

# EML Summer School - Day 1

07/07/19

## Intro to deep learning - Razvan Riscan (DeepMind)

- SGD intuition from Taylor expansion (1st order):

$$\underset{\Delta \theta}{\text{argmin}} L(\theta + \Delta \theta) = \underset{\Delta \theta}{\text{argmin}} \left[ L(\theta) + \Delta \theta \frac{\partial L}{\partial \theta} \right] \text{ s.t. } \|\Delta \theta\| \leq \epsilon$$

$$\Leftrightarrow \underset{\Delta \theta}{\text{argmin}} L(\theta) + \Delta \theta \frac{\partial L}{\partial \theta} + \frac{1}{2} (\Delta \theta)^T \nabla^2 L(\theta) \Delta \theta \quad \begin{array}{l} \xrightarrow{\text{Lagrange H.}} \\ \text{constr.} \\ \text{(step of learning rate!)} \end{array}$$

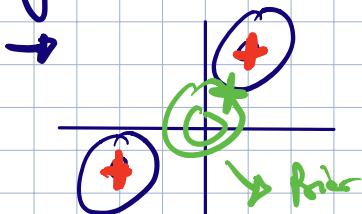
- Black-Box via Evolutionary algs  $\Rightarrow$  Mutation  $\oplus$  Next Generation
  - $\rightarrow$  allows for potentially global optimizers  $\rightarrow$  problem of many peaks
  - $\rightarrow$  More useful for hyperparams
- Population-based Theory  $\Rightarrow$  Train models in parallel / stop + evaluate afterwards only first to develop best performing model
  - $\rightarrow$  Between parametric and non-parametric / optimizable!
  - $\rightarrow$  slow learning
  - $\rightarrow$  Comparison / fast learning
  - $\rightarrow$  Non-parametric
  - $\rightarrow$  Hyper-Parametric
  - $\rightarrow$  bad early / slow tree
  - $\rightarrow$  fast learning
- $h: X \rightarrow Y \Rightarrow$  Predictor  $\Leftrightarrow$  Hypothesis  
(Classification / Regression / Structured Pred. (graph))
  - $\hookrightarrow h: \Theta \times X \rightarrow Y \Rightarrow$  fit  $\Theta$  via Loss fct.
  - $\rightarrow$  ML / MLE perspective!

$$\underset{\Theta}{\text{argmax}} P(\Theta | D) = P(D | \Theta) P(\Theta) / P(D) \rightarrow \underset{\Theta}{\text{argmax}} P(D | \Theta) P(\Theta)$$

$\rightarrow$  Uniform prior in MLE formulation results in Max. Likelihood!  
together with i.i.d assumption  $\rightarrow$  split data into  $\prod_{i=1}^n$   $\oplus$  log fully

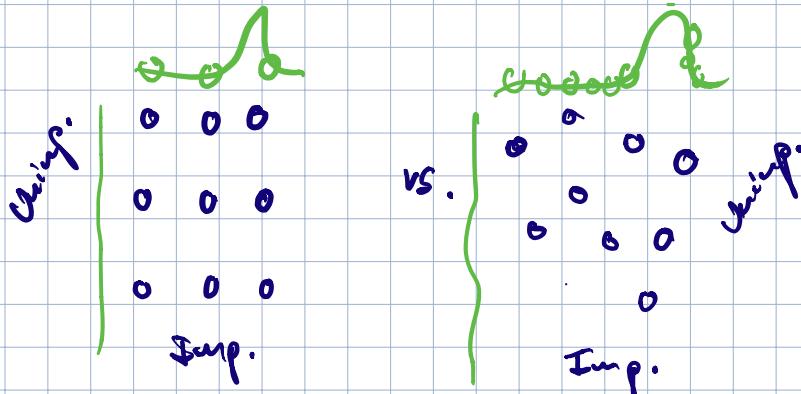
- Loss fct often times results from Bayesian / MLE perspective  $\oplus$  assumptions on the data  $\rightarrow$  E.g. MSE from MLE + Uniform Prior + Gaussian Data!
- $\rightarrow$  Multi-label classification: Multinomial  $\Leftrightarrow$  log. loss LLL.

- Regularize w/o not assuming  $P(\Theta)$  being uniform! [Together w. Optimizer.]



