# Dual–SPMA Framework for Convex MDPs

## Speaker Script with Q&A Bank

### Two-Column Format: Script + Anticipated Questions

Pegah Aryadoost  Danielle Nguyen  Shervin Khamooshian  Ahmed Magd

December 2025

## Speaker Assignments

| Speaker | Section | Slides |
|---------|---------|--------|
| Pegah | Background & Motivation | 1–9 |
| Danielle | Problem Formulation & Fenchel Duality | 10–17 |
| Shervin | Related Work & Our Method | 18–27 |
| Ahmed | Experiments & Conclusion | 28–36 |

## Notation Pronunciation Guide

> **Pronunciation Guide**
>
> **Greek Letters:** $\pi$ [pie] (policy) — $\gamma$ [gamma] (discount) — $\rho$ [rho] (initial dist.) — $\alpha$ [alpha] (step size) — $\beta$ [beta] (dual step) — $\lambda$ [lambda] (Lagrange mult.) — $\mu$ [mu] (penalty weight) — $\tau$ [tau] (threshold) — $\eta$ [eta] (learning rate) — $\omega$ [omega] (dual params) — $\phi$ [phi] (features)
>
> **Calligraphic:** $\mathcal{S}$ [script-S] (states) — $\mathcal{A}$ [script-A] (actions) — $\mathcal{D}$ [script-D] (occupancy set)
>
> **Functions:** $f^*$ [f-star] (conjugate) — $\nabla$ [nabla/gradient] — sup [soup/supremum] — $\langle \cdot, \cdot \rangle$ [inner product] — $[\cdot]_+$ [positive part]
>
> **Key Terms:** $d_\pi$ [d-pi] (occupancy) — $J(\pi)$ [J of pi] (return) — $L(\pi, y)$ [Lagrangian] — $r_y$ [r-sub-y] (shaped reward) — $\hat{d}_\pi$ [d-hat-pi] (estimated occupancy)

# Section 1: Background & Motivation (Pegah, Slides 1–9)

## Slide 1 — *A Dual–SPMA Framework for Convex MDPs*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Hi everyone, thanks for being here. I'm Pegah, and this project is joint work with Ahmed, Shervin, and Danielle from SFU's School of Computing Science."</li><li>"Our title is 'A Dual–SPMA Framework for Convex MDPs'."</li><li>"The one-sentence claim: **combining Fenchel duality with a fast policy optimizer, SPMA, gives a simple and fairly competitive way to solve convex MDPs**."</li><li>"We'll compare this 'Dual–SPMA' recipe against a strong baseline called NPG–PD."</li></ul> | **Q: What's the main takeaway?** <br> A: If you take the convex-MDP framework of Zahavy et al. and plug in SPMA as the policy optimizer, you get a simple "Dual–SPMA" recipe that (i) is easy to implement as shaped-reward RL, and (ii) performs competitively with NPG–PD on constrained safety tasks. <br><br> **Q: Is Dual–SPMA supposed to *beat* NPG–PD?** <br> A: No, we show SPMA is a **viable policy player** inside the convex-MDP dual framework: it's stable, respects constraints, and gives comparable reward. |

## Slide 2 — *Outline*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Here's the plan for the next 20–25 minutes."</li><li>"First, I'll quickly review basic RL background."</li><li>"Then I'll motivate **convex MDPs**—why we care about them beyond standard RL."</li><li>"Danielle: core math (Fenchel duality → saddle-point)."</li><li>"Shervin: related work and Dual–SPMA method."</li><li>"Ahmed: experiments and results."</li><li>"Finally: takeaways and future work."</li></ul> | **Q: Why compare against NPG–PD specifically? Why not PPO or TRPO?** <br> A: Our work is about **convex objectives & dual variables**, not generic performance on Atari. NPG–PD is a principled *primal–dual* algorithm specifically designed for constrained MDPs, with guarantees on both optimality and constraint violation. It also uses a geometry-aware policy update (natural gradient), which makes the comparison with SPMA conceptually clean. PPO/TRPO are great empirically, but they're not formulated directly as primal–dual CMDP algorithms. |

## Slide 3 — *What is Reinforcement Learning?*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Let's start with a quick refresh on RL."</li><li>"In RL, an **agent** interacts with an **environment** repeatedly over time."</li><li>"At each time step $t$ [t]: (1) observe state $s_t$ [s-sub-t], (2) choose action $a_t$ [a-sub-t], (3) receive reward $r_t$ [r-sub-t], (4) transition to $s_{t+1}$ [s-sub-t-plus-one]."</li><li>"Goal: learn a policy that maximizes long-term reward."</li></ul> | *(Basic RL background—unlikely to get technical questions here.)* <br><br> **Q: What's the difference between $y$ and $\lambda$?** <br> A: $y$ [y] is the general Fenchel dual variable over state–action pairs. $\lambda$ [lambda] is the scalar Lagrange multiplier for CMDP safety constraints. In the constrained case: $y_\lambda(s,a) = \lambda c(s,a) - r(s,a)$. |

## Slide 4 — *Markov Decision Process (MDP): Formal Definition*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Formally, the environment is a **Markov Decision Process**."</li><li>"MDP: $(\mathcal{S}, \mathcal{A}, P, r, \gamma, \rho)$ [script-S, script-A, P, r, gamma, rho]"</li><li>"$\mathcal{S}$: state space; $\mathcal{A}$: action space"</li><li>"$P(s' \mid s,a)$ [P of s-prime given s, a]: transition probability"</li><li>"$r(s,a)$: reward; $\gamma$ [gamma]: discount factor (0 to 1)"</li><li>"$\rho(s)$: initial state distribution"</li><li>"Policy $\pi(a \mid s)$ [pi of a given s]: probability distribution over actions"</li></ul> | **Q: What assumptions do you need for your convex MDP formulation?** <br> A: We assume a discounted MDP with $\gamma < 1$, bounded reward and cost, and that we restrict to stationary policies. For convexity, we require $f$ to be proper, closed, and convex on the occupancy polytope $\mathcal{D}$. For the Fenchel reduction, we also need standard conditions that guarantee the biconjugate equality and existence of a saddle point—essentially the Fenchel–Moreau assumptions. |

## Slide 5 — *Trajectory and Return*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| • "Running a policy gives a **trajectory**: sequence of states, actions, rewards." <br><br> • "$s_0 \to a_0 \to r_0 \to s_1 \to \ldots$" <br><br> • "The **discounted return**: $J(\pi) = \mathbb{E}_\pi \left[ \sum_{t=0}^\infty \gamma^t r(s_t, a_t) \right]$ [J of pi equals E-sub-pi of the sum...]" <br><br> • "Classical RL objective: find $\pi^*$ [pi-star] that maximizes this." | **Q: Can you restate the key identity relating $J(\pi)$ and $d_\pi$?** <br> A: Yes: we define <br><br> $$d_\pi(s,a) = (1-\gamma) \sum_{t=0}^\infty \gamma^t \Pr_\pi(s_t = s, a_t = a).$$ <br><br> Then <br><br> $$\langle r, d_\pi \rangle = \sum_{s,a} r(s,a)\, d_\pi(s,a) = (1-\gamma)J(\pi).$$ <br><br> So up to the constant factor $1 - \gamma$, the RL objective is just a linear function of the occupancy measure. |

