

Intro to deep learning - Razvan Pascanu (DeepMind)

SGD intuitions from Taylor expansion (1st order):

argmin_{Δθ} L(θ + Δθ) ≈ argmin_{Δθ} [L(θ) + Δθ $\frac{\partial L}{\partial \theta}$] s.t. ||Δθ|| ≤ ε

⇔ argmin_{Δθ} L(θ) + Δθ $\frac{\partial L}{\partial \theta}$ + $\frac{1}{2} (\Delta\theta)^T I (\Delta\theta)$

Lagrange H. multiplier
↳ form of loss req!

Black-Box via Evolutionary algos ⇒ Multi-modal Next Generation
 → allows for potentially global optimizers ⇒ problem of many peaks
 → More useful for hypothesis

Population-based Training ⇒ Train models in parallel / step + evaluate afterwards only first develop best performing model
 → Between parametric and evolutionary optimization!

slow learning
compression / fast Eval

h: X → Y ⇒ Predictor ⇒ Hypothesis
 (Classification / Regression / Structured Pred. (graphs))

Parametric
 Non-Parametric
 ↳ bad scaling / slow time
 ↳ Fast Learning

↳ h: θ x X → Y ⇒ fit θ via loss fun.

↳ ML / MXP perspective!

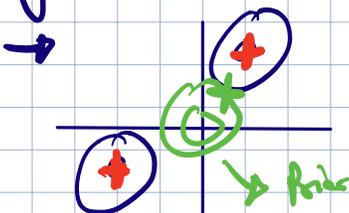
argmax_θ P(D|θ) = P(D|θ) P(θ) / P(D) ⇒ argmax_θ P(D|θ) P(θ)

→ Uniform prior in MXP formulation results in Max. Likelihood!
 ↳ together with iid assumption → split data into $\frac{N}{k} \oplus$ log trials

Loss fun often times results from Bayesian / MXP perspective ⊕ assumptions on the data → E.g. MSE from MXP + Uniform Prior + Gaussian Data!

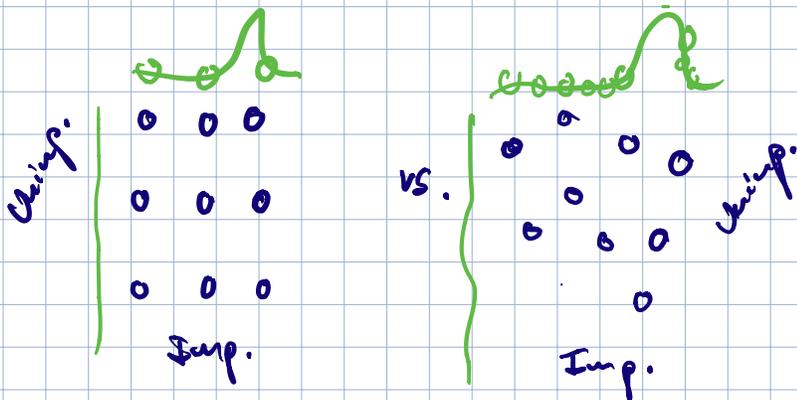
↳ Multi-Label Classification: Multivariate ⇔ reg. by Ltt.

Regularize via not assuming P(θ) being uniform! [Together w. outward approx.]



