

## Chapt. 7 Stochastic Methods

We've considered:

analytic methods       $\xrightarrow{\text{computing derivatives}}$        $\xrightarrow{\text{solving eqns}}$        $\Rightarrow$  optimal model parameters

for complicated models — local derivatives  
 $\rightarrow$  gradient methods —  $m_2$   
 $\rightarrow m_N$

These approach don't scale well

Direct search, i.e., exhaustive search is impractical

$\Rightarrow$  stochastic search

idea: 1) bias search towards promising regions  
2) let randomness drive exploration

Text discusses 2 classes of methods:

- 1) ~~Boltzmann~~ Simulated annealing / Boltzmann learning
- 2) evolutionary alg. i.e. genetic alg. & genetic programming

$\Rightarrow$  we'll start with Genetic Algorithms

## 7.2 Stochastic Search

7.2.1

Assume variables  $s_i \cdot i \geq 1, \dots, N$

$s_i$  are discrete values i.e.  $\pm 1$

$$E = -\frac{1}{2} \sum_{i,j=1}^N w_{ij} s_i s_j \quad \text{is the energy in this system}$$

Goal: choose values of  $s_i$  so as to minimize  $E$

Note:  $w_{ii} = 0$   $w_{ij}$  are bidirectional weighted edges  
symmetric i.e. interaction between  $s_i$

Obs:  $2^N$  possible configuration

⇒ not practical to solve directly

7.2.1 Introduction of Simulated annealing (Need to refresh classes memory?)

7.2.2 Boltzmann factor

Each configuration (choice of  $s_i$  values) is indexed by  $\gamma$

and has energy  $E_\gamma$  has the probability

$$P(\gamma) = \frac{e^{-E_\gamma/T}}{Z(T)}$$

where  $Z$  is a normalization constant and  $T$  is temperature.

The numerator:  $e^{-E_\delta/T}$  is the Boltzmann factor.

The denominator is the partition function,

$$Z(T) = \sum_{\delta} e^{-E_\delta/T}$$

i.e. the sum over all possible configurations

$\Rightarrow P(\delta)$  is a true probability

Note  $\delta$  ranges from 1 to  $2^N$

$\Rightarrow$  not practical to compute  $Z(T)$  in general

### Simulated Annealing Alg

Start with:

- High  $T(1)$
- random states  $s_i$

choose node  $i$  randomly suppose  $s_i = +1$

- calculate  $E_a$  for  $s_i$  in state  $+1$

- calculate  $E_b$  for  $s_i$  in state  $-1$

if  $E_b < E_a$  accept  $s_i = -1$

Else accept  $s_i = -1$  with probability

$$e^{-\Delta E_{ab}/T} \quad \text{where } \Delta E_{ab} = E_b - E_a$$

Note

with high  $T$  energetically less favored state will be accepted more often

$\Rightarrow$  local energy minima ~~can't~~ possible to escape

Or with large enough  $T$  all configurations are  $\sim e^0$  prob

~~repeat~~  $\Rightarrow$  continue randomly selecting + testing nodes

$\Rightarrow$  gradually decrease  $T$

Stop when  $T$  is very low

Show Algorithm 1 (Stochastic Simulated Annealing)

Note 1  $N_i$  denotes set of nodes connected with nonzero weights  $w_{ij}$  to node  $i$

Note 2 when comparing  $E_a + E_b \Rightarrow$  only consider nodes directly connected to  $s_i$  since  $s_i$  only affects its immediate neighbors in the ~~connected~~ network.

7 3 4

Issues: starting value of  $T$   
ending value  
cooling schedule possibility  $T(k+1) = cT(k)$   
stopping criterion  $c < < 1$

$T(1)$  should be high enough that all configurations are roughly equally likely.

$\Rightarrow T(1) \gg \text{any } \Delta E_{ab}$

, in practice

$\Rightarrow$  cooling must be gradual  $\Rightarrow 0.8 < c < 0.99$   
we need at <sup>allow</sup> least  $N/2$  transitions

why?  $\Rightarrow$  # transitions from any configuration do  
the global optimum is at most  $N/2$  units

Show fig 7.3 explain

Show fig 7.4 focus on probability of state  $\delta(\delta)$   
and expected energy  $E[E]$   
as a function of  $T$

7.2.5

## Deterministic Annealing

~~Solve~~ idea: let  $s_i$  be continuous and dependent on  $T$

$$\text{i.e. } s_i = f(l_i, T) = \tanh [l_i/T]$$

where  $l_i = \sum_j w_{ij} s_j$  i.e. effect of nodes connected to  $s_i$

Show Fig 7.5 explain

$\Rightarrow$  with large  $T$  every  $l_i$  does not force discrete  $s_i$

$\Rightarrow$  small  $T \Rightarrow$  discrete  $s_i$  with smaller  $l_i$

Show Algorithm 2 (Deterministic Simulated Annealing)

## Boltzmann Learning

7.3.1

Note: recall concept of opposing magnets connected in a network

modality: separate magnets/nodes into two sets

$\alpha$  - input/output nodes  $\xrightarrow{\alpha^o + \alpha^i}$

$\beta$  - internal/hidden nodes

Concept: clamp - hold input node fixed to the input value during classification

Probability of a visible configuration

$$\Rightarrow P(\alpha) = \sum_{\beta} P(\alpha, \beta)$$
$$= \frac{\sum_{\beta} e^{-E_{\alpha\beta}/T}}{Z}$$

where  $E_{\alpha\beta}$  is the energy defined of the entire sys

$$\text{i.e. } -\frac{1}{2} \sum_{i,j=1}^N w_{ij} s_i s_j \quad (\text{from last class})$$

and  $Z$  is again the full partition function  $\sum_{\gamma} e^{-E_{\gamma}/T}$

