

Classify New Instance (x)

$$v_{NB} = \underset{v_j \in V}{\operatorname{argmax}} \hat{P}(v_j) \pi \hat{P}(a_i | v_j)$$

$|V|$ should be small.
(# of classification states)

Play Tennis

$$P(Y) P(\text{sun}/Y) P(\text{cool}/Y) P(\text{high}/Y) P(\text{strong}/Y) = 0.005$$

$$P(n) P(\text{sun}/n) P(\text{cool}/n) P(\text{high}/n) P(\text{strong}/n) = 0.021$$

$$\hookrightarrow V_{AB} = n.$$

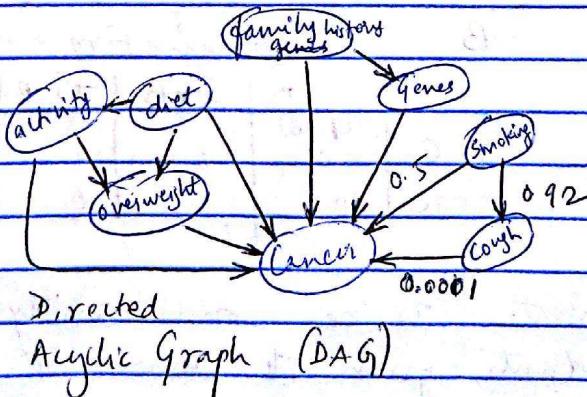
even if one probability is 0, then the whole probability = 0.

so use lots of data so that no probability is 0.

If it is 0, then we use an m -estimator to find some value
 $m \rightarrow$ estimator $P(a_i | v_j) = \frac{n_c + mp}{n+m}$ to estimate a non zero
value when you encounter it in N.B.

March 14th, 2017

= 0 ?



Belief network
how to infer?

Bayesian Belief Network

- Naive Bayes is too restrictive
- The problem may prove to be intractable without such assumptions
- Bayesian Belief Networks describe conditional independencies among subsets of variables.
- allows combining the prior knowledge about (in)dependencies

among variables observed in training data.

"Bayesian Belief Networks" = "Bayes Nets"

Conditional Independence

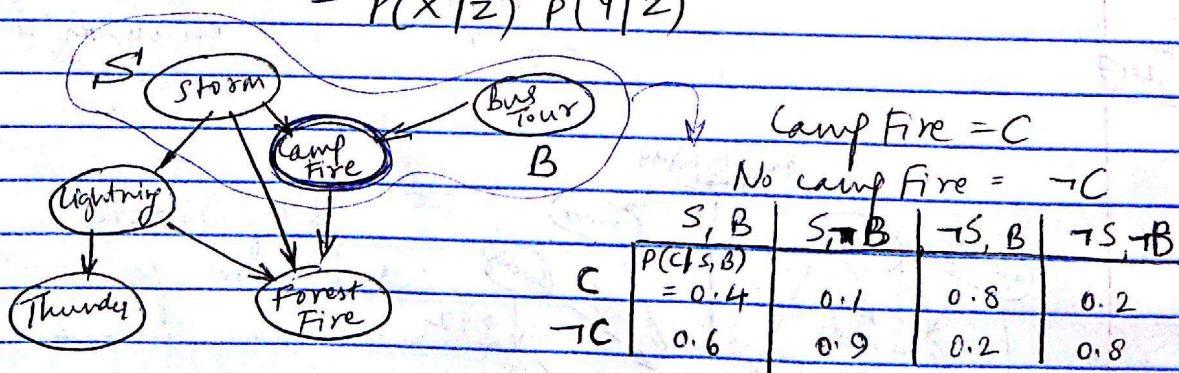
X is conditionally independent of Y given Z if the probability distribution governing X is independent of the value of Y given the value of Z ; i.e., if

$$P(X=x_i | Y=y_j | Z=z_k) = P(X=x_i | Z=z_k)$$

$$(or) P(X|Y, Z) = P(X|Z)$$

e.g. $P(\text{Thunder} | \text{Rain, Lightning}) = P(\text{Thunder} | \text{Lightning})$ \times is not dependent on Y

$$\begin{aligned} P(X, Y|Z) &= P(X|Y, Z) P(Y|Z) \\ &= P(X|Z) P(Y|Z) \end{aligned}$$



Each node is asserted to be conditionally independent of its non-descendants, given its immediate predecessors.