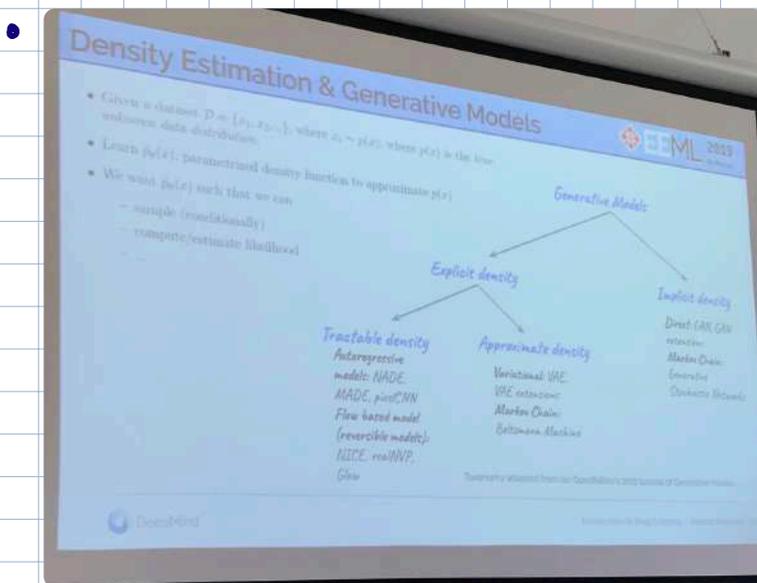
⇒ Reg. get rid of local min problem + by preferring one that is closer to 0 → +

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



\rightarrow Be careful not to miss best cluster!
 \rightarrow Unless what global level concept is!

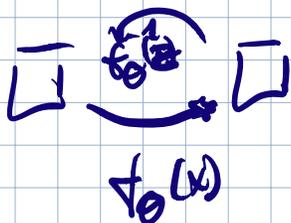
- Early-stopping as a prior \rightarrow Form of L_2 \Rightarrow See Bishop! Trust Region?!
- \rightarrow Duvenaud et al 2016 - Early Stopping as Hypothesis.



Generative Models

\rightarrow autoregressive: decompose joint as product of conditionals
 \rightarrow Not necessarily clear how to do so if there is no temporal chain.
 Eg. when to start sampling image
 \rightarrow also: slow sampling!

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



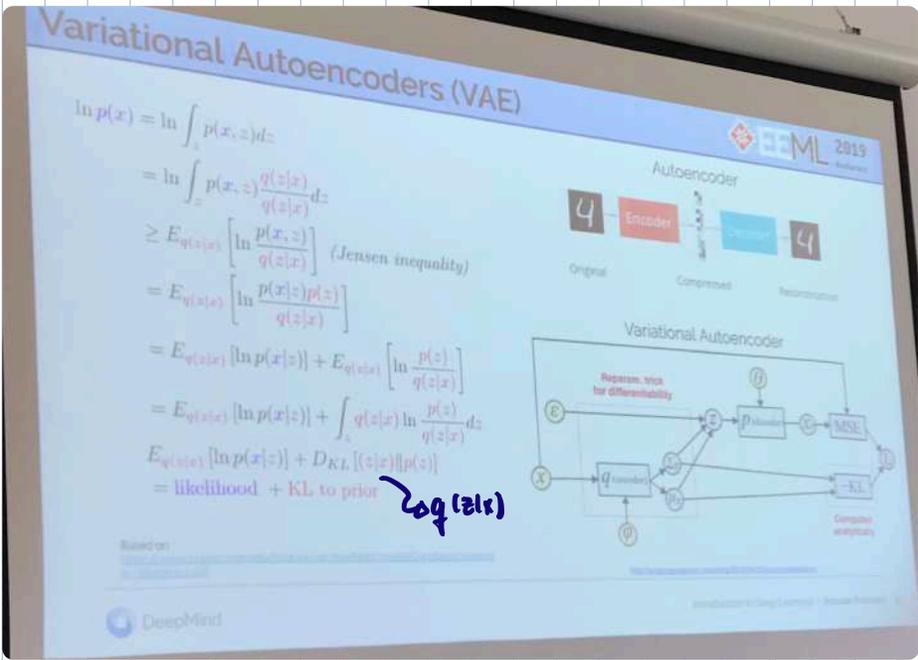
\rightarrow Change of variables \rightarrow Gaussian in latent space
 \rightarrow optimize jointly and invert

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

\rightarrow Not clear how to parametrize f_{θ} s.t. it is invertible!

