

## March 16, 2017 Evolutionary Algorithm

Not considered as "learning". Mainly "optimization"

Inspired by natural evolution

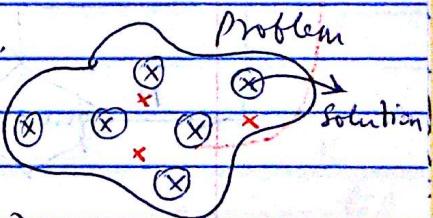
↳ Survival of the fittest (Natural selection)

- EAs are function optimizers

- Introduced by John Holland. (1975) "Adaptation in Natural and Artificial Systems!"

- GAs are EA<sup>s</sup> which are basically stochastic search methods (Genetic algorithms)

Sparcity is an enemy of stochastic algorithm.



Terminology

① population (candidate solutions = guesses)

② Generation (i.e., "iteration"; "epochs"; "episodes")

(in clustering) (Neural networks) → (in Reinforcement Learning)

③ Fitness → driving force every member of population that  
is fit is allowed to survive and reproduce itself.  
→ = objective function

35 million years

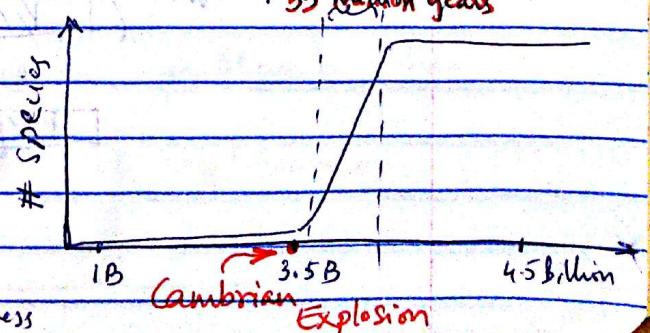
④ Crossover (getting offsprings)

