

S M U

- MODELY NAUČITELNOSTI PAC A ONLINE
 - NAUČITELNOST KONJUNKCÍ A DISJUNKCÍ
 - BAYESOVSKÉ SÍŤE
 - REINFORCEMENT LEARNING
- SEQUENTIAL
- $$p^u(x_{t+1}) = p^u(x_1) p^u(a_1|x_1) \cdot p^u(x_2|x_1, a_1) \cdot \dots \cdot p^u(x_{t+1}|x_t, a_t) p^u(a_t|x_t, x_{t-1})$$
- POLICY
- $$a_t = \pi(x_{t-1}, x_t)$$
- DETERMINISTIC
- $$p^u(a_t|x_{t-1}, x_t) = 1 \quad \text{FOR } a_t = \pi(x_{t-1}, x_t)$$
- $$= 0 \quad \text{OTHERWISE}$$
- OBSERVATION, REWARD
- $$x_t = (o_t, r_t)$$
- $$p^u(x_t|x_{t-1}) = p^u(o_t, r_t|x_{t-1}) = p^u(o_t|r_t, x_{t-1}) \cdot p^u(r_t|x_{t-1}) =$$
- $$= p^u(r_t|o_t, x_{t-1}) \cdot p^u(o_t|x_{t-1})$$
- NON-SEQUENTIAL
- $$p^u(o_t|r_t, x_{t-1}) = p^u(o_t)$$
- REWARDS OR DEPENDENT DEPEND ON PREVIOUS OBSERVATION AND ACTION ONLY
- $$p^u(r_t|o_t, x_{t-1}) = p^u(r_t|o_{t-1}, a_{t-1})$$

- BATCH LEARNING

- NONSEQUENTIAL

- LEARNING PHASE $\lambda = 1, 2, \dots, K$

- ACTION PHASE $\lambda = K+1, K+2, \dots$

- AGENT DOES NOT CHANGE ITS DECISION POLICY IN ACTION PHASE

IF $\lambda, \lambda' > K$ AND $x_\lambda = x_{\lambda'}$ THEN $\sigma_\lambda = \sigma_{\lambda'}$

AND IGNORES REWARDS

- FOR $\lambda > K$

$$\sigma_\lambda = \pi(x | \sigma_{\lambda \leq K}, \theta_\lambda)$$

- POLICY DEPENDS ON ACTUAL OBSERVATION AND STEPS OF LEARNING PHASE

- REWARDS

- WE WANT TO MAXIMIZE SUM OF REWARDS UP TO FINITE HORIZON $m \in \mathbb{N}$

$$r_1 + r_2 + \dots + r_m$$

- OR DISCOUNTED REWARD

$$\sum_{\lambda=1}^{\infty} r_\lambda \gamma^\lambda \quad \text{WHERE } \forall \lambda: \gamma^\lambda \geq 0 \quad \text{AND} \quad \sum_{\lambda=1}^{\infty} \gamma^\lambda < \infty$$

- BUT REWARDS ARE PROBABILISTIC

- WE WANT TO MAXIMIZE EXPECTED CUMULATIVE REWARD

$$\sum_{r \in \mathcal{R}} \gamma^W(r_{\leq m} | \sigma_{\leq m}) \cdot (r_1 + r_2 + \dots + r_m)$$

- OR DISCOUNTED REWARD

$$\lim_{m \rightarrow \infty} \sum_{r \in \mathcal{R}} \gamma^W(r_{\leq m} | \sigma_{\leq m}) \sum_{\lambda=1}^m r_\lambda \gamma^\lambda$$

- IN BATCH LEARNING WE DON'T CARE ABOUT REWARDS IN LEARNING PHASE

$$\sum_{r_2 \in \mathcal{R}} \gamma^W(r_2 | x | \sigma_{\leq K}) \cdot r_2$$

- ENVIRONMENT STATE

- WE DON'T HAVE TO LOOK AT WHOLE SEQUENCE OF HISTORY

- WE CAN LOOK ONLY AT STATE OF ENVIRONMENT

$$g^u(x_2 | s_2)$$

- INITIAL STATE IS "ПОЧМНУ"