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

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



- ↳ Gaussian cores on lucky for the capabilities as well as simple
- ↳ Gaussian + beta //
- also new with on Global Path all Multidimensional
- ↳ Problem to get steep color days!
- ↳ want VAE for latent interpolation!

→ Interpressive models vs. GMMs → Under GMM better, under BR

- ReLU activations: Splits space into linear regions → Can be a lot #
 - ↳ Montufar et al 2014 → NIPS → Learning
 - ↳ Form of optimal efficient policy of the space!
 - ↳ optimization such that decision boundary becomes simple!
 - ↳ No as fitting some regions are restricted to close parameters!
 - Need high-dimensional data → Not necessarily best
 - ↳ 1d mixture of sinusoids hard to fit!

- Loss Opt. Surfaces: Not necessarily crazy!
 - ↳ Yann Lecun: Simple random heuristic of complex prob. of all dir. putting up → very messy
 - ↳ Sharp vs. flat minima! → What do they mean? local depression
 - ↳ Can switch eigenvalue profile without changing features!

but a lot lower about restrictive parameters!

- RMS prop → Runif sup, Mean/Vos. → larger step if small var
 - ↳ Approx. of Natural Gradient
 - ↳ Better when → More accurate closer to identity → Fish & NG!

- Convolution = restriction of net \rightarrow need more inductive biases!
- Graph Nets = do conv not on locally defined pixels but on other local structure
- RNN \rightarrow Exploding gradient from product of Jacobians \rightarrow Gradient clipping
 - \hookrightarrow LSTM \rightarrow varying gradient constant via linear cell state but need to learn gates
 - \hookrightarrow Interpret gates more as a form of low pass filters ^{passes 0 at 1!}

Intro to RL - John Breck (DeepMind / McGill) \rightarrow TD

- How to overcome sample efficiency problem with sparse rewards
 - \rightarrow self-play \rightarrow easy acquisition // intrinsic motivation \rightarrow surrogate reward
- Key features:
 - * Trial-and-Error Search
 - * Stochastic Env
 - * Delayed reward \rightarrow Temp. Credit Assignment
 - * $E-E$ Trade-off
- heavily with multiple output heads / auxiliary heads
 - \hookrightarrow Relationship to Ego / Bounded Rationality?!
- TD Ganner \rightarrow (Tesauro et al, 1995) \Rightarrow "Early Alpha-Go"
- Intrinsic learning / behavioral cloning \rightarrow form of bootstrapping
 - \rightarrow important to train in simulator with exploration to robustify!
- Different interpretations: Survival prob., bias-variance tradeoff, self-heals
- ML methods \approx supervised $\rightarrow V(S_t)_{t+1} \leftarrow V(S_t)_{t_2} + \alpha [R_t - V(S_t)_{t_2}]$
 - \rightarrow Update for each transition \Rightarrow not iid! sequentially correlated
 - \hookrightarrow Overcome via ER suffers for example!
- Bias-Variance Trade-Offs:
 - * Function approx. \rightarrow Gradient updates
 - * Generalization across environments
 - * Discount \rightarrow Estimation of future rewards
 - \hookrightarrow schedule of discount params \rightarrow relate to TD as in PER?!

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

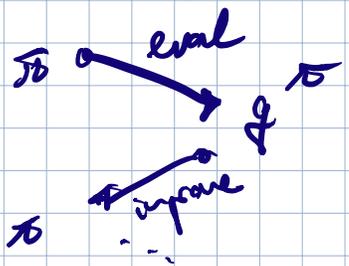
$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_{\pi}(s')]$$

- ↳ But $v_{\pi}(s)$ is a non-linear set of equations due to argmax!
- ↳ Can only be done if full MDP is known! \Rightarrow Estimate dynamics
- ↳ Importance $\gamma < 1$: 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 exploitable heuristic



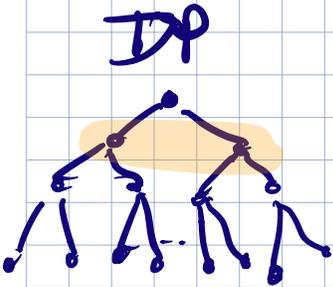
• Dyn-Pr: Full-Width Backup

↳ TD: In between MC and DP

• Crucial assumption of Markovity for bootstrapping \Rightarrow BIAS

↳ Introduce bias to reduce variance

↳ MC not biased! \Rightarrow k -step TD introduces!



