Deep Q-Life: Markov, Bellman, and Habits

Dr. Yves J. Hilpisch¹ with GPT-5

November 8, 2025 (preliminary draft)

Abstract

We recast everyday habit building as a small, friendly reinforcement learning problem. Using the Markov decision process (MDP) lens to define state, action, reward, and transitions, the Bellman principle provides the north star: choose the option that is good now and sets up good options next. Deep Q-Learning (DQL) contributes a simple learning loop—try a bit, remember what worked, and stabilize goals with weekly reviews. We keep the math precise but lightweight, pair each concept with relatable examples (cake vs. walk, saving vs. spending, texting vs. sleep), and add visual intuitions and checklists. The result is a practical toolkit for designing compact policies you can practice, review, and improve—light-hearted in tone, formal enough to be science.

¹Get in touch: https://linktr.ee/dyjh. Web page https://hilpisch.com.

Contents

1	The Big Idea	1
2	Life Is (Mostly) Markovian	2
3	Foundations: MDPs and Bellman	3
4	From Equations to Everyday Decisions	5
5	Deep Q-Learning for Humans	6
6	Building Habits via the Bellman Principle	7
7	Case Study: 9-Ball and the Bellman Principle	9
8	Case Study: A Week of Decisions	11
9	Case Study: Quitting Drinking with DQL	12
10	Practical Toolkit	14
11	Limits, Ethics, and Human Factors	15
12	Related Work	16
13	Conclusion	16

1 The Big Idea

This paper makes a simple, slightly cheeky promise: you can borrow serious mathematical ideas—Markov decision processes, the Bellman principle, and Deep Q-Learning (DQL)—to design better everyday decisions and habits. We keep the tone light-hearted, but the foundations formal enough to be science, not vibes.

Why this matters: life is a stream of states (your health, finances, relationships, energy) and actions (eat cake, go for a walk, save or spend, send the text, hit snooze). Rewards show up now (yum, cake) or later (stable weight, rainy-day savings). The Markov lens says: what matters for the next move is mostly your current state; the Bellman lens says: decide today as if you will decide optimally tomorrow; DQL adds: learn good choices by mixing exploration and exploitation, remembering what worked, and stabilizing your learning.

What you will get from this paper:

- A crisp mapping from life to an MDP: states, actions, rewards, transitions, and discounting.
- An intuitive explanation of the Bellman optimality principle and how it becomes a habit-building rule of thumb.
- A human-friendly take on DQL: exploration (try things), exploitation (use what works), replay (reflect), and target networks (keep goals stable).
- A practical scaffolding: checklists, templates, and a week-long case study to try at home.
- A balanced view of limits: emotions, uncertainty, non-Markovian baggage, and ethical guardrails.

Who this is for: curious practitioners who like both rigor and relatable examples. No background in RL required; we will separate formal statements from analogies and provide an appendix for the math-inclined.

How to use this: skim the narrative sections for ideas, then adopt the "toolkit" section to structure one habit you care about (health, money, relationships, focus). Iterate weekly, treating your life like a friendly learning system rather than a pass/fail test.

Running example (to keep us honest): "The Cake, the Coach, and the Credit Card." We will revisit three micro-scenarios—dessert vs. gym, saving for a treat next week, and texting before bed—to show how the same Bellman logic and DQL patterns travel across domains.

At a Glance

Life as MDP (state, action, reward); act today for tomorrow's best (Bellman); learn by trying, remembering, and stabilizing (DQL). Build habits by extracting simple policies from your learned values.

2 Life Is (Mostly) Markovian

We will model everyday decision-making as a *Markov decision process (MDP)*. The Markov idea is humble and practical: once you summarize your current situation well ("the state"), the distant past adds little extra predictive power for what happens next. Breakfast five years ago matters less than your sleep, hunger, and schedule right now.

Definition 1 (Markov Property). Let $(S_t)_{t\geq 0}$ be states and $(A_t)_{t\geq 0}$ actions. The process is Markov if for all t and states s'

$$\Pr(S_{t+1} = s' \mid S_0, \dots, S_t, A_0, \dots, A_t) = \Pr(S_{t+1} = s' \mid S_t, A_t).$$

In words: given where you are and what you do, the full history need not be consulted.

Definition 2 (Markov Decision Process (MDP)). An MDP is a tuple (S, A, P, R, γ) with state space S, action space A, transition kernel $P(s' \mid s, a)$, expected reward R(s, a) (or R(s, a, s')), and discount factor $\gamma \in [0, 1)$. A policy chooses actions from states.