$$g^u(x_1 | s^*) = g^u(x_1)$$

- UPDATE DISTRIBUTION

$$E(s_2 | s_{2-1}, x_{2-1}, \sigma_{2-1})$$

- ~~THE~~ $g^u(x_2 | s_2) = g^u(x_2 | x_2, \sigma_2)$

- ASSUMPTION

$$|E| < E_{\max}$$

- STATES ARE BOUNDED BY E_{\max}

- THERE ARE LESS STATES THAN ALL POSSIBLE HISTORIES

- PERCEPTS DEPEND ON STATE, BUT STATE DOESN'T DEPEND ON PERCEPTS

$$E(x_2 | s_{2-1}, x_{2-1}, \sigma_{2-1}) = E(s_2 | s_{2-1}, \sigma_{2-1})$$

- AGENT STATE

$$|A| < A_{\max}$$

- ACTION IS DETERMINED BY AGENT'S STATE

$$\sigma_2 = \pi(a_2)$$

$$\sigma_2 = \bar{\pi}(a_2, s_2)$$

- STATE IS UPDATED DETERMINISTICALLY

$$s_2 = A(s_{2-1}, x_2)$$

- FOR BATCH LEARNING, THE AGENT DOESN'T UPDATE ITS STATE IN ACTION

PHASE $s_2 = s_k \quad \forall 2 \geq k$

- ON-LINE CONCEPT LEARNING

- CONCEPT OF ENVIRONMENT

$$- C = \{o \in O \mid \mathcal{N}_o(o|1) > 0\}$$

- SET OF OBSERVATIONS WHICH CAN BE GENERATED FROM STATE 1

- OBSERVATIONS COMING FROM THIS CONCEPT ARE POSITIVE OBSERVATIONS

- $O \setminus C$ ARE NEGATIVE OBSERVATIONS

- AGENT PERFORMS CLASSIFICATION LEARNING

- AGENT TRIES TO PREDICT THE STATE OF ENVIRONMENT

$$Y = E = \{0, 1\}$$

- REWARD

$$r_{t+1} = \begin{cases} 0 & \text{if } \mathcal{N}_x = \mathcal{D}_x \\ -1 & \text{OTHERWISE} \end{cases}$$

- AGENT'S (HYPOTHESIZED) CONCEPT

$$- C(a_x) = \{o \in O \mid \pi(a_x, o) = 1\}$$

- AGENT WANT TO IDENTIFY ITS CONCEPT WITH THE CONCEPT OF ENVIRONMENT

$$C(a_x) = C$$

- AGENT'S HYPOTHESIS

$$\mathcal{N}_x = a_x$$

$$\mathcal{D}_x = \pi(\mathcal{N}_x, o_x)$$

$$C(\mathcal{N}_x) = \{o \in O \mid \pi(\mathcal{N}_x, o) = 1\}$$

$$C(\mathcal{N}_x) = C$$

- UPDATE RULE

$$w_x = (o_x, h_x)$$

$$A(w_{x-1}, x_x) = A((o_{x-1}, h_{x-1}), (o_x, h_x)) = (o_x, h_x)$$

- GENERALIZING AGENT

- USES CONJUNCTIONS

$$r_x = \pi(h_x, o_x) = \begin{cases} 1 & \text{IF } o_x \models h_x \\ 0 & \text{OTHERWISE} \end{cases}$$

- STARTS WITH THE MOST SPECIFIC HYPOTHESIS

- CONJUNCTION OF ALL PROPOSITIONAL VARIABLES AND THEIR NEGATION

- INITIAL HYPOTHESIS IS GENERALIZED TOWARDS THE CORRECT ONE

$$- h_1 = p_1 \wedge \neg p_1 \wedge p_2 \wedge \neg p_2 \wedge p_3 \wedge \neg p_3$$

- UPDATE RULE

$$h_x = \begin{cases} h_{x-1} & \text{IF } r_x = 0 \\ \text{DELETE } (h_{x-1}, o_{x-1}) & \text{OTHERWISE} \end{cases}$$

$$\text{DELETE } \left(\bigwedge_{i \in I} p_i \wedge \bigwedge_{j \in J} \neg p_j, (o^1, o^2, \dots, o^m) \right) = \bigwedge_{i \in I} p_i \wedge \bigwedge_{j \in J} \neg p_j \\ o^{i=1} \quad o^{j=0}$$

- DELETE KEEPS LITERALS FROM h_{x-1} WHICH ARE CONSISTENT WITH o_{x-1}

- WE NEED TO ASSUME THAT THERE EXISTS "CORRECT" CONJUNCTION h^*

$$\pi(h^*, o_x) = r_x \quad \forall o_x \in O$$

- AGENT MAKES MISTAKES ONLY ON POSITIVE EXAMPLES

$$h_x \supseteq h^*$$

- MISTAKES ARE CORRECTED BY REMOVING AT LEAST ONE INCONSISTENT

LITERAL

- GENERALIZING AGENT MAKES AT MOST $2m$ MISTAKES

- $2m$ LITERALS (POSITIVE AND NEGATIVE) (REMOVES AT LEAST ONE FOR EACH SAMPLE)