In General,

$$P(Y_1, Y_2, \dots, Y_n) = \prod_{i=1}^n P(Y_i | \text{Parents}(Y_i))$$

where $\text{Parents}(Y_i)$ denote immediate predecessors of Y_i in the graph.

(Engagement)

The joint distribution is fully defined by the graph plus the probability of: $(y_i \mid \text{parents}(y_i))$

Inference in Belief Networks

How can we infer the values (probabilities) of one or more network variables given observed values of others?

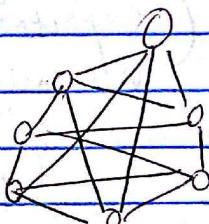
The Bayes Net contains all information needed for inference. If only one variable with unknown value, that would be easy to infer.

Generally, this is NP-hard. (= intractable)

NP-hard: e.g. Travelling

Salesman

Problem



minimizing the gas mileage,
how can I travel to all
cities?

Learning in Belief Nets

- Structure may be known/unknown.
- Training examples might provide values of all variables or just some.

If structure is known and we observe all variables, then it's easy. Just use Naive Bayes classifier!!!

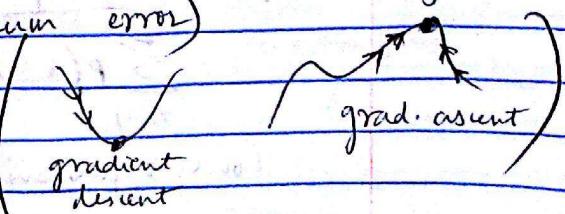
Learning: Gradient Ascent for Bayes Nets

(we are looking to maximize something, not to get the minimum error)

Let W_{ijk} denote one entry in the

cond. probability table for y_i in the network.

$W_{ijk} = P(y_i = y_{ij} \mid \text{parents}(y_i) = \text{the list } u_{ik} \text{ of values})$



E.g. if $y_i = \text{Campfire}$, then w_{ik} might be $\langle \text{Storm} = \text{True}, \text{BusTour} = \text{False} \rangle$

Perform gradient ascent repeatedly.

① update all weights w_{ijk} using training data D .

$$w_{ijk} \leftarrow w_{ijk} + \eta \sum_{d \in D} \left[\frac{p_h(y_{ij}, u_{ik} | d)}{w_{ijk}} \right]$$

② Renormalize the w_{ijk} to assume

$$\sum_j w_{ijk} = 1 \quad 0 \leq w_{ijk} \leq 1$$

One may also use EM algorithm for learning

(Expectation

Maximization)

① calculate probabilities of unobserved variables assuming h .

② calculate new w_{ijk} to maximize $E[\ln p(D|h)]$ where D now includes both observed & unobserved (calculated) variables.

(hypothesis, like "Yes there is cancer" or "no cancer")

EM algorithm:

Pick random initial $h = \langle \mu_1, \mu_2 \rangle$

(mixture of Gaussians)

E-step : $E[z_{ij}]$ (estimate the unobserved value based on n distributions)

$$E[z_{ij}] = \frac{P(x=x_i | \mu=\mu_j)}{\sum_{n=1}^2 P(x=x_i | \mu=\mu_n)}$$

$$\sum_{n=1}^2 P(x=x_i | \mu=\mu_n)$$

(only 2 gaussians to build the mixture)

M-step calculate a new max likelihood hypothesis

$$h' = \langle \mu'_1, \mu'_2 \rangle, \text{ assuming the value taken on}$$

by each hidden variable z_{ij} is its expected value $E[z_{ij}]$.