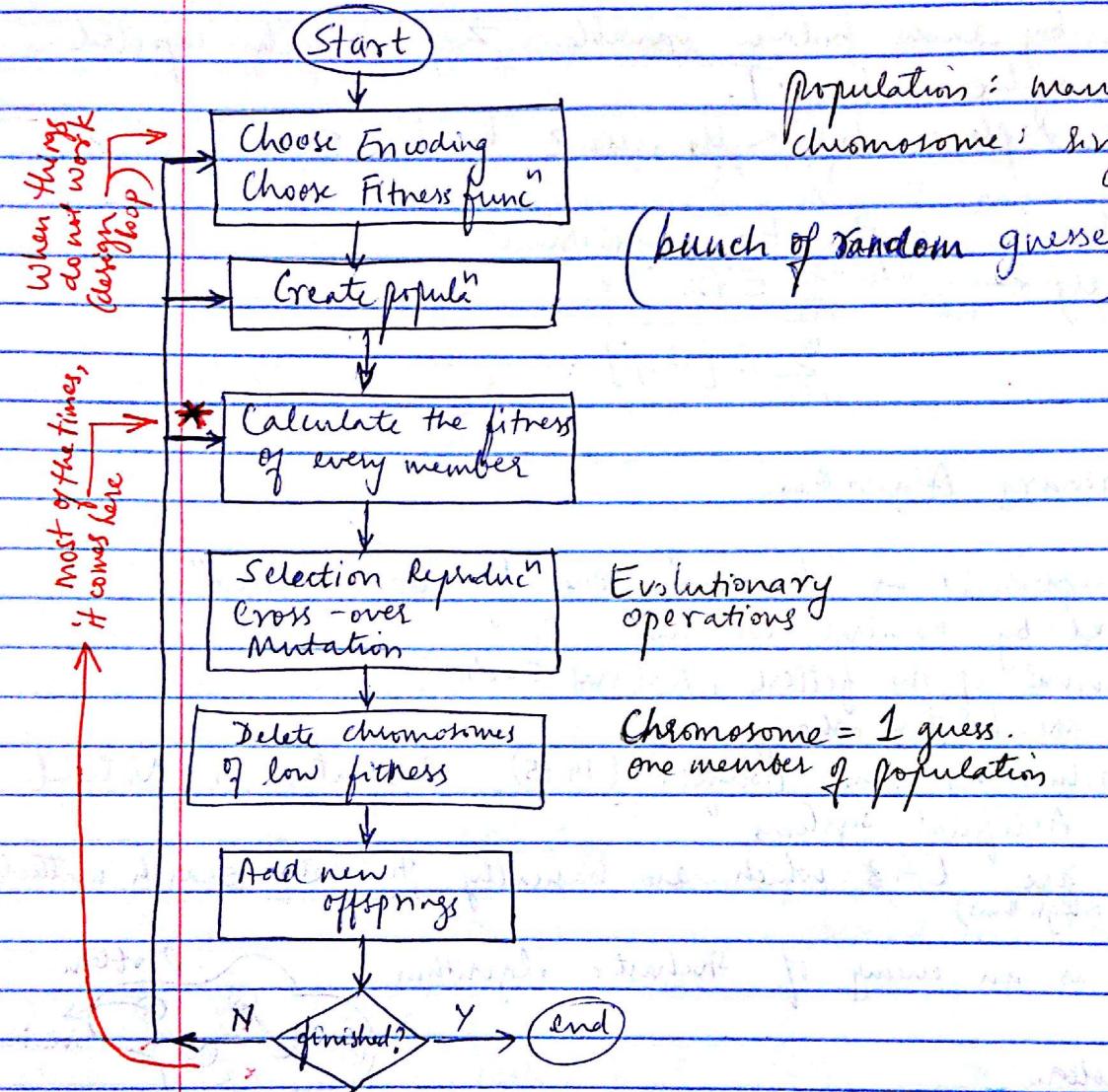
⑤ Reproduction

⑥ Mutation

⑦ Encoding (e.g. DNA → Binary code)

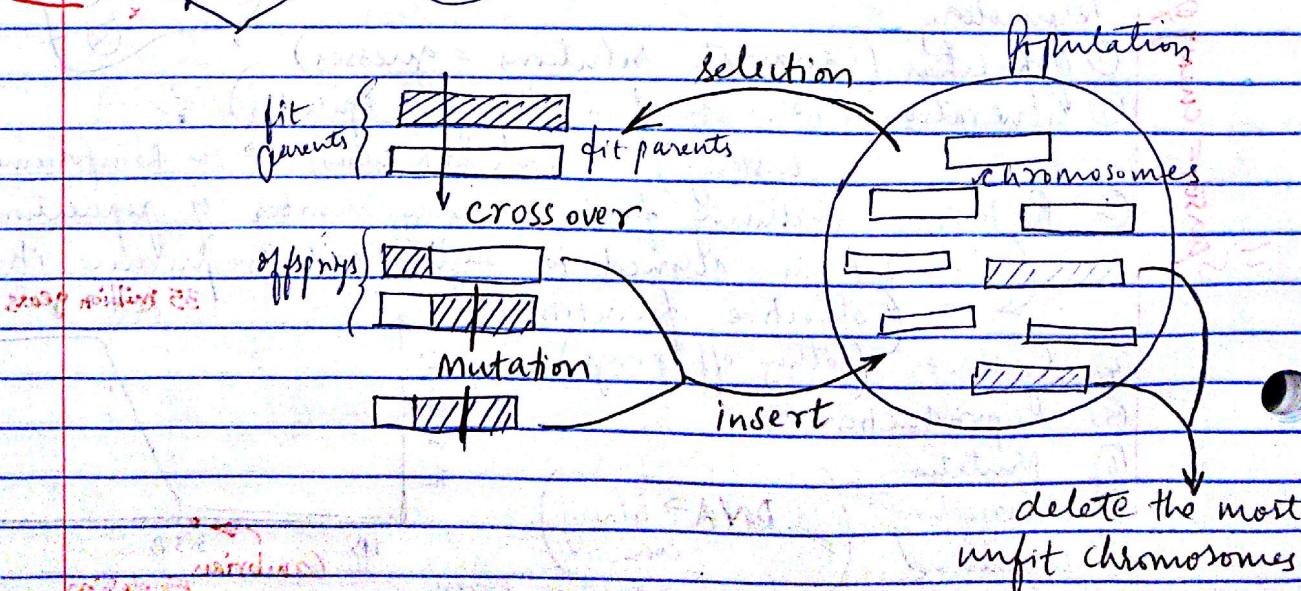
⑧ Chromosome (1 individual member in population. 1 guess)