- CUMULATIVE REWARD

$$\sum_{i=1}^m r_i \geq -2m$$

- CAN ALSO LEARN DISJUNCTIONS

$$\neg (P_1 \vee P_2 \vee \dots \vee P_m) = \neg P_1 \wedge \neg P_2 \wedge \dots \wedge \neg P_m$$

- WE HAVE TO CHANGE

$$\bar{\sigma}_i = (1 - \sigma_{i1}^+, 1 - \sigma_{i1}^-, \dots, 1 - \sigma_{i2}^+) = \neg P_1 \wedge \neg P_2 \wedge \dots \wedge \neg P_m$$

$$\bar{\sigma}_i = 1 - \sigma_i$$

- SUBSUMPTION

- SUBSUMPTION LATTICE

- PARTIALLY ORDERED

- EACH TWO ELEMENTS HAVE THEIR UNIQUE LEAST UPPER BOUND AND THE GREATEST LOWER BOUND

- $\mathcal{A}_1 \subseteq \mathcal{A}_2$ IMPLIES $\mathcal{A}_2 \vdash \mathcal{A}_1$ IF \mathcal{A}_1 AND \mathcal{A}_2 ARE CONJUNCTIONS

- \mathcal{A}_1 ENTAILS \mathcal{A}_2 IF ANY MODEL OF \mathcal{A}_1 IS ALSO MODEL OF \mathcal{A}_2

$$\mathcal{A}_1 \vdash \mathcal{A}_2$$

- FOR DISJUNCTION

$$\mathcal{A}_1 \subseteq \mathcal{A}_2 \text{ IMPLIES } \mathcal{A}_1 \vdash \mathcal{A}_2$$

- SEPARATING AGENT

- \mathbf{h}_2 WILL BE REPRESENTED BY NON-LOGICAL MEANS

- DEFINES HYPERPLANE IN $\mathcal{O} = \{0, 1\}^n$

- $C(\mathbf{h}_2)$ INCLUDE OBSERVATIONS LYING ABOVE HYPERPLANE

- POLICY

$$r_2 = \pi(\mathbf{h}_2, \mathbf{o}_2) = \begin{cases} 1 & \text{IF } \mathbf{h}_2 \cdot \mathbf{o}_2 > \frac{n}{2} \\ 0 & \text{OTHERWISE} \end{cases}$$

- \mathbf{h}_2 IS n -TUPLE OF INTEGER VALUES BOUNDED BY CONSTANT $q \in \mathbb{N}$

- INITIAL HYPOTHESIS

$$\mathbf{h}_1 = (1, 1, \dots, 1)$$

- UPDATE RULE

$$\mathbf{h}_2 = \begin{cases} \mathbf{h}_{2-1} & \text{IF } r_2 = 0 \\ \text{UPDATE}(2, \mathbf{h}_{2-1}, \mathbf{o}_{2-1}) & \text{IF } \mathbf{h}_{2-1} \cdot \mathbf{o}_{2-1} \leq \frac{n}{2} \\ \text{UPDATE}(0, \mathbf{h}_{2-1}, \mathbf{o}_{2-1}) & \text{IF } \mathbf{h}_{2-1} \cdot \mathbf{o}_{2-1} > \frac{n}{2} \end{cases}$$

$$\text{UPDATE} : \mathbf{h}_2^i = \begin{cases} \alpha \cdot \mathbf{h}_{2-1}^i & \text{IF } \mathbf{o}_{2-1}^i = 1 \\ \mathbf{h}_{2-1}^i & \text{OTHERWISE} \end{cases}$$

- INTEGER COUNTERPART OF PERCEPTRON ALGORITHM

- MAKES AT MOST $2 + 2s \lg n$ MISTAKES

- CUMULATIVE REWARD IS

$$\sum_{\lambda=1}^m r_\lambda \geq -2 - 2s \lg n$$

- PERFORMS BETTER IF NUMBER OF VARIABLES n IS LARGER THAN NUMBER OF RELEVANT VARIABLES

- HYPOTHESIS CLASS

- SET OF ALL HYPOTHESES AN AGENT CAN EXPRESS
- FINITE SET
- GENERALIZING AGENT
 - ALL CONJUNCTIONS MADE OF AT MOST n VARIABLES
- SEPARATING AGENT
 - SET OF q -BOUNDED n -TUPLES OF INTEGERS
- EACH HYPOTHESIS CLASS INDUCES SET OF CONCEPTS

$$C(H) = \{c(h) \mid h \in H\}$$

- CONCEPT CLASS
- WE WOULD LIKE HUGE HYPOTHESIS SPACE FOR AGENT, WHICH INCREASES THE CHANCE THAT $C \in C(H)$, BUT IT INCREASES NUMBER OF MAXIMAL ERRORS $|H|-1$

