

Bayes Theorem

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

\nearrow Prior probability of hypothesis
 \searrow Evidence

Find best hypothesis given data

$D \rightarrow$ training data
 $h \rightarrow$ function mappings for target

From candidate hypotheses, iterate all of them with same data & find maximum a posteriori (MAP) hypothesis.

$$h_{MAP} = \operatorname{argmax} P(h|D)$$

\nearrow most probable hypothesis = highest a posteriori

$$= \operatorname{argmax} \frac{P(D|h)P(h)}{P(D)} = \operatorname{argmax} P(D|h)P(h)$$

\searrow This doesn't change from one hypothesis to another

\searrow we are left with this!

★ If we don't have any prior or for some cases all $P(h)$ might be equal.

$$h_{ML} = \operatorname{argmax} P(D|h)$$

\searrow max likelihood hypothesis

$$P(h) = \frac{1}{|H|}$$

ex 11

	$P(\text{cancer}) = 0.008$	$P(+ \text{cancer}) = 0.98$	
priors \leftarrow	$P(\neg\text{cancer}) = 0.992$ (dun)	$P(- \text{cancer}) = 0.02$	\rightarrow likelihoods
	$P(+ \neg\text{cancer}) = 0.03$	$P(- \neg\text{cancer}) = 0.97$	\nearrow

Calculating MAPs:

You had test result (+), do you have cancer?

$$P(\text{cancer} | +) = P(+|\text{cancer})P(\text{cancer}) = 0.0078$$

Probability of not having cancer given your test is (+):

$$P(\neg\text{cancer} | +) = P(+|\neg\text{cancer})P(\neg\text{cancer}) = 0.0298$$

If A is $n \times n$ is there a nonzero vector x such that Ax is a scalar multiple of x ?

$Ax = \lambda x$

Bayesian Inference

Likelihood: Probability of observing the data that has been observed assuming that the data came from a specific scenario.

ex II $P(\text{test} = \text{positive} | \text{disease})$

Prior: Previous knowledge on the subject.

ex II Disease is very rare.

Posterior Probability: What is the probability of you having the disease given that the test is positive.

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

$$P(\text{disease} | \text{test} = \text{positive}) = \frac{P(\text{test} = \text{positive} | \text{disease}) \times P(\text{disease})}{P(\text{test} = \text{positive})}$$

$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \rightsquigarrow P(A \& B)$

Labels: Posterior (points to the result), Likelihood (points to $P(\text{test} = \text{positive} | \text{disease})$), Prior (points to $P(\text{disease})$), Evidence (points to $P(\text{test} = \text{positive})$)

ex II $P(\text{disease}) = 0.00148$

$$P(\text{test} = (+) | \text{disease}) = 0.93$$

$$P(\text{test} = (-) | \text{disease}) = 0.99$$

$$P(\text{disease} | \text{test} = (+)) = \frac{P(\text{Person has disease} \& \text{test} (+))}{P(\text{test} (+))}$$

$$P(\text{disease} \& \text{test} (+)) = \underbrace{P(\text{disease})}_{0.00148} \times \underbrace{P(\text{test} = (+) | P(\text{disease}))}_{0.93}$$

$$P(A \& B) = P(A|B)P(B) = \text{Likelihood} \times \text{Prior}$$

$$P(\text{disease} | \text{test} (+)) = \frac{0.0013}{0.011} \approx 0.12 \rightarrow \text{Posterior}$$

$\frac{0.0013}{0.011}$

 Evidence (points to 0.011) $\rightarrow P(\text{test} = \text{positive})$

$$= P(\text{disease} | (+)) + P(\text{no disease} | (+))$$

Bayes Updating

What if you take a 2nd test if first one is (+)?

Assume: Correctness of 2nd is not influenced by first.

$P(\text{disease}) = 0.12 \rightarrow$ having disease given that first test = (+)

$$P(\text{disease} | 2^{\text{nd}} \text{ test} = (+)) = \frac{P(\text{disease}) P(2^{\text{nd}} = (+) | P(\text{disease}))}{P(2^{\text{nd}} = (+))}$$

$$= 0.93$$

\rightarrow 2 tests makes it much more probable than one (+) test

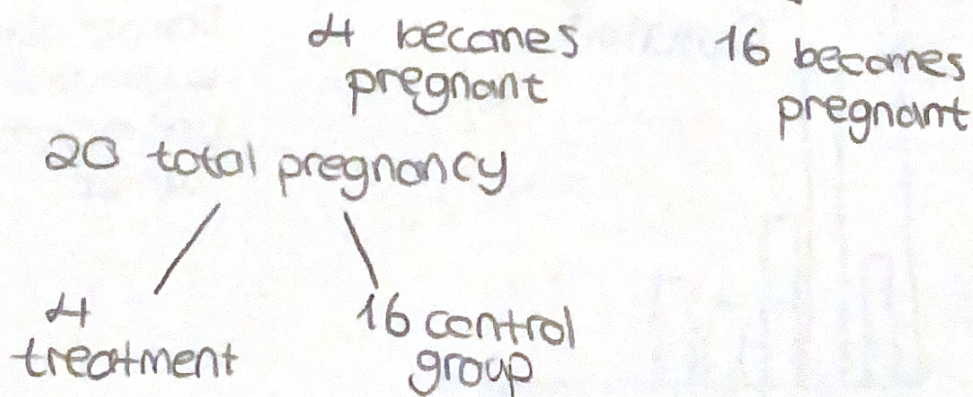
Bayes updating \rightarrow updating a probability based on an event affecting it.

Update prior to get posterior

prior $\xrightarrow[\text{Updating}]{\text{Bayesian}}$ Posterior

Bayesian Inference

ex // 40 women, 20 asking protection, 20 doesn't.

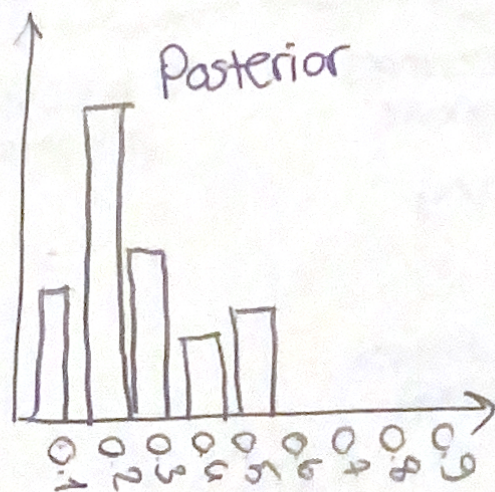
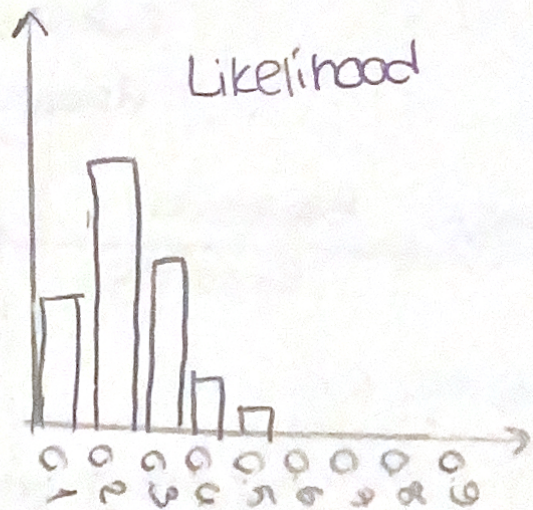
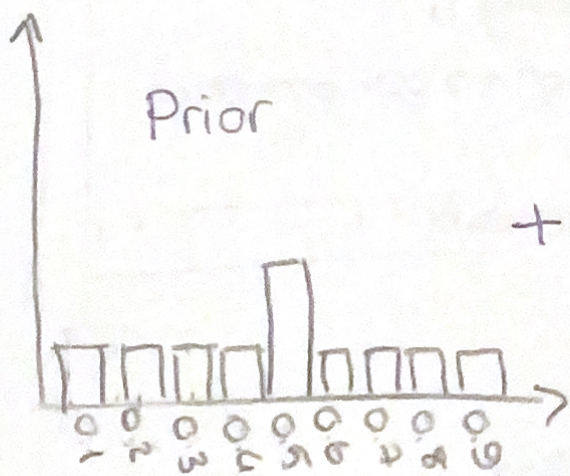