## Slide 6 — *Occupancy Measure: Where the Policy Spends Time*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| • "Now we introduce the **discounted occupancy measure**." <br><br> • "$d_\pi(s,a) = (1-\gamma)\sum_{t=0}^\infty \gamma^t \Pr_\pi(s_t = s, a_t = a)$ [d-sub-pi of s, a...]" <br><br> • "Intuitively: **normalized discounted frequency** of visiting $(s,a)$ under $\pi$." <br><br> • "$(1-\gamma)$ normalizes so $\sum_{s,a} d_\pi(s,a) = 1$." <br><br> • "Think of $d_\pi$ as a heatmap over states and actions." | **Q: Why introduce occupancy measures at all? Why not just work with $J(\pi)$?** <br> A: Because in terms of policy parameters, $J(\pi)$ is generally non-convex. In terms of occupancy, return is **linear**: $\langle r, d_\pi \rangle = (1-\gamma)J(\pi)$. Many interesting extras (entropy, penalties, constraints) are convex in $d$. The set of valid occupancies is a convex polytope defined by linear flow constraints. So in $d$-space we get a clean convex program: minimize a convex $f(d)$ over a convex set $\mathcal{D}$. <br><br> **Q: What exactly is $\mathcal{D}$, the set of feasible occupancies?** <br> A: $\mathcal{D}$ is the set of $d \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ that satisfy the Bellman flow constraints plus non-negativity: for every state $s$, <br><br> $$\sum_a d(s,a) = (1-\gamma)\rho(s) + \gamma \sum_{s',a'} P(s|s',a')\, d(s',a'), \quad d(s,a) \geq 0.$$ <br><br> This says "discounted occupancy flowing *into* a state equals initial occupancy plus discounted flow from other states." It's a convex polytope. |

## Slide 7 — *Key Identity: RL is Linear in Occupancy*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| • "The **fundamental identity**:" <br><br> • "$\langle r, d_\pi \rangle = \sum_{s,a} r(s,a)\, d_\pi(s,a) = (1-\gamma)J(\pi)$" <br><br> • "So the RL objective is simply a linear function of $d_\pi$." <br><br> • "However, many interesting objectives are **not** linear in $d_\pi$:" <br><br> • "Safety constraints: $\langle c, d_\pi \rangle \leq \tau$" <br><br> • "Imitation learning (match expert occupancy)" <br><br> • "Exploration bonuses (entropy of $d_\pi$)" <br><br> • "That's where **convex MDPs** come in." | **Q: Is there a 1–1 correspondence between policies and occupancy measures?** <br> A: Under mild regularity conditions (e.g., communicating MDP, $\gamma < 1$), each stationary policy has a unique discounted occupancy measure, and each valid occupancy in $\mathcal{D}$ corresponds to at least one stationary policy (you can recover $\pi(a|s) = d(s,a)/\sum_{a'} d(s,a')$ where the denominator is nonzero). So working in $d$-space is equivalent, up to degeneracies. <br><br> **Q: Why do you call $f(d)$ "convex" when there's a minus reward term $-\langle r, d \rangle$?** <br> A: Because $-\langle r, d \rangle$ is linear in $d$, and linear functions are both convex and concave. When we add convex terms like entropy or penalties, the sum stays convex. So all the examples on the "Convex MDP: Examples" slide fit into "convex $f$" without breaking anything. |

## Slide 8 — *Why Convex MDPs?*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Linear RL is great if all we care about is maximizing reward."</li><li>"But in practice we often care about more:"</li><li>"**Enforce safety** constraints"</li><li>"**Match expert behaviour** in imitation learning"</li><li>"**Encourage exploration** or control risk"</li><li>"Convex MDPs: instead of maximizing linear $\langle r, d_\pi \rangle$, we **minimize convex** $f(d_\pi)$."</li><li>"Challenge: optimizing over $d$ is hard—dimension $|\mathcal{S}||\mathcal{A}|$, complicated polytope."</li><li>"Our approach: **Fenchel duality → min–max game →** shaped-reward RL."</li></ul> | **Q: Why do we even care about *convex* MDPs instead of just standard RL?** <br> A: Standard RL maximizes a linear function of the occupancy measure—expected return—so it's great when you only care about reward. Many real problems add structure: safety constraints, exploration terms, risk penalties, imitation losses, etc. Many of those can be written as convex functionals of the occupancy $d_\pi$. Convex MDPs give a unified way to handle these goals while preserving nice optimization properties like existence of a saddle point. <br><br> **Q: How is this different from just "constrained RL with a Lagrange multiplier"?** <br> A: Classic constrained RL focuses on a particular structure: reward plus linear constraints. Convex MDPs are more general: you minimize an arbitrary convex function $f(d_\pi)$ over the convex occupancy polytope. Constrained RL is one example where $f$ adds a penalty or indicator for constraint violations; entropy-regularized RL and some imitation objectives are *other* examples. Our framework—and the Fenchel duality step—work for all of them, not just linear constraints. |

## Slide 9 — *Convex MDP: Examples*

| Speaker Script (Pegah) | Related Q&A |
|---|---|
| <ul><li>"Examples of convex MDP objectives:"</li><li>"Standard RL: $f(d) = -\langle r, d \rangle$ (linear, hence convex)"</li><li>"**Entropy-regularized RL**: <br> $f(d) = -\langle r, d \rangle + \alpha \sum_{s,a} d(s,a) \log d(s,a)$"</li><li>"**Constrained safety**: <br> $f(d) = -\langle r, d \rangle + \mu \max\{0, \langle c, d \rangle - \tau\}$"</li><li>"Here $c$: cost, $\tau$ [tau]: threshold, $\mu$ [mu]: penalty weight"</li><li>"All three are **convex functions of occupancy**."</li><li>"Next: Danielle shows **Fenchel duality**."</li></ul> | **Q: Why use $\mu$ for penalty weight instead of $\lambda$?** <br> A: We intentionally use $\mu$ there to distinguish a **fixed penalty weight** in a penalized objective from the **Lagrange multiplier** $\lambda$ that evolves as a dual variable. This avoids confusion between "tunable hyperparameter" and "optimization variable." <br><br> **Q: What's the difference between $y$ and $\lambda$?** <br> A: $y$: general Fenchel dual variable over state–action pairs. $\lambda$: scalar Lagrange multiplier for CMDP constraints. In constrained case: $y_\lambda(s,a) = \lambda c(s,a) - r(s,a)$. |

# Section 2: Problem Formulation & Fenchel Duality (Danielle, Slides 10–17)

## Slide 10 — *Roadmap (checkpoint)*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| • "Thanks, Pegah. So far: RL basics and convex MDPs motivation."<br><br>• "Next: I'll show how **Fenchel duality** turns this into a saddle-point formulation." | *(Transition slide—take a breath, check audience engagement.)* |

## Slide 11 — *Fenchel Conjugate: Definition*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| • "The **Fenchel conjugate** of convex $f$:"<br><br>• "$f^*(y) = \sup_x\{\langle y, x\rangle - f(x)\}$ [f-star of y equals supremum...]"<br><br>• "$y$: dual variable; $\langle y, x\rangle$: dot product; sup [soup]: supremum"<br><br>• "Intuitively: how large can $\langle y, x\rangle$ be relative to $f(x)$?"<br><br>• "We'll use this to expose a min-max structure." | **Q: Why bring in Fenchel conjugates at all?**<br>A: Fenchel conjugates give us an identity:<br><br>$$f(d) = \sup_y\{\langle y, d\rangle - f^*(y)\}.$$<br><br>Plugging this into $\min_{d\in\mathcal{D}} f(d)$ lets us rewrite the convex MDP as<br><br>$$\min_{d\in\mathcal{D}}\max_y\{\langle y, d\rangle - f^*(y)\}.$$<br><br>So we get a saddle-point game with a natural dual variable $y$. This structure is what allows us to run two-player algorithms: one player updates $y$, the other updates the policy.<br><br>**Q: How do you get the conjugate $f^*(y)$ in the entropy-regularized case?**<br>A: For<br><br>$$f(d) = -\langle r, d\rangle + \alpha\sum d\log d,$$<br><br>you can complete the square in log-space / use standard convex analysis to show<br><br>$$f^*(y) = \alpha\log\sum_{s,a}\exp((y(s,a) + r(s,a))/\alpha).$$<br><br>Its gradient is a softmax: $\nabla f^*(y) = \text{softmax}((y + r)/\alpha)$. That's the expression we use in the dual update. |

