

B. Schölkopf (MPI) : 'Symbolic, Statistical & Causal AI'

FEYNMAN APPROACH TO INTELLIGENCE CREATE UNDERSTAND

▫ N. WIENER → Cybernetics → FROM ENERGY TO INFO-PROCESSING

▫ F. ROSENBLATT → PERCEPTRON - CONNECTIONISM \Rightarrow LEARNING!

→ XOR problem: COULD NOT BE LEARNED AT THE TIME!

→ NOVIKOFF (1962): PERCEPTRON CONVERGE THEOREM

▫ McCARTHY ET AL → INTELLIGENCE = PROCESS OF MANIPULATING
DISCRETE SYMBOLS

▫ MINSKY & PAPERT (1969): PARITY PROBLEM



↳ PERCEPTRONS FELL OUT OF FAVOR

↳ SYMBOLIC METHODS GAINED MOMENTUM [FROM CONNECTIONIST FAILURE]

▫ MORAVEL'S PARADOX (1980s):

↳ WELL SUITED FOR COMPUTERS!

- PROBLEMS THAT APPEAR HARD \Rightarrow SIMPLE
- PROBLEMS THAT APPEAR SIMPLE \Rightarrow HARD
- SYMBOLIC METHODS COULDN'T SCALE
- BOLTZMANN MACHINES, BACKPROP
- PROBABILISTIC EXPERT SYSTEMS

} BIRTH OF THE FIELD OF MACHINE LEARNING

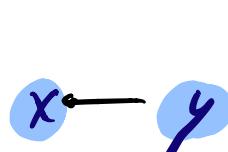
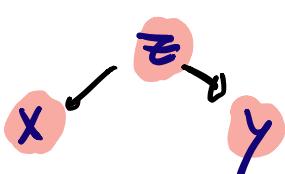
▫ DEEP MIND DQN - PROBLEM OF IID SETTING \Rightarrow RE-TAKIN + TRICKS

→ WE ARE INTERESTED IN MORE THAN STATISTICAL REGULARITIES

→ CAUSAL RELATIONSHIPS VS. DEPENDENCE

↳ REICHENBACH'S COMMON CAUSE PRINCIPLE

SYMMETRY
↳ ASYMMETRY



→ ALL THREE OPTIONS HAVE SAME OBSERVED DATA

"RUCKSTÜCK PROBLEM" // CARGO CULT \rightarrow NOT PROBLEM FOR IID

INDUSTRIAL REVOLUTION \rightarrow PROCESS ENERGY

} GENERATE / PROCESS,
DIGITAL REVOLUTION \rightarrow PROCESS INFORMATION
CONVERT!

B. Schölkopf + S. Bauer (MPI) : Causality

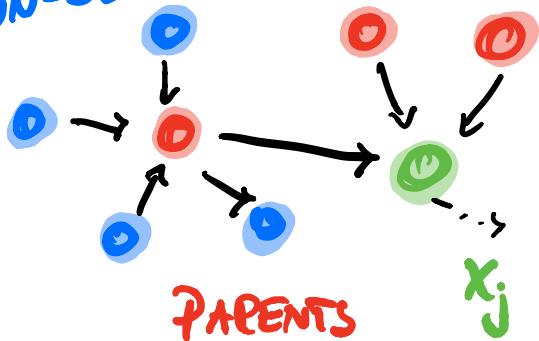
- INDEPENDENCE OF RVs: $X \perp\!\!\!\perp Y : p(x,y) = p(x) \cdot p(y)$ $\forall x,y$
 $\hookrightarrow X \perp\!\!\!\perp Y \Rightarrow E[X \cdot Y] = E[X] \cdot E[Y], \text{cov}[X,Y] = 0$
- COND. INDEP. OF RVs: $X \perp\!\!\!\perp Y | Z : p(x,y|z) = p(x|z) \cdot p(y|z), p(z) > 0$

INDEPENDENCE OF MECHANISMS ASSUMPTION: GENERATIVE PROCESS OF SYSTEM'S VARS. = COMPOSITION OF AUTONOMOUS MODULES THAT DO NOT INFLUENCE EACH OTHER

→ STRUCTURAL / FUNCTIONAL CAUSAL MODEL:

- DAG G w. VERTICES $x_1, \dots, x_n \Rightarrow V$: OBS., E : DIRECT CAUSATION
- $X_i := f_i(PA_i, U_i) \rightarrow$ NOISE VARS. / UNMODELED \Rightarrow INDEPENDENT
 \hookrightarrow WHY INDEP. U_i ? \Rightarrow REICHENBACH! \rightarrow CAUSAL SUFFICIENCY
 \hookrightarrow REICHENBACH'S COMMON CAUSE PRINCIPLE \Rightarrow ALL THREE OPTIONS HAVE SAME OBS. DATA!
- ENTAILED DISTRIBUTION $\rightarrow P(X_1, \dots, X_n) \Rightarrow$ OBSERV. + NOISE

NON-DESC.