$p \rightarrow$ pregnancy comes from treatment group 0.5
 (assume it can take 0.2, 0.5, 0.8)

p	0.2	0.5	0.8
prior $P(\text{model})$	0.06	0.52	0.06
likelihood $P(\text{data} \text{model})$	0.21	0.004	0.003
likelihood x prior	0.03	0.002	0
posterior $P(\text{model} \text{data})$	0.42	0.98	0

$$P(\text{data}|\text{model}) = \underbrace{P(k=4 | n=20, p)}_{P(n, p)} \text{ Binomial}$$

Choose the model with highest posterior



Sample size büyüdükçe likelihood'un posterior'a etkisi artar.

Probability formulas

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A) \rightarrow \text{Product Rule}$$

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} \rightarrow \text{Posterior probability}$$

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i) \rightarrow \text{Law of Total Probability}$$

Brute Force Max a Posteriori + Concept Learning

Assumptions:

1. Training data is noise free. \rightarrow hyp. space
2. Target concept c is in the H .
3. Every hypothesis has equal priors

Then: $P(h) = \frac{1}{|H|}$ $P(D|h)$: probability of data D given hypothesis h is 1 if data is consistent with h , 0 otherwise.

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)} \rightarrow P(h|D) = \frac{1 \cdot \frac{1}{|H|}}{P(D)} = \frac{1 \cdot \frac{1}{|H|}}{\frac{|V_{S_{H,D}}|}{|H|}}$$

$$P(D) = \frac{|V_{S_{H,D}}|}{|H|} = \frac{1}{|V_{S_{H,D}}|}$$

subset of hypotheses from H consistent with data.

$$P(h|D) = \begin{cases} \frac{1}{|V_{S_{H,D}}|} & \text{if } h \text{ is consistent with } D \\ 0 & \text{otherwise} \end{cases}$$

Initially all $h = \frac{1}{|H|}$ as we add data posterior of inconsistent hypotheses become zero while consistent ones have a probability of $\frac{1}{\# \text{consistent ones}}$

Bayes Optimal Classifier

What if we have following hypotheses classifying a new instance as follows:

$h_1 \rightarrow$ classifies x as (+) with 0.4 \rightarrow MAP!

$h_2 \rightarrow$ (-) with 0.3 \rightarrow summing up (-) seems

$h_3 \rightarrow$ (-) with 0.3 \rightarrow more probable

* If new example can only take on any value v_j from V , then $P(v_j | D)$ - correct classification - is

$$P(v_j | D) = \sum P(v_j | h_i) P(h_i | D)$$

Look for the maximum

\hookrightarrow how hypothesis

\hookrightarrow prob of hypothesis given data

classifies instance in each label & their probabilities

$$\operatorname{argmax} P(v_j | D)$$

$$= \operatorname{argmax} \sum P(v_j | h_i) P(h_i | D)$$

\hookrightarrow Bayes optimal classification

ex 11

$$\begin{aligned} P(h_1 | D) &= 0.4 & P(- | h_1) &= 0 & P(+ | h_1) &= 1 \\ P(h_2 | D) &= 0.3 & P(- | h_2) &= 1 & P(+ | h_2) &= 0 \\ P(h_3 | D) &= 0.3 & P(- | h_3) &= 1 & P(+ | h_3) &= 0 \end{aligned}$$

$$\sum_h P(+ | h_i) P(h_i | D) = 0.4$$

$$\sum_h P(- | h_i) P(h_i | D) = 0.6 \rightarrow \text{classified as (-)!}$$

* If we're learning boolean concepts we can take weighted vote among all hypotheses in version space with weights = posterior probabilities.

Naive Bayes Classifier

NB classifier applies to learning tasks where each instance x is described by a conjunction of attribute values \mathcal{A} where target function $f(x)$ can only take any value from finite set V .

$$u_{MAP} = \operatorname{argmax}_{v_j} P(v_j | a_1, a_2, \dots, a_n)$$

\downarrow class label \rightarrow attributes (data)

rewrite as

$$u_{MAP} = \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, a_3, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

\rightarrow this is same across classes

$$= \operatorname{argmax}_V P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

- $P(v_j)$: class proportions in training data
- \rightarrow prob of attribute values given class \times Prob of class

- Assuming all attributes are independent of each other:

$$P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$$

plugging these in

$$u_{NB} = \operatorname{argmax}_{v_j} P(v_j) P(a_i | v_j) \quad \left. \vphantom{u_{NB}} \right\} \text{Naive Bayes Classifier}$$

\rightarrow target value output by NB classifier

ex// Instance to predict:

	Outlook	Temp	Humid	Wind
Target = Play Tennis (yes/no)	Sunny	cool	High	Strong

- Estimate $P(v_j)$ from class proportions

$$P(\text{Play Tennis} = \text{Yes}) = 0.64$$

$$P(\text{Play Tennis} = \text{No}) = 0.36$$

- Estimate conditional probabilities.

$$P(\text{Wind} = \text{Strong} | \text{Play Tennis} = \text{yes}) = 0.33$$

$$P(W = \text{strong} | PT = \text{no}) = 0.6$$

$$P(\text{yes} | \text{inst}) = P(\text{yes}) \times P(\text{sunny} | \text{yes}) \times P(\text{cool} | \text{yes}) \times P(\text{high} | \text{yes}) \times P(\text{str} | \text{yes}) = 0.0206$$

Calculate same for other labels:

$$P(\text{no} | \text{inst}) = P(\text{no}) P(\text{sumy} | \text{no}) P(\text{cool} | \text{no}) P(\text{high} | \text{no}) P(\text{strong} | \text{no})$$
$$P(\text{w} | \text{j}) \prod_i P(\text{a}_i | \text{v}_j)$$

= 0.0206 \rightarrow Our prediction is larger than "yes"

Normalizing these:

$$\frac{0.0206}{0.0206 + 0.0053} = \underline{\underline{0.795}}$$

Bayesian Belief Networks

"Conditional independence is overrated"