Replace $h = \langle \mu_1, \mu_2 \rangle$ by $h' = \langle \mu'_1, \mu'_2 \rangle$

Then we apply the adjustment:

$$\mu_j' \leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]}$$

March 16th, 2017 Evolutionary Algorithm

Not considered as "learning". Mainly "optimization"

Inspired by natural evolution

↳ Survival of the fittest (Natural selection)

- EA's are function optimizers

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

- GA's are EA's 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", "epoch", "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

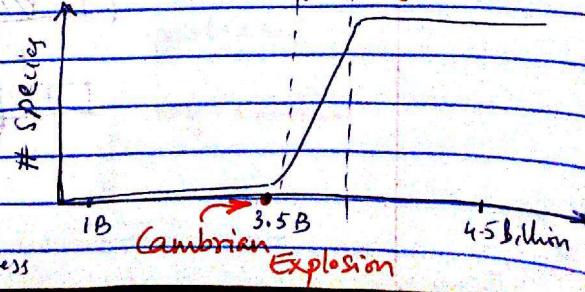
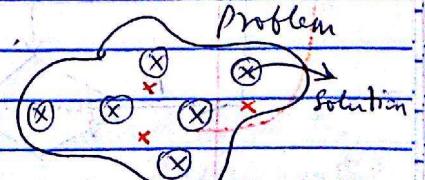
④ Crossover (getting off springs)