## Slide 12 — *Fenchel–Moreau Theorem*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| • "Key result: **Fenchel–Moreau theorem**."<br><br>• "For proper, closed, convex $f$:"<br><br>• "$f(d) = \sup_y\{\langle y, d\rangle - f^*(y)\}$"<br><br>• "So $f = f^{**}$ [f equals f-double-star]."<br><br>• "This allows us to turn **minimizing** $f(d)$ into a **min–max optimization** over $d$ and $y$." | **Q: What assumptions for convex MDP formulation?**<br>A: Discounted MDP ($\gamma < 1$), bounded reward/cost, stationary policies. $f$ must be proper, closed, convex on $\mathcal{D}$. Standard Fenchel–Moreau conditions for biconjugate equality and saddle point existence. |

## Slide 13 — *Applying Fenchel Duality to Convex MDPs*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| <ul><li>"Start with: $\min_{d \in \mathcal{D}} f(d)$"</li><li>"Apply Fenchel–Moreau: $f(d) = \sup_y \{\langle y, d \rangle - f^*(y)\}$"</li><li>"Plug in: $\min_{d \in \mathcal{D}} \sup_y \{\langle y, d \rangle - f^*(y)\}$"</li><li>"This is a **convex-concave saddle-point problem**."</li><li>"Under standard assumptions, solving this equals solving original convex MDP."</li><li>"Replace $d$ by $d_\pi$: $\min_\pi \max_y L(\pi, y)$ where $L(\pi, y) = \langle y, d_\pi \rangle - f^*(y)$"</li></ul> | **Q: Why do you write $\min_d \max_y$ instead of $\max_y \min_d$? Are you allowed to swap min and max?**<br>A: We derive the saddle-point form by directly applying Fenchel–Moreau:<br><br>$$f(d) = \sup_y \{\langle y, d \rangle - f^*(y)\},$$<br><br>then taking $\min_d$ outside. Under standard convexity/closedness assumptions, there is no duality gap and the saddle-point value equals both the minimax and maximin. So while we don't literally swap orders in the algebra, we rely on these conditions to interpret the solution as an equilibrium of the game. |

## Slide 14 — *Two-Player Game Interpretation*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| <ul><li>"Saddle-point as a **two-player game**:"</li><li>"$\min_\pi \max_y L(\pi, y) = \langle y, d_\pi \rangle - f^*(y)$"</li><li>"**Policy player**: chooses $\pi$, wants to **minimize** $L$ (RL agent)"</li><li>"**Dual player**: chooses $y$, wants to **maximize** $L$ (shapes reward)"</li><li>"At equilibrium (saddle point): neither can unilaterally improve"</li><li>"Policy player has $\pi^*$, dual player has $y^*$"</li></ul> | **Q: Why is the solution of the min–max at a saddle point?**<br>A: Because in convex-concave problems, under suitable conditions, the minimax and maximin values coincide and are achieved at a pair $(\pi^*, y^*)$ satisfying<br><br>$$L(\pi^*, y) \le L(\pi^*, y^*) \le L(\pi, y^*)$$<br><br>for all $\pi, y$. That point is simultaneously optimal for the min player (policy) and the max player (dual), and is called a saddle point. Solving the min–max is thus equivalent to finding such a saddle point. |

## Slide 15 — *Visualizing the Saddle Point*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| <ul><li>"3D visualization: $L(d, y)$ as a surface"</li><li>"$d$ (horizontal): policy-side variable"</li><li>"$y$ (other axis): dual variable"</li><li>"**Red point**: saddle point"</li><li>"Along $d$-direction (fix $y$): red point is a **minimum**"</li><li>"Along $y$-direction (fix $d$): red point is a **maximum**"</li><li>"It's 'best response' for both players $\rightarrow$ solution of min–max game"</li></ul> | **Q: Does your algorithm *guarantee* convergence to the saddle point?**<br>A: In the idealized tabular, exact-oracle setting, existing theory (Zahavy et al. + SPMA paper) shows that if the dual player and the policy player both use low-regret algorithms, the average iterates converge to a saddle point. In our implementation we have function approximation, noisy occupancy estimates, and finite inner loops, so we don't claim a formal convergence proof. Instead we check empirically that the saddle objective and constraints behave sensibly. |

## Slide 16 — *From Saddle Point to Shaped Reward*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| <ul><li>"How to minimize over policies in saddle-point?"</li><li>"For fixed $y$, policy subproblem: $\min_\pi \langle y, d_\pi \rangle$"</li><li>"Using occupancy definition: $\langle y, d_\pi \rangle = (1 - \gamma)\mathbb{E}_\pi[\sum_t \gamma^t y(s_t, a_t)]$"</li><li>"So minimizing $\langle y, d_\pi \rangle = $ **maximizing** return with **shaped reward**: $r_y(s, a) = -y(s, a)$"</li><li>"**Key insight**: policy player can use **standard RL** on shaped reward $r_y = -y$"</li></ul> | **Q: What does the dual variable $y$ represent intuitively?**<br>A: It's a vector of the same dimension as $d$—one coordinate per state–action pair. In the Fenchel dual, it's the variable that "linearizes" the convex objective. Operationally, for the policy player, $y$ is just a shaped reward field: we treat $r_y(s, a) = -y(s, a)$ as the reward. Different choices of $y$ encourage the policy to put occupancy where $f$ wants it.<br><br>**Q: Does the shaped reward $r_y = -y$ break the Markov property or anything?**<br>A: No, because $y$ is a function of the current state–action pair $(s, a)$ only; it doesn't introduce dependence on the entire trajectory. The environment dynamics $P(s'|s, a)$ are unchanged, and the reward at each step is just a different function of $(s, a)$. So the process is still a valid MDP, just with a different reward. |

## Slide 17 — *Summary: The Fenchel Dual Reduction*

| Speaker Script (Danielle) | Related Q&A |
|---|---|
| <ul><li>"Summary of reduction:"</li><li>"Started with convex MDP: $\min_\pi f(d_\pi)$"</li><li>"Fenchel duality $\rightarrow$ saddle-point: $\min_\pi \max_y L(\pi, y)$"</li><li>"Fixing $y$: minimizing over $\pi$ = RL with shaped reward $r_y = -y$"</li><li>"Simple algorithmic structure:"</li><li>"1. **Dual step**: $y_{k+1} = y_k + \alpha(\hat{d}_{\pi_k} - \nabla f^*(y_k))$"</li><li>"2. **Policy step**: run RL on $r_{y_k} = -y_k$ to get $\pi_{k+1}$"</li><li>"Reduced to: **alternating convex optimization (dual) and standard RL (policy)**"</li></ul> | **Q: Why not just solve the convex program over $d$ with CVX or some LP solver?**<br>A: You *could* if the MDP is small and fully known. In our setting we assume a **model-free** RL scenario: we don't know the transition probabilities, we only interact with the environment by sampling trajectories. That means we can't write down $\mathcal{D}$ explicitly or compute exact expectations. RL gives us a way to approximate $\min_\pi \langle y, d_\pi \rangle$ from samples by treating $-y$ as a reward and running a policy optimizer.<br><br>**Q: Why is alternating updates (policy step then dual step) reasonable, given that vanilla Gradient Descent–Ascent can diverge?**<br>A: It's true that naive GDA can oscillate. We're not using plain GDA: we use specific online-learning style updates. The dual player uses a simple FTL/gradient step in a convex Euclidean space; the policy player uses SPMA, which is a mirror-descent method in logit space with convergence guarantees. Online saddle-point theory says that when both players have low regret, the average of their strategies converges to a saddle point, so alternating these particular updates is justified. |

# Section 3: Related Work & Our Method (Shervin, Slides 18–27)

## Slide 18 — *Roadmap (checkpoint)*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"We've completed the core math using Fenchel duality."</li><li>"Next: related work that inspired this framework and the specific policy optimizer, SPMA."</li></ul> | *(Transition slide—Shervin takes over.)* |