$\Rightarrow$  Reg. gets rid of local min problem + by preferring one that is closer to 0  $\rightarrow$  +

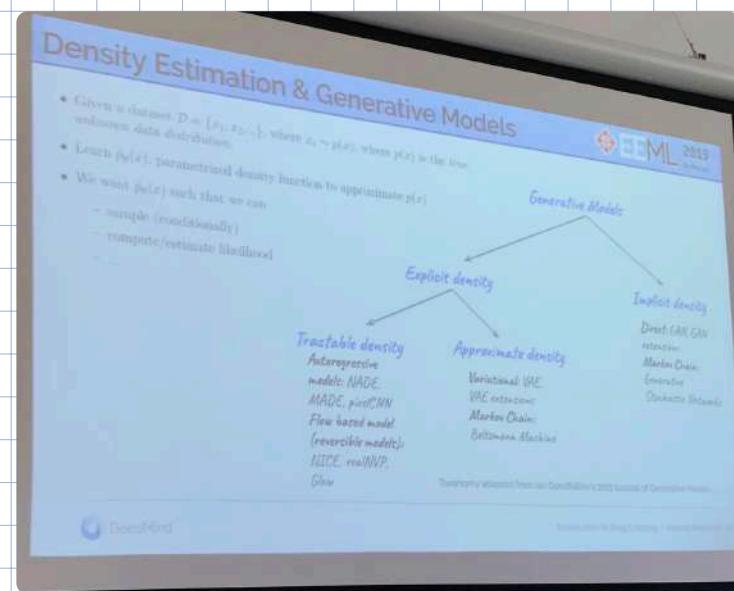
- Hyperparameter Optimization: Random vs. Grid  $\rightarrow$  Random covers more space!



↳ Be careful not to output next cluster!

↳ Unclear what global level complete is!

- Early-stopping as a prior  $\rightarrow$  Form of  $L_2$   $\Rightarrow$  See Bishop! that's why?
- ↳ Duvenaud et al 2016 - Early Stopping as Hyperprior.



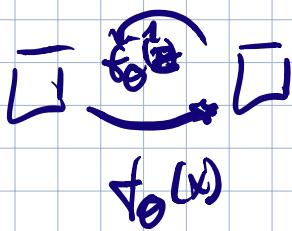
$\rightarrow$  Autoregressive: Decompose joint as product of conditionals

↳ Not necessary clear how to do so if there is no temporal chn.

E.g. when to start sampling lange

↳ also: slow Sampling!

$\rightarrow$  Flow-based / Reversible Models: Find  $f_\theta^{-1}(z)$



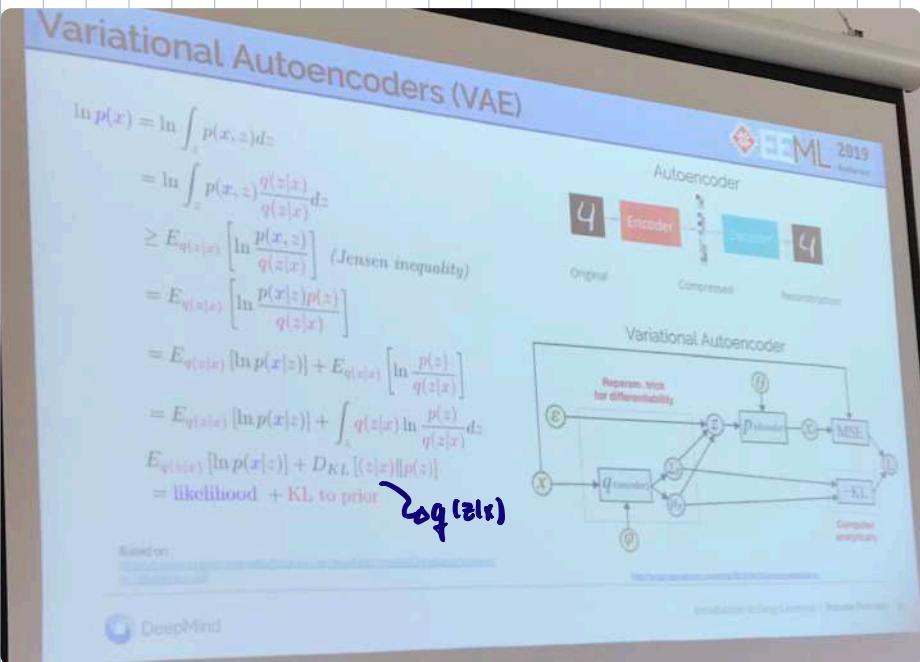
$\rightarrow$  Change of variables  $\rightarrow$  Gaussian in latent space to optimize easily and fast

$$f_x(x) = p_z(z) \left| \det \left( \frac{\partial x}{\partial z} \right) \right|^{-1}$$

↳ Not clear how to parameterize  $f_\theta$  s.t. it is invertible!