Naive Bayes → All of the variables are conditionally independent

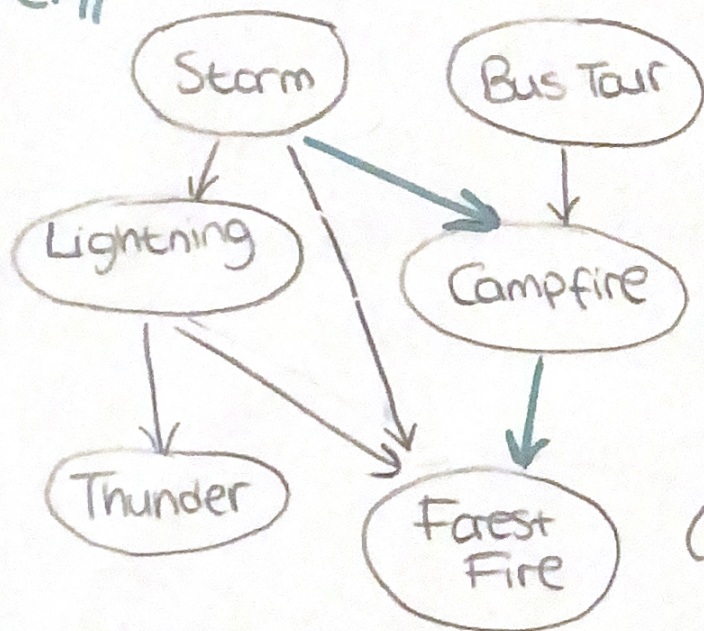
Belief Nets → Some of the variables are conditionally independent

ex // $P(x=x_i | y=y_j, z=z_k) = P(x=x_i | z=z_k)$ if x is independent of y then ignore y !

x is conditionally independent of y given a value of z

What NB does → $P(A_1, A_2 | V) = P(A_1 | A_2, V) P(A_2 | V)$
 features target = $P(A_1 | V) P(A_2 | V)$

ex //



if all the features are independent, simply multiply all the likelihoods

What BBN does:

Features are conditionally independent of their non-descendants.

(x is descendant of y if there's directed path from x to y)

$$P(y_1, \dots, y_n) = \prod P(y_i | \text{Parents}(y_i))$$

→ should be immediate parent $P \rightarrow C$

	S, B	S, !B	!S, B	!S, !B
C	0.4	0.1	0.8	0.2
!C	0.6	0.9	0.2	0.8

$$P(C = \text{True} | S = \text{True}, B = \text{True}) = 0.4$$

$$P(T = \text{True} | L = \text{True}, C = \text{True}) = P(T | L)$$

$$P(L = \text{True} | S = \text{True}, B = \text{True}) = P(L | S)$$

→ Look at immediate parents only!

ex // Thunder is conditionally independent of other variables, given Lightning is True. (we calculate Lightning by looking at its parents and go backwards in the graph for all parents)

Bayesian Belief Nets

Conditional independence assumptions only apply to a subset of variables.

ex 11 $P(X=x_i | Y=y_j, Z=z_k) = P(X=x_i | Z=z_k)$

↳ we can get rid of this

- Assumption: X is conditionally independent of Y given Z if distribution governing X is independent of value of Y given a value for Z, if above equation holds.

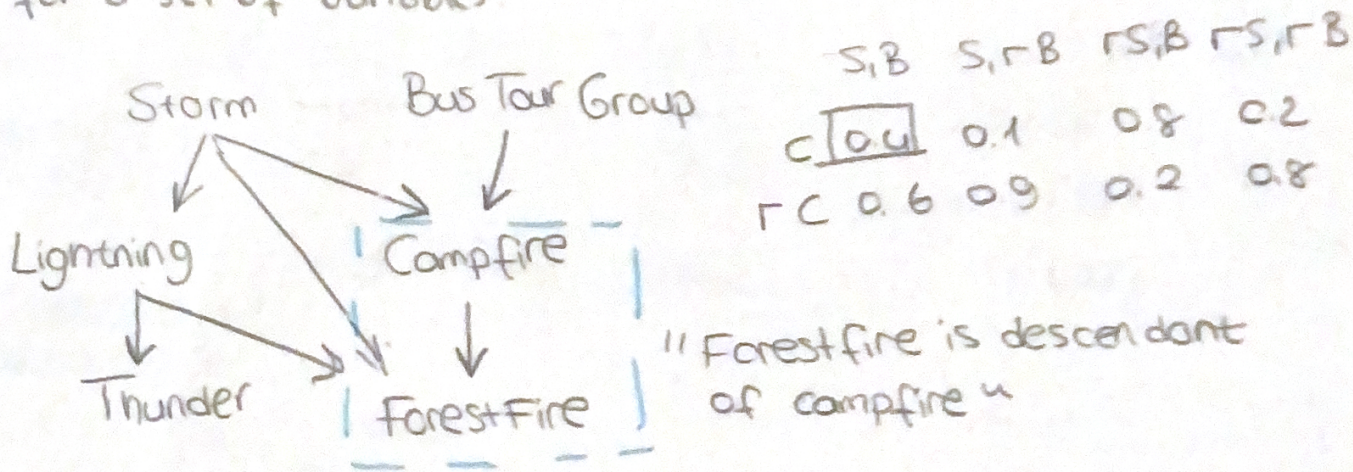
$$P(X|Y,Z) = P(X|Z)$$

↳ X ve Y independent'sa y den kurtulabilirsin

- Naive Bayes'teki independence assumption: sayesinde:

$$P(A_1, A_2 | V) = P(A_1 | A_2, V) P(A_2 | V) \\ = P(A_1 | V) P(A_2 | V)$$

A Belief Network represents joint probability distribution for a set of variables.



$$P(y_1, \dots, y_n) = \prod_{i=1}^n P(y_i | \text{Parents}(y_i))$$

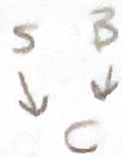
↳ predecessors of y_i
graph'ta y den önceki olaylar

Campfire is independent of it's nandescendants Lightning & Thunder.

$$P(\text{Campfire} = \text{True} | \text{Storm} = \text{True}, \text{Bus Tour G} = \text{True}) = \underline{\underline{0.4}}$$

In a Bayesian Net, we can number each node E_i parents will have smaller indexes than their children.

$$S \rightarrow L \rightarrow T$$



- $P(\text{Campfire} | S, \neg B) = 0.1$
- $P(\text{Thunder, Lightning, Storm}) = P(T|L) P(L|S) P(S)$
- $P(\text{Campfire}) = P(C | S, B) P(S) P(B) + P(C | S, \neg B) P(S) P(\neg B) + P(C | \neg S, B) P(\neg S) P(B) + P(C | \neg S, \neg B) P(\neg S) P(\neg B)$

- $P(B) = \sum_{i=1}^n P(B|A_i)$ where A_i add up to 1

- $P(C|S) = ?$ (we don't know about bus tour group?)

$$P(C|S) = P(C|S, B) P(B) + P(C|S, \neg B) P(\neg B) \rightarrow \begin{array}{l} \text{Storm} = \text{True} \\ \text{iterate over} \\ \text{Bus tour group} \end{array}$$

- If most of the variables are unknown the complexity becomes exponential (2^n)

- We can use Monte Carlo simulation for this.

ex 11	$P(S)$	$P(B)$	$P(C S, B)$	} Generate cases
$P(S) = 0.9$	$S \rightarrow T$	$B \rightarrow T$	$C \rightarrow F$	
$P(B) = 0.5$	$S \rightarrow T$	$B \rightarrow F$	$B \rightarrow F$	
$P(C B, S) = 0.4$	---	---	---	
$P(C) = ?$	$S \rightarrow F$	$B \rightarrow T$	$C \rightarrow T$	