## Slide 19 — *Related Work: "Reward Is Enough"*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"First key paper: '*Reward Is Enough for Convex MDPs*' by Zahavy et al."</li><li>"Foundational work for convex-MDP perspective."</li><li>"Three main contributions:"</li><li>"1. **Fenchel dual reduction**: convex MDP $\rightarrow$ saddle-point"</li><li>"2. **Meta-algorithm**: alternate policy & dual updates"</li><li>"3. Many RL paradigms (imitation, constrained, entropy-reg.) are special cases"</li><li>"Our contribution: **implement this with SPMA** and compare to NPG–PD."</li></ul> | **Q: If you can plug "any RL algorithm" into the policy player, why specifically SPMA?** <br> A: The convex-MDP theory says any RL algorithm could in principle be used as the policy player. We choose SPMA because (i) it has **strong convergence results** in tabular MDPs, (ii) it's geometry-aware, operating in logit space with a softmax mirror map, (iii) its loss looks similar to PPO/MDPO, so it's practical to implement, and (iv) the SPMA paper already shows it's competitive with PPO/TRPO in continuous control. |

## Slide 20 — *Related Work: Softmax Policy Mirror Ascent (SPMA)*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"Second key ingredient: **SPMA** from Asad et al."</li><li>"**Mirror ascent in logit space** using log-sum-exp mirror map."</li><li>"In tabular MDPs: achieves **linear convergence**."</li><li>"Tabular update: $\pi_{t+1}(a\|s) = \pi_t(a\|s)[1 + \eta A_{\pi_t}(s,a)]$"</li><li>"**Geometry-aware**: respects probability simplex"</li><li>"No explicit per-state normalization"</li><li>"Clean extension to function approximation as convex classification"</li></ul> | **Q: How is SPMA different from vanilla policy gradient?** <br> A: Vanilla PG minimizes $L_{\mathrm{PG}} = -\mathbb{E}[\log \pi(a\|s) A(s,a)]$ and takes gradient steps in parameter space. SPMA instead uses an **exponential-gradient / mirror-descent update in logit space**. The actor loss includes a regularizer $$\frac{1}{\eta}(\exp(\Delta \log \pi) - 1 - \Delta \log \pi)$$ which penalizes large changes in log-probabilities. This gives better control over how far the policy moves per step and leads to faster convergence in the tabular theory. <br> **Q: Why is SPMA "geometry-aware"?** <br> A: Because it does mirror ascent with the log-sum-exp mirror map, which is tailored to the simplex geometry. Instead of taking Euclidean steps in parameter space and then projecting back onto the simplex, SPMA treats the logits as dual variables and the softmax as the mirror map. This respects the probability constraints and tends to produce well-behaved updates when probabilities are near 0 or 1. |

## Slide 21 — *Related Work: NPG–PD (Our Baseline)*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"Main baseline: **NPG–PD** (Ding et al. 2020)"</li><li>"Designed for constrained MDPs with Lagrangian:"</li><li>"$L(\pi, \lambda) = J_r(\pi) + \lambda(J_c(\pi) - \tau)$, $\lambda \geq 0$"</li><li>"**Primal**: natural policy gradient ascent (Fisher geometry)"</li><li>"**Dual**: $\lambda_{k+1} = [\lambda_k + \beta(J_c(\pi_k) - \tau)]_+$"</li><li>"Has $\mathcal{O}(1/\sqrt{T})$ bounds on optimality gap & constraint violation"</li><li>"We use NPG–PD as our **comparison point**."</li></ul> | **Q: How does SPMA compare to natural policy gradient conceptually?**<br>A: Both are trying to account for the geometry of the policy space. Natural PG uses the Fisher information matrix to modify the gradient direction and step size; SPMA uses a Bregman divergence in logit space. In small tabular problems you can show they are closely related; in our setting we treat them as two reasonable but different choices of geometry-aware update and compare them empirically. |

## Slide 22–23 — *Roadmap & Our Contributions*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| **Slide 22 (Roadmap):**<br><ul><li>"So far: convex-MDP theory, Fenchel dual framework, SPMA optimizer."</li><li>"Next: what we actually built."</li></ul>**Slide 23 (Our Contributions):**<br><ul><li>"1. **Complete Dual–SPMA framework**: outer dual loop + SPMA policy oracle"</li><li>"Supports entropy-regularized RL and constrained safety CMDPs"</li><li>"2. **Three occupancy estimators**: tabular MC, feature-based MC, MLE-style"</li><li>"3. **Faithful NPG–PD baseline**"</li><li>"4. **Empirical comparison** of SPMA vs NPG–PD on constrained safety"</li></ul> | **Q: Could you use other policy optimizers in place of SPMA?**<br>A: Yes. The meta-algorithm from Zahavy et al. lets you plug in any RL method as the policy player. We chose SPMA because of its theory and similarity to modern methods, but in principle you could try PPO, TRPO, or MDPO. An interesting piece of future work would be a systematic empirical comparison of different policy players inside the same convex-MDP dual framework. |

## Slide 24 — *Dual–SPMA Loop: High-Level View*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"High-level Dual–SPMA structure:"</li><li>"Saddle-point objective: $L(\pi, y) = \langle y, d_\pi \rangle - f^*(y)$"</li><li>"**Outer loop** (dual): $y_{k+1} = y_k + \alpha(\hat{d}_{\pi_k} - \nabla f^*(y_k))$"</li><li>"Use estimated occupancy $\hat{d}_{\pi_k}$ from rollouts"</li><li>"**Inner loop** (policy): run $K_{\text{inner}}$ SPMA steps on shaped reward $r_y = -y$"</li><li>"After fixed outer iterations, return final policy $\pi_K$"</li></ul> | **Q: Why do you use a simple gradient/FTL step for the dual variable $y$?**<br>A: The dual variable lives in a straightforward Euclidean space and the dual objective is convex with a simple gradient $\hat{d}_\pi - \nabla f^*(y_k)$. For such problems, basic online convex optimization algorithms like gradient ascent or FTL are known to have sublinear regret. They're easy to implement and come with clear theoretical guarantees, so there's no need for more complex methods.<br><br>**Q: How do you choose the SPMA step size $\eta$?**<br>A: We use a simple Armijo line search over a small discrete set of candidate $\eta$ values. For each candidate we compute the SPMA loss on a minibatch and pick the largest $\eta$ that produces a sufficient decrease. This is inspired by the SPMA paper and by MDPO/TRPO line-search practices. |

## Slide 25 — *Policy Player: SPMA Actor Loss*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| • "SPMA policy update via custom **actor loss**." <br><br> • "Standard PG: $L_{\mathrm{PG}} = -\mathbb{E}[\log \pi(a\|s) A(s,a)]$" <br><br> • "SPMA adds 'stay close to old policy' regularizer in logit space:" <br><br> • "$L_{\mathrm{SPMA}} = \mathbb{E}[-\Delta \log \pi \cdot A + \frac{1}{\eta}(\exp(\Delta \log \pi) - 1 - \Delta \log \pi)]$" <br><br> • "$\Delta \log \pi = \log \pi_{\mathrm{new}} - \log \pi_{\mathrm{old}}$" <br><br> • "Exponential term: KL-like regularizer penalizing large log-prob changes" | **Q: Where does your SPMA loss come from?** <br> A: It's the convex surrogate corresponding to the tabular SPMA update. If you linearize the value function and consider an exponential-gradient step in logit space, the associated Bregman divergence gives rise to exactly that $\exp -1 - \Delta \log \pi$ term. The SPMA paper derives it formally; we implement it as a differentiable loss so we can use standard automatic differentiation. <br><br> **Q: Do you enforce the $[1 + \eta A]_+$ and normalization from the tabular update?** <br> A: In the tabular theory, the closed-form update is <br><br> $$\pi_{t+1}(a\|s) \propto \pi_t(a\|s)[1 + \eta A(s,a)]_+.$$ <br><br> In our neural approximation we don't directly enforce this per-state formula. Instead we optimize the SPMA loss over minibatches of trajectories; the softmax layer plus the KL-like regularizer implicitly enforce a similar behaviour. So we're using the **function-approximation version** of SPMA instead of the exact tabular update. |