Mapping life into this structure is more art than surgery. A useful state captures what matters *now*: energy, mood, calendar commitments, bank balance, and social context. Actions are the options on your next step: go for a walk, eat cake, save or spend, send the message, snooze or get up, focus or multitask. Rewards can be multi-dimensional (health, money, joy), but we will often combine them into a single score to keep learning simple.

Transitions encode how today's choices shape tomorrow's situation: sleep affects energy; spending affects bank balance; late-night texting affects tomorrow's focus and relationships. The discount factor γ tunes patience: smaller γ favors immediate comfort; larger γ emphasizes compounding benefits. We visualize immediate versus delayed rewards in Figure 1.

In reinforcement learning (RL) practice, we frequently do not know P or R exactly; we learn from experience. The MDP framing still helps because it tells us what to track, what to try, and how to update our expectations.

Approximate Markovization: A Mini-Checklist

To make the Markov assumption useful in daily life:

- Include observable signals with real leverage (sleep hours, meeting density, step count, cash-on-hand).
- Track slowly changing resources (energy, time budget, money) and immediate context (location, people).
- Keep the state small; if you can't remember or measure a feature easily, drop it.
- Encode commitments and hard constraints (deadlines, no-spend categories) explicitly.
- Standardize scales (e.g., 0–10 energy) so rules of thumb transfer across days.

Running example: dessert vs. gym. Define a compact state

s = (energy, hunger, time-left, steps-today)

Two actions: $a \in \{\text{cake, walk}\}$. Rewards: $r_{\text{cake}} = \text{taste-future-weight-cost}$; $r_{\text{walk}} = \text{mood-boost+health-credit}$. The transition increases steps if you walk and may reduce energy; cake increases satiety now but nudges health downward tomorrow. With this framing, writing and improving a policy ("if low energy and little time, take a 10-minute walk; otherwise skip dessert") becomes natural.

Visualization: Immediate vs. Delayed Rewards

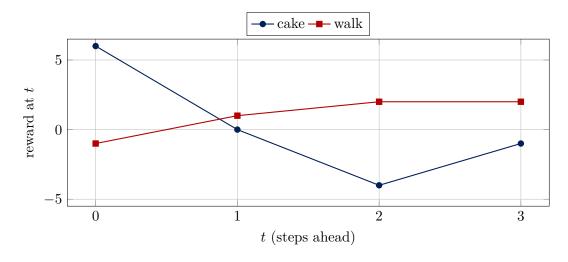


Figure 1: Reward timing matters: cake pays now, costs later; a short walk can be the reverse. Discounting in Section 3 weights these series into comparable values.

3 Foundations: MDPs and Bellman

We briefly collect the formal pieces that power the rest of the paper. An MDP (S, A, P, R, γ) describes states, actions, dynamics P, rewards R, and a discount factor $\gamma \in [0, 1)$. A policy $\pi(a \mid s)$ selects actions from states.

Value functions. For a policy π , the state-value and action-value functions are

$$V^{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^{t} R(S_{t}, A_{t}) \mid S_{0} = s \right],$$

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^{t} R(S_{t}, A_{t}) \mid S_{0} = s, A_{0} = a \right].$$

Plain Words

 $V^{\pi}(s)$: the expected total, discounted reward you collect when you start in state s and then follow policy π forever. Near-term rewards count more because of the powers of γ . $Q^{\pi}(s,a)$: the expected total, discounted reward when you start in state s, take action a now, and then follow policy π thereafter.

Symbols: \mathbb{E}_{π} averages over the randomness from transitions P and any randomness in π ; $R(S_t, A_t)$ is the immediate reward at step t; $\gamma \in [0, 1)$ discounts step-t rewards by γ^t ; the conditions $S_0 = s$ and $A_0 = a$ fix the starting point.

These satisfy the Bellman expectation equations:

$$V^{\pi}(s) = \sum_{a} \pi(a \mid s) \Big[R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) V^{\pi}(s') \Big],$$
$$Q^{\pi}(s, a) = R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) \sum_{a'} \pi(a' \mid s') Q^{\pi}(s', a').$$

Plain Words

 $V^{\pi}(s)$ recursion: average over actions the policy might take in s; for each action, add the *immediate reward* R(s,a) to the *discounted expected value* of the next state, $V^{\pi}(s')$, weighted by transition probabilities $P(s' \mid s, a)$.