- VERSION SPACE AGENT

- MISTAKE BOUND IMPROVED FROM $|H|-1$ TO $\lg |C|$
- ON EACH OBSERVATION AGENT DISCARDS ALL HYPOTHESES FROM HYPOTHESIS CLASS WHICH ARE INCONSISTENT WITH THE OBSERVATION
- VERSION SPACE AGENT'S STATE IS

$$s_t = (o_t, H_t)$$

- o_t IS MEMORIZED OBSERVATION AND H_t IS SET OF HYPOTHESES
- DECISIONS ARE DETERMINED BY VOTING AMONG ALL HYPOTHESES IN H_t

$$y_t = \pi(H_t, o_t) = \begin{cases} 1 & \text{IF } |\{h \in H_t \mid o_t \models h\}| > |H_t|/2 \\ 0 & \text{OTHERWISE} \end{cases}$$

- AGENT STARTS WITH FULL HYPOTHESIS STATE
- WE REMOVE ALL HYPOTHESIS INCONSISTENT WITH o_{t-1}
- AGENT MAKES AT MOST $\log |H|$ MISTAKES
- CUMULATIVE REWARD IS

$$\sum_{t=1}^m r_t \geq -\log |H|$$

- IF MISTAKE IS MADE, AT LEAST HALF OF HYPOTHESIS ARE REMOVED.
- IN WORST CASE, THE LAST HYPOTHESIS IS CORRECT
- IT NEEDS A LOT OF MEMORY

- AGENT LEARNS CLASS H ON-LINE IF IT MAKES AT MOST $\text{poly}(m)$ IN THE ON-LINE SCENARIO IF $C \in C(H)$. m IS SIZE OF OBSERVATIONS.

- AGENT THAT LEARNS HYPOTHESIS CLASS H ON-LINE IS SAID TO LEARN IT EFFICIENTLY IF IT SPENDS AT MOST POLYNOMIAL TIME (IN m) BETWEEN THE RECEIPT OF A PERCEPT AND GENERATION OF THE NEXT ACTION

- IF $|H|$ IS AT MOST EXPONENTIAL IN m , THEN VERSION SPACE AGENT LEARNS CLASS H ON-LINE

- VERSION SPACE AGENT DOESN'T LEARN H EFFICIENTLY

- VAPNIK - CHERVONENKIS DIMENSION

- CARDINALITY OF LARGEST SET $O' \subseteq O$ THAT IS SHATTERED BY H

- NO UPPER BOUND ON THE NUMBER OF MISTAKES MADE BY AGENT IS SMALLER BY $V_C(H)$

- BATCH LEARNING WITH GENERAL ON-LINE AGENTS

- STANDARD ON-LINE AGENT

- USES SINGLE HYPOTHESIS

- CHANGES IT IF AND ONLY IF A MISTAKE HAS BEEN MADE

- IF STANDARD ON-LINE AGENT RETAINS A HYPOTHESIS h_x FOR q STEPS ($h_x = h_{x+1} = \dots = h_{x+q}$) THEN $\text{ERR}(h_x) \leq \epsilon$ WITH PROBABILITY AT LEAST $1 - 2^{-q\epsilon}$