## Slide 26 — *Occupancy Estimation: MC vs MLE*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| • "Critical ingredient: occupancy estimate $\hat{d}_\pi$" <br><br> • "**Default**: Monte Carlo estimator:" <br><br> • "$\hat{d}_\pi(s,a) = \frac{1-\gamma}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \gamma^t \mathbf{1}\{s_t^{(i)} = s, a_t^{(i)} = a\}$" <br><br> • "Simple: count discounted visits. Variance grows with $\|\mathcal{S}\|\|\mathcal{A}\|$." <br><br> • "Also explored **MLE-based estimator** (Barakat et al.): log-linear $\lambda_\omega(s,a) \propto \exp(\omega^\top \phi(s,a))$" <br><br> • "Error depends on feature dimension, not $\|\mathcal{S}\|\|\mathcal{A}\|$" <br><br> • "Sanity check: verified $\sum_{s,a} \hat{d}_\pi(s,a) \approx 1$" | **Q: Why not use exact occupancies in FrozenLake instead of Monte Carlo?** <br> A: We intentionally use Monte Carlo to simulate the **model-free** setting most RL algorithms operate in—where you don't have access to $P(s'\|s,a)$. It also keeps the pipeline more consistent with what you'd do in larger environments. Using the exact occupancy is possible in FrozenLake, but we'd then be solving a much easier problem that doesn't stress the estimator. <br><br> **Q: How do you ensure your Monte Carlo occupancy estimate is valid?** <br> A: We compute <br><br> $$\hat{d}_\pi(s,a) = \frac{1-\gamma}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \gamma^t \mathbf{1}\{s_t^{(i)} = s, a_t^{(i)} = a\}.$$ <br><br> By construction this is non-negative. As a sanity check, we track $\sum_{s,a} \hat{d}_\pi(s,a)$ over training and verify it stays very close to 1; that's the plot on our entropy-regularized results slide. Small deviations come from truncating trajectories at finite $T$. <br><br> **Q: What's the point of the MLE-based occupancy estimator?** <br> A: The tabular MC estimator's variance scales with $\|\mathcal{S}\|\|\mathcal{A}\|$; in large MDPs that's bad. The MLE estimator instead fits a log-linear model $\lambda_\omega(s,a) \propto \exp(\omega^\top \phi(s,a))$ by maximum likelihood on samples. Its error depends mostly on the feature dimension, not the raw state–action count. So it's a more scalable idea we wanted to prototype, although we only briefly tested it here. |

## Slide 27 — *Example: Constrained Safety CMDP*

| Speaker Script (Shervin) | Related Q&A |
|---|---|
| <ul><li>"How constrained MDP fits our framework:"</li><li>"Problem: maximize $J_r(\pi)$ s.t. $J_c(\pi) \leq \tau$"</li><li>"In Dual–SPMA:"</li><li>"1. Build dual: $y_\lambda(s,a) = \lambda c(s,a) - r(s,a)$"</li><li>"2. Policy sees: $r_y = r - \lambda c$"</li><li>"3. Run SPMA inner loop on $r_y$"</li><li>"4. Update dual: $\lambda_{k+1} = [\lambda_k + \beta(J_c(\pi_k) - \tau)]_+$"</li><li>"SPMA & NPG–PD have **same dual update**—difference is policy optimizer."</li><li>"Now Ahmed presents experimental results."</li></ul> | **Q: Are there any constraints on $y$? Do you project it?**<br>A: In the entropy-regularized case, $y$ is unconstrained—we don't project. For CMDPs we sometimes parameterize the dual as $y_\lambda = \lambda c - r$ with $\lambda \geq 0$, so the only constraint is non-negativity on $\lambda$, which we enforce with a $[\cdot]_+$ projection (same as NPG–PD). More sophisticated bounds or regularizers on $y$ are possible but we didn't explore them in this project.<br><br>**Q: Is the dual step size $\alpha$ sensitive?**<br>A: Yes, like most gradient methods it needs to be in a reasonable range: too small and progress is slow; too large and you can oscillate. In practice we tuned $\alpha$ on a log scale and picked values where the dual objective and the constraint violation decreased smoothly. Because our setting is tabular and relatively small, this wasn't too painful. |

# Section 4: Experiments & Conclusion (Ahmed, Slides 28–36)

## Slide 28–30 — *CMDP Example, Roadmap, Experimental Setup*

| Speaker Script (Ahmed) | Related Q&A |
|---|---|
| **Slide 28**: CMDP example (same as Slide 27, concrete walkthrough).<br>**Slide 29**: "We've described Dual–SPMA, occupancy estimators, NPG–PD baseline. Next: experimental setup and results."<br>**Slide 30 (Experimental Setup):**<br>• "Tabular setting for visualizing occupancies."<br>• "Environment: **FrozenLake 4×4** gridworld, deterministic transitions."<br>• "Unsafe states (holes): cost $c = 1$."<br>• "Compare: **Dual–SPMA** vs **NPG–PD**"<br>• "Track: $J_r(\pi)$, $J_c(\pi)$, $J_c - \tau$ (violation), occupancy sum $\approx 1$"<br>• "Settings: $\gamma = 0.99$, $\tau = 0.1$, 30 outer iterations, 2048 steps/rollout" | **Q: Why FrozenLake 4×4?**<br>A: Wanted **tabular** CMDP to (i) compute/visualize occupancies, (ii) understand safety constraint intuitively (holes vs safe), (iii) debug without function approximation error. Standard simple gridworld.<br><br>**Q: How many trajectories/steps per outer iteration?**<br>A: 2048 environment steps per rollout per iteration. Enough for stable occupancy estimation and critic training in this small environment. Didn't push sample efficiency; systematic study would vary this.<br><br>**Q: How do you measure "success" in your experiments?**<br>A: Three criteria: (1) **constraint satisfaction**—$J_c(\pi) \leq \tau$; (2) **reward**—$J_r(\pi)$ as high as possible; (3) **stability**—no blow-up or wild oscillations. On these, Dual–SPMA $\approx$ NPG–PD. |

## Slide 31–32 — *Results: Entropy-Regularized RL & Constrained CMDP*

| Speaker Script (Ahmed) | Related Q&A |
|---|---|
| **Slide 31 (Entropy-Regularized):**<br>• "No constraints, just entropy bonus."<br>• "Left: $L$ (Lagrangian) vs iteration—both methods decrease it."<br>• "Right: occupancy sum $\approx 1$ (sanity check)."<br>• "SPMA and NPG–PD behave similarly."<br>**Slide 32 (Constrained CMDP):**<br>• "Left: reward $J_r(\pi_k)$ increases over iterations."<br>• "Right: violation $J_c(\pi_k) - \tau$ goes below zero (constraint satisfied)."<br>• "Both methods reach similar final performance." | **Q: Which method converges faster?**<br>A: Depends on hyperparameters. Typically SPMA takes larger, more "aggressive" steps, can reduce violation faster initially, but slightly more sensitive to hyperparameters. NPG–PD smoother but needs careful tuning of $\beta$. Final performance broadly comparable.<br><br>**Q: How do you handle exploration?**<br>A: Stochastic policy (softmax) naturally explores early. Entropy-regularized experiments explicitly add entropy term. No extra exploration bonuses. |

## Slide 33–34 — *Results: SPMA vs NPG–PD & Occupancy Heatmaps*

| Speaker Script (Ahmed) | Related Q&A |
|---|---|
| **Slide 33 (Comparison):**<br>• "Left: reward $J_r$ vs iterations—similar once constraint satisfied."<br>• "Right: violation—both bring close to zero, different curves."<br>• "SPMA: larger steps, sometimes faster, more hyperparameter-sensitive."<br>• "NPG–PD: smoother, needs careful $\beta$ tuning."<br>• "Conclusion: Dual–SPMA is competitive, SPMA is viable policy optimizer."<br>**Slide 34 (Heatmaps):**<br>• "**Left**: unconstrained policy—explores broadly, includes unsafe states."<br>• "**Right**: safety-constrained—avoids unsafe states (occupancy $\approx 0$)."<br>• "Constraint reflected in agent behaviour, not just numerically." | **Q: Do both methods satisfy the constraint?**<br>A: Yes. Both reduce $J_c(\pi) - \tau$ toward zero. Violation plots approach/oscillate around zero. Learned policies' heatmaps clearly avoid unsafe states.<br><br>**Q: Are advantages estimated with critic or MC?**<br>A: Standard actor–critic: critic is value network trained with TD, advantage computed as $A = \hat{Q} - V$ via GAE or simple TD returns. Both SPMA and NPG–PD share same critic architecture.<br><br>**Q: Do you use MLE estimator for main results?**<br>A: No, we use MC estimator for main FrozenLake experiments. MLE implemented and tested briefly; full comparison is future work. |