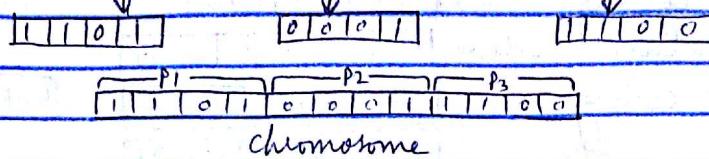
Evolutionary operations

Chromosome = 1 guess.  
one member of population



Function to be optimized

$$f(x) = P_1 x^2 + P_2 \log(x) + P_3 \sin(x)$$

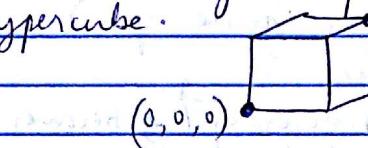


## Coding of the GA's secret

(Stochastic optimization)

we look at bit strings of size  $L_F = \boxed{01111111} \underbrace{\dots}_{L_F}$   
Size of search space =  $2^{L_F}$

GAs basically sample the corners of this  $L$  dimensional hypercube.



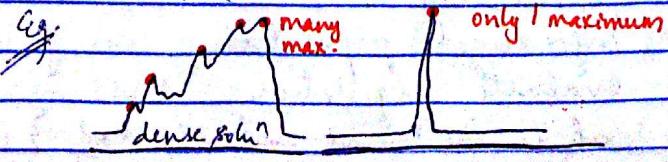
this analogy is  
only for  
binary  
encoding.

But the problems you can solve using ANN is  
something that cannot be solved using  
GA & vice versa.

Most (test functions =) problems can be encoded with  
max 30 bits. Anything less can be enumerated

Eg. chess game  $\rightarrow$  branching factor = 16  
 $\rightarrow$  a game is made of 100 moves (on average)  
 $\therefore 16^{100}$  possibilities.  $= (2^4)^{100} = 2^{400}$

If # of solutions is small (sparse problem) then  
stochastic search will not work.



fitness: driving force for evolutionary algo.  
 error: driving force for neural network.  
 reward/punishment: for learning algo.

## Encoding

① Binary coding

1|1|0|0|1|0

② Permutation coding

1|3|7|5|9|2|8|6|4

③ Value coding.

1.234|9.687|0.1501

back back left

ACWYBXYNIOJDI

## Diversity & GAs

global maximum

Converge towards  
actual maxima.

(but we may never  
get to exact maxima)

after many generations

$| \text{population} | = 5$  (5 chromosomes. i.e., 5 guesses randomly  
in the space)

→ Size of population?

- large, accuracy ↑, becomes sluggish
- small, less accurate, but fast.

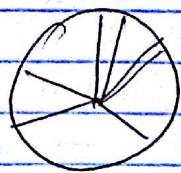
→ Cross over probability = 100%

→ Mutation probability = 1/1000 → in 1000 of flipper; you perform mutation. i.e., just flip a bit.