- SO WE CAN WAIT FOR $h_x = h_{x+1} = \dots = h_{x+q}$ TO HAPPEN AND THEN KEEP h_x WITH PROBABILISTIC ERROR BOUND

- CONSISTENT AGENT

- EQUIVALENT OF VERSION SPACE AGENT, BUT FOR BATCH LEARNING

- COLLECT ALL OBSERVATIONS

- AT THE END OF TRAINING PHASE SELECTS ~~OR~~ RANDOM HYPOTHESIS CONSISTENT WITH ALL OBSERVATIONS

- PAC LEARNING MODEL

- AGENT PROBABLY APPROXIMATELY LEARNS CLASS H (IN BATCH SETTING) IF AT THE END OF THE TRAINING PHASE IT PRODUCES h_k SUCH THAT $\text{ERR}(h_k) \leq \epsilon$ WITH PROBABILITY AT LEAST $1 - \delta$ AND $k - 1 = m \leq \text{POLY}(n, \frac{1}{\delta}, \frac{1}{\epsilon})$
- IT LEARNS THE CLASS EFFICIENTLY IF IT SPENDS AT MOST POLYNOMIAL (IN THE SAME VARIABLES) TIME BETWEEN THE RECEIPT OF PERCEPT AND THE GENERATION OF THE NEXT ACTION IN THE TRAINING PHASE
- GENERALIZING AGENT EFFICIENTLY PAC-LEARNS CONJUNCTIONS
- ANY STANDARD AGENT LEARNING (EFFICIENTLY) A CONCEPT CLASS C ON-LINE, HAS A COUNTERPART WHICH (EFFICIENTLY) PAC-LEARNS C .
- IF $C \in \text{CLH}$ AND $|H|$ IS AT MOST EXPONENTIAL IN m THEN THE CONSISTENT AGENT PAC-LEARNS H .

- REINFORCEMENT LEARNING

- AGENT INTERACT WITH ENVIRONMENT
- LEARNS HOW TO ACT OPTIMALLY
- WORLD IS MARKOV DECISION PROCESS
- TRANSITION AND REWARD IS NOT KNOWN
- GOAL IS TO LEARN OPTIMAL POLICY

- DIFFICULTIES

- ACTION EFFECTS ARE NON-DETERMINISTIC
- REWARDS ARE SPARSE AND DELAYED
- GREEDY DOESN'T WORK
- LARGE AND COMPLEX WORLD

- PASSIVE LEARNING

- FIXED POLICY π
- IT LEARNS HOW GOOD THE POLICY IS BY INTERACTING WITH THE ENVIRONMENT (LEARNS $U^\pi(s)$)

- ACTIVE LEARNING

- AGENT SEARCHES FOR OPTIMAL POLICY
- IT EXPLORES MANY DIFFERENT ACTIONS IN MANY DIFFERENT STATES

- MODEL BASED LEARNING

- LEARNS MDP MODEL
 - S - TRANSITION FUNCTION
 - R - REWARD FUNCTION

- MODEL FREE LEARNING

- DOESN'T LEARN MODEL DIRECTLY
- LEARNS A POLICY

- PASSIVE REINFORCEMENT LEARNING

- DIRECT UTILITY ESTIMATION (DUE)

- MODEL FREE

- LEARN EXPECTED REWARD-TO-GO FROM OBSERVED REWARDS-TO-GO

- REWARD-TO-GO OF STATE S IS THE SUM OF THE DISCOUNTED REWARDS FROM STATE S UNTIL A TERMINAL STATE IS REACHED

- EXPECTED REWARD-TO-GO MATCHES THE TRUE STATE UTILITY GIVEN THE POLICY

- ESTIMATE $U^\pi(s)$ AS AVERAGE TOTAL REWARD OF EPOCHS CONTAINING S (IF SAMPLE GOES THROUGH S , WE TAKE IT FOR CALCULATION OF $U^\pi(s)$)

- REDUCES RL TO SUPERVISED LEARNING

- GUARANTEED CONVERGENCE, BUT SLOW

- DOES NOT EMPLOY DEPENDENCE AMONG STATE UTILITIES

- ADAPTIVE DYNAMIC PROGRAMMING (ADP)

- COMPLEX MODEL BASED METHOD

- STEP 1

- LEARN TRANSITION MODEL S AND REWARD FUNCTION R

- STORE EACH $R(S)$

- $S(S'/S, a)$ - KEEP TRACK HOW OFTEN WE ENTER S' FROM S BY ACTION a

- STEP 2

- PERFORM POLICY EVALUATION BASED ON UNDERLYING MDP

- SOLVES n BELLMAN EQUATIONS $U^\pi(s) = R(s) + \gamma \sum_{s'} S(s'/s, \pi(s)) U^\pi(s')$

- TEMPORAL DIFFERENCE LEARNING (TD)

$$- U^\pi(s) = R(s) + \gamma \sum_{s'} S(s'|s, \pi(s)) U^\pi(s')$$

- INSTEAD OF SUMMING OVER ALL SUCCESSORS, ONLY ADJUST THE UTILITY OF THE STATE BASED ON THE SUCCESSOR OBSERVED IN THE TRIAL