## Slide 35–36 — *Takeaways, Limitations & Future Work*

| Speaker Script (Ahmed) | Related Q&A |
|---|---|
| **Slide 35 (Takeaways & Limitations):** | **Q: What are the main limitations of your project?** |
| • "**Recipe**: Convex MDP → Fenchel dual → SPMA-based shaped-reward RL" | A: The big three are: (1) we only test in small tabular environments, so we don't claim scalability yet; (2) our MLE estimator is implemented but not fully stress-tested on large problems; and (3) we don't provide a formal convergence-rate analysis for the full Dual–SPMA algorithm with function approximation—only for the inner SPMA and the baseline, via existing theory. |
| • "**Built**: Dual–SPMA loops, NPG–PD baseline, occupancy estimators" | |
| • "**Result**: Dual–SPMA is stable, satisfies constraints, competitive with NPG–PD" | |
| • "**Limitations**: only tabular settings; MLE estimator needs stress-testing; hyperparameter sensitivity not fully explored" | **Q: How hard would it be to extend this to MuJoCo or other continuous-control tasks?** |
| **Slide 36 (Future Work):** | A: Conceptually it's straightforward: we keep the same dual update and shaped-reward view, but replace the tabular policy with a neural policy over continuous actions (e.g., Gaussian). The main challenges are practical: occupancy estimation in continuous spaces (we'd rely on feature-based or MLE estimators), tuning SPMA or NPG–PD in higher dimensions, and computational cost. It's doable but beyond the scope of this course project. |
| • "**Scaling up**: larger CMDPs (MuJoCo), continuous states/actions" | |
| • "**Better occupancy estimation**: MLE in high dimensions" | |
| • "**More convex objectives**: risk-sensitive, imitation learning" | **Q: What would a convergence-rate analysis of Dual–SPMA look like?** |
| • "**Theory**: convergence rates and sample complexity bounds" | A: It would likely combine: (1) regret bounds for the dual player (FTL/gradient); (2) convergence or regret bounds for SPMA as a policy optimizer under shaped rewards; and (3) error terms due to approximate occupancy estimation and finite inner loops. Those could be turned into an $\mathcal{O}(1/\sqrt{T})$ or better bound on the primal–dual gap of averaged iterates. We didn't attempt this full derivation, so we list it explicitly as future work. |

## Slide 37–38 — *Q&A & References*

| Speaker Script | Related Q&A |
|---|---|
| **Slide 37 (Q&A):** | *(Use the complete Q&A bank on the next page for any questions during Q&A session.)* |
| • "Thanks for listening—that concludes our presentation." | |
| • "We're happy to take questions." | |
| • (Remind: "Pegah: background; Danielle: Fenchel dual; Shervin: related work & method; Ahmed: experiments.") | |
| **Slide 38 (References):** | |
| • "Main references: Zahavy et al. (convex MDPs), Asad et al. (SPMA), Ding et al. (NPG–PD), Barakat et al. (MLE occupancy)." | |

# Complete Q&A Bank by Category

## 1. Big-picture & motivation (Q1–Q5)

**Q1. What's the main takeaway of your project in one sentence?**
A1. We show that if you take the convex-MDP framework of Zahavy et al. and plug in SPMA as the policy optimizer, you get a simple "Dual–SPMA" recipe that (i) is easy to implement as shaped-reward RL, and (ii) performs competitively with a strong primal–dual baseline, NPG–PD, on constrained safety tasks.

**Q2. Why do we even care about *convex* MDPs instead of just standard RL?**
A2. Standard RL maximizes a linear function of the occupancy measure—expected return—so it's great when you only care about reward. Many real problems add structure: safety constraints, exploration terms, risk penalties, imitation losses, etc. Many of those can be written as convex functionals of the occupancy $d_\pi$. Convex MDPs give a unified way to handle these goals while preserving nice optimization properties like existence of a saddle point.

**Q3. How is this different from just "constrained RL with a Lagrange multiplier"?**
A3. Classic constrained RL focuses on a particular structure: reward plus linear constraints. Convex MDPs are more general: you minimize an arbitrary convex function $f(d_\pi)$ over the convex occupancy polytope. Constrained RL is one example where $f$ adds a penalty or indicator for constraint violations; entropy-regularized RL and some imitation objectives are *other* examples. Our framework—and the Fenchel duality step—work for all of them, not just linear constraints.

**Q4. Why do you compare against NPG–PD in particular? Why not PPO or TRPO?**
A4. The core of our work is about **convex objectives & dual variables**, not generic performance on Atari. NPG–PD is a principled *primal–dual* algorithm specifically designed for constrained MDPs, with guarantees on both optimality and constraint violation. It also uses a geometry-aware policy update (natural gradient), which makes the comparison with SPMA conceptually clean. PPO/TRPO are great empirically, but they're not formulated directly as primal–dual CMDP algorithms.

**Q5. Is Dual–SPMA supposed to *beat* NPG–PD?**
A5. For a term project we don't claim to beat state of the art. What we want to show is that SPMA is a **viable policy player** inside the convex-MDP dual framework: it's stable, respects constraints, and gives reward comparable to NPG–PD. Beating a heavily tuned algorithm like PPO or NPG–PD consistently would require much more engineering and experimentation than we can do here.

## 2. Occupancy measures & convex MDP formulation (Q6–Q11)

**Q6. Why introduce occupancy measures at all? Why not just work with $J(\pi)$?**
A6. Because in terms of policy parameters, $J(\pi)$ is generally non-convex. In terms of occupancy, return is **linear**: $\langle r, d_\pi \rangle = (1-\gamma)J(\pi)$. Many interesting extras (entropy, penalties, constraints) are convex in $d$. The set of valid occupancies is a convex polytope defined by linear flow constraints. So in $d$-space we get a clean convex program: minimize a convex $f(d)$ over a convex set $\mathcal{D}$.

**Q7. Can you restate the key identity relating $J(\pi)$ and $d_\pi$?**
A7. Yes: we define

$$d_\pi(s,a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr_\pi(s_t = s, a_t = a).$$

Then

$$\langle r, d_\pi \rangle = \sum_{s,a} r(s,a)\, d_\pi(s,a) = (1-\gamma)J(\pi).$$

So up to the constant factor $1-\gamma$, the RL objective is just a linear function of the occupancy measure.

**Q8. What exactly is $\mathcal{D}$, the set of feasible occupancies?**
A8. $\mathcal{D}$ is the set of $d \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ that satisfy the Bellman flow constraints plus non-negativity: for every state $s$,

$$\sum_a d(s,a) = (1-\gamma)\rho(s) + \gamma \sum_{s',a'} P(s|s',a')\, d(s',a'), \quad d(s,a) \geq 0.$$

This says "discounted occupancy flowing *into* a state equals initial occupancy plus discounted flow from other states." It's a convex polytope.

**Q9. Is there a 1–1 correspondence between policies and occupancy measures?**
A9. Under mild regularity conditions (e.g., communicating MDP, $\gamma < 1$), each stationary policy has a unique discounted occupancy measure, and each valid occupancy in $\mathcal{D}$ corresponds to at least one stationary policy (you can recover $\pi(a|s) = d(s,a)/\sum_{a'} d(s,a')$ where the denominator is nonzero). So working in $d$-space is equivalent, up to degeneracies.

**Q10. Why do you call $f(d)$ "convex" when there's a minus reward term $-\langle r, d \rangle$?**
A10. Because $-\langle r, d \rangle$ is linear in $d$, and linear functions are both convex and concave. When we add convex terms like entropy or penalties, the sum stays convex. So all the examples on the "Convex MDP: Examples" slide fit into "convex $f$" without breaking anything.