March 21<sup>st</sup>, 2017

## Selection of Chromosomes

① Roulette wheel selection



② Rank Selection

Sort chromosomes based on fitness, then rank them

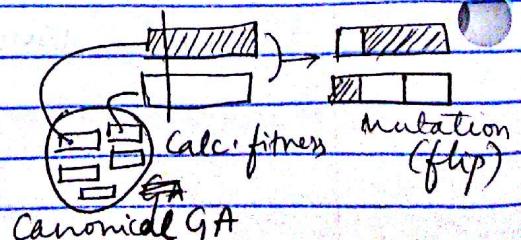
pseudo code

Genetic Alg () {

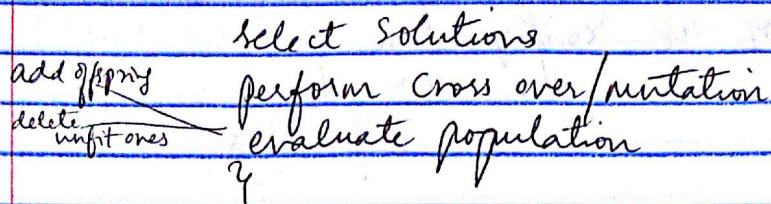
  initialize population,

  evaluate population

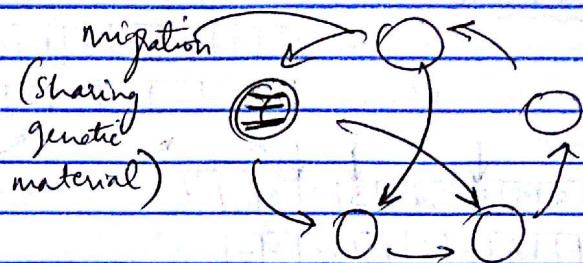
  while the end not reached {



- evolu" tech. are very slow  
! done offline, get solution  
go online if we at  
or application is not time critical
- low diversity  $\Rightarrow$  you are going to solution - stop.
- evolutionary alg. does not require data. No huge to save / store. You only obtain a set of solutions.
- diversity of solution over time



Many GA Models  
 $\hookrightarrow$  Island mode



Generally:

- Initialization - usually done randomly  
- can integrate existing solutions
- Termination - check it in every generation
  - by reaching a fitness
  - by arriving at max. # of generations
  - by achieving a minimum level of diversity.

Problem: 8- Queens Problem

The problem is to place 8 queens on chessboard such that they do not check each other.

Question: What is the fitness function?

fitness:  $\frac{1}{\# \text{check.}}$

fitness function: (This has to be maximized)  
 $\# \text{ of checks have to be min.}$

Manually: 92 solutions.

Phenotype:

	1	2	3	4	5	6	7	8
1								
2								
3	•					•		
4		•		•	•		•	
5								
6								
7								
8								

chromosome is based on permutations

Chromosome:  $[3 | 4 | 6 | 4 | 4 | 3 | 4 | 4]$

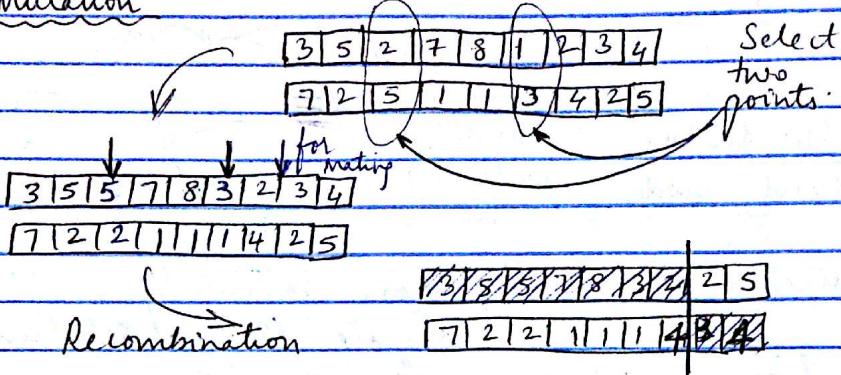
- Fitness:
- Penalty of one queen is the # of queens she can check.
  - Penalty for configuration is the sum of all penalties.
  - Penalty has to be minimised

④ Fitness of the configuration :

$$\max \frac{1}{\sum \text{penalties}}$$

Elements  
(has to be 8 in every row)  
in all chromosomes (eg not 9.)

Mutation



Selection

Pick 5 parents and take the best 2 to undergo cross-over.

Recombination probability = 100%

Mutation → swap (generates more diversity)

Mutation probability = 80%

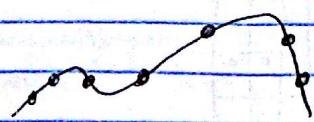
Parent selection → best 2 out of 5.

Survival selection → replace the worst.