\rightarrow MARKOV CONDITIONS \Rightarrow EQUIVALENCE:

EXISTENCE OF SCM

\Leftrightarrow LOCAL CAUSAL MARKOV

\Leftrightarrow GLOBAL CAUSAL MARKOV

$\Leftrightarrow p(X_1, \dots, X_n) = \prod p(X_i | PA_i)$ causal Markov kernels

- INTERVENTION: REPLACING $X_i := f_i(PA_i, U_i)$ w. ANOTHER ASSIGNM.
 \hookrightarrow RESULTING ENTAILED DISTR. = INTERVENTIONAL DISTR.

- MECHANISM INDEPENDENCE: $p(X_i | PA_i)$ CHANGE DOESN'T CHANGE $p(X_j | PA_j)$ for $j \neq i$. \rightarrow INVARIANCE \Rightarrow DISENTANGLING FACTORIZATION

□ COUNTERFACTUAL \Rightarrow 'WOULD BE' \rightarrow CAUSAL INF. = MISSING DATA P.

\rightarrow DO-CALCULUS: $p(Y | \text{do } X)$

"DOING VS.
SEEING"
CAN BE DIFFERENT
FROM $p(Y|X)$

DISTR. IF WE INTERVENE ON X
AND ASSIGN IT TO x

NEED $\& p$!

$p(X_1, \dots, X_n | \text{do } x_i) = \prod_{j \neq i} p(X_j | p_{X_j}) \delta_{x_i x_j}$

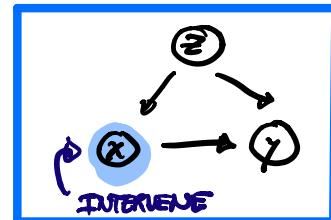
[CAUSAL FACT.] $p(X, Y, Z) = p(Z) p(X|Z) p(Y|X, Z)$

[INTERVENE] $p(X, Y, Z | \text{do } x) = p(Z) \delta_{Xx} p(Y|X, Z)$

[INTV. OUT X] $p(Y, Z | \text{do } x) = p(Z) p(Y|X, Z)$

[MARG. Z] $p(Y | \text{do } x) = \sum_z p(z) p(Y|X, z) \Rightarrow$ ADJUSTMENT FORMULA

[SIMPSON'S P.] $\neq p(Y|X)$



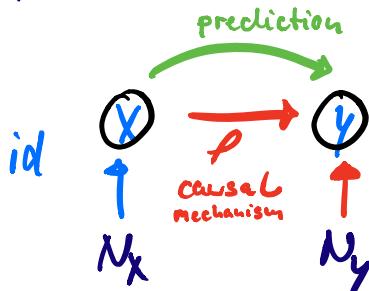
□ ADDITIONAL NOISE MODEL FOR 2 VARIABLE - CASE

$X \rightarrow Y$ ASSUME TO BE ADDITIVE: $f(X) = Y + N$ w. $X \perp\!\!\!\perp N$
 N [NOISE]

\hookrightarrow NON-LINEAR REG. + CHECK IF RESIDUALS ARE CORREL. W. X ?!

□ CAUSAL LEARNING:

$p(X), p(Y|X)$ INDEP.

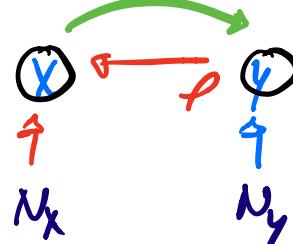


\Rightarrow SEMI-SUPERV. L. IMPOSSIBLE
 $\Rightarrow p(Y|X)$ INVARIANT UNDER
CHANGE IN $p(X)$

\hookrightarrow COVARIATE SHIFT

ANTIC. LEARNING:

$p(Y), p(X|Y)$ INDEP.