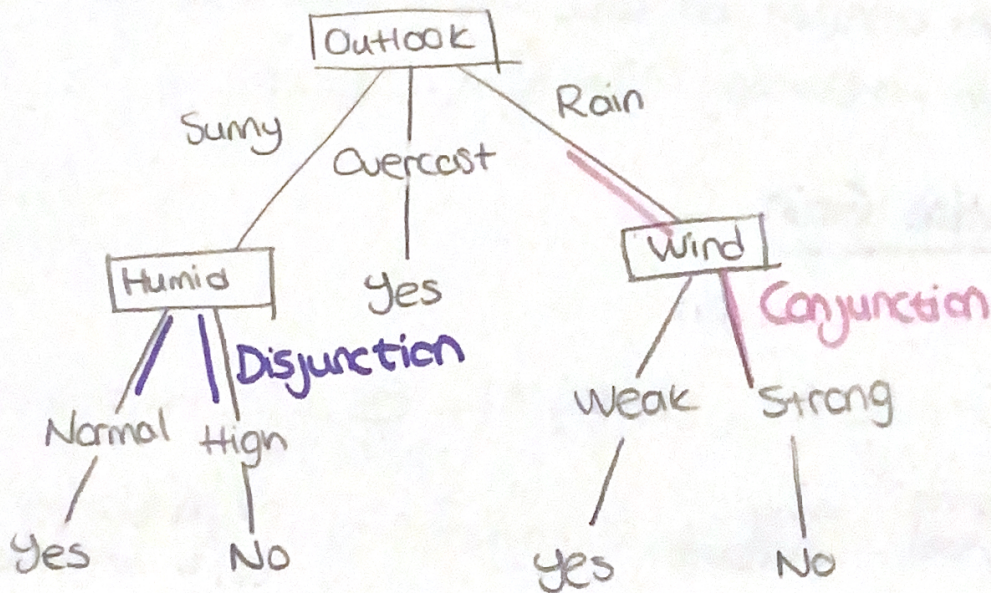
↘ consider these cases

ratio obtained will be approximation to the campfire

Decision Tree Learning

Represents disjunction of conjunction of constraints on values attributes can get.

$(\text{Outlook} = \text{Sunny} \wedge \text{Humidity} = \text{Normal}) \vee (\text{Outlook} = \text{Overcast})$
 $\vee (\text{Outlook} = \text{Rain} \wedge \text{Wind} = \text{Weak})$



ID3

Algorithm

1. Create a root node for the tree
2. If all examples are (+) return single node tree $\boxed{+}$
3. Do the same if all examples are (-)
4. If attributes are empty return tree root with most common value of target attribute.

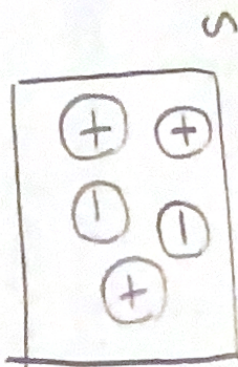
Otherwise begin

1. Find attribute with best information gain = A
 2. Put that attribute to root node (A)
 3. For each possible value of A (v_i) do:
 1. Add a new tree branch below root, $A = v_i$
 2. If examples of $A = v_i$ empty, put the most common target label
 3. If examples of $A = v_i$ exist, add the subtree with ID3 (Examples v_i , Target attribute, Attributes - A)
- ↪ recursion

↪ Datasets entropy - Entropy of attribute

Attribute Selection

Entropy



$$\text{Entropy}(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

0.6 0.4

→ $-0.6 \log_2 0.6 - 0.4 \log_2 0.4$

Target attribute can take values more than 2!

Then → $\text{Entropy}(S) = \sum_{i=1}^n -P_i \log_2 P_i$

Information Gain

For entropy reduction

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v_i} \frac{|S_{v_i}|}{|S|} \text{Entropy}(S_{v_i})$$

anceki
feature
(parent)

aday
attribute

→ parent 40
0 degeri
alan example sayisi
parent site

Choosing attributes

Plurality value \rightarrow return the most frequent value.

Entropy of dataset $\rightarrow H(\text{Goal}) = B\left(\frac{P}{P+N}\right)$

feature $\rightarrow \sum_{e} P(C_{ve}) \log_2 \frac{1}{P_{ve}}$

$$= - \sum_{e} P(C_{ve}) \log_2 P(C_{ve})$$

Information

gain $\rightarrow G(Y, X) = E(Y) - E(Y|X)$

\rightarrow target given X
 \rightarrow test for X

$$\text{Remainder}(A) = \sum_{e=1}^D \frac{p_{e+n}}{p+n}$$

\rightarrow proportion

$$B\left(\frac{P_e}{P_e+n}\right) \rightarrow \text{entropy}$$

\hookrightarrow A'nin aldig' aegerer \rightarrow ver value iain
he saplayip tcapa