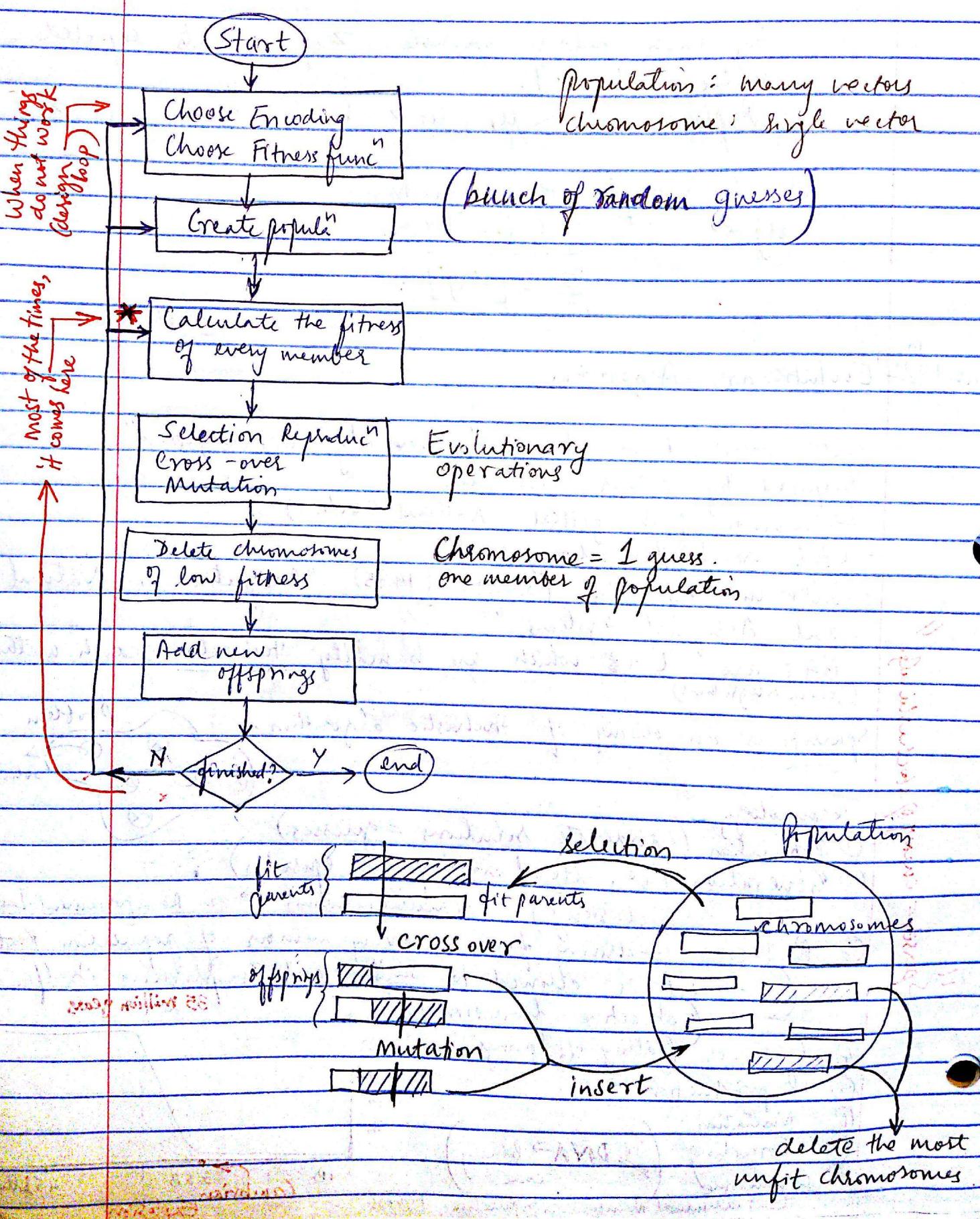
⑤ Reproduction

⑥ Mutation

⑦ Encoding (e.g. DNA → Binary code)

⑧ Chromosome ↗ 1 individual member in population. 1 guess





Function to be optimized

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

\downarrow \downarrow \downarrow

1	1	1	0	1
0	0	0	1	1

$\underbrace{\quad\quad\quad}_{P_1}$ $\underbrace{\quad\quad\quad}_{P_2}$ $\underbrace{\quad\quad\quad}_{P_3}$

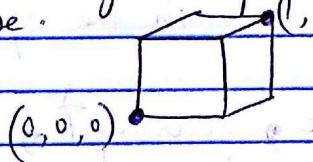
chromosome

Coding of the GA's secret

(Stochastic optimization)

we look at bit strings of size $L_E = \underbrace{[0111111111]}_{Z}$
size of search space = 2^L

GAs basically sample the corners of this L dimensional hypercube.



(for neural network, you go inside, the continuum of all answers) \therefore more accurate answers using NN,

analogy is only for binary encoding.

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

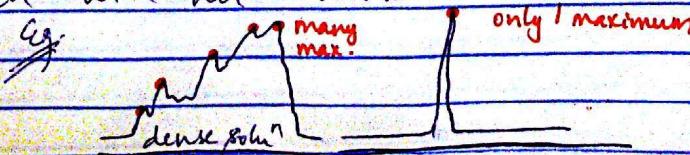
Most (test functions =) problems can be encoded with max 30 bits. Any thing less can be enumerated

Eg. chess game \rightarrow branching factor = 16

\rightarrow a game is made of 100 moves (on average)

$$\therefore 16^{100} \text{ possibilities. } = (2^4)^{100} = 2^{400}$$

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



Encoding

① Binary coding

1	1	0	1	1	0
---	---	---	---	---	---

② Permutation coding

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

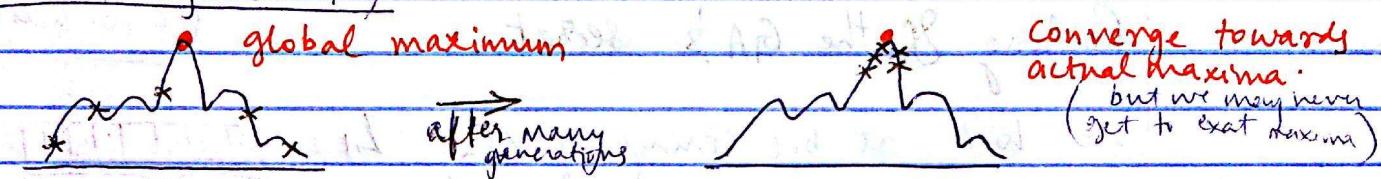
③ Value coding.

1.234	9.687	0.1501
-------	-------	--------

A	C	W	Y	B	X	N	I	O	J	D	I
---	---	---	---	---	---	---	---	---	---	---	---

back back left

Diversity & GAs



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

→ Size of population?

- large, accuracy \uparrow , becomes sluggish
- Small, less accurate, but fast

→ Cross over probability = 100%

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