-layer pretraining

- "Deep Boltzmann Machine", AISTATS 2009
  - The very first deep generative model
  - Ruslan Salakhutdinov & Geoffrey Hinton



deep belief net

Deep Boltzmann Machine

# Deep Boltzmann Machine

- Can we have a ***deep*** version of RBM?
  - Deep Belief Net (2006)
  - Deep Boltzmann Machine (2009)

- Sampling?
  - F
  - B
  - P

- "Dee
  - T
  - R

**Deep Boltzmann Machine**

4000 units

4000 units

4000 units

Preprocessed transformation

Stereo pair

Gaussian visible units (raw pixel data)

Hidden layer 3

Hidden layer 2

**Training Samples**

**Generated Samples**

Boltzmann Machine

# Summary

- Hopfield Network
  - The very first generative neural network

- Boltzmann Machine
  - A stochastic version of Hopfield network
  - An undirected probabilistic model

- Restricted Boltzmann Machine
  - Layered structure for fast inference
- Next
  - General formulation of energy-based models

# Energy-Based Model

- Goal of generative model
  - A probability distribution of "patterns" $P(x)$

- Requirement
  - $P(x) \geq 0$ (non-negative)
  - $\int_x P(x)dx = 1$ (sum to 1)

- Energy-Based Model
  - Energy function: $E(x; \theta)$ parameterized by $\theta$
  - $P(x) = \frac{1}{Z}\exp(-E(x; \theta))$
  - $Z = \int_x \exp\big(-E(x; \theta)\big)\,dx$   *partition function*

Why use exp() function?
e.g. $|x|$ or $|x|^2$

# Energy-Based Model

- A particular class of density function

$$P(x) = \frac{1}{Z}\exp(-E(x;\theta))$$

- Pros
  - Compatible with log-probability measure to capture large variations
  - Exponential family (e.g., Gaussian)
  - Common in statistical physics
  - Extremely flexible, i.e., use any $E(x)$ you like (e.g., any $f(x): \mathbb{R}^d \to \mathbb{R}$, even CNNs)
- Cons
  - Non-trivial to sample and train due to the partition function $Z$
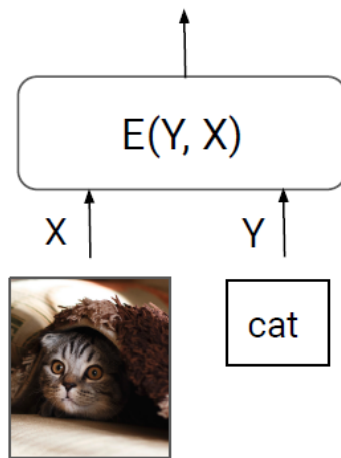  - Is it possible to avoid $Z$?

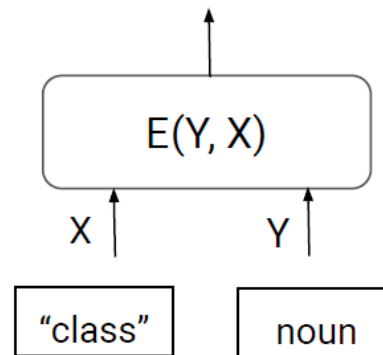# Energy-Based Model

- A particular class of density function

$$P(x) = \frac{1}{Z}\exp(-E(x;\theta))$$

- The ratio of two samples does not require $Z$!
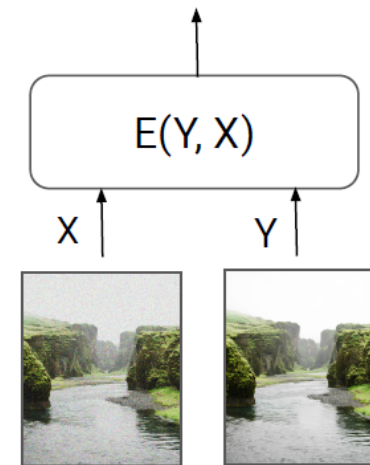
$$\frac{P(x)}{P(x')} = \exp(-E(x;\theta) + E(x';\theta))$$



| object recognition | sequence labeling | image restoration |

# Energy-Based Model: Training

- A particular class of density function
$$P(x) = \frac{1}{Z} \exp(-E(x; \theta))$$

- Maximum Likelihood Training
  - $L(\theta) = \log P(x) = -E(x; \theta) - \log Z(\theta)$
  - use Monte-Carlo Estimates for $Z(\theta)$

- Contrastive Divergence Algorithm
  - $\nabla_\theta L(\theta) \approx \nabla_\theta \big(-E(x_{train}; \theta) + E(x_{sample}; \theta)\big)$
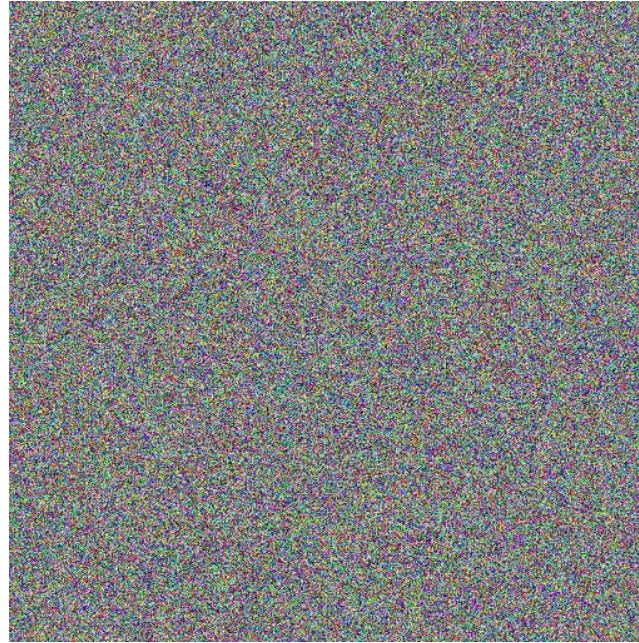
- How to sample from an energy-based model?

# A Generic Solution

- Sampling from the energy-based model $p(s) = \frac{1}{Z}\exp(-E(s))$
  - Random initialize $s^0$
  - $s' \leftarrow s^t + \text{noise}$
  - If $E(s') < E(s^t)$; then accept $s^{t+1} \leftarrow s'$
  - Else accept $s'$ with probability $\exp(E(s^t) - E(s'))$
  - Repeat
- Then after enough iterations, we get samples from $p(s)$

- *Details to be explained in the next lecture!* ☺

# Modern Energy-Based Model Examples



Du & Mordatch, 2019

Song et. al., 2021

**Additional Readings**
OpenAI Blog: https://openai.com/blog/energy-based-models/
A nice overview from Yang Song & Diederik Kingma: https://arxiv.org/abs/2101.03288

# Summary

- Hopfield Network
  - The first generative neural network
  - Undirected complete graph

- Boltzmann Machine
  - A probabilistic interpretation of Hopfield Network
  - The first deep generative model

- Energy-Based
  - Extremely flexible and powerful, designed to be multi-modal
  - Hard to sample and learn
  - Closely related to probabilistic inference and Bayesian methods

# Thanks!