 $Q^{\pi}(s, a)$ recursion: take the *immediate reward* R(s, a) and add the *discounted expectation*, over next states s' and the policy's next action a', of $Q^{\pi}(s', a')$. First you move to s' according to P, then the policy chooses a' according to π .

Optimality. Define $V^*(s) = \sup_{\pi} V^{\pi}(s)$ and $Q^*(s, a) = \sup_{\pi} Q^{\pi}(s, a)$. They solve the Bellman *optimality* equations:

$$V^*(s) = \max_{a} \left[R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) V^*(s') \right],$$
$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) \max_{a'} Q^*(s', a').$$

Any greedy policy with respect to Q^* is optimal.

Plain Words

 $V^*(s)$: the best possible total, discounted reward from state s. The recursion says: pick the action that maximizes *immediate reward plus discounted best continuation* $V^*(s')$. $Q^*(s,a)$: the best possible total, discounted reward if you take a now in s and act optimally thereafter. The recursion keeps the immediate reward, then adds the discounted best next-step value $\max_{a'} Q^*(s',a')$ averaged over next states.

Greedy optimality: if, in every state, you choose the action with the largest $Q^*(s, a)$, that policy is optimal.

Everyday Intuition

Bellman says: make today's choice by combining today's payoff with tomorrow's best continuation. In practice: pick the action that is good now and sets you up for good options next.

Visualization: Discounting Emphasizes the Near Future

We visualize how different discount factors weight the near future in Figure 2.

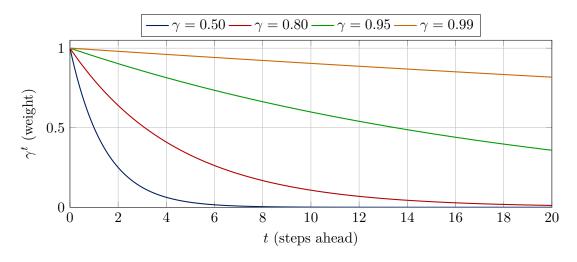


Figure 2: Heavier discounting (smaller γ) focuses learning and decisions on the near future. Larger γ values spread credit further forward in time.

4 From Equations to Everyday Decisions

We now translate the formal ingredients into a practical, human-sized workflow. The aim is not to capture every nuance of life, but to choose a compact state, a sane reward, and a few actionable policies you can practice.

State design. Keep it small, stable, and observable. Favor features you can sense or measure in under 10 seconds (e.g., energy 0–10, hunger 0–10, calendar load 0–10, steps today, cash-on-hand). If a feature rarely changes your choice, drop it.

Reward design. Combine immediate and delayed consequences into a single score that reflects your values. Use units you care about (wellbeing points, money, time saved). Be wary of "reward hacking" where a proxy crowds out the real goal (e.g., step count without sleep quality). See Figure 3 for a simple decomposition.

Constraints and context. Encode hard constraints (budgets, time windows, commitments) and soft norms (e.g., no screens after 22:00). Policies become simpler when constraints rule out brittle choices.

Friction and cues. Many transitions are habit loops in disguise: $cue \rightarrow routine \rightarrow reward$. Reduce friction on good routines (lay out shoes) and add friction to risky ones (hide the sweets).

Choosing γ . Your discount factor expresses patience: higher γ values prefer compounding habits; lower values prefer quick comfort. You can even set domain-specific γ (e.g., higher for health than for entertainment).

Visualization: Reward Shaping vs. Hacking (Waterfall)

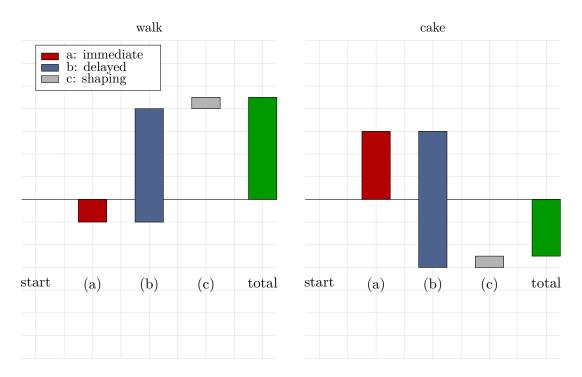


Figure 3: The figure shows how immediate and future rewards and penalties lead to net total rewards or penalties. Or as they say, that "The costs of your good habits are in the present. The costs of your bad habits are in the future.", see [6]. Left: walk (a: -1, b: +5, c: +0.5, total = 4.5). Right: cake (a: +3, b: -6, c: +0.5, total = -2.5). Colors: a red, b blue, c grey, total green.