↳ Easily sample based on sampling Gaussian  $z$  and then apply  $f_\theta^{-1}(z)$

$\rightarrow$  VAE  $\Rightarrow$  Restrict latent codes by having encoder output Bernoulli parametrization  $\rightarrow$  Optimize both Reconstruction as well as KL objective!  $\rightarrow$  Restriction allows for efficient sampling



↳ Gaussian cores on lucy  
for KL computation  
as well as sampling

↳ Gaussian + Beta //

also new one on Gumbel with  
all multivariates

↳ problem to get sleep over  
days!

↳ want VAE for latent interpolation!

→ compressive models vs. GANs  $\Rightarrow$  VAEs GANs better, didn't DR

• ReLU activations: Splits space into linear regions  $\rightarrow$  can backprop

$\hookrightarrow$  Montufar et al 2014  $\rightarrow$  NIPS

↳ form of optimal efficient fully of the space!

$\hookrightarrow$  optimizable such that decision boundary becomes simpler!

↳ No a<sup>ffine</sup> some regions are constrained to else regions!

$\Rightarrow$  Need high-dimensional data  $\rightarrow$  unnecessarily bad

$\hookrightarrow$  1d mixture of sinusoids will be fit!

but a lot

more about  
non-affine regions!

• Loss Opt. Surfaces: Not necessarily crazy!

$\hookrightarrow$  Yann LeCun: Simple random forest of capture pods.  
of all dir. pointing up  $\rightarrow$  very nicely

$\hookrightarrow$  sleep vs. flat minimum!  $\rightarrow$  what do they mean? algorithm

$\hookrightarrow$  Can switch eigenvalue profile without sleep problem!

• RMSprop  $\Rightarrow$  found by Muen/Ves.  $\rightarrow$  larger step if small var

$\hookrightarrow$  approx. of Natural Gradient

$\hookrightarrow$  Batch mean  $\rightarrow$  more accurate  
closer to identity  $\rightarrow$  Fisher  $\approx$  NG!

- Convolutions = restrictions of net  $\rightarrow$  need more trainable biases!
- Graph Nets = do conv over locally defined pixels but a other local structure
- RNN  $\rightarrow$  Exploding gradient from product of Jacobians  $\rightarrow$  Gradient Clipping
  - $\hookrightarrow$  LSTM  $\rightarrow$  vanishing gradient combat via linear cell state  
But need to learn gates
  - $\hookrightarrow$  Interpret gates more as a form of low pass filters  $\xrightarrow{\text{gates}} \text{at } 1$ !

## Intro to RL - David Silver (DeepMind / McGill) $\rightarrow$ TD

- How to overcome sample efficiency problem with sparse rewards
  - $\rightarrow$  Self-play or very cooperative / intrinsic motivation  $\rightarrow$  surrogate reward
- Key features:
  - \* Trial-and-Error Search
  - \* Stochastic Env
  - \* Delayed reward  $\rightarrow$  Temp. Credit Assignment
  - \* E-E Trade-off
- Learning with multiple output heads / auxiliary heads
  - $\hookrightarrow$  Relationship to ECE / Generalizability?
- TD Gamma  $\rightarrow$  (Tsetsas et al., 1995)  $\Rightarrow$  "Early Alpha-Go"
- Function learning / function approximation  $\rightarrow$  form of poaching!
  - $\rightarrow$  important to train in tandem with explorable vs robustness!
- Different of interpretations: Survival prob., bias-variance tradeoff, confidence
- MC methods  $\approx$  Supervised  $\rightarrow$   $V(S_t)_{\text{MC}} \leftarrow V(S_t)_{\text{MC}} + \alpha [G_t - V(S_t)_{\text{MC}}]$ 
  - $\circlearrowleft$   $\rightarrow$  Update for each transition  $\Rightarrow$  not iid! Sequentially correlated
    - $\hookrightarrow$  Observe via ER buffer for example!
- Bias-Variance Trade-Offs:
  - \* Function approx.  $\rightarrow$  Gradient Updates
  - \* Generalization across environments
  - \* Disent  $\rightarrow$  Epsilon-like off pure rewards

$\hookrightarrow$  schedule of disent parents  $\rightarrow$  relate to TD as in PER?