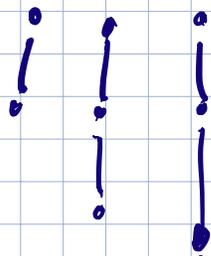
TD(1)



k -step TD



TD(2)



• Epilepsy project: Use SDS to simulate
↳ likelihood-free methods!!

exp. weighted updates!

EEML Summer School - Day 2

Intro to RL II - Policy Search → Control

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

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

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

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

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

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

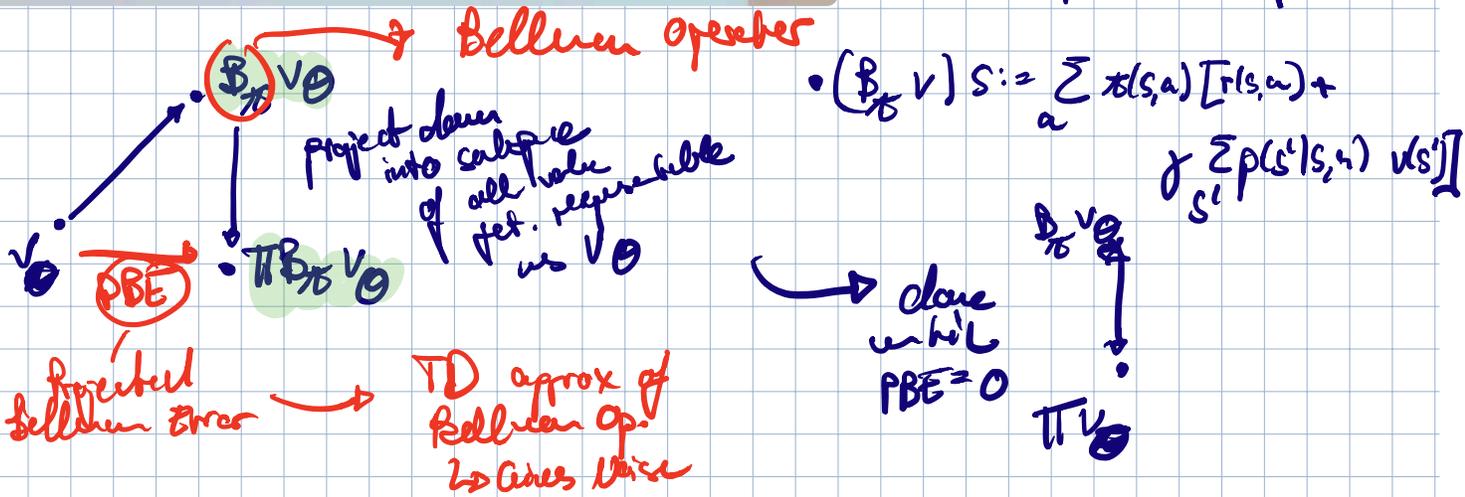
Action-value form: $\theta \leftarrow \theta + \alpha [\text{Target}_t - q(S_t, A_t, \theta)] \phi(S_t, A_t)$

↳ Transform into form of supervised learning

↳ Assume that target does not depend on θ

↳ Stability problem ⇒ target updates

↳ Would not be a problem if full MC



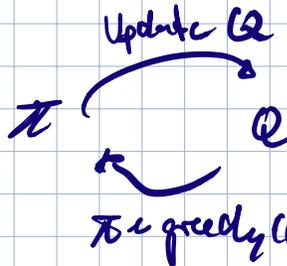
↳ Policy shifts due policy shifts! ⇒ need to constrain that shift

E.g. introducing trust region! → clipping / MC

↳ Linear fct. approx: $MSVE(\theta_{TD}) \leq \frac{1}{1-\gamma} MSVE(\theta_{MC})$ → noisy since in con-linear!