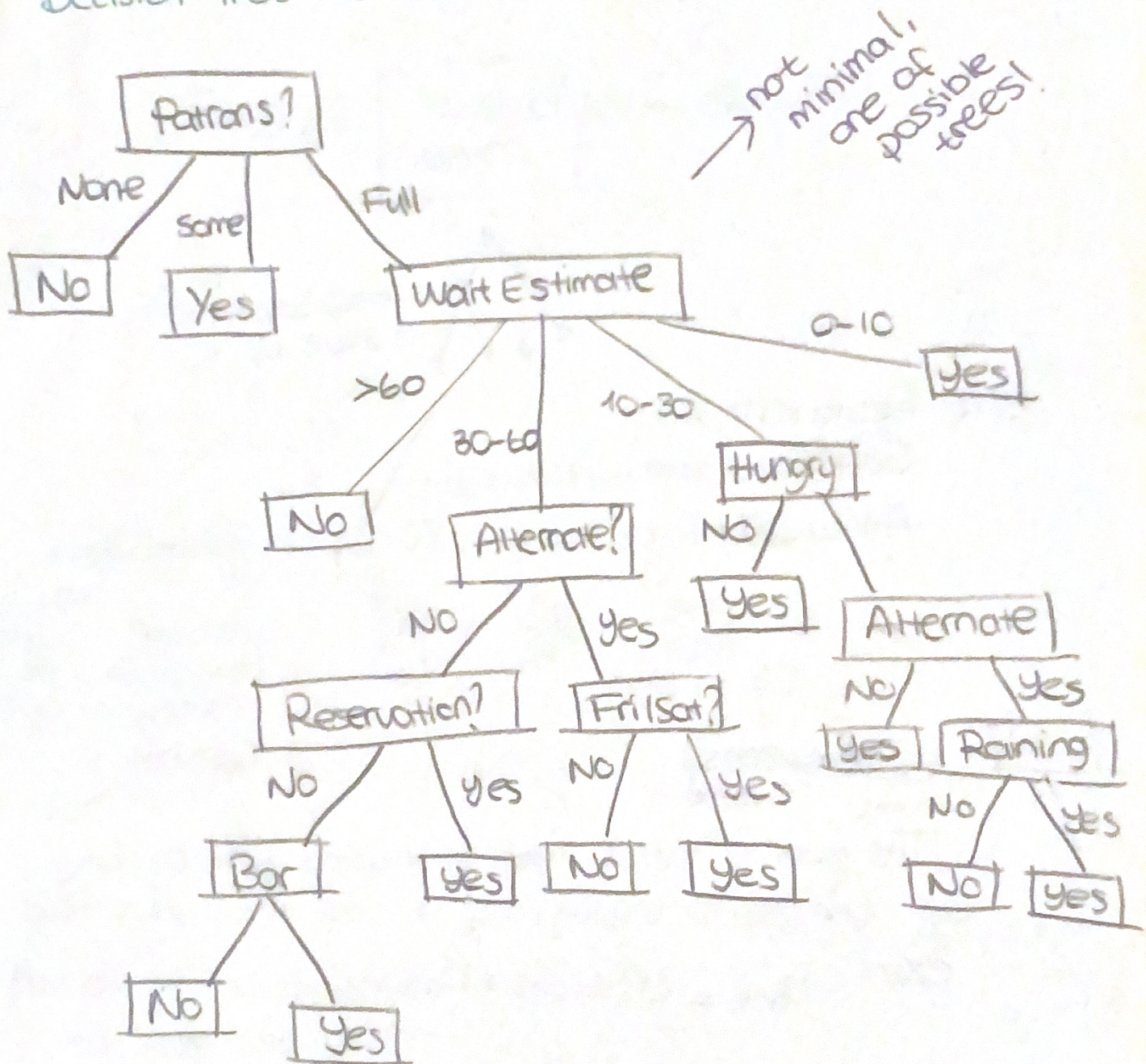
Root node 'da \rightarrow dataset 'in

entropy 'in - feature 'in

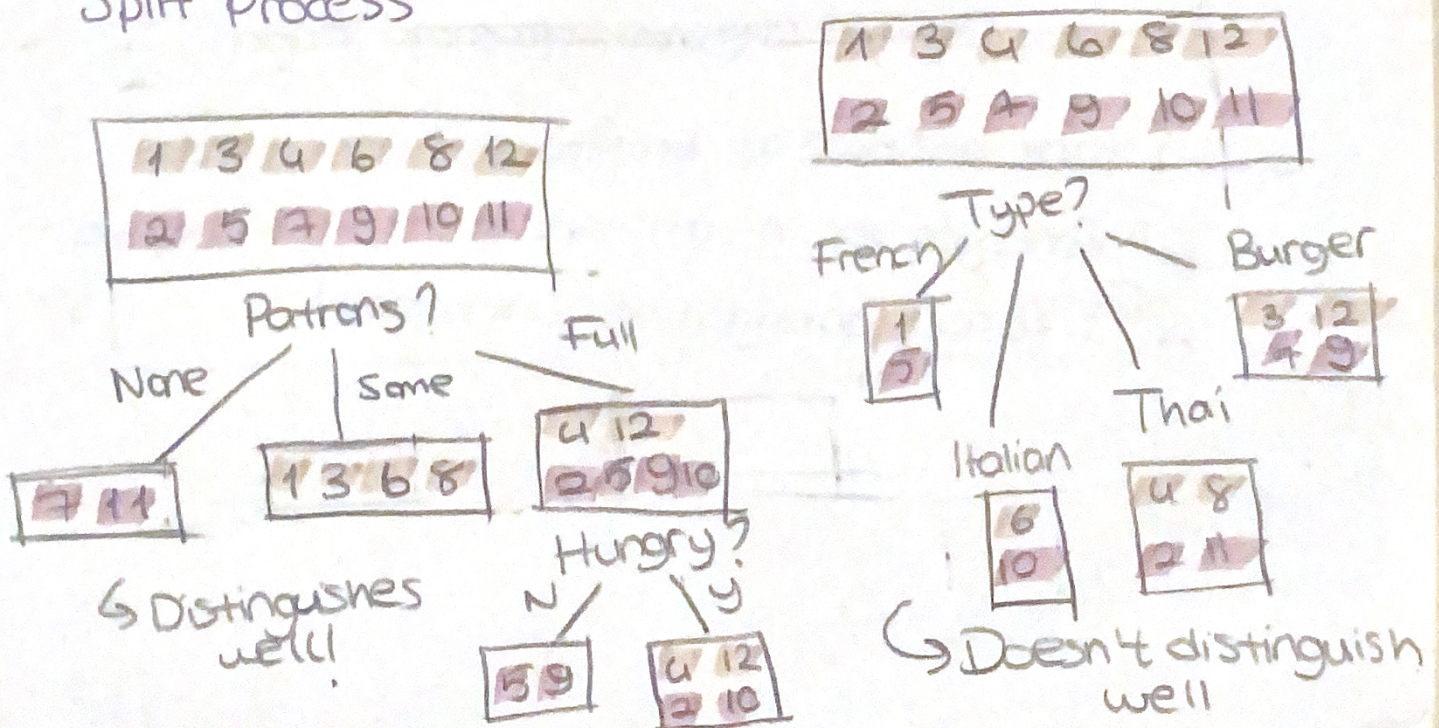
IG = Parent - Feature 'in

entropy entropy

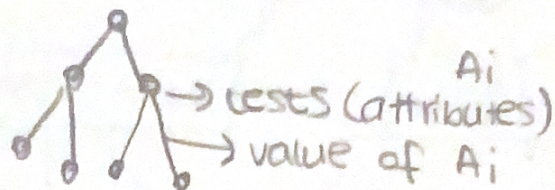
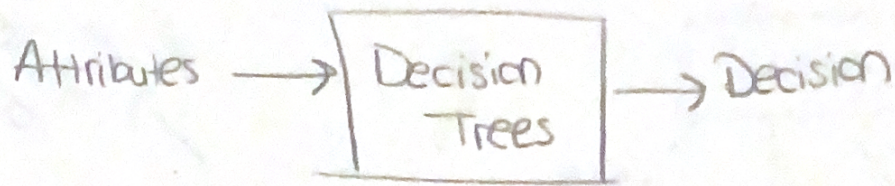
Decision Tree From Examples



Split Process



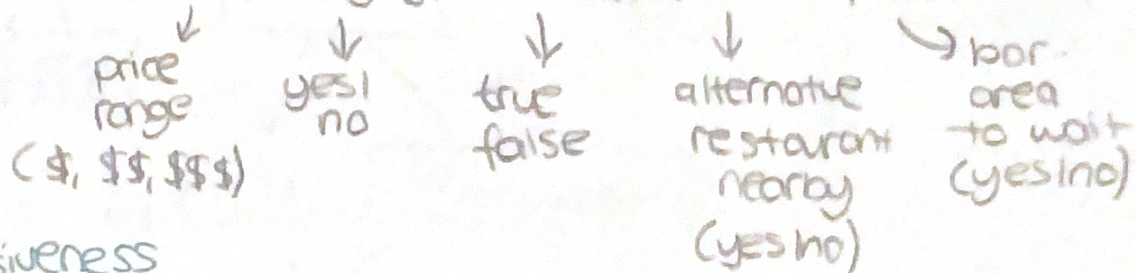
Decision Trees



ex I Restaurant Wait

Goal Predicate \rightarrow will wait

Attributes \rightarrow Price, Hungry, Fri/Sat, Alternative, Bar..