5 Deep Q-Learning for Humans

Deep Q-Learning (DQL) extends tabular Q-learning by using a function approximator (e.g., a neural network) to estimate Q(s,a) from features of s. You do not need a literal network; a compact set of rules of thumb plays the same role for everyday decisions.

Update rule. At each experience (s, a, r, s'), form a temporal-difference (TD) target

$$y = r + \gamma \max_{a'} Q_{\text{target}}(s', a')$$

and then adjust the online estimator Q_{online} to reduce the squared error $(Q_{\text{online}}(s, a) - y)^2$ (using a learning rate). The *target network* is a slowly updated copy of the online estimator, which stabilizes learning by keeping the target y steadier. We visualize this stability in Figure 4.

Visualization: Stable Targets Tame Oscillation

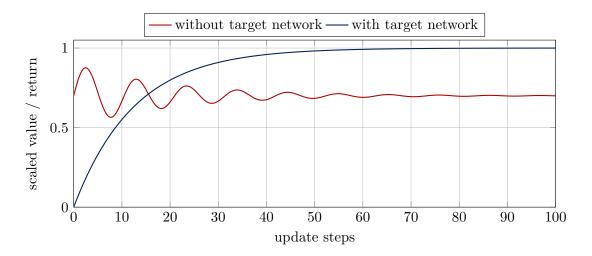


Figure 4: A slowly updated target makes the TD target steadier, reducing oscillations and divergence in value estimates. Synthetic illustration.

Plain Words

TD target: today's payoff r plus a discounted best guess of tomorrow's value.

Online vs. target: your current heuristics learn from a slowly moving goalpost (weekly targets), not from themselves instant-to-instant.

Experience replay: learn from a shuffled batch of recent situations (your journal), which breaks up streaks and reduces overfitting to today's mood.

Epsilon-greedy: most of the time do the current best action; sometimes deliberately try something else to discover better options.

Mini-protocol for a week.

- 1. Define a small state and actions for one domain (Section 2).
- 2. Set a reward rule that reflects values (Section 4).
- 3. Each day: $\log 3-5$ decisions (s, a, r, s') with one-sentence notes.
- 4. Nightly: sample 5–10 past entries (replay), update heuristics toward the TD target.
- 5. Weekly: update goals (target network), review one new experiment, lower ε slightly.

6 Building Habits via the Bellman Principle

Bellman gives us a crisp rule for habits: act today as if you will keep acting well tomorrow. In learning terms, once your action-values Q(s,a) are reasonably shaped, you can extract a simple policy and make it automatic.

From Q to policy. A pragmatic approach:

- 1. Pick a compact state (Section 2) and a small action set (e.g., walk, cake, stretch, water).
- 2. Track a week of experiences and update a human-sized Q estimate (Section 5).
- 3. For common states, write *if*-then rules that choose $\arg \max_a Q(s, a)$.

4. Add one sentence of why (links to values) to aid memory and motivation.

Example: "If energy ≥ 4 and time-left ≥ 10 minutes, take a 10-minute walk (maximizes mood/health value). Otherwise, drink water and stretch for 2 minutes." We visualize the idea of "choose the higher Q" in Figure 5.

Visualization: Policy from Q

We compare two simple action-value curves Q_{walk} and Q_{cake} as a function of energy. The greedy policy picks the higher line at each energy level; see Figure 5.

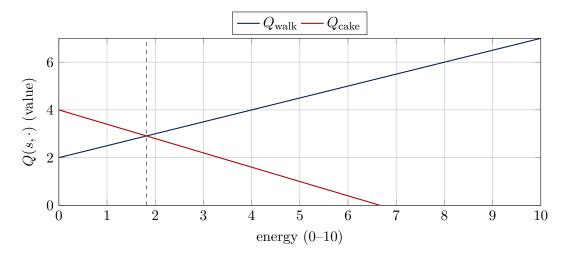


Figure 5: Policy from Q: choose the action with the larger predicted value. For this toy model, below energy ≈ 1.8 the cake's immediate payoff dominates; above it, walking wins on total value.