↳ u-step backups: Trade-off bias-variance ⇒ Full MC has large variance!
↳ of controls for TD

• Monte-Carlo Control:



⇒ Random exploring shifts!
↳ on vs off-policy

↳ Exploration: ϵ -greedy → distorting ⇒ problem: indep. of actual values

↳ In practice: exploration schedules are hard to implement

- **Targets:** $y_t = r_{t+1} + \gamma Q(s_t, a_t; \theta_t)$ [SARSA]
- $y_t = r_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a; \theta_t)$ [Exp. SARSA]
- $y_t = \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_a \pi(a|s_{t+1}) Q(s', a; \theta_t)]$
- ↳ SARSA: on-policy → next action is chosen at t to update
↳ will explore!

- **Q-learning:** **Exploring** → since **stochastically** converge from completely random state!
- ↳ SARSA vs. Q-L: **Safety vs. Optimality**

• **Maximization bias:**

$N(-0.1, 1)$ → $(b) \leftarrow (a) \rightarrow 0$ → **Table PDN**

↳ **variance reduction?** → Takes a long time to learn small values

↳ Use **softmax** instead ⇒ **form of ensemble!**

↳ **Stochastic policy** or **stochasticity** in environment

↳ **Tradeoff** **time shifts** → # iterations between target and update
more iterations ⇒ larger shift but smaller variance!

- **Reinforcement Learning:** Off-policy requirements ⇒ but maybe there is **safety belts!** → **META RL!**
- **Used for different set of optimizers in RL** → **different targets than supervised learning**

• **Policy gradient Methods:**

$\pi(a|s, \theta)$ → Objective: $J(\theta) = y_{\theta}(s_0)$ → **initial state distr.** → **policy**

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

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

Deriving REINFORCE from the PGT

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

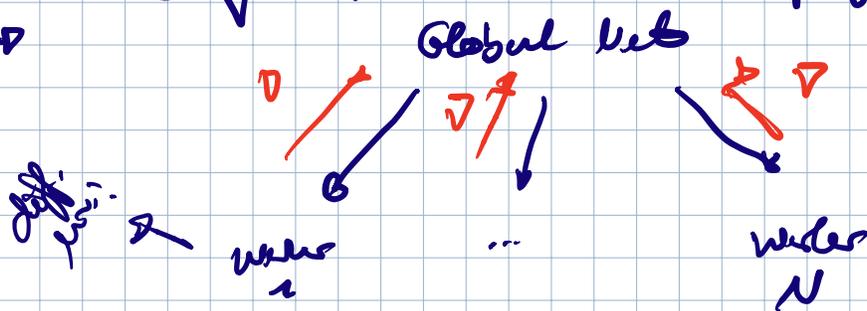
Thus

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

- No parameterization of value fun. needed
- ↳ Use returns
- ↳ MC estimation
- ↳ off policy!
- ↳ bad to stabilize!
- AC \rightarrow parameterize value fun!

- *3C \rightarrow burden exploration by parallelization!
- ↳ Pooling of experiences \rightarrow Mean of gradients

↳



- \rightarrow have set of gradients and update themselves
- ↳ you use average gradient!

- TRPO / PPO \rightarrow before KL constraint as large!