Expressiveness

The goal is true if input attributes satisfy one of the paths leading to a leaf with value true.

ex II Path = (Patrons = Full \wedge WaitEstimate = 0-10)

\rightarrow bu true leaf'ine giden pathlerden

Goal \leftrightarrow (Path1 \vee Path2 ..) biri

n tane attribute'la boolean function setinde

2^n fonksiyon var. n attribute'la truth table'da

2^{2^n} farklı fonksiyon var. (intractable)

Concept Learning

Problem of searching through a predefined space of potential hypotheses for the hypothesis that fits the data best.

- Each constraint can be
 - a value (water = warm)
 - don't care (water = "?")
 - no value allowed (water = "∅")

? : Feature has no predictive power. (almost redundant)

★ We're looking for a conjunction of constraints that describe (+) examples the best.

★ most general hypothesis → $\langle ?, ?, ?, ?, ?, ? \rangle$

Every day (instance)
is a positive example,
can't do classification

most specific hypothesis → $\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

No day (instance)
is positive example

- Try to induce general functions from specific training examples.

Version Spaces & Candidate Elimination Alg.

Candidate Elimination Algorithm: Output a description of the set of all hypotheses consistent with the training examples

Version space: Subset of hypotheses in H that is consistent with training data.

- General boundary G is the set of maximally general members of H consistent with data.
- Specific boundary S is minimally general members of H consistent with data.

Version Space = General boundary + Specific Boundary + Everything in Between

Candidate Elimination Learning Algorithm

1. G ve S 'yi tanımla $G_0 = \langle ?, ?, ?, ?, ?, ? \rangle$
 $S_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

2. Datayı tara. (for d in D)

- Eğer d pozitif örneğe:

- G 'nin içinde d ile çelişen hipotezleri sil.

- S 'deki hipotezi minimal şekilde genelleştir.

- S 'de daha genel bir hipotez varsa kaldır.

- S ile d çelişmiyorsa güncelleme yapma.

- Eğer d negatif örneğe

- S 'de çelişen örnekleri sil.

- G 'de d ile çelişen hipotez varsa (g)

- g 'yi sil.

- G 'yi d 'ye göre genelle.

- G 'de daha spesifik olan hipotezleri sil.

S 'de d ile çelişen hipotezler varsa



Candidate Elimination Example

Training Examples:

$\langle \text{Sunny, Warm, Normal, Strong, Warm, Same} \rangle \rightarrow \text{yes}$

$\langle \text{Sunny, Warm, High, Strong, Warm, Same} \rangle \rightarrow \text{yes}$

$\langle \text{Rainy, Cold, High, Strong, Warm, Change} \rangle \rightarrow \text{No}$

1. $S_0 = \langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

$G_0 = \langle ?, ?, ?, ?, ?, ? \rangle$

2. Training Ex #1: (+) example, inconsistent with S

$G_1 = \langle ?, ?, ?, ?, ?, ? \rangle$

$S_1 = \langle S, W, N, S, W, S \rangle$

3. Ex #2: (+) example, inconsistent with S

$G_2 = \langle ?, ?, ?, ?, ?, ? \rangle$

$S_2 = \langle S, W, ?, S, W, S \rangle$

↳ generalized

4. Ex #3: (-) example, inconsistent with G

$S_3 = \langle S, W, ?, S, W, S \rangle$

$G_3 = \langle \text{Sunny}, ?, ?, ?, ?, ? \rangle, \langle ?, W, ?, ?, ?, ? \rangle$

$\langle ?, ?, ?, ?, ?, S \rangle \rightarrow$ Alternative maximally general hypothesis

5. Ex #4: $\langle S, W, H, S, C, C \rangle \rightarrow (+)$

Inconsistent with S

$S_4 = \langle S, W, ?, S, ?, ? \rangle \rightarrow$ Son 2 si S_3 ile gelişiyor

$G_4 = \langle S, ?, ?, ?, ?, ? \rangle, \langle ?, W, ?, ?, ?, ? \rangle$

↳ $\langle ?, ?, ?, ?, ?, S \rangle$ yi ekledik sonunda C olduğu için (inconsistent)

Concept Learning

Problem of searching through a predefined space of potential hypotheses for the hypothesis that fits the data best.

- Each constraint can be
 - a value (water = warm)
 - don't care (water = "?")
 - no value allowed (water = "∅")

? : Feature has no predictive power. (almost redundant)

★ We're looking for a conjunction of constraints that describe (+) examples the best.

★ most general hypothesis → $\langle ?, ?, ?, ?, ?, ? \rangle$

Every day (instance)
is a positive example,
can't do classification

most specific hypothesis → $\langle \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle$

No day (instance)
is positive example

- Try to induce general functions from specific training examples.

Inductive Bias

Some of the hypotheses are impossible to exist.

ex//

1	S	W	N	S	C	C	(+)
2	C	W	N	S	C	C	(+)
3	R	W	N	S	C	C	(-)

↳ sadece bu farklı

elimizdeki en

spec set $\rightarrow \langle ?, W, N, S, C, C \rangle$

↳ çok genel ve (-) i de kapsuyor.

An Unbiased Learner

↓
biased learner

Inductive Bias: Set of all assumptions made by an ML algorithm to perform induction.

↳ generalization

Inductive bias of candidate elimination: target concept c is contained in given hypothesis space H .

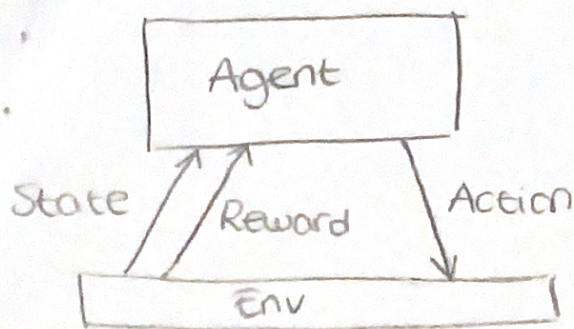
Rule Learner: Store examples, classify x iff it matches previously observed example (training data)

Candidate elimination: Assumes target concept can be explained in hypothesis space.

Find-S: Assumes all instances are negative unless the opposite is given. (Datadaki (+) örnekler dışında hepsi (-))

↓
weakest
to
strongest
bias

Reinforcement Learning



$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots$$

Agent chooses action a_i in state s_i and gets reward r_i . But goal is to maximize:

$$r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \quad 0 \leq \gamma < 1$$

Reward on the long term

γ : discount factor
(a parameter on how much we care about immediate & future rewards)

- Learn a control policy

$$\pi: S \rightarrow A$$

set of states

set of actions

- Take a from A given current state s from S

Problems

- Delayed Reward & Temporal Credit Assignment
Determine which of the actions in its sequence are to be credited for eventual rewards.
- Exploration vs Exploitation
Trade-off in choosing exploration of unknown states & actions (more info) or exploitation of states & actions that are known to yield reward.

The Learning Task

Markov Decision Process

$S \rightarrow$ set of states
 $A \rightarrow$ set of actions

At each t , in s_t , perform a_t , get $r(s_t, a_t)$, produce $s_{t+1} = \delta(s_t, a_t)$

- $\pi: S \rightarrow A$

learn this

- $V^\pi(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$

cumulative reward

$$= \sum_{i=0}^{\infty} \gamma^i r_{t+i}$$

$$0 \leq \gamma < 1$$

- When $\gamma = 0 \rightarrow$ only immediate reward is considered
- When $\gamma \approx 1 \rightarrow$ future rewards are more important

These only depend on current state.