• Bellman Self-Consistency Equations: Set of linear equations!

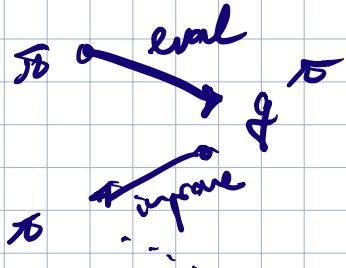
$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_{\pi}(s')]$$

- ↳ But  $v^*(s)$  is a non-linear set of equations due to optimality
- ↳ Can only be done if full MDP is known!  $\Rightarrow$  expensive dynamics
- ↳ Impose  $\gamma \leq 1$ : Contraction  $\Rightarrow$  Banach Fixed Point Theorem
- ↳ Problem also with large state spaces  $\Rightarrow$  storage is expensive!

• Contraction is what makes bootstrapping work!

↳ Would it make sense to start off with MC return backup and slowly move towards full bootstrap?  $k$ -step where  $k=1 \rightarrow k=1$

• Policy Improvement: Assumption of no local optima otherwise explorative needed



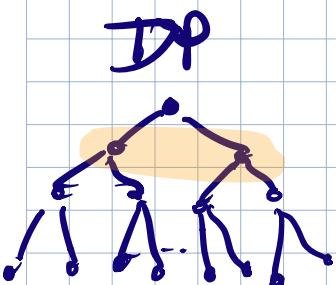
• Dyn.-fr.: Full-Wieliczka Backup

↳ TD: In between MC and DP

• Crucial assumption of stationarity for bootstrapping  $\Rightarrow$  STAS

↳ Introduce bias to reduce variance

↳ MC not biased!  $\Rightarrow$   $k$ -step TD interpolates!



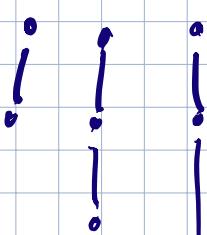
TD(1)



$k$ -step TD



TD(2)



• Epilepsy project: Use SDS  
↳ simulate

↳ likelihood-free methods!!

exp.  
weightly  
update!

# EML Summer School - Day 1

## Intro to RL II - Action Space → Control

Stochastic Gradient Descent (SGD) is the idea behind most approximate learning

General SGD:  $\theta \leftarrow \theta - \alpha \nabla_{\theta} \text{Error}^2$

For VFA:  $\leftarrow \theta - \alpha \nabla_{\theta} [\text{Target}_t - \hat{v}(S_t, \theta)]^2$

Chain rule:  $\leftarrow \theta - 2\alpha [\text{Target}_t - \hat{v}(S_t, \theta)] \nabla_{\theta} [\text{Target}_t - \hat{v}(S_t, \theta)]$

Semi-gradient:  $\leftarrow \theta + \alpha [\text{Target}_t - \hat{v}(S_t, \theta)] \nabla_{\theta} \hat{v}(S_t, \theta)$

Linear case:  $\leftarrow \theta + \alpha [\text{Target}_t - \hat{v}(S_t, \theta)] \phi(S_t)$

Action-value form:  $\theta \leftarrow \theta + \alpha [\text{Target}_t - \hat{q}(S_t, A_t, \theta)] \phi(S_t, A_t)$

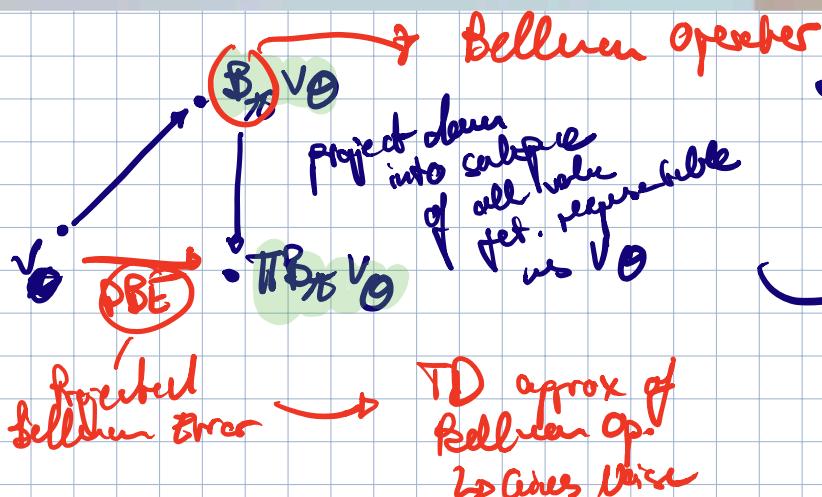
→ Transform into form of supervised learning