Making habits stick. Prediction is glue. When your brain can predict a small, reliable reward from a routine, it becomes easier to start. Design cues (time, place, preceding action), reduce friction (shoes out), and protect the first minute (start tiny). Use precommitments (calendar blocks, snack-free zone) to keep policies deterministic when willpower is low.

Habit Recipe

When *state cue* occurs, I *do action* because it scores well on *my values*. Then I record a 1-line note. Weekly, I replay a few notes, nudge the rule, and keep going.

We also visualize the cue-routine-reward loop that underlies many habits in Figure 6.

Visualization: Cue-Routine-Reward Loop

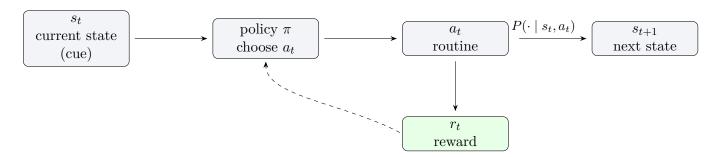


Figure 6: Habit loop as an MDP: the current state (cue) flows through a policy to an action (routine), yielding a reward and a new state. Rewards and reflections nudge future policy updates.

Domain mini-policies.

- Health: sleep-first; 10-minute walk if energy ≥ 4 ; sweets after protein.
- Money: 24-hour delay on discretionary buys; auto-transfer to savings on payday.
- Relationships: assume positive intent; weekly check-in; no late-night texting.
- Productivity: 5-minute start; single-task sprint; phone in another room.

7 Case Study: 9-Ball and the Bellman Principle

In 9-ball pool, you must pocket balls in order (1 through 9). Great players don't just pocket the current ball—they *position the cue ball* to make the *next* shot easy. That's Bellman in the wild: choose the shot that maximizes "good now + sets up good next."

Mapping to an MDP.

- State s: positions of the cue ball and legal object ball, plus simple features (angle, distance, rail clearance).
- Action a: choose aim angle and stroke speed (coarsened into a few options: soft stun, follow, draw, two-rail).
- Reward r: immediate pocket success (1 if potted, else 0), plus a small bonus for landing the cue ball inside a "position region" for the next ball.
- Transition: deterministic-ish physics with noise (cloth/friction, tip error) mapping to a new cue-ball position and next legal object ball.

We sketch a toy table with a target pocket and a shaded position region in Figure 7.

Visualization: Position Play on the Table

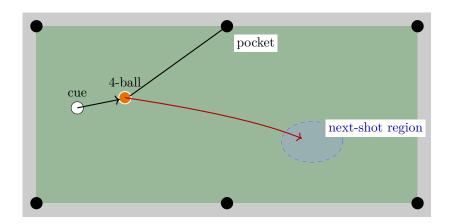


Figure 7: Bellman on the felt: pocket the current ball while positioning the cue ball inside a high-value region for the next shot. The best shot maximizes immediate success and good continuation.

A second plot shows the trade-off between pot difficulty and cue-ball position quality as stroke power increases, and their Bellman-style combination in Figure 8.

Visualization: Pot vs. Position Trade-off

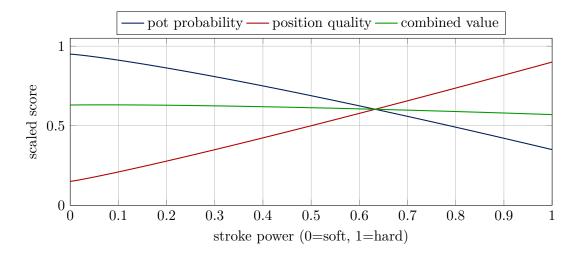


Figure 8: Trade-off: a harder stroke may reduce pot probability but improve cue position for the next shot. The Bellman view favors the stroke that optimizes the combined value.

We also include a classic two-rail position pattern in Figure 9 to illustrate how intentional cue-ball routes set up the next shot.

Visualization: Two-Rail Position Variant

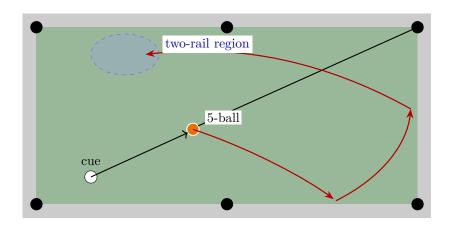


Figure 9: Two-rail position: choose spin and speed to send the cue ball via two cushions into a high-value region for the next shot. This often sacrifices a little pot margin now to make the next shot trivial—pure Bellman.