Termination → after 100,000 iterations

Typical EAs behaviour

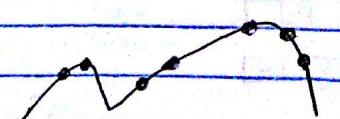
early



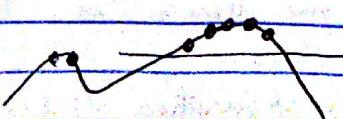
high diversity

(required to be high otherwise it is a killer, if you have few diversity)

mid



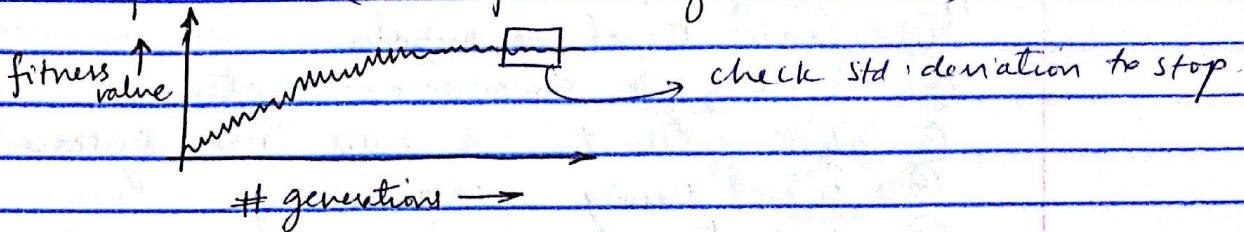
end



low diversity

(means converged to a solution)

Fitness function (the way it changes over time)

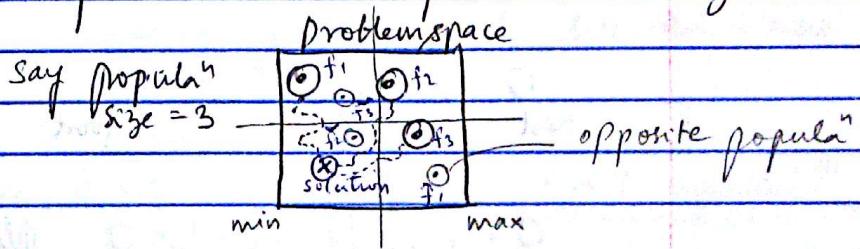


Is it worth it to spend time to improve initialization?  
→ It depends

Given a problem,  
we make  
random guesses

$r_1$  and  $r_2$   
to solve the  
problem. Compared

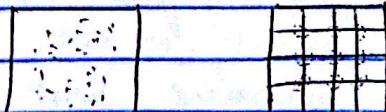
to when we get  $r_1$  and its opposite  $\tilde{r}_1$



$$\bar{x} = \min + \max - x$$

random trajectory  
towards solution

$$P(g(r_1, r_2) < g(r_1, \tilde{r}_1)) \approx 12.5\%$$



March 23<sup>rd</sup>, 2017

## Swarm Intelligence

Inspired by social behaviours of animals

Ant colonies, Bird flocking, fish schooling

### Ant colonies (AC)

• current population =  $10^{16}$  ants on this planet.

• been there since 100s of millions of years.