↳ assume that target does not depend on previous

↳ L2-Square loss problem

⇒ separable datasets

↳ Weight not be a problem if full MC



$$\bullet (\frac{\partial}{\partial t} v) S := \sum_a \pi(s, a) [r(s, a) + \gamma \sum_{s'} p(s'|s, a) v(s')]$$

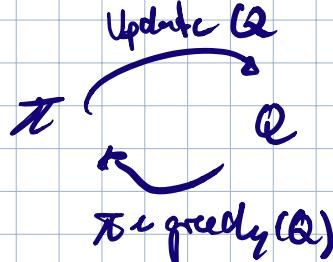
$\xrightarrow{\frac{\partial}{\partial t} v = 0}$   
done until  
 $PBE = 0$   
 $\Pi v$

↳ place shifts due policy shifts! → need to constrain that shift  
E.g. stochastic trust region! → clipping / KL

↳ Linear feature space:  $MSVE(\theta_{id}) \leq \frac{1}{1-\gamma} MSVE(\theta_{uc}) \rightarrow$  linear gives in-con-linear!

↳ n-step backups: Trade-Off bias-variance  $\Rightarrow$  Full MC has large variance!  
↳ f controls for bias

• Monte-Carlo Control:



⇒ Random exploring starts!

↳ on vs off-policy

↳ Exploration:  $\epsilon$ -greedy  $\rightarrow$  different  $\Rightarrow$  problem: indep. of action values  
↳ In practice: exploration schedules are hard to implement

- Targets:  $y_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}; \theta_t)$  [SARSA]

$$y_t = r_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a; \theta_t) \quad [\text{SARSA}]$$

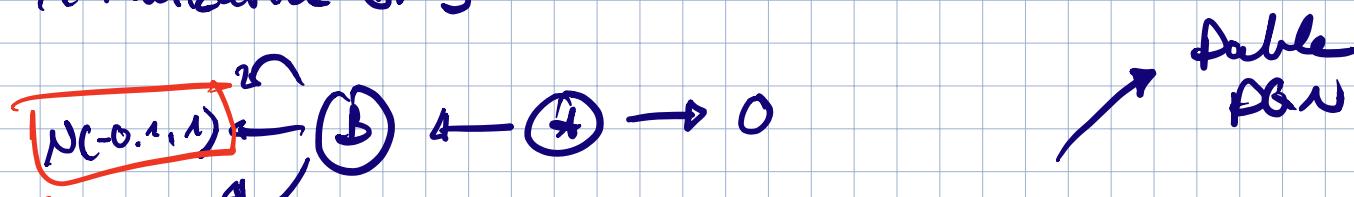
$$y_t = \sum_{s', r} \rho(s', r | s, a) \left[ r + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a) \right]$$

↳ SARSA: on-policy  $\rightarrow$  next action is chosen at  $t$  to update  
↳ still explore?

- Q-Learning: Bootstrapping  $\rightarrow$  does not necessarily converge from completely random start!

↳ SARSA vs. Q-L: Safety vs. Optimality

- Maximization Policy:



Covering radius?  
 ↳ Takes a long time to learn small variance  
 ↳ Use separate rewards  $\Rightarrow$  form of variable!  
 ↳ Stochastic policy or stochasticity in environment  
 ↳ Trade-off fine shift  $\rightarrow$  # iterations between target net prob  
 more iterations  $\Rightarrow$  bigger shift but smaller variance!