Optimal policy is:

$$\pi \in \arg \max V^\pi(s)$$

pick the policy with the most value

Definitions

Action (A): All the possible moves agent can make.

State (S): Current situation returned by environment

Reward (R): An immediate return sent back from the environment to evaluate last action.

Policy (π): Strategy that agent employs to determine next action based on current state.

Value (V): Expected long-term return with discount, as opposed to short term reward.

$V^\pi(s) \rightarrow$ long term return of current state under policy π

Q-value: $Q^\pi(s,a) \rightarrow$ long term return of current state s , taking action a under policy π .

Q-Learning

Reward = $Q(s,a) \rightarrow$ probability of a reward in action-state pair } Q-value

$$Q(s,a) = r + \gamma \max_{a'} Q(s',a')$$

Diagram illustrating the components of the Q-learning equation:

- r : reward by making action a
- γ : next reward r'
- s' : next state
- a' : action in next state
- The entire right-hand side is labeled as "update".

$$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))$$

Cartpole Example

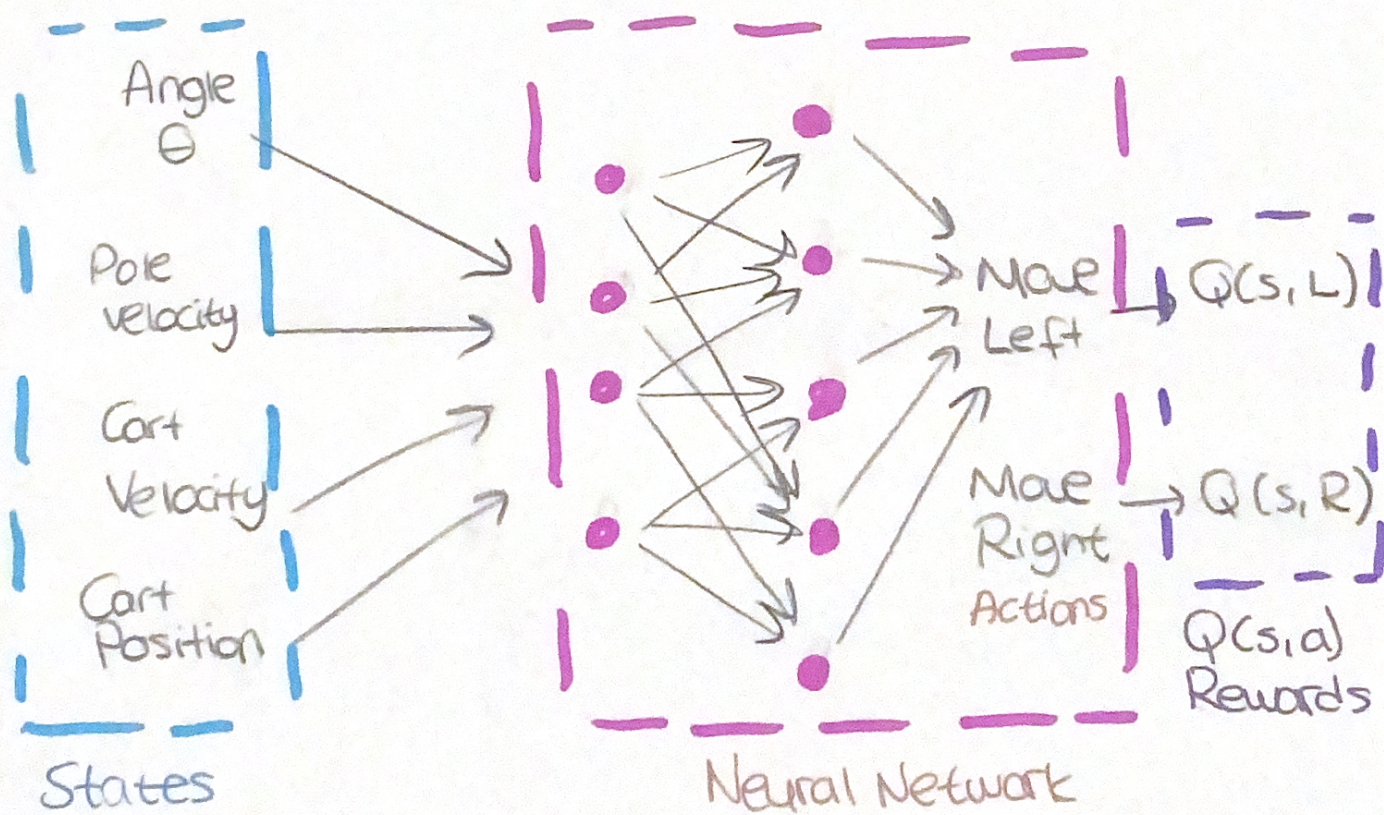
Actions \rightarrow move L & move R (discrete)

State space $\rightarrow (\theta, \dot{\theta}, x, \dot{x})$ (continuous)

\hookrightarrow we must shift state in every slight change in angle, not feasible to hold in Q-table

Deep Q-Network

Approximate $Q(s,a)$



★ It is essential to have a loss function that minimizes error between approximation and true $Q(s,a)$ calculated from equation

(S, A, R, P, P)
 \downarrow state \downarrow action \downarrow reward \downarrow transition probs \rightarrow initial state distribution

Reward Function $\rightarrow r_{t+1} = R(s_t, a_t, s_{t+1})$

\hookrightarrow This formulation is called a Markov Decision Process

MDP: A method to select an action a given state s .
Then observe $o' \in s'$ based on transition probs P .

★ Most DQNs are flat CNNs + batch normalization

★ DQN is the agent's learning in the env.

Markov Decision Processes

Special stochastic time control process for decision making which assumes random probability \in a decision maker having complete control.

S: Set of states. At each time step, the state of the environment is an element $s \in S$.

A: Set of actions. At each time step, agent chooses an action $a \in A$ to perform.

$P(s_{t+1} | s_t, a_t)$: State transition model that describes how the env state changes when user performs an action a depending on current state s .

$P(r_{t+1} | s_t, a_t)$: Reward model that describes the real valued reward value that the agent receives from the env after performing an action. (Depends on state \in action)

γ : Discount factor that controls importance of future rewards.

— In MDP, we search for a policy function that the agent or decision maker will choose in next state s .