Q: How can we measure the difference between the desired probability of visible state  $\alpha$  & actual  $\alpha_2$

A: Kullback-Leibler divergence i.e. relative entropy

$$D_{KL}(Q(\alpha), P(\alpha)) = \sum_{\alpha} Q(\alpha) \log \frac{Q(\alpha)}{P(\alpha)}$$

note:  $D_{KL}$  is zero only if  $P(\alpha) = Q(\alpha)$  for all  $\alpha$

also - depends only on visible units  $\alpha$

learning: idea: adjust weights based on difference between desired and actual visible states

$$\Rightarrow \Delta w_{ij} = \eta \frac{\partial D_{KL}}{\partial w_{ij}}$$

↑ learning rate

⋮

$$> \frac{n}{T} \left[ E_Q[s_i s_j]_{\text{clamped}} - E[s_i s_j]_{\text{free}} \right]$$

Learning component aka teacher component

unlearning component aka student component

## Stochastic Learning of Input-Output Association

idea: clamp  $\alpha^i + \alpha^o$  in learning component  
clamp only  $\alpha^i$  in unlearning component

$$\Rightarrow \Delta w_{ij} = \frac{\eta}{T} [E_Q[s_i s_j] \alpha^o \text{clamped} - E[s_i s_j] \alpha^i \text{clamped}]$$

show fig 7.8.

## Evolutionary Methods

7.5.1

Inspiration: biological evolution

Note: biological evolution is massively parallel  
→ These algorithms are naturally parallel

basic idea: population of individuals are instances of a solution to an optimization problem

Fitness function: score individual  
→ fit individuals reproduce  
→ offspring inherit traits from parents

Example: 3-D conformation of a molecule

goal: → find most probable conformers  
→ lowest energy state solutions  
→ may be several stable states

reproduction would ~~combine~~ combine partial solutions  
from best/fittest parents

Text focuses on population of classifiers

7.5.2

⇒ optimization goal: find a good classifier

notion: chromosome - binary string: encodes classifier

note: Two requirements for using G.A. approach!

- 1) binary string encoding of ~~sln~~ soln instance
- 2) fitness fun for evaluating individual solns.

Show fig Alg 4

$\theta$  - fitness function

$p_{co}$  - probability of crossover

$p_{mut}$  - probability of site mutation

(Show fig 7.13)

Genetic Operators!

Replication: copy chromosome → no change

(Show fig 7.14)

Crossover: 2 chromosomes are "mated"

- 1) randomly select a split point
- 2) cut both chromosomes at same point
- 3) recombine to get 2 new chromosomes

Mutation : bit flip

All bits in a chromosome are evaluated for mutation

Issue: Representation : How is the classifier encoded in a bit string?

Singlet approach: bits represent features

Example on text: bits specify feature in 2-layer perceptron  
ie pixel features

①  $\Rightarrow$  maps out certain parts of string for certain features ~~with~~

Cross-over will mostly preserve segments

This allows good ~~and~~ partial solns to be combined

or

②

chromosome segments represent  $w$  in Multilayer NN.

$\Rightarrow$  need fixed topology (for mapping ~~to~~ purpose)

or

③

Show Fig 7.15 chromosome segments map to parts of a decision tree.

# Selection:

7.5.4

Fitness function: to choose <sup>ancestors</sup> parents of next generation

Pitfall: Q: What happens if we select only the most fit individuals

A: Approach a monoculture

⇒ inbreeding (:( )) → premature convergence!!

Note: progress depends on variance in a population

Sol 1. Always keep a few low scoring individuals

Soln 2. Use a probabilistic scoring function

→ occasionally low fitness individual selected

idea: probability of selection  $\propto$  fitness

Example: Show slide:  $P(i) = \frac{e^{f_i/T}}{\sum e^{f_i/T}}$

- probability of selecting chromosome  $i$
- with fitness  $f_i$
- and temperature  $T$  (control parameter)
- note: expectation is over current generation

temp is high early  $\Rightarrow$  all chromosomes have good chance for selection  
(later  $\rightarrow$  lower temperature)

# Genetic Programming

7.6.1

like GA, but different representation

GA → chromosomes are bit strings

GP → chromosomes are of computer programs

Show Fig 7.16

As with GA, we have:

replication - as before

cross-over: restrictions on valid snipping points

mutation: Element may mutate

i.e. number,  $\rightarrow$  number,

operator 1  $\rightarrow$  operator 2 (must have same arity)

Insertion: randomly insert code snippet from a set

Representation: must use some language

Spotted! Rich languages are difficult to use ex C++

$\Rightarrow$  LISP : syntactically very simple  $\Rightarrow$  already parsed  
by parentheses

7.6.2

LISP expressions are of the form

(Operator) <operand<sub>1</sub>> ... ) typically 1 - 2 operands

show by 7.17

notion: wrapper - decides if code snippet is meaningful  
→ catch ungrammatical code snippets