- Reinforcement Optimization Requirements  $\Rightarrow$  but maybe there is  
safety better!  $\rightarrow$  META RL!
- Need for different sets of optimizers in RL  $\rightarrow$  different objectives than  
optimization learning

- Policy gradient Methods:

$$\pi(a|s, \theta) \rightarrow \text{Objective: } \gamma(\theta) = y_{\theta}(s_0) \quad \begin{array}{l} \text{↳ initial state dist.} \\ \text{↳ learning of} \\ \text{policy} \end{array}$$

$$\text{↳ } \theta_{\theta+1} := \theta_t + \alpha \nabla_{\theta} \gamma(\theta_t)$$

$$\text{↳ } \nabla \gamma(\theta) = \sum_s d_{\theta}(s) \underbrace{\sum_a q_{\theta}(s, a)}_{\text{↳ baseline constraint}} \nabla_{\theta} \pi(a|s, \theta)$$

## Deriving REINFORCE from the PGT

$$\begin{aligned}
 \nabla \eta(\theta) &= \sum_s d_\pi(s) \sum_a q_\pi(s, a) \nabla_\theta \pi(a|s, \theta), \\
 &= \mathbb{E}_\pi \left[ \gamma^t \sum_a q_\pi(S_t, a) \nabla_\theta \pi(a|S_t, \theta) \right] \\
 &= \mathbb{E}_\pi \left[ \gamma^t \sum_a \pi(a|S_t, \theta) q_\pi(S_t, a) \frac{\nabla_\theta \pi(a|S_t, \theta)}{\pi(a|S_t, \theta)} \right] \\
 &= \mathbb{E}_\pi \left[ \gamma^t q_\pi(S_t, A_t) \frac{\nabla_\theta \pi(A_t|S_t, \theta)}{\pi(A_t|S_t, \theta)} \right] \quad (\text{replacing } a \text{ by the sample } A_t \sim \pi) \\
 &= \mathbb{E}_\pi \left[ \gamma^t G_t \frac{\nabla_\theta \pi(A_t|S_t, \theta)}{\pi(A_t|S_t, \theta)} \right] \quad (\text{because } \mathbb{E}_\pi[G_t|S_t, A_t] = q_\pi(S_t, A_t))
 \end{aligned}$$

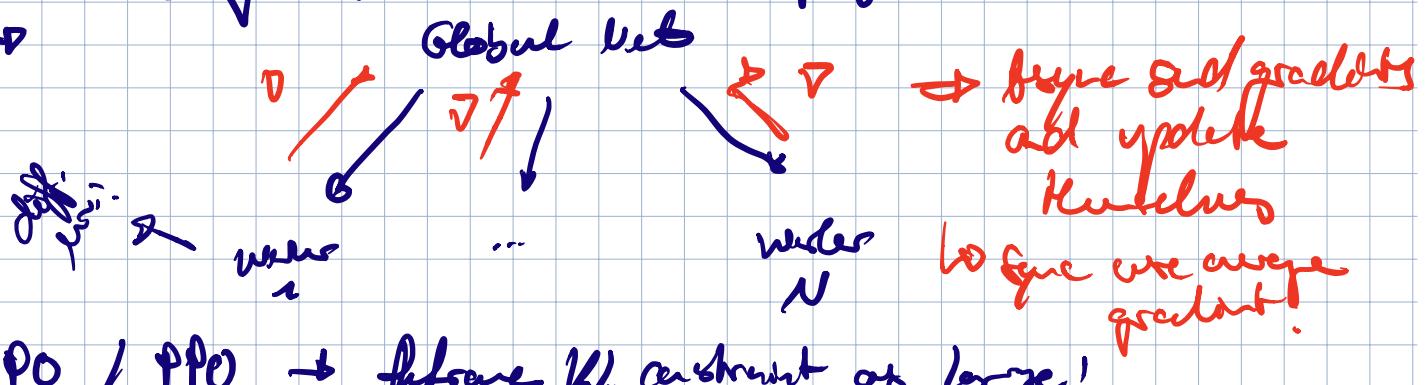
Thus

$$\theta_{t+1} \triangleq \theta_t + \alpha \widehat{\nabla \eta(\theta_t)} \triangleq \theta_t + \alpha \gamma^t G_t \frac{\nabla_\theta \pi(A_t|S_t, \theta)}{\pi(A_t|S_t, \theta)}$$

- No parametrization of value fn. needed  
↳ Use returns
- MC estimation  
↳ Avg. Value!  
Need to stabilize!
- AC  $\rightarrow$  parametrize value fn!

- \*3C  $\rightarrow$  harder implement by parallelization!  
↳ Pooling of explores  $\rightarrow$  Mean of gradients

↳



- TRPO / PPO  $\rightarrow$  before KL constraint of large!