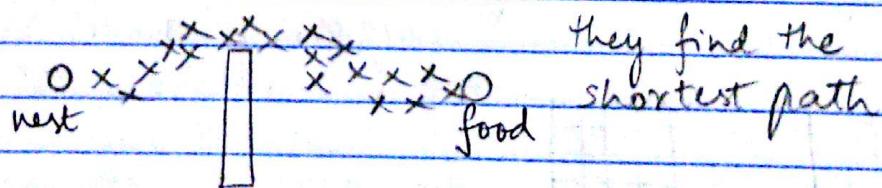
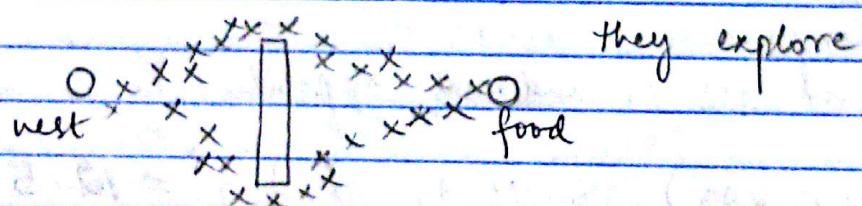
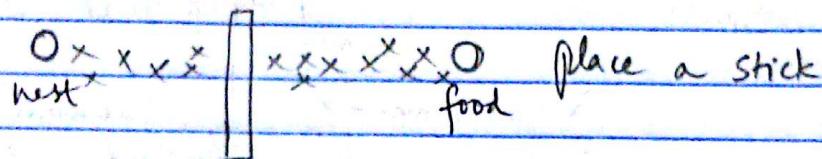
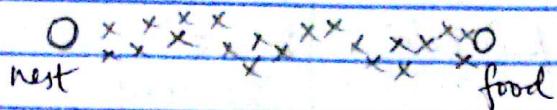
• one of the most successful species on the planet.

Concept of stigmergy: indirect communication used by social insects to co-ordinate activities.

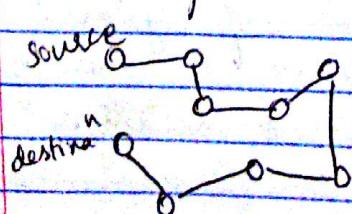
## Meta heuristic

### Applications of ACO

- ① Combinatorial optimization
- ② Routing for communication systems.
- ③ Task allocation in multi-robot systems.
- ④ Graph drawing and partitioning
- ⑤ Exploratory data analysis.



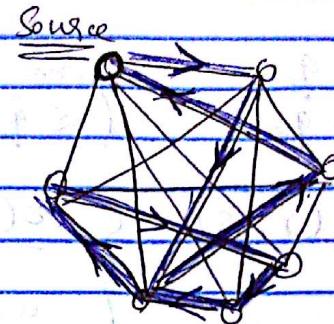
- Ants use pheromone trail
- Pheromone is a diffuse chemical substance
- Auto catalytic process.
- shortest path = higher pheromone concentration .
- Trail laying / trail - following behaviour.
- Artificial ants = piece of software.



(an agent moving from point to point on a TSP (Travelling Salesman Problem) graph.)

TSP:

find a path with  
minimum cost. Simple one!



Choose the cities that are  
connected by edges with  
a lot of pheromone trail.

Initially  $m$  ants are placed on randomly selected cities.

At each step, the ants move to new cities and  
modify the pheromone trail. → "local trail update".

When all ants have completed a tour (meaning cycle of  
actions → (source → destination → source)) the ant that  
made the shortest tour modifies the edges of its tour.  
"global trail update"

### Ideas from nature

- ① Preference for paths with high pheromone concentration
- ② Higher rate of growth of the amount of pheromone  
on shorter paths
- ③ The trail mediated communication among ants.

The ant  $k$  in city  $q$  chooses the city  $s$  not in  
the working memory  $M_k$  by applying a  
probabilistic decision.

$$s = \left\{ \begin{array}{l} \text{argmax } \{ \tau(r, u) * [\eta(r, u)^\beta] \} \text{ if } q \leq q_0 \\ s^* \text{ else} \end{array} \right.$$

$\tau(r, u)$  = pheromone on the edge  $(r, u)$

$\eta(r, u)$  = heuristic function (= inverse of distance)

$\beta$  = weighs the relative importance of pheromone trail & closeness.

$q = \text{random number in } [0, 1] \quad (\text{changes for a tour})$   
 $q_0 = (0 \leq q_0 \leq 1) \quad (\text{doesn't change for a tour})$

$$P_k(r, s) = \begin{cases} \frac{\tau(r, s) \eta(r, s)^{\beta}}{\sum_{u \in M_k} \tau(r, u) \eta(r, u)^{\beta}} & \text{if } s \in M_k \\ 0 & \text{else} \end{cases}$$

### Global trail update

$$\tau(r, s) \leftarrow (1-\alpha) \tau(r, s) + \Delta \tau(r, s)$$

$$\Delta \tau(r, s) = \frac{1}{\text{length of shortest tour.}}$$

Similar to Reinforcement Learning  $\Rightarrow$  better solutions  
get a higher reward (= pheromones in our case)

### Local Trail Update

$$\tau(r, s) \leftarrow (1-\alpha) \tau(r, s) + \alpha \tau_0$$

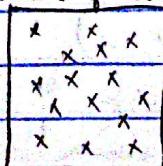
Motivated by trail evaporation.