(Contd)

- Once we specify policy & fix actions for each state,
- state depends on $\Pi(s)$ & $\Pr(s_{t+1} = s' | s_t = s, a_t = a)$
 - ↓
policy on state s
 - transition probability

Reward = total discounted reward

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

↳ discounted sum of all rewards obtained from time t .

prevents reward from going to infinity

Value Iteration VS Policy Iteration

- Constantly refines value function v (or Q)

Find $Q(s, a)$

$a = \operatorname{argmax} Q(s, a)$

- Defines policy function that converges to most optimal policy (through policy gradient)

Find $\Pi(s)$

Sample $a \sim \Pi(s)$

} both uses MDP

Other RL algorithms

SARSA: (State-Action-Reward-State-Action) uses MDP to adjust value of Q -function based on next state (modified Q -learning that uses extra action E_t state)

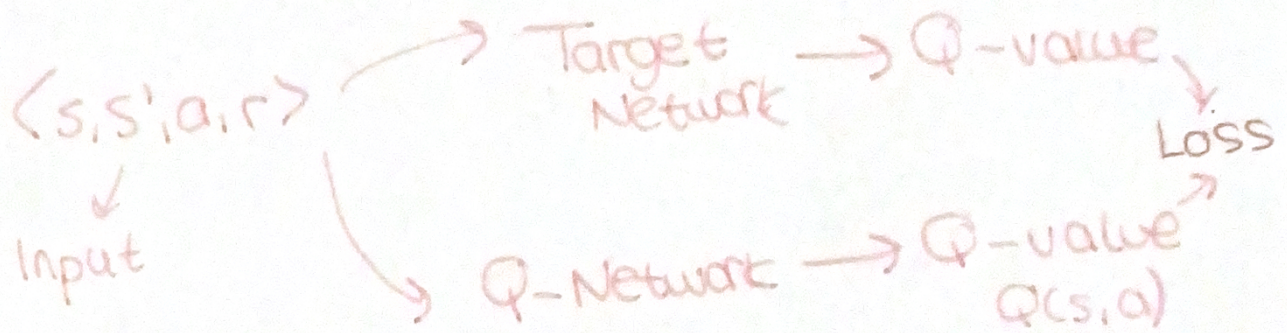
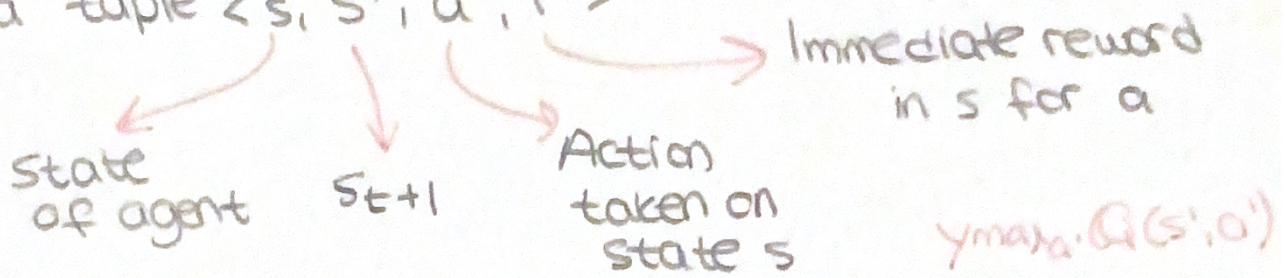
Monte Carlo Methods: Directly learns from experience E_t past $a-s$ pairs without any prior knowledge of MDP probs. MC uses policy iteration.

Experience Replay

Neural network training process becomes more stable when training is done on random batch of previous experiences

Experience Replay: Memory that stores experiences

as a tuple $\langle s, s', a, r \rangle$



$$\text{Loss} = (r + \gamma \max_{a'} Q(s', a') - Q(s, a))^2 \quad \leftarrow \text{Backprop}$$

Deep Q-Learning Steps

1. Provide state of environment to agent. Get Q-values from Target Net & Q-Net.
2. Pick action a , based on epsilon value.
3. Perform action a .
4. Observe reward & next state s' . Put them in $\langle s, s', a, r \rangle$
5. Sample random batches from experience replay memory & perform training on Q-Net.
6. Each n^{th} iteration, copy weights from Q to Target.
7. Repeat 2-7 for each episode.

Episode \rightarrow all states between initial & terminal states.

Epoch \rightarrow Forward & backward pass

Genetic Algorithms

- Genetic algorithms are used to exploit a search space for the best solution.
- 1. Pool of hypotheses is iterated.
- 2. Every hypothesis is evaluated on a fitness function
- 3. A new population is generated with most fit individuals from current population. Some of them are passed intact & some of them are used for crossover & mutation.

Basic Algorithm

Chromosomes

Chromosomes are solutions encoded in strings.

ex: 1 0 0 1 1

Crossover

11100101 > 11111101
01011101 > 01000101

Mutation

11100101
↓
turn this into a 0

Selection Strategies

- Roulette wheel

A	B	C	D
0.4	0.3	0.2	0.1

Total Fitness

} Every chromosome has chance of being picked according to its fitness value

- Tournament Selection

Choose k individuals at random
Then choose the best individual from the pool.

Schema Theorem

Characterizes evolution of population (collectively)

Schemas \rightarrow 1*0*
↓
don't care
→ 1101
→ 1001
→ 1100
→ 1000

Population Evaluation

Schema Theorem

Given: Schema $\rightarrow 1*1**$ fitness

Chromosomes \rightarrow

$C_1: 00101$	10	
$C_2: 11101$	25*	fit chromosomes
$C_3: 00000$	15	
$C_4: 10010$	20	
$C_5: 11111$	30*	

Calculate: $m(s,t) \rightarrow$ # of instances of schema s in the population at time t

- $f(n) \rightarrow$ fitness of hypothesis n
- $\bar{f}(t) \rightarrow$ average fitness of all individuals in pop at t
- $\hat{u}(s,t) \rightarrow$ average fitness of instances of s .

$Pr(n) = \frac{f(n)}{\sum f(n_i)}$ } fitness of n

$Pr(n \in s) = \frac{\sum f(n)}{n \bar{f}(t)} = \frac{\hat{u}(s,t) m(s,t)}{n \bar{f}(t)}$

↓
hyp. in schema

$\hat{u}(s,t) = \frac{\sum_{n \in s} f(n)}{m(s,t)}$ \rightarrow schemaya uyantarin fitness toplami

(schemaya uyantarin sayisi)

prob that we'll pick a hypothesis in the schema

Solution:

- $\hat{u}(s,t) = \frac{25+30}{2} = 27.5$
- $\bar{f}(t) = \frac{100}{5} = 20$ $\rightarrow 27.5 > 20$
elements of schema expected to increase
- $P_c = 1$
 $d(s) = 3-1$
 $L = 5$

$E[m(s,t+1)] \geq \frac{\hat{u}(s,t)}{\bar{f}(t)} m(s,t) \cdot (1 - P_c \frac{d(s)}{L-1}) \cdot (1 - P_m)$

bits ↑

$E[m(s,t+1)] \geq 1.375$

\hookrightarrow Expected # is decreasing due to crossover

Genetic Algorithms

Expressing hypotheses as logic functions & genomes

IF $a_1 = T \wedge a_2 = F$ THEN $c = T$

a_1	a_2	c
10	01	1

Modelling Problems & Defining Fitness Functions

A company owns 3 buses & employs 5 drivers. Only one driver can drive any bus on a single day. Each driver cannot work for more than two consecutive days in a week. Company uses buses every day. Make a weekly schedule (7 days) denoting drivers assigned to each of three buses. Provide encoding & fitness function.

Genes (Drivers) $\rightarrow A, B, C, D, E$

Occupied drivers for a given day \rightarrow

Assuming bus order is ignored

- ABC
 - ABD
 - ABE
 - ACD
 - ACE
 - ADE
 - BCD
 - BCE
 - BDE
 - CDE
- } 10

ABC CDE ABD ABE CDE ABC CDE

ABC
ABD

⋮
CDE

\rightarrow can be modeled with idleness

ABC \rightarrow 11100

CDE

Fitness: How many days each crew has left before a day off.

$$F = d_1 + d_2 + \dots + d_m = \sum_{i=1}^m d_i$$

$d =$ # Days an employee has before a day off

$m =$ # employees