8 Case Study: A Week of Decisions

We run a seven-day loop on one domain (health) with the dessert-vs-walk decision as our focal action. State features: energy (0-10), hunger (0-10), time-left (minutes), steps-today. Actions: cake, walk, stretch. Reward: mood + health credit - future weight cost (scaled to 0-6).

Protocol. Monday and Tuesday emphasize exploration (higher ε); midweek consolidates what worked; the weekend reviews and updates the target goals (Section 5). We summarize daily total returns and the exploration schedule in Figure 10 and show explore vs. exploit counts in Figure 11.

Visualization: Week at a Glance

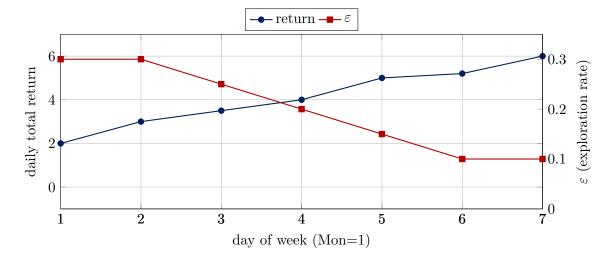


Figure 10: Daily total returns trend upward as exploration ε decays through the week. Values are illustrative but consistent with the protocol.

Visualization: Explore vs. Exploit Counts

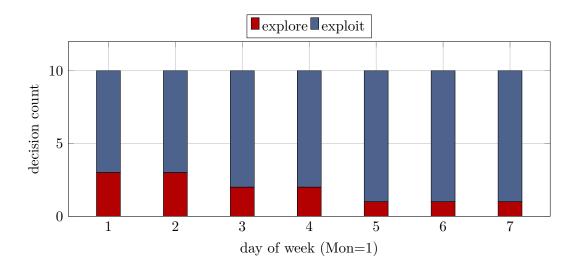


Figure 11: Exploration shrinks through the week (about 10 decisions/day). Counts are illustrative and consistent with the ε schedule.

What happened. As exploration decays, returns trend upward. A small shaping reward for daylight walks nudged choices without flipping preferences. The week ends with a compact policy and a slightly lower exploration rate for the next week.

Practical Outcomes

By Sunday you should have: (1) one or two confident if—then rules, (2) a realistic exploration rate for next week, and (3) a short list of cues and friction tweaks to make the good action the easy action.

9 Case Study: Quitting Drinking with DQL

This case study shows how to apply the same math-based habit philosophy to reducing or quitting alcohol when it is a social/psychological habit rather than a medical dependence.²

The pattern. Drinking often yields high immediate rewards (relaxation, social ease) with meaningful delayed costs (sleep disruption, low next-day energy, mood dips, money, strained commitments). Bellman's guidance remains the same: pick the option that is good now and sets you up for good options next—often an alcohol-free alternative plus a friendly exit plan.

Mapping to an MDP.

- State s: time of day; day of week; location; social context (alone/friends/work); stress 0–10; energy 0–10; sleep score; cash-on-hand; planned early obligation (boolean).
- Actions a: abstain; alcohol-free drink; one drink then water; leave early; text a buddy; breathing or short walk; order food first.

²Alcohol use can be serious. If you suspect dependence, withdrawal risk, or alcohol use disorder, seek medical advice. Abrupt cessation can be dangerous for some people; talk to a clinician or addiction professional. Community and peer support can help (e.g., local counseling services or support groups).

- Reward r: immediate social/mood benefit minus delayed penalties (sleep debt, money spent, next-day energy/mood, missed workout). Add small shaping for aligned identity (e.g., "kept promise to self").
- Transitions P: late-night drinking usually reduces sleep and next-day energy; AF choices often preserve sleep and tomorrow's options (gym, early meeting, family plans).
- Discount γ : for health and relationships, a higher γ emphasizes compounding benefits from many sober nights.

Policy in practice (narrative). Friday, 19:00. A message pops up: "Happy hour?" You feel the familiar mix: a little stress from the week, a little thrill at the first sip, and a small voice that knows tomorrow's run will be better if you skip. You reply "See you there," and decide to position the next shot—Bellman on the dance floor, not the pool table. At the bar you order an AF drink you actually like, tell yourself, "I can always change my mind in 20 minutes," and text a friend, "Hold me to one AF round." The first urge passes; the room is still friendly. When someone jokes, "Come on, live a little," you smile: "I'm on a sleep streak." You are not arguing, just narrating your plan. Later, when boredom flickers, you step outside for two minutes, breathe, and check HALT (hungry, angry, lonely, tired). "Hungry" lights up; a snack helps more than a buzz. At 22:00 you leave with Future You in mind—tomorrow's shot is a hanger.