**Q11. What assumptions do you need for your convex MDP formulation?**
A11. We assume a discounted MDP with $\gamma < 1$, bounded reward and cost, and that we restrict to stationary policies. For convexity, we require $f$ to be proper, closed, and convex on the occupancy polytope $\mathcal{D}$. For the Fenchel reduction, we also need

standard conditions that guarantee the biconjugate equality and existence of a saddle point—essentially the Fenchel–Moreau assumptions.

## 3. Fenchel duality & saddle-point formulation (Q12–Q18)

**Q12. Why bring in Fenchel conjugates at all?**
    A12. Fenchel conjugates give us an identity:

$$f(d) = \sup_y \{\langle y, d \rangle - f^*(y)\}.$$

Plugging this into $\min_{d \in \mathcal{D}} f(d)$ lets us rewrite the convex MDP as

$$\min_{d \in \mathcal{D}} \max_y \{\langle y, d \rangle - f^*(y)\}.$$

So we get a saddle-point game with a natural dual variable $y$. This structure is what allows us to run two-player algorithms: one player updates $y$, the other updates the policy.

**Q13. What does the dual variable $y$ represent intuitively?**
    A13. It's a vector of the same dimension as $d$—one coordinate per state–action pair. In the Fenchel dual, it's the variable that "linearizes" the convex objective. Operationally, for the policy player, $y$ is just a shaped reward field: we treat $r_y(s, a) = -y(s, a)$ as the reward. Different choices of $y$ encourage the policy to put occupancy where $f$ wants it.

**Q14. Why is the solution of the min–max at a saddle point?**
    A14. Because in convex-concave problems, under suitable conditions, the minimax and maximin values coincide and are achieved at a pair $(\pi^*, y^*)$ satisfying

$$L(\pi^*, y) \le L(\pi^*, y^*) \le L(\pi, y^*)$$

for all $\pi, y$. That point is simultaneously optimal for the min player (policy) and the max player (dual), and is called a saddle point. Solving the min–max is thus equivalent to finding such a saddle point.

**Q15. Does your algorithm *guarantee* convergence to the saddle point?**
    A15. In the idealized tabular, exact-oracle setting, existing theory (Zahavy et al. + SPMA paper) shows that if the dual player and the policy player both use low-regret algorithms, the average iterates converge to a saddle point. In our implementation we have function approximation, noisy occupancy estimates, and finite inner loops, so we don't claim a formal convergence proof. Instead we check empirically that the saddle objective and constraints behave sensibly.

**Q16. Why is alternating updates (policy step then dual step) reasonable, given that vanilla Gradient Descent–Ascent can diverge?**
    A16. It's true that naive GDA can oscillate. We're not using plain GDA: we use specific online-learning style updates. The dual player uses a simple FTL/gradient step in a convex Euclidean space; the policy player uses SPMA, which is a mirror-descent method in logit space with convergence guarantees. Online saddle-point theory says that when both players have low regret, the average of their strategies converges to a saddle point, so alternating these particular updates is justified.

**Q17. How do you get the conjugate $f^*(y)$ in the entropy-regularized case?**
    A17. For

$$f(d) = -\langle r, d \rangle + \alpha \sum d \log d,$$

you can complete the square in log-space / use standard convex analysis to show

$$f^*(y) = \alpha \log \sum_{s,a} \exp((y(s, a) + r(s, a))/\alpha).$$

Its gradient is a softmax: $\nabla f^*(y) = \text{softmax}((y + r)/\alpha)$. That's the expression we use in the dual update.

**Q18. Why do you write $\min_d \max_y$ instead of $\max_y \min_d$? Are you allowed to swap min and max?**
    A18. We derive the saddle-point form by directly applying Fenchel–Moreau:

$$f(d) = \sup_y \{\langle y, d \rangle - f^*(y)\},$$

then taking $\min_d$ outside. Under standard convexity/closedness assumptions, there is no duality gap and the saddle-point value equals both the minimax and maximin. So while we don't literally swap orders in the algebra, we rely on these conditions to interpret the solution as an equilibrium of the game.

## 4. Shaped-reward RL and "why RL" (Q19–Q21)

**Q19. Why not just solve the convex program over $d$ with CVX or some LP solver?**
    A19. You *could* if the MDP is small and fully known. In our setting we assume a **model-free** RL scenario: we don't know the transition probabilities, we only interact with the environment by sampling trajectories. That means we can't write down $\mathcal{D}$ explicitly or compute exact expectations. RL gives us a way to approximate $\min_\pi \langle y, d_\pi \rangle$ from samples by treating $-y$ as a reward and running a policy optimizer.

**Q20. Does the shaped reward $r_y = -y$ break the Markov property or anything?**

A20. No, because $y$ is a function of the current state–action pair $(s, a)$ only; it doesn't introduce dependence on the entire trajectory. The environment dynamics $P(s'|s, a)$ are unchanged, and the reward at each step is just a different function of $(s, a)$. So the process is still a valid MDP, just with a different reward.

**Q21. If you can plug "any RL algorithm" into the policy player, why specifically SPMA?**

A21. The convex-MDP theory says any RL algorithm could in principle be used as the policy player. We choose SPMA because (i) it has **strong convergence results** in tabular MDPs, (ii) it's geometry-aware, operating in logit space with a softmax mirror map, (iii) its loss looks similar to PPO/MDPO, so it's practical to implement, and (iv) the SPMA paper already shows it's competitive with PPO/TRPO in continuous control.

## 5. SPMA details & comparison to other policy gradients (Q22–Q27)

**Q22. How is SPMA different from vanilla policy gradient?**

A22. Vanilla PG minimizes $L_{\text{PG}} = -\mathbb{E}[\log \pi(a|s) A(s, a)]$ and takes gradient steps in parameter space. SPMA instead uses an **exponential-gradient / mirror-descent update in logit space**. The actor loss includes a regularizer

$$\frac{1}{\eta}(\exp(\Delta \log \pi) - 1 - \Delta \log \pi)$$

which penalizes large changes in log-probabilities. This gives better control over how far the policy moves per step and leads to faster convergence in the tabular theory.

**Q23. Where does your SPMA loss come from?**

A23. It's the convex surrogate corresponding to the tabular SPMA update. If you linearize the value function and consider an exponential-gradient step in logit space, the associated Bregman divergence gives rise to exactly that $\exp -1 - \Delta \log \pi$ term. The SPMA paper derives it formally; we implement it as a differentiable loss so we can use standard automatic differentiation.

**Q24. Do you enforce the $[1 + \eta A]_+$ and normalization from the tabular update?**

A24. In the tabular theory, the closed-form update is

$$\pi_{t+1}(a|s) \propto \pi_t(a|s)[1 + \eta A(s, a)]_+.$$

In our neural approximation we don't directly enforce this per-state formula. Instead we optimize the SPMA loss over minibatches of trajectories; the softmax layer plus the KL-like regularizer implicitly enforce a similar behaviour. So we're using the **function-approximation version** of SPMA instead of the exact tabular update.

**Q25. How do you choose the SPMA step size $\eta$?**

A25. We use a simple Armijo line search over a small discrete set of candidate $\eta$ values. For each candidate we compute the SPMA loss on a minibatch and pick the largest $\eta$ that produces a sufficient decrease. This is inspired by the SPMA paper and by MDPO/TRPO line-search practices.

**Q26. Why is SPMA "geometry-aware"?**

A26. Because it does mirror ascent with the log-sum-exp mirror map, which is tailored to the simplex geometry. Instead of taking Euclidean steps in parameter space and then projecting back onto the simplex, SPMA treats the logits as dual variables and the softmax as the mirror map. This respects the probability constraints and tends to produce well-behaved updates when probabilities are near 0 or 1.

**Q27. How does SPMA compare to natural policy gradient conceptually?**

A27. Both are trying to account for the geometry of the policy space. Natural PG uses the Fisher information matrix to modify the gradient direction and step size; SPMA uses a Bregman divergence in logit space. In small tabular problems you can show they are closely related; in our setting we treat them as two reasonable but different choices of geometry-aware update and compare them empirically.