Multi-agent RL - Simon Whiteson (Oxford)

- **Robustness**: classical RL abstracts many aspects
- ↳ World is full of multi-agent systems

COOPERATIVE

- ↳ Shared reward
- ↳ Coordination

COMPETITIVE

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

MIXED

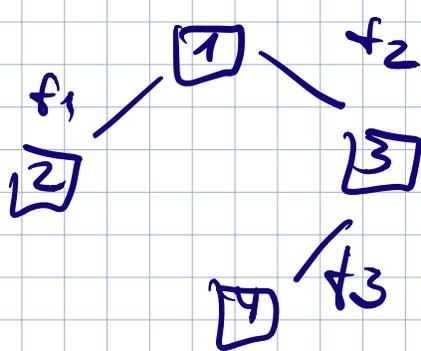
- ↳ General Sum
- ↳ Nash Eq.
- ↳ Q? Etc.

• Cooperative setting: Absence of social contract \Rightarrow risk of coordination

\rightarrow Exponential joint-action space \Rightarrow grows very fast

\rightarrow Question: what does \Rightarrow Coordination Graphs

$$Q(u) = f_1(u^1, u^2) + f_2(u^1, u^3) + \dots$$



↳ Probabilistic graphical Model

\Rightarrow Solvers: MDP estimation!

↳ Encoding of cond. independence

\Rightarrow No separability!

\rightarrow Agent 1's best response to get 2.5

↳ Variable elimination:

$$\max_u Q(u) = \max_{u_2, u_3, u_4} [f_3(u^3, u^4) + \max_{u_1} [f_1(u^1, u^2) + f_2(u^1, u^3)]]$$

↳ form of cooperative Markov!

↳ Computationally expensive the more connected the graph

↳ Computes only optimal value of joint action \Rightarrow No balanced paths to actually derive that value!

↳ Need to know the graph / cond. indep. conditions

\rightarrow Max-Plus \Rightarrow Vlassis et al (2004) \rightarrow Message Passing

↳ Two agent local payoff assumption

↳ Improve messages themselves

↳ Convergence guarantees in acyclic graphs

→ MDP (+ Synchronicity!) : Centralized Multi-Agent MDP

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

→ Independent Q-learning: No attempt to model $Q(s, u)$

↳ each agent learns $Q(s_i, u_i^a)$! \Rightarrow consistency if all agents learn \rightarrow no convergence guarantee! \rightarrow Tan et al (1988)

→ Correlated Q-learning: Busdrian et al (2002)

$$Q^{tot}(s, u) = \sum_{e \in \mathcal{E}} Q_e(s^e, u^e)$$

↳ extremely strong assumption! \rightarrow Transitions might be coupled

↳ Update each factor \Rightarrow graph is stationary!

• Partial observability \rightarrow requires decentralized execution!

↳ learning should be centralized! \rightarrow share params / gradients

↳ Minimal vs. right assumptions!

• Dec-POMDP

$Q(s, a): S \times \mathcal{A} \rightarrow \mathbb{R}$ \rightarrow believe the hist: $\tau^a \in T \subseteq (\mathbb{R} \times U)^T$

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

↳ dilemma: Exploit private info vs. being predictable \rightarrow ^{depends on} reward fun

• Policy Gradient Methods for MxRL:

\rightarrow PG useful if gradient is hard!

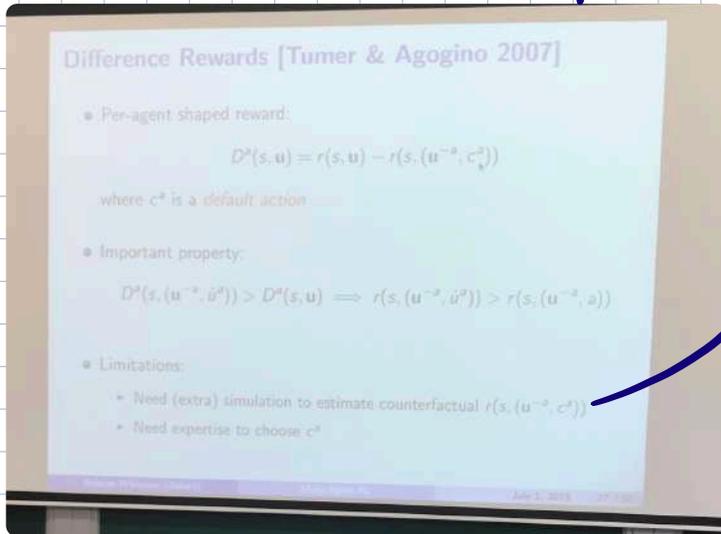
\rightarrow MxRL independent: share parameters - still indep. inputs + joint value
↳ consistency, hard to learn correlated, multi-agent coord. envs.