\hookrightarrow $p(X), p(Y|X)$ dependent!
 \Rightarrow SEMI-SUPERV. L. POSSIBLE
 $\Rightarrow p(Y|X)$ CHANGES W. $p(X)$

| | RODENG-TAKOUNGY | STAT | CAUSAL | DE |
|---------|-----------------|------|--------|----|
| IID | ✓ | ✓ | ✓ | ✓ |
| SHIFT | ✗ | ✓ | ✓ | ✓ |
| PHYSICS | ✗ | (✓) | ✓ | ✓ |
| PERSON | ✗ | ? | ? | ? |
| LEARN | ✓ | (✓) | ✗ | |

□ ALGORITHMIC CAUSAL MODELS \rightarrow KOLMOGOROV COMPLEXITY

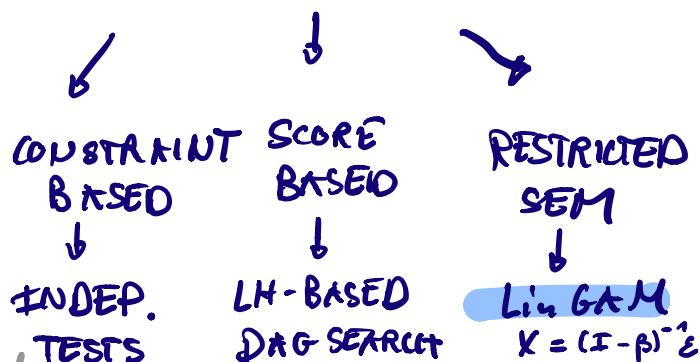
\hookrightarrow THERMODYNAMIC ARROW OF TIME / ENTROPY PRESERVATION

□ CAUSAL MODEL =
POSETS OF DISTRIB.

↳ CLEARLY STATE
ASSUMPTIONS!
CAN ALLOW FOR
CAUSAL CONCLUSION

⇒ DATA IS NOT ENOUGH! → NEED INDUCTIVE BIASES! ⇒ ASSUMPTIONS

□ CAUSAL STRUCTURE LEARNING [HEINZE - DEML ET AL. 17']



CAUSAL DIAGRAM ⇒ GRAPH!

GAIN !!!
ASSUMPTIONS

□ THEOREM: IDENTIFIABILITY

OF NON-GAUSSIAN
LINEAR MODELS

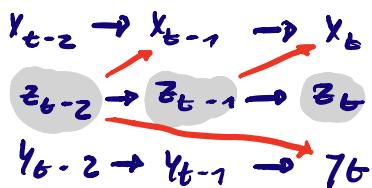
↳ RELATION TO IGT
⇒ INDEP. COMPONENTS!



⇒ OPEN REQUIRES LARGE SEARCHES OVER DIFF. GRAPHS!

□ TIME SERIES → GRANGER "CAUSALITY" ⇒ PREDICTIVE POWER

↳ IF $y_{\text{present}} \perp\!\!\!\perp x_{\text{past}} | y_{\text{past}}$ THERE MUST BE ARROW FROM x TO y



$y_{\text{present}} \perp\!\!\!\perp x_{\text{past}} | y_{\text{past}}$
 $x_{\text{present}} \perp\!\!\!\perp y_{\text{past}} | x_{\text{past}}$

⇒ CONFOUNDED
PROBLEM!

Problem: GRANGER INFERS $x \rightarrow y$

□ ODE INVARIANCE ACROSS ENVIRONMENTS → FORM OF CROSS-VAL.

↳ PRISTER ET AL. ⇒ STABILIZED REGRESSION → MODEL AVERAGING HELPS WHEN FUNCTIONAL FORM OF ODE NOT KNOWN / REJECTED

⇒ KEY PROBLEM: CAUSAL VARS ARE ASSUMED TO BE FULLY OBSERVED

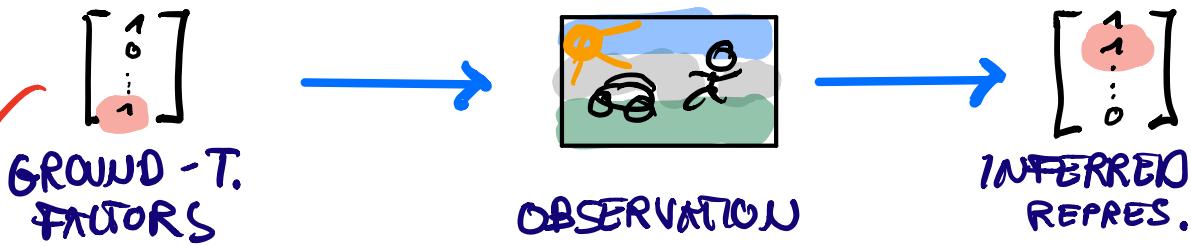
↳ 'WHAT ARE CAUSAL VARIABLES OF AN IMAGE?')

↳ DEEP REPRESENTATION LEARNING OF INDEP. COMPONENTS!

↳ GOAL: TRANSFER OF MECHANISMS ⇒ OUT-OF-SAMPLE!