$$U^\pi(s) \sim R(s) + \gamma U^\pi(s')$$

- TRANSITION MODEL IS NOT NEEDED

$$- U^\pi(s) = U^\pi(s) + \alpha \cdot (R(s) + \gamma U^\pi(s') - U^\pi(s))$$

- α - LEARNING RATE

- IF LEARNING RATE DECREASES, IT WILL CONVERGE TO TRUE UTILITY VALUE

- ACTIVE RL

- FIND OPTIMAL POLICY

- ACTIVE ADP

- DO NOT KEEP POLICY FIXED

- ALWAYS TAKE ACTION THAT MAXIMIZES EXPECTED REWARD

- GIVEN EXISTING ENVIRONMENT MODEL AND UTILITY ESTIMATES

$$U(s) = R(s) + \gamma \max_a \sum_{s'} S(s'|s, a) U(s')$$

- BUT CHOOSING OPTIMAL ACTION LEAD TO SUBOPTIMAL RESULTS

- WE NEED MORE EXPLORATION

- ϵ - GREEDY STRATEGY

- ϵ - RANDOM CHOICE

- REST GREEDY

- SOFTMAX STRATEGY

- ACTION a PROBABILITY AT TIME t IS A FUNCTION TO ITS RELATIVE VALUE

$$P(a) = \frac{e^{Q_t(a)/\tau}}{\sum_{a=1}^n e^{Q_t(a)/\tau}}$$

- τ IS TEMPERATURE

- $\tau \rightarrow 0$ CONVERGES TO GREEDY STRATEGY

- SHOULD DECREASE WITH TIME

- OPTIMISTIC UTILITIES

- ACTIONS NOT TRIED OFTEN GETS HIGHER WEIGHTS

$$U^+(s) = R(s) + \gamma \max_a \left(\sum_{s'} S(s'|s, a) U^+(s'), N(s, a) \right)$$

\uparrow NUMBER OF TIMES
TO TAKE IN s

- GLIE

- MUST TRY EACH ACTION IN EACH STATE AN UNBOUNDED NUMBER OF TIMES SO THAT IT DOES NOT MISS ANY OPTIMAL ACTION

- MUST EVENTUALLY BECOME GREEDY

- EXAMPLE: ϵ -GREEDY, OPTIMISTIC UTILITY

Q - LEARNING

- WE WILL LEARN $Q(s, a)$ INSTEAD OF $U(s)$

$$U(s) = \max_a Q(s, a)$$

$$a_{\text{GREEDY NEXT}} = \text{ARG MAX}_a Q(s, a)$$

$$U(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) U(s')$$

↓

$$Q(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) U(s')$$

↓

$$Q(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$$

- UPDATE STEP FROM TD

$$U^\pi(s) = U^\pi(s) + \alpha (R(s) + \gamma U^\pi(s') - U^\pi(s))$$

$$\downarrow Q^\pi(s, a) = Q^\pi(s, a) + \alpha (R(s) + \gamma \max_{a'} Q^\pi(s', a') - Q^\pi(s, a))$$

- OFF-POLICY UPDATE

- SARSA

- SLIGHTLY DIFFERENT UPDATE

$$Q^\pi(s, a) = Q^\pi(s, a) + \alpha (R(s) + \gamma Q^\pi(s', a') - Q^\pi(s, a))$$

- ON-POLICY UPDATE

- GUARANTEED CONVERGENCE

- SLOWER THAN ADP IN TERMS OF EPOCHS, BECAUSE IT DOESN'T ENFORCE CONSISTENCY AMONG VALUES THROUGH THE MODEL

- PROBABILISTIC GRAPHICAL MODELS
- FULL JOINT MODEL
- MARGINALIZATION

$$Pr(X) = \sum_{y \in Y} Pr(X, y)$$

- NORMALIZATION

$$Pr(X|Y) = \frac{Pr(X, Y)}{Pr(Y)}$$

- OR WITH NORMALIZATION CONSTANT α

$$Pr(X|Y) = \alpha \cdot Pr(X, Y), \quad \alpha \text{ IS SET SO THAT } \sum_{x \in X} Pr(x, Y) = 1$$