• Counterfactual MDPs → Foerster et al 2018 (COMA)

① Centralize Critic

② Wolpert & Tumer (2000) ⇒ Wondrous Life Utility

↳ reason about counterfactual reward if you would not have participated



↳ need to go back into simulator!
↳ Counterfactual Baseline

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

↳ problem: as π becomes deterministic we can't get the counterfactual aspect

↳ COMA critic: output kind gives all values

↳ still non-stochastic joint value fun.

• Value Decomposition Networks - Sunehag et al 2017

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

↳ decentralizes the agents ⇒ no lower level gradification

• QMIX: Not only summation but mixing which

⇒ constraint to non-neg weights!

↳ can blow away after centralized learning!

SMAC framework → StarCraft baseline!

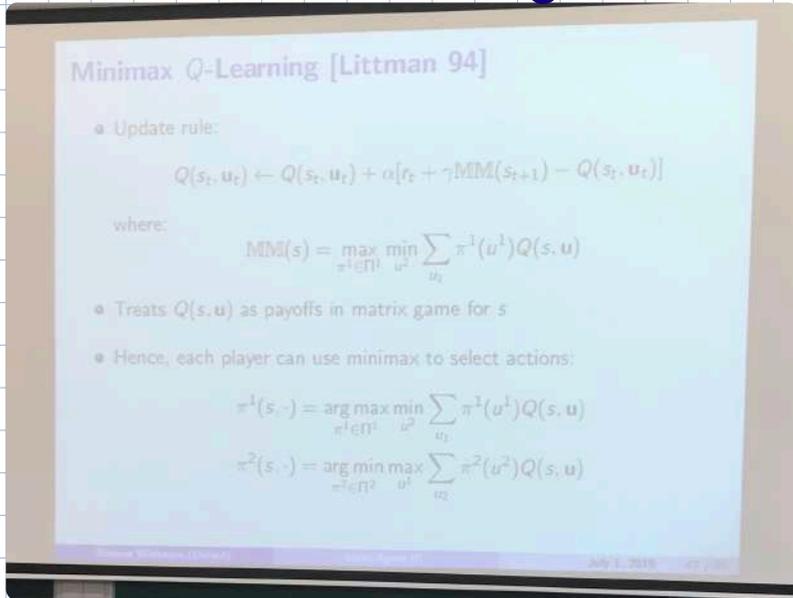
Competitive setting: In competitive setting one optimal policy exists → low: stochasticity

→ mixed strategies ↔ Minimax Theorem

↳ Rational vs. Strategic agents → Equilibrial ordering does not matter

↳ inner player does not have to believe stochastically!

→ Minimax Q-Learning → Littman (1994)



→ Training approach: Learn about the opponent from data
↳ Fwd. Q-Learning → implicit
↳ Explicit: Fictitious play
Rousselle
→ Opponent Modeling

→ Self-Play: Checkers → Arthur Samuel (1956)

↳ natural curriculum ⇒ on full if game is not transitive
↳ leave yourself vulnerable to bad opponents while beating stronger ones ⇒ help star → League / Hall of Fame!

• Fixed: Nash Equilibria → Existence!

↳ Nash Equilibrium Q-Learning → Hu & Wellia (1998)
→ almost never guaranteed to converge!

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

↳ Heuristics help! → Examples

↳ Fictitious-Q Q-Learning