## 6. Dual-update & FTL choice (Q28–Q30)

**Q28. Why do you use a simple gradient/FTL step for the dual variable $y$?**

A28. The dual variable lives in a straightforward Euclidean space and the dual objective is convex with a simple gradient $\hat{d}_\pi - \nabla f^*(y_k)$. For such problems, basic online convex optimization algorithms like gradient ascent or FTL are known to have sublinear regret. They're easy to implement and come with clear theoretical guarantees, so there's no need for more complex methods.

**Q29. Are there any constraints on $y$? Do you project it?**

A29. In the entropy-regularized case, $y$ is unconstrained—we don't project. For CMDPs we sometimes parameterize the dual as $y_\lambda = \lambda c - r$ with $\lambda \geq 0$, so the only constraint is non-negativity on $\lambda$, which we enforce with a $[\cdot]_+$ projection (same as NPG–PD). More sophisticated bounds or regularizers on $y$ are possible but we didn't explore them in this project.

**Q30. Is the dual step size $\alpha$ sensitive?**

A30. Yes, like most gradient methods it needs to be in a reasonable range: too small and progress is slow; too large and you can oscillate. In practice we tuned $\alpha$ on a log scale and picked values where the dual objective and the constraint violation decreased smoothly. Because our setting is tabular and relatively small, this wasn't too painful.

# 7. Occupancy estimation (Q31–Q34)

**Q31. Why not use exact occupancies in FrozenLake instead of Monte Carlo?**
A31. We intentionally use Monte Carlo to simulate the **model-free** setting most RL algorithms operate in—where you don't have access to $P(s'|s,a)$. It also keeps the pipeline more consistent with what you'd do in larger environments. Using the exact occupancy is possible in FrozenLake, but we'd then be solving a much easier problem that doesn't stress the estimator.

**Q32. How do you ensure your Monte Carlo occupancy estimate is valid?**
A32. We compute

$$\hat{d}_\pi(s,a) = \frac{1-\gamma}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \gamma^t \mathbf{1}\{s_t^{(i)} = s, a_t^{(i)} = a\}.$$

By construction this is non-negative. As a sanity check, we track $\sum_{s,a} \hat{d}_\pi(s,a)$ over training and verify it stays very close to 1; that's the plot on our entropy-regularized results slide. Small deviations come from truncating trajectories at finite $T$.

**Q33. What's the point of the MLE-based occupancy estimator?**
A33. The tabular MC estimator's variance scales with $|\mathcal{S}||\mathcal{A}|$; in large MDPs that's bad. The MLE estimator instead fits a log-linear model $\lambda_\omega(s,a) \propto \exp(\omega^\top \phi(s,a))$ by maximum likelihood on samples. Its error depends mostly on the feature dimension, not the raw state–action count. So it's a more scalable idea we wanted to prototype, although we only briefly tested it here.

**Q34. Do you actually use the MLE estimator for the main results, or just MC?**
A34. No, we use MC estimator for main FrozenLake experiments. MLE implemented and tested briefly; full comparison is future work.

# 8. Experiments & empirical behaviour (Q35–Q40)

**Q35. Why did you choose FrozenLake 4×4 as your main environment?**
A35. Wanted **tabular** CMDP to (i) compute/visualize occupancies, (ii) understand safety constraint intuitively (holes vs safe), (iii) debug without function approximation error. Standard simple gridworld.

**Q36. Do both methods actually satisfy the safety constraint?**
A36. Yes. Both reduce $J_c(\pi) - \tau$ toward zero. Violation plots approach/oscillate around zero. Learned policies' heatmaps clearly avoid unsafe states.

**Q37. Which method converges faster in your experiments?**
A37. Depends on hyperparameters. Typically SPMA takes larger, more "aggressive" steps, can reduce violation faster initially, but slightly more sensitive to hyperparameters. NPG–PD smoother but needs careful tuning of $\beta$. Final performance broadly comparable.

**Q38. How many trajectories / steps do you need per outer iteration?**
A38. 2048 environment steps per rollout per iteration. Enough for stable occupancy estimation and critic training in this small environment. Didn't push sample efficiency; systematic study would vary this.

**Q39. Are the advantages estimated with a critic network or Monte Carlo?**
A39. Standard actor–critic: critic is value network trained with TD, advantage computed as $A = \hat{Q} - V$ via GAE or simple TD returns. Both SPMA and NPG–PD share same critic architecture.

**Q40. How do you handle exploration in these experiments?**
A40. Stochastic policy (softmax) naturally explores early. Entropy-regularized experiments explicitly add entropy term. No extra exploration bonuses.

# 9. Limitations & future work (Q41–Q44)

**Q41. What are the main limitations of your project?**
A41. The big three are: (1) we only test in small tabular environments, so we don't claim scalability yet; (2) our MLE estimator is implemented but not fully stress-tested on large problems; and (3) we don't provide a formal convergence-rate analysis for the full Dual–SPMA algorithm with function approximation—only for the inner SPMA and the baseline, via existing theory.

**Q42. How hard would it be to extend this to MuJoCo or other continuous-control tasks?**
A42. Conceptually it's straightforward: we keep the same dual update and shaped-reward view, but replace the tabular policy with a neural policy over continuous actions (e.g., Gaussian). The main challenges are practical: occupancy estimation in continuous spaces (we'd rely on feature-based or MLE estimators), tuning SPMA or NPG–PD in higher dimensions, and computational cost. It's doable but beyond the scope of this course project.

**Q43. Could you use other policy optimizers in place of SPMA?**
A43. Yes. The meta-algorithm from Zahavy et al. lets you plug in any RL method as the policy player. We chose SPMA because of its theory and similarity to modern methods, but in principle you could try PPO, TRPO, or MDPO. An interesting piece of future work would be a systematic empirical comparison of different policy players inside the same convex-MDP dual framework.

**Q44. What would a convergence-rate analysis of Dual–SPMA look like?**

A44. It would likely combine: (1) regret bounds for the dual player (FTL/gradient); (2) convergence or regret bounds for SPMA as a policy optimizer under shaped rewards; and (3) error terms due to approximate occupancy estimation and finite inner loops. Those could be turned into an $\mathcal{O}(1/\sqrt{T})$ or better bound on the primal–dual gap of averaged iterates. We didn't attempt this full derivation, so we list it explicitly as future work.

## 10. Misc / notation and implementation details (Q45–Q48)

**Q45. What's the difference between $y$ and $\lambda$ in your notation?**

A45. $y$ is the general Fenchel dual variable over state–action pairs in the convex-MDP formulation. $\lambda$ is the scalar Lagrange multiplier for CMDP safety constraints. In the constrained case we sometimes define a structured dual variable $y_\lambda(s, a) = \lambda c(s, a) - r(s, a)$, so $y$ is parameterized by $\lambda$.

**Q46. Why does your "Convex MDP examples" slide use $\mu$ as a penalty weight instead of $\lambda$?**

A46. We intentionally use $\mu$ there to distinguish a **fixed penalty weight** in a penalized objective from the **Lagrange multiplier** $\lambda$ that evolves as a dual variable. This avoids confusion between "tunable hyperparameter" and "optimization variable".

**Q47. What exactly are $J_r(\pi)$ and $J_c(\pi)$?**

A47. They are just discounted returns with respect to reward and cost: $J_r(\pi) = \mathbb{E}_\pi[\sum_t \gamma^t r(s_t, a_t)]$, $J_c(\pi) = \mathbb{E}_\pi[\sum_t \gamma^t c(s_t, a_t)]$. We estimate them by Monte Carlo averaging over rollouts.

**Q48. How do you measure "success" in your experiments?**

A48. We consider three main criteria: (1) **constraint satisfaction**—$J_c(\pi)$ should be at or below $\tau$; (2) **reward**—$J_r(\pi)$ should be as high as possible given the constraint; and (3) **stability**—dual variables and occupancies shouldn't blow up or oscillate wildly. On those metrics, Dual–SPMA behaves about as well as NPG–PD in our tests.