Common Thoughts and Helpful Counters

- "If I don't drink, the night will be boring." \rightarrow "Boredom is a cue, not a verdict; try a micro-challenge: new conversation, AF taste-test, quick dance, or step outside and re-evaluate in 10 minutes." Decision rule: start AF, set a 20-minute recheck.
- "They'll think I'm judging them." → "People care more about their glass than mine." Decision rule: hold an AF drink, change the subject, offer a toast. One-liner: "I'm chasing great sleep tonight."
- "One won't hurt." → "Future me pays compound interest." Decision rule: if early obligation is true or stress > 5, choose AF; otherwise delay 10 minutes and eat first.
- "I already had one; might as well keep going." \rightarrow "Reset now beats reset tomorrow." Decision rule: water + food, text a buddy, switch context, and set a leave time.
- "I'm more fun with a buzz." \rightarrow "Relaxation is trainable." Decision rule: aim for two AF social reps this week and rate fun the next morning (reward comes in the review).
- "I don't want to explain." \rightarrow "Pre-commit a one-liner." Decision rule: rehearse a phrase that fits your identity (athlete, early riser, driver).

Friction and cues. Treat the scene like gentle choice architecture. Arrange an AF option you genuinely like, bring it if needed, and stand where ordering is slower. Put spending cash in a small separate wallet, set a rideshare cutoff, and tell one supportive friend your plan. Reduce strong cues (bar counter, tequila round), amplify identity cues ("driver," "runner," "early meeting").

A light DQL loop. Keep a tiny replay buffer (3–5 lines/night: state, action, reward, next-morning feel, notable thought). In the weekly review, sample 10 entries: What surprised me? Which counter-thought worked? Update one rule and stabilize goals (target network). Decay

exploration (fewer "just testing" sips) as confidence grows. Name urges as passing states ("a wave") and practice *urge surfing*: notice, breathe, choose; the value target lives in tomorrow morning's mood and energy.

Visualization: Abstain vs. Drink (Waterfall)

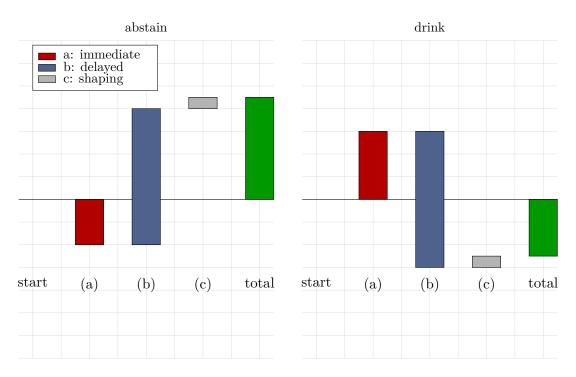


Figure 12: Immediate versus delayed consequences for alcohol choices. Left (abstain): a = immediate unease/effort to say no and not-yet-relaxed feeling; b = great sleep, little spend, clear morning, better training/fitness; c = small identity reward ("kept promise to self"). Right (drink): a = immediate buzz and social ease; b = fragmented sleep, lower next-day energy and mood, extra spend, higher chance to derail plans; c = small ritual reward ("Friday vibe"). Numeric example: abstain a: -2, b: +6, c: $+0.5 \Rightarrow$ total 4.5; drink a: +3, b: -6, c: $+0.5 \Rightarrow$ total -2.5. Echoing [6]: the costs of good habits are in the present; the costs of bad habits are in the future.

Takeaways

Treat evenings as a sequence of states and small actions. Make the good action easy (AF first, exit plan), keep a gentle replay buffer, and let Bellman's "good now + good next" steer you—all while keeping compassion and medical prudence in the loop.

10 Practical Toolkit

This section provides simple artifacts that make the framework usable in everyday life.

One-page template. Capture: (1) 4–6 state features with scales, (2) 3–5 available actions, (3) reward rule, (4) constraints/guardrails, (5) discount γ , and (6) one sentence of purpose.

Daily scorecard. For 3–5 key decisions per day, record (s, a, r