# Multi-agent RL - Shimon Whiteson (Oxford)

- Motivation: Classical RL also treats many aspects  
however is full of multi-agent systems

## COOPERATIVE

- ↳ Shared reward
- ↳ Coordination

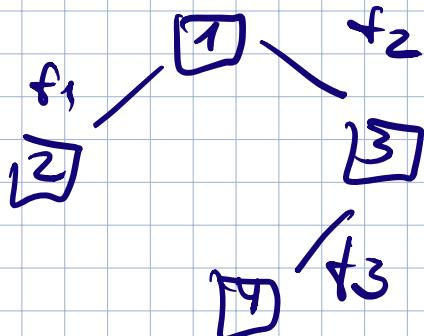
## COMPETITIVE

- ↳ Zero-Sum Game
- ↳ Opposing rewards
- ↳ Minimax

## MIXED

- ↳ General Sum
- ↳ Nash Eq.
- ↳ QZ. Econ

- Cooperative setting: absence of social conflict  $\Rightarrow$  risk of coordination
  - $\rightarrow$  Exponential joint-action space.  $\Rightarrow$  grows very fast
  - $\rightarrow$  Anstreicher et al. 2002  $\Rightarrow$  Coordination Graphs  $\xrightarrow{\text{+ add. eqpt 3}}$
$$Q(u) = f_1(u^1, u^2) + f_2(u^1, u^3) + \dots$$



$\hookrightarrow$  Probabilistic graphical model  
 $\Rightarrow$  Solutions: Most estimable!  
 ↳ Encoding of cond. independence  
 $\Rightarrow$  No sequentiality!

$\hookrightarrow$  Variable elimination:

$$\max_u Q(u) = \max_{u^1, u^2, u^3, u^4} [f_3(u^3, u^4) + \underbrace{\min_{u^1} [f_1(u^1, u^2) + f_2(u^1, u^3)]}_{f_2(u^2, u^3)}]$$

$\rightarrow$  first is best response  
 to get 2.5

$\hookrightarrow$  form of cooperative Nash!

$\hookrightarrow$  Computationally expensive if the were connected the graph

$\hookrightarrow$  Computes only optimal value of joint action  $\Rightarrow$  do not need to pass to actually desire that action!

$\hookrightarrow$  Need to know the graph / Cond. indeps. conditions

$\rightarrow$  Max-Plus  $\Rightarrow$  Vlassis et al. (2004)  $\rightarrow$  Message passing

$\hookrightarrow$  Two agent local payoff assumption

↳ Improving messages iteratively

↳ Converge guarantee in acyclic graphs

→ MDP (+ Symmetry!): Centralized Multi-agent MDP

↳ All agents see state  $\rightarrow$  Not really multi-agent: Simply large action/state space  $\rightarrow$  simply sub. optimization

→ Independent Q-Learning: No attempt to model  $Q(s,a)$

↳ each agent learns  $Q(s_a, a^a)$ !  $\Rightarrow$  last iteration of all agents learning  $\rightarrow$  no convergence guarantee!  $\rightarrow$  Tari et al (1993)

→ Coordinated Q-Learning: Gaudreault et al (2002)

$$Q^{tot}(s, a) = \sum_{c=1}^C Q_c(s^c, a^c)$$

↳ Subjects our agents / ad parties

↳ extremely strong assumption!  $\rightarrow$  Transitions might be coupled

↳ Update each peer  $\Rightarrow$  graph is shiftable!

• Partial observability  $\rightarrow$  requires decentralized execution!

↳ Learning should be centralized  $\rightarrow$  flex. power / creativity

↳ Minimal vs. right assumptions?

• Dec-POMDP

↳  $O(s, a)$ :  $S \times A \rightarrow Z$   $\rightarrow$  Action-Obs hist:  $\tau^a \in T \subseteq (Z \times U)^*$

↳ Decentralized policies:  $\pi^a(u^a | \tau^a)$ :  $T \times U \rightarrow [0, 1]$

↳ Dilemma: Exploit private info vs. being predictable  $\rightarrow$  switch fast

• Policy Gradient Methods for MTRL:

→ PGM useful if generalization is hard!