Goal: avoid a very strong edge being chosen by all ants.

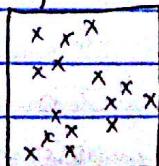
(Especially in the beginning, we can get into local maxima f get stuck here).

### Geometry Organization

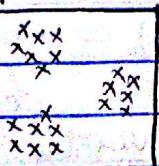
dead bodies of ants:



After 3 hours:



After 5 hours:



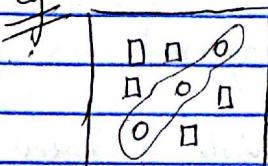
Exploratory State Analysis

March 28<sup>th</sup>, 2017 Boosting

Is it worth it to design a complex solution for a problem at hand?

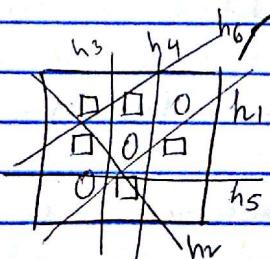
- \* Yes: we need complex structures to solve difficult problems.
- \* No: we can reach the same level of generalization with simple solutions.

Eg:



non linear, perhaps SVM / MLP / deep net.  
(if higher dimension)

Yes can be solved but requires lot of effort.



No

Instead of a powerful classifier, use several weak classifiers.

What is a weak classifier?

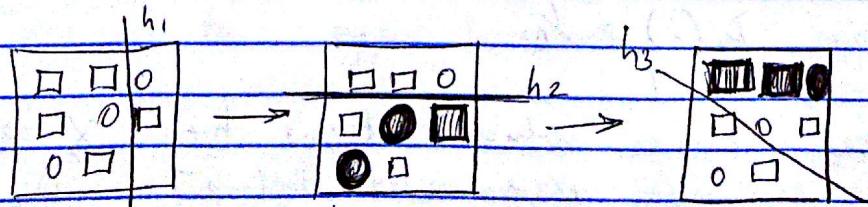
It can classify data

It has to be better than random ( $> 50\%$  accuracy)

ideally, the # of hypothesis should be  $\rightarrow \infty$

|H|  $\rightarrow \infty$

Boosting parts come from manipulation of weights.



change weights

of those that have been  
misclassified.

$$Y = \text{sign}(\sum x_i h_i)$$

For  $h_i \in H$ , we can approach high accuracy if

- ① weights of misclassified instances are manipulated.
- ② each classifier is better than random (accuracy  $> 50\%$ )
- ③ no outliers/noise  
(Boosting will collapse if we have outliers)

Most common out of the box classifier is boosting multiple small decision trees.

1990 - Boost by majority algorithm

1995 - AdaBoost (Freund and Shapire)

1997 - Generalized AdaBoost

2001 - AdaBoost for face detection (Viola & Jones approach)

2006 - RAMO Boost (Boosting for imbalanced data)

- AdaBoost is a linear classifier.

- AdaBoost can be used as feature selector.

- Sequential decision making.

- AdaBoost delivers a "strong" classifier as

$$f(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

weak classifier,  $y \in \{-1, +1\}$   
weight we try to induce

$$H(x) = \text{sign}(f(x))$$

Say,

Given  $(x_1, y_1), \dots, (x_m, y_m)$ ;  $x_i \in X$ ;  $y_i \in \{-1, +1\}$

Initialize  $D_1(i) = 1/m$

For  $t = 1, 2, \dots, T$

① Weak learner which returns  $h_t : X \rightarrow \{-1, +1\}$

with minimum error w.r.t distribution  $D_t$

② choose  $\alpha_t \in \mathbb{R}$

③ Update

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t Y_i h_t(x_i))}{Z_t}$$

$Z_t$  is a normalization factor so that  $D_t$  is

a distribution to normalize for each instance

from training data

from test data

from validation data

from output

from hypothesis

Strong classifier:

$$H(x) = \text{sign} \left( \sum_{t=1}^T \alpha_t h_t(x) \right)$$

Complexity  $\propto f(t)$

All information is stored in  $D_t$ .  
When we call "weak learner"  
given the distribution  $D_t$ ,  
it will return a weak  
classifier  $h_t : \mathcal{X} \rightarrow \{-1, +1\}$  from  $H = \{h(x)\}$   
(all decision trees we can  
generate)

Select a weak classifier

with the smallest weighted error.

$$h_t = \operatorname{argmin}_{h_j \in H} \underbrace{\varepsilon_j}_{\text{error}} = D_t(x_i) [y_i \neq h_j(x_i)].$$

find distribution value when

Condition:  $E_t < 0.5$

Most common decision tree builder, perceptron learning rule.

Main goal: minimize  $E_t = \frac{1}{m} \left| \{x_i : H(x_i) \neq y_i\} \right|$

(Count when  $x_i$  (how many times) has been misclassified)

This can be upper bounded by  $E_t(H) \leq \sum_{t=1}^T Z_t$

How to select  $x_t$ :

Select such that we minimize  $Z_t(\alpha)$  in each step.

If  $Z_t(\alpha)$  is a convex differentiable function with one extremum, then  $h_t(x) \in \{-1, +1\}$  then optimum ~~exists~~

$$\text{then optimal } \alpha_t = \frac{1}{2} \log \left( \frac{1+r_t}{1-r_t} \right)$$

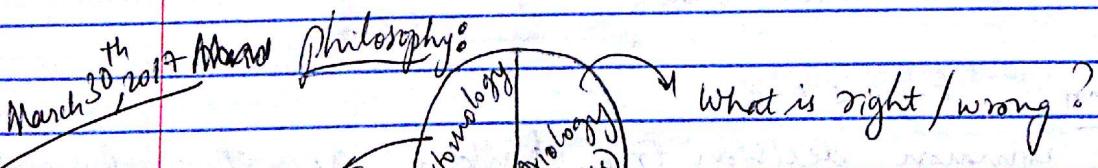
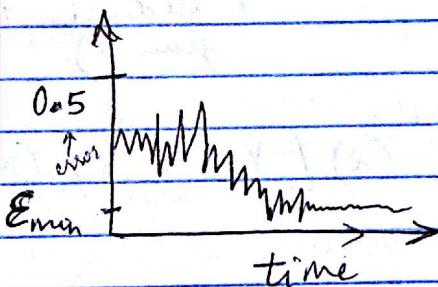
$$\text{where } Z_t = \sum_{i=1}^m D_t(i) h_t(x_i) y_i$$

If we have that, we can define:

$$Z_t = 2 \sqrt{\varepsilon_t (1-\varepsilon_t)} \leq 1 \text{ for optimal } \alpha_t$$

$$\Rightarrow D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

new value



"Does god exist"? religion... some driving force.

① Socrates : Inquiry, (find it yourself)

② Plato : Ideas are the source of everything  
: There is such a thing as "perfect state"  
when enlightened people rule.

③ Thomas Aquinas : Faith & reason will converge  
(must)

④ Thomas Hobbes: Social Atomism: Selfish interest drives everything.

To end the war of all against all, he proposes a "Social Contract".

⑤ John Locke (one of the initiators of Enlightenment).  
• He is strongly a protective of private property.  
• There are god given rights.

⑥ Immanuel Kant

- Happiness cannot be the motivation of our deeds.
- Ethics is simply there → There is god since we know ethics. We know naturally what is right/wrong.

⑦ Jeremy Bentham: All actions are motivated by pleasure.  
• Utilitarianism: whatever is useful, useful is right.  
Thinking is not the main measure for protection.  
• Law: If you suffer, you should have rights.

⑧ J. S. Mill "On Liberty".

- Harm principle. You are free to do whatever you like unless you are harming others or yourself.