- UNIVERSAL

- FOR SUFFICIENT SAMPLE SIZE ITS LEARNING CONVERGES

- INTRACTABLE FOR REAL PROBLEMS

- $2^n - 1$ PROBABILITIES FOR n PROPOSITIONS

- COURSE OF DIMENSIONALITY

- IMPRETRABLE FOR REAL TASKS

- MODEL GIVES NO EXPLICIT KNOWLEDGE ABOUT THE DOMAIN

- CONDITIONAL INDEPENDENCE

- A AND B ARE CONDITIONALLY INDEPENDENT GIVEN C IF

$$Pr(A, B|C) = Pr(A|C) \cdot Pr(B|C) \quad \forall A, B, C \quad Pr(C) \neq 0$$

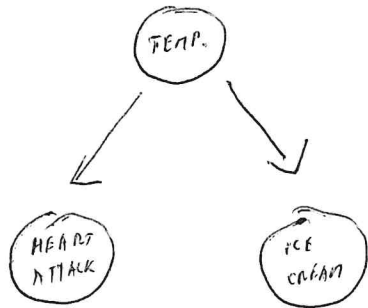
- DENOTES

$$A \perp\!\!\!\perp B | C$$

- SOME OBSERVATIONS BECOMES REDUNDANT

$$Pr(B|C) = Pr(B|A, C) \quad Pr(A|C) = Pr(A|B, C)$$

- ~~CON~~ DIVERGING CONNECTION



HEART - HEART ATTACK AND ICE CREAM INDEPENDENT GIVEN TEMPERATURE

$$Pr(H|I, T) = Pr(H|T)$$

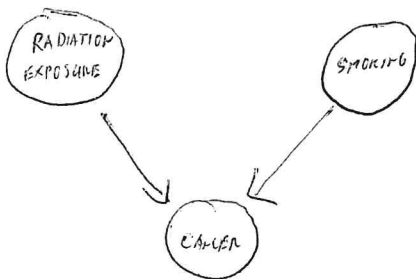
- LINEAR CONNECTION



- PHD_g AND PHD INDEPENDENT GIVEN PHD_p

$$Pr(PHD | PHD_p, PHD_g) = Pr(PHD | PHD_p)$$

- CONVERGING CONNECTION



- SMOKING AND RADIATION EXPOSURE BECAME DEPENDENT GIVEN CANCER

$$Pr(RADIATION EXPOSURE | CANCER, SMOKING) \neq Pr(RADIATION EXPOSURE | CANCER)$$

- D-SEPARATION

- EQUIVALENT TO CONDITIONAL INDEPENDENCE BETWEEN X AND Y GIVEN Z

- BAYESIAN NETWORK

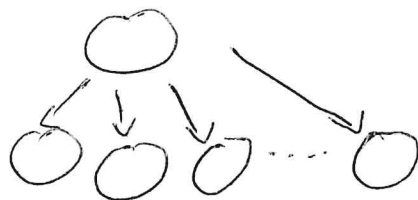
- DIRECTED ACYCLIC GRAPH
- NODES REPRESENTS RANDOM VARIABLES
- EDGES REPRESENTS DIRECT DEPENDENCE
- NODES ANNOTATED BY PROBABILITIES
 - NODE CONDITIONED BY CONJUNCTION OF ALL ITS PARENTS
 - ROOT NODES ANNOTATED BY PRIOR
- REDUCE COMPLEXITY BY FOCUSING ONLY ON FUNDAMENTAL RELATIONSHIPS

- NAIVE

- NO EDGES
- COMPLETELY INDEPENDENT
- $P_r(O_1, O_2, \dots, O_n) = P_r(O_1) \cdot P_r(O_2) \cdot \dots \cdot P_r(O_n)$

OR

- CI OF VARIABLES GIVEN DIAGNOSIS



- FULLY CONNECTED GRAPH
 - DEPENDENCE OF ALL VARIABLES
 - COMPLEXITY AS JOINT MODEL
- REASONABLE MODEL LIES IN BETWEEN

- FUNDAMENTAL TASK

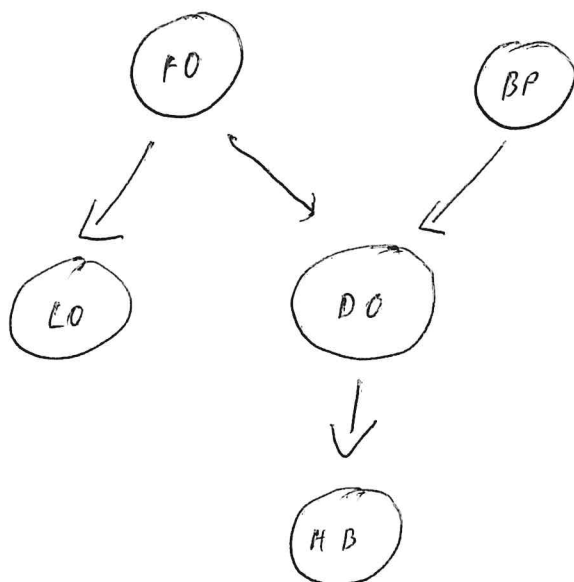
- INFERENCE

- FROM OBSERVED EVENTS ASSUME PROBABILITY OF OTHER EVENTS

- LEARNING MODEL PARAMETERS FROM DATA

- LEARNING NETWORK STRUCTURE

- WHICH EDGES TO USE



$$- Pr(FO, BP, DO, LO, HB) = Pr(FO) \cdot Pr(BP) \cdot Pr(LO|FO) \cdot Pr(DO|FO, BP) \cdot Pr(HB|DO)$$