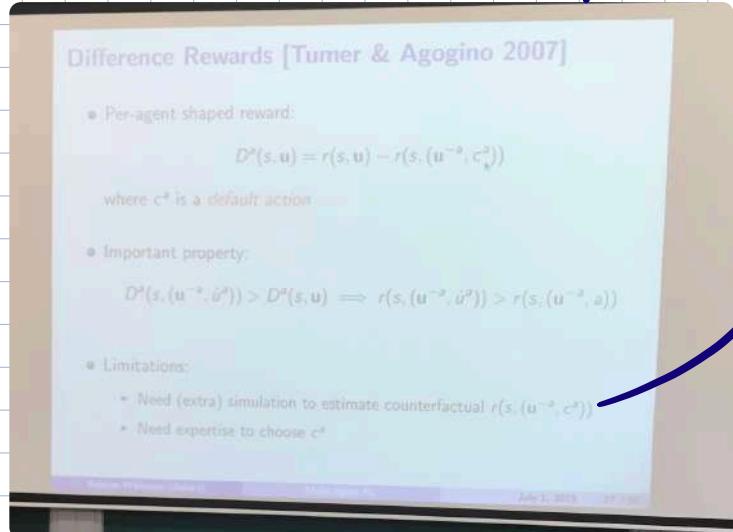
→ Not MTRL independent: Some parameters - still help. inputs + output later  
↳ sensitivity, hard to learn coordinate, multi-agent tasks vs.

- Counterfactual MDPs  $\rightarrow$  Foerster et al 2018 (COM\*)

① Centralize Critic

② Wolpert & Tumer (2000)  $\Rightarrow$  Weakform Life Utility

$\hookrightarrow$  reason about counterfactual reward if you would not have participated



$\hookrightarrow$  need to go back into simulation!  
 $\hookrightarrow$  Counterfactual Baseline

$$A^a(s, u) = Q(s, u) - \sum_{u'} D^a(u', \pi^a)$$

$$Q(s, (u^{-a}, u^a))$$

$\hookrightarrow$  problem: as  $\pi^a$  becomes deterministic  
 we can't get the counterfactual aspect

$\hookrightarrow$  critic: output kind gives all values

$\hookrightarrow$  still monotonic w.r.t. value function.

- Value Decomposition Methods - Guehag et al 2017

$$\rightarrow \text{per agent: } Q_{\text{tot}}(\pi, u) = \sum_{a=1}^N Q_a(\pi^a, u^a; \theta^a)$$

$\hookrightarrow$  decentralizes the agents  $\Rightarrow$  no layer level greedification

- QMIX: not only summable but nicely rebalanced

$\hookrightarrow$  constraint to non-negative weights!

$\hookrightarrow$  can throw away off-counterfactual terms!

SMAC framework  $\rightarrow$  StarCraft baseline!

- Competitive setting: for cooperative setting  $\Rightarrow$  one optimal policy exists  $\rightarrow$  how: stochasticity

$\rightarrow$  Mixed strategies  $\leftrightarrow$  Monotone Measures

$\hookrightarrow$  Rational vs. Strategic agents  $\rightarrow$  symmetries of others not matter

$\hookrightarrow$  lower player does not have to believe stochastically!

→ Minimax Q-Learning → Littman (1994)

Minimax Q-Learning [Littman 94]

- Update rule:
$$Q(s_t, u_t) \leftarrow Q(s_t, u_t) + \alpha[r_t + \gamma \text{MM}(s_{t+1}) - Q(s_t, u_t)]$$

where:

$$\text{MM}(s) = \max_{\pi^1 \in \Pi^1} \min_{u_1} \sum_{u_2} \pi^1(u^1) Q(s, u)$$

- Treats  $Q(s, u)$  as payoffs in matrix game for  $s$
- Hence, each player can use minimax to select actions:

$$\pi^1(s, \cdot) = \arg \max_{\pi^1 \in \Pi^1} \min_{u_1} \sum_{u_2} \pi^1(u^1) Q(s, u)$$

$$\pi^2(s, \cdot) = \arg \min_{\pi^2 \in \Pi^2} \max_{u_2} \sum_{u_1} \pi^2(u^2) Q(s, u)$$

→ Tricky approach: Learns about the opponent from data to find Q-Learning → Implicit

↳ Explicit: Fictitious play framework

→ Opponent Modeling

→ Self-Play: Checkers → Alter Served (1996)

↳ natural combinator → can fail if game is not transitive

↳ leave yourself beatable to beat opponents while beating stronger ones ⇒ triple star → League/Tournament!

• Fixed: Nash Equilibrium → Fixbee!

↳ Nash Equilibrium Q-Learning → Hu & Wellha (1998)

→ almost never guaranteed to converge!

↳ "If multi-agent learning is the answer, what is the question?"

↳ Meelissen design! → Fixbees

↳ Friend-or-Foe Q-Learning