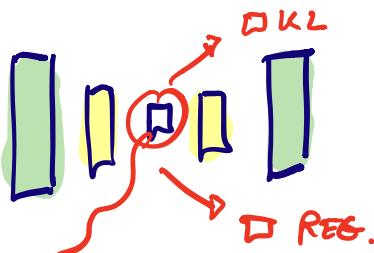
D CAUSAL PERSPECTIVE ON DEEP REPRESENTATION LEARNING

→ BENOU, COURVILLE, VINCENT 12' → MULTI-TASK MOTIVATION



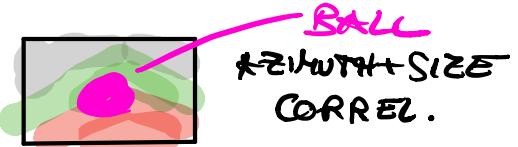
→ DISENTANGLEMENT: SINCE CHANGE FACTOR
⇒ SINGLE CHANGE REPRES. → ENCOURAGE FACTORIZATION

→ DISENTANGLEMENT VIA VAE + REGULARIZER



RESULT: FOR ARBITRARY UNSUPERVISED LEARNING OF DISENTANGLING REPRESENTATIONS IMPOSSIBLE!
↳ NEED FOR INDUCTIVE BIASES!
[FOR WEAK SUPERVISION]

REPRES. ⇒ SPLIT SOURCES OF VAR.!

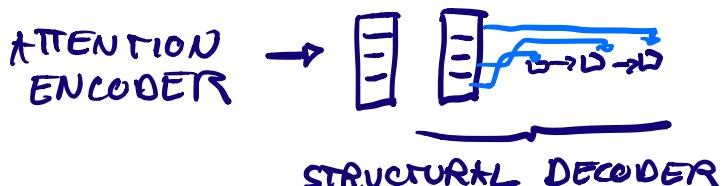


→ SCALING TO REAL-WORLD DATA?

↳ TRÜMBLE ET AL 20' → INDEP. OF MECHANISMS VS. VARS

↳ DISENTANGLEMENT ← DOWNSTREAM PERFORMANCE

⇒ ENCODE CAUSAL STRUCTURE INTO AUTOENC. → LEEB ET AL. 20'



→ LEARN TO DISTINGUISH AND ORDER MECHANISMS
↳ STRUCTURE ARCHITECTURE VS. GAUSSIAN REGULARIZER

D TOWARDS CAUSAL WORLD MODEL

→ LEARNING INDEP. MECHANISMS



LEARN DIFF.
DE-NOISING
MECHANISM

CAUS. → DISCRIMINATOR

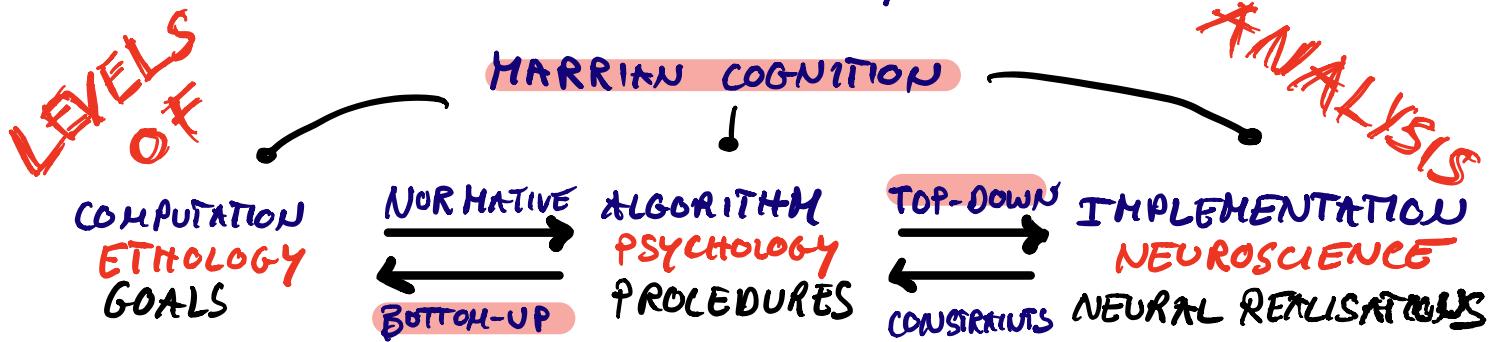
→ INDEP. COMPONENTS ⇒ SOCIAL INTELLIGENCE

→ RECURRENT IKS
↳ GOYAL ET AL. 18'



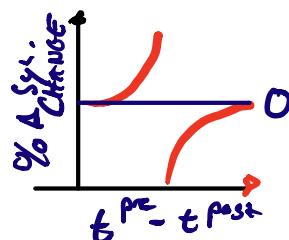
ATTENTION MODULE
ON TOP OF LSTM

Peter Dayan (MPI): Comp. Neuro in KL



SPIKE-TIME DEP. PLASTICITY

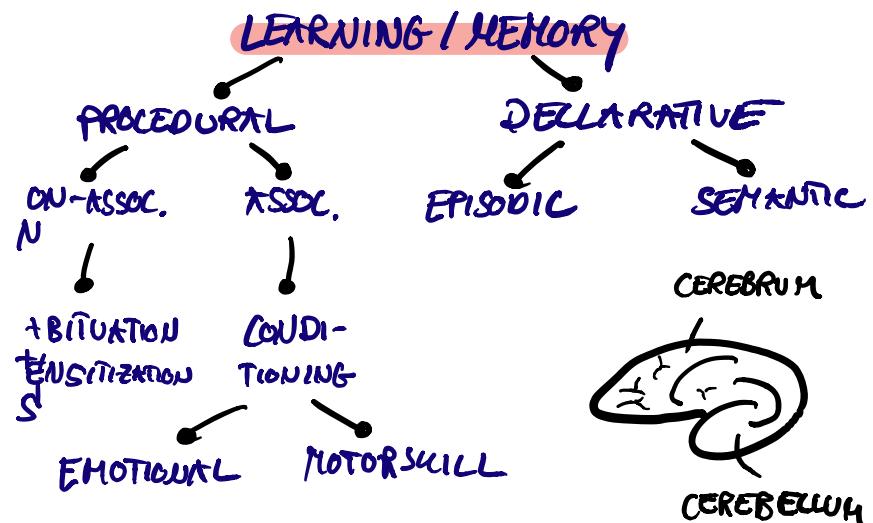
[STOP] → H. MARKRAM



↳ HEBBIAN INTUITION
STRENGTHEN / WEAKEN BASED ON TIMING!

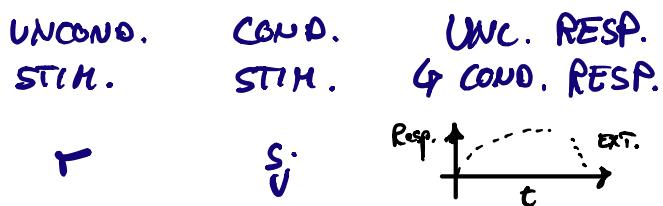
↳ "COINCIDENCE DETECTION"

⇒ MODULATION: DOPAMINE + ACETYL CHOLINE
↳ LOCAL VS. GLOBAL



⇒ MISSING: * LEARNING OF PREDICTIONS
* FORWARD / INVERSE MODELS
* REPRESENTATION LEARNING

PAVLOVIAN CONDITIONING



↳ BLOCKING: LEARNING OF ADD. STIMULUS CAN BE BLOCKED
⇒ CONTINENCY IS NOT ENOUGH
↳ NEED SURPRISE TO LEARN!

↳ 2ND ORDER: TIMING IS CRUCIAL
⇒ LEARN PREDICTOR OF PREDICTOR

$$V_t = \mathbb{E}_\pi \left[\sum_{t'=t}^T r_{t'} \right]$$

↳ EXPEC. FUTURE REWARD SEQ.

RESCORLA - WAGNER LEARNING

$$V_{ut+1} = \alpha \sum_{v=1}^u (1-\alpha)^{u-v} r_v + (1-\alpha)^u V_t$$

CASE FOR: $s_j \in \text{Eq. 13}$

↳ LRATE → BUT ALSO: FORGET RATE?
⇒ SPEED OF ENVIRONMENT CHANGE

↳ CHOICE VIA PROPENSITIES:

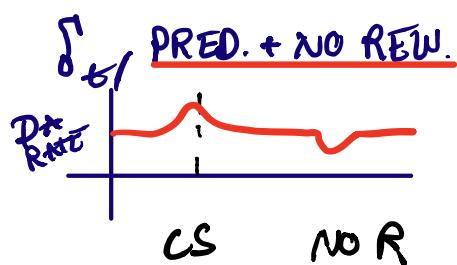
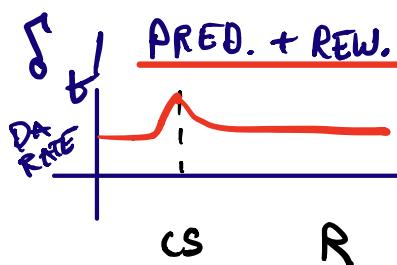
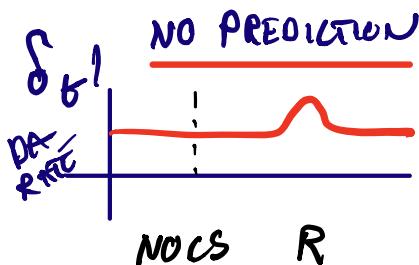
- 1 $\pi(s, a) = Q(s, a)$
- + SOFTMAX POLICY w. β TEMP.
- 2 POLICY GRAD. / "DIRECT" ACTOR
- ↳ REINFORCE (Williams)

$$\frac{\partial \ln \pi(s, a)}{\partial \theta} \cdot [r(s, a) - V(s)]$$

→ KEY: BELLMAN CONSISTENCY \Rightarrow USE INCONSISTENCY
ERROR FOR LEARNING 1 → TEMP. DIFFERENCE ERROR

$$V_t \leftarrow V_t + \alpha (\tau_t + V_{t+1} - V_t) \stackrel{\delta_t}{\rightarrow} \text{TD}(0)$$

\Rightarrow DOPAMINE REWARD PRED. ERROR HYPOTHESIS



↳ MODEL-BASED (PLANNING) VS. MODEL-FREE RL

ACTOR-CRITIC METHODS:

LEARN PROPENSITIES + VALUES SIMULTANEOUSLY!

↳ $\delta > 0$: - V IS TOO PESSIMISTIC
- a IS BETTER THAN AVG.

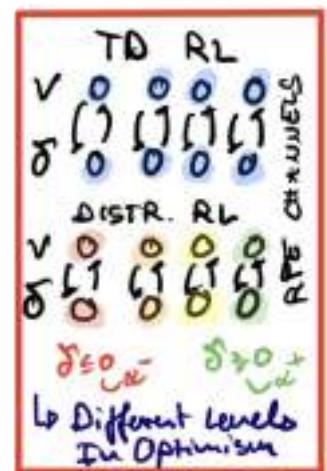
↳ SAME SIGNAL \Rightarrow VENTRAL
→ DORSAL STRIATUM!

↳ NEURO-WORK BY Yael Niv

DABNEY ET AL. 20' - DISTR. TD

=> DISTRIBUTIONAL
→ NEURON-SPEC α_i !

↳ ENCODE RISK SENSITIVITY /
OPTIMISM!
→ CAPTURES DATA NEURON HETEROGEN.



DIFFERENT SORTS OF LEARNING



LECUN CAKE



DIFFERENT ALGORITHMS



Q&t:

TABULAR VS. FCT. APPROX.

↳ Q OF STATE REPRES.

Q GENERALISATION

DIFFERENT NEURAL SUBSTRATES



Costis Daskalakis (MIT): Game Theory & ML

UTILITIES \leftrightarrow LEARN + DECIDE



\Rightarrow STRATEGIC REASONING

\square SYMMETRY VS. ASYMMETRY
 \Rightarrow OPTIMIZATION NOT ENOUGH!

MINIMIZATION \rightarrow CURRENT AI

MIN-MAX OPTIM. \rightarrow FUTURE AI

$$\inf_{\theta} \sup_w f(\theta, w) \xrightarrow{\text{HIGH-DIM. + CONSTRAINT}}$$

BEST CASE: f CONVEX IN θ , CONCAVE IN w

GANS [GOODFELLOW ET AL. 14', ARJOSKY ET AL. 17']

GENERATOR \leftrightarrow DISCRIMINATOR

$$\min_{\theta} \max_w \mathbb{E}_{x \sim \text{Real}} [D_\theta(x)] - \mathbb{E}_{z \sim \text{HAWLICKED}} [D_\theta(G_\theta(z))]$$

\rightarrow TRAINING OSCILLATIONS [1ST ORDER]

\rightarrow HIGH-DIM. STATISTICS PROBLEM

\rightarrow COMMON: MORE FREQ. UPDATES FOR GEN.

TRAINING OSCILLATIONS

\square METZ ET AL. 17' \Rightarrow NODE COLLAPSE

\hookrightarrow SGD CYCLES BETWEEN

\square DASKALAKIS ET AL. 18'

\hookrightarrow EVEN SIMPLER 2D GAN'S!

\hookrightarrow PURE OPTIMIZ. PROBLEM

\Rightarrow EVEN FOR CONVEX-CONCAVE OBJECTIVES!

\Rightarrow EVEN WHEN FCT. IS PERFECTLY KNOWN

\Rightarrow EVEN WHEN THERE ARE NO STATES

\square PLATFORM DESIGN

\hookrightarrow OPTIM. NOT ENOUGH \rightarrow NEED BEHAV. MODEL \rightarrow COUNTERFACTUAL

\square PHYSICAL RECOMMENDER SYSTEM

\hookrightarrow COMBINE LEARNING + MARKET DESIGN

\Rightarrow STATIONARY WORLD

\Rightarrow NON-STAT. WORLD

PART I: CONVEX-CONCAVE OBJECTIVES

\square MIN-MAX TH.: IF $X \subset \mathbb{R}^n, Y \subset \mathbb{R}^m$ ARE COMPACT & CONVEX, AND $f: X \times Y \rightarrow \mathbb{R}$ IS CONVEX-CONCAVE, THEN

$$\min_{x \in X} \max_{y \in Y} f(x, y) = \max_{y \in Y} \min_{x \in X} f(x, y)$$

\hookrightarrow OPTIMAL POINT IS UNIQUE [IF STRICT] + VALUE IS ALWAYS UNIQUE

POINTS = EQUILIBRIA OF ZERO-SUM GAME

\square BROWN 49' \rightarrow FICTITIOUS PLAY!

\square NO-REGRET-LEARNING

$$\frac{1}{T} \sum_t \ell_t(\pi_t) - \frac{1}{T} \min \sum_t \ell_t(.) \rightarrow 0$$

\hookrightarrow FOLLOW-THE-REGULARIZED LEADER

$$\pi_t \in \arg \min \left[\sum_{l=1}^T \ell_l(\cdot) + \frac{1}{n} \| R(\cdot) \|_2 \right]$$

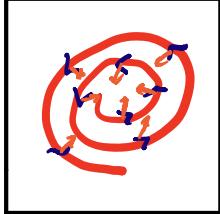
\hookrightarrow GAME VERSION X, Y PLAYERS

$$f(\cdot, \pi_t) \quad -f(x_{t-1})$$

\Rightarrow NO-REGRET LEARNING \rightarrow TWO-STRATEGY APPROX. NASH

REMOVING OSCILLATIONS VIA NEGATIVE MOMENTUM

- OPTIMISTIC GD [POPOV 80'] : $x_{t+1} = x_t - \gamma \nabla f(x_t) + \frac{\gamma}{2} \nabla f(x_{t-1})$
 - FTRL: UNIFORM POLICY HISTORY \rightarrow UPWEIGHT MOST IMMEDIATE HISTORY
 - EXTRA-GRAD. [KORPELEVICH 76'] : $x_{t+\frac{1}{2}} = x_t - \gamma \nabla f(x_t)$
 $f(x, y) = x \cdot y$
 $x_{t+1} = x_t - \gamma \nabla f(x_{t+\frac{1}{2}})$



NEG. MOM.
"GRAVITY PULL"

\Rightarrow SPINNING OUT
- \hookrightarrow OGA: LAST-ITERATE CONV. + LINEAR RATES FOR UNCONSTRAINED BILINEAR GAMES W. WELL. COND. PAYOFF MATRIX
- \hookrightarrow EG: SAME FOR UNCONSTRAINED FOR STRONG CASE!
- CONSTRAINED CASE \Rightarrow PROJECTED VERSIONS!

PART II: NONCONVEX - NONCONCAVE OBJECTIVES

- DASKALAKIS ET AL. 18': OPTIMISTIC ADAM \rightarrow CLEAR TO GAN
- \hookrightarrow DOES NOT NEED DIFFERENT # STEPS GENERATOR VS. DISCRIMINATOR?
- $f(x, y)$ NOT CONVEX-CONCAVE \Rightarrow MIN-MAX THEOREM BREAKS! \Rightarrow GAN?
- \hookrightarrow MORE GENERAL CASE: [GENERAL CONSTRAINED MIN-MAX PROBLEM]

GIVEN CONTINUOUS f & CONVEX, COMPACT SETS:

$$\min_x \max_y f(x, y) \text{ s.t. } (x, y) \in S$$

$$\text{DEF.: } S(x, \cdot) = \{y | (x, y) \in S\}$$

$$S_x = \{x | \exists y \text{ s.t. } (x, y) \in S\}$$

① LOCAL MIN - GLOBAL MAX SOLUTION

$(x^*, y^*) \in S$ s.t., for some $\delta > 0$,

$$f(x^*, y) \leq f(x^*, y^*) \leq \max_{y' \in S(x^*)} f(x^*, y')$$

$\forall y \in S(x^*)$ and $x \in N_\delta(x^*) \cap S_x$

\hookrightarrow ASYMMETRY INNER/OUTER PLAYER

① + ② : SEQ. MOVES vs. ③ : SIMULT.

①: GUARANTEED TO EXIST \rightarrow BUT NP-HARD TO FIND vs. ② + ③: NOT ALWAYS EXISTING

② LOCAL MINIMIX

\hookrightarrow Y PLAYER IS ALSO CONSTRAINED TO MOVE LOCALLY

③ LOCAL MIN-MAX EQ.

$(x^*, y^*) \in S$ s.t., for some $\delta > 0$,

$$f(x^*, y) \leq f(x^*, y^*) \leq f(x, y^*)$$

$\forall y \in N_\delta(y^*) \cap S(x^*)$ & $x \in N_\delta(x^*) \cap S(\cdot, y^*)$

STABILITY?

FIXED POINTS
OF GDA

LOCAL MIN-MAX
EQUILIBRIUM

LOCAL MINIMAX
SOLUTION

NON-EMPTY, IF f IS
CONT. DIFFERENTIABLE
 \Rightarrow BROUWER'S FP TH.

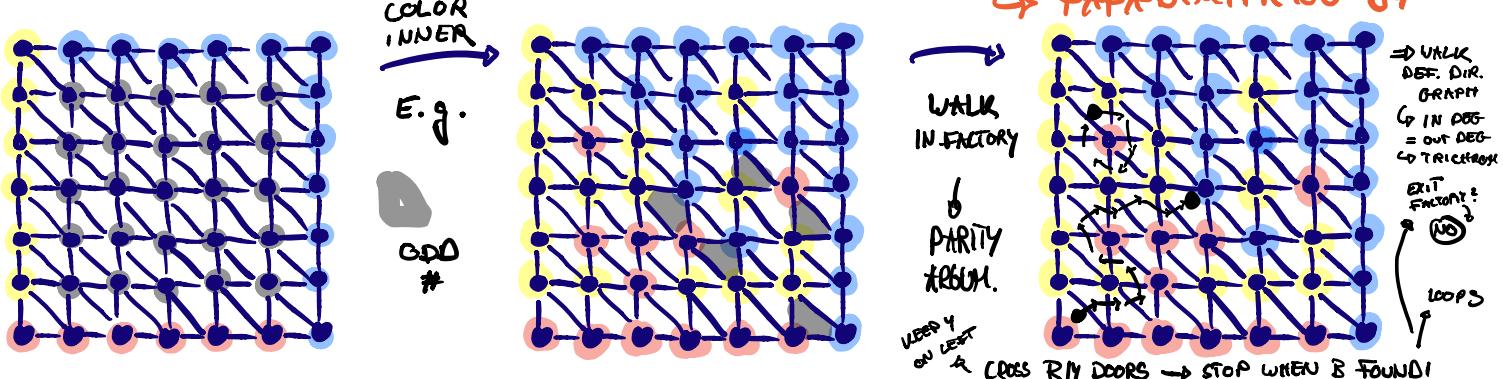
CAN BE
EMPTY

■ GOAL: FIND APPROX. FP: $\|(\mathbf{x}^*, \mathbf{y}^*) - F_{\text{GDA}}(\mathbf{x}^*, \mathbf{y}^*)\|_2 \leq \varepsilon$

\hookrightarrow APPROX. BROUWER FP \Rightarrow PPAD COMPLETE \rightarrow BETWEEN P & PPAD

PPAD: THE CLASS OF ALL PROBLEMS IN NP W. GUARANTEED SOLUTION
BY GRAPH-THEORETIC LEMMA: IF \exists UNBALANCED NODE, THEN \exists KNOWNER

\hookrightarrow PAPADIMITRIOU '94!



SPERNER'S LEMMA: NO MATTER HOW THE INTERNAL NODES ARE COLORED
THERE EXISTS A TRI-CROMATIC TRIANGLE. \rightarrow AN ODD # OF THEM.

\hookrightarrow BOILS DOWN TO DIRECTED PARITY ARGUMENT \rightarrow LIKE PPAD?

\hookrightarrow WHAT MAKES PROBLEM HARD?

\Rightarrow COLORING CIRCUIT MAY BE SMALL BUT IMPLIED GRAPH EXP. LARGE!

\hookrightarrow CAN BE USED TO PROOF BROUWER!

DSA-KALAKIS ET AL. '20': COMPUTING LOCAL MIN-MAX EQ. IN ZERO-SUM GAMES W.
NONCONVEX-NONCONCAVE OBS. \Rightarrow AS HARD AS NASH-EQ. IN GENERAL SUM GAMES

■ MULTI-AGENT RL \Rightarrow STRUCTURE ALLOWS NASH IN NONCONVEX-NONCONCAVE CASE

$$\max_{\pi_i} V_{\pi_i, \pi_{-i}}^i = \mathbb{E} \left[\sum_{t=0}^T R^i(s_t, a_t, s_{t+1}) \right]$$

\Rightarrow NASH EQ.: $\pi^{1, \dots, N}$ s.t. $\forall i: V_{\pi_i, \pi_{-i}}^i \geq V_{\pi_i^*, \pi_{-i}}^i$

\hookrightarrow LOTS OF OPEN QUESTIONS!

\Rightarrow DECENTRALIZED CONTROL
 \Rightarrow TOWARDS GENERAL SUM!

SHAPLEY
SS'