- INFERENCE IS NP-HARD

- COMPLEXITY GROWS WITH NUMBER OF EDGES

- APPROXIMATE INFERENCE BY STOCHASTIC SAMPLING

- MONTE CARLO

- SAMPLES FROM JOINT PROBABILITY DISTRIBUTION

- ESTIMATE TARGET CONDITIONAL PROBABILITY FROM SAMPLE SET

- FORWARD SAMPLING

- TOPOLOGICALLY SORT THE NETWORK NODES
 - PARENT COMES BEFORE CHILDREN
- INSTANTIATE VARIABLES ALONG TOPOLOGICAL ORDERING
 - TAKE $P_r(O_j | \text{PARENTS}(O_j))$ AND SAMPLE O_j
- $P_r(q|e) \approx \frac{N(q, e)}{N(e)}$
- SAMPLES CONTRADICTING EVIDENCE ARE NOT USED

- REJECTION SAMPLING

- REJECT PARTIALLY GENERATED SAMPLES AS SOON AS THEY VIOLATES EVIDENCE e

- LIKELIHOOD WEIGHTING

- FIX VALUES OF e , WE SAMPLE THE REST
- BUT SAMPLES MUST BE WEIGHTED
- WEIGHT EQUALS TO LIKELIHOOD OF EVENT GIVEN THE EVIDENCE
- INITIALIZE SAMPLE WEIGHT

$$w^m = 1$$

- INSTANTIATE VARIABLES

- IF $O_j \in E$ THEN TAKE O_j^m FROM e AND $w^m \leftarrow w^m \cdot P_r(O_j | \text{PARENTS}(O_j))$
- $P_r(O_i | e) \approx \frac{\sum_{m=1}^n w^m \delta(O_i^m, O_i)}{\sum_{m=1}^n w^m}$ $\delta(\tilde{i}, j) = \begin{cases} 1 & \text{FOR } \tilde{i} = j \\ 0 & \text{FOR } \tilde{i} \neq j \end{cases}$

- GIBBS SAMPLING

- LEARNING BAYESIAN NETWORK FROM DATA

- FREQUENCY TABLE

- GIVES NUMBER OF SAMPLES WITH PARTICULAR CONFIGURATION

- 2^m ENTRIES

- WE KNOW THE STRUCTURE, WE SEARCH FOR CONDITIONAL PROBABILITY TABLES IN EACH NODE

- MAXIMUM LIKELIHOOD ESTIMATE

$$\begin{aligned} L(\theta; D) &= \prod_{m=1}^n Pr(d_m; \theta) = \prod_{m=1}^n Pr(O_{1m} O_{2m} \dots O_{r_m m}; \theta) \\ &= \prod_{j=1}^m \prod_{m=1}^n Pr(O_j; \text{PARENTS}(O_j); \theta_j) = \prod_{j=1}^m L_j(\theta_j; D) \end{aligned}$$

- CONTRIBUTION OF EACH NETWORK NODE $L_j(\theta_j; D)$ IS DETERMINED (MAXIMIZED) INDEPENDENTLY

$$\hat{\theta}_j = \underset{\theta}{\text{ARG MAX}} L_j(\theta_j; D)$$

- GENERALIZED FORMULA FOR ML IS

$$\hat{\theta}_{O_j | \text{PARENTS}(O_j)} = \frac{N(O_j, \text{PARENTS}(O_j))}{N(\text{PARENTS}(O_j))} \approx Pr(O_j | \text{PARENTS}(P_j))$$

- PROBLEM OF INCOMPLETE DATA

- EM ALGORITHM

- INITIALIZE NETWORK

- E-STEP

- TAKE EXISTING NETWORK AND CALCULATE MISSING VALUES (INFERENCE)

- M-STEP

- MODIFY THE NETWORK ACCORDING TO THE CURRENT COMPLETE OBSERVATION

- REPEAT UNTIL CHANGE

- STRUCTURE LEARNING

- CAN'T SEARCH WHOLE SPACE OF POSSIBILITIES

- NAIVE SEARCH CAN'T BE USED

- K2, MCMC ALGORITHMS

- SEVERAL SCORES HOW TO EVALUATE HOW WELL A STRUCTURE MATCHES THE DATA

- LOG LIKELIHOOD DOESN'T WORK, IT OVERFITS

- BAYESIAN SCORE

- LOCAL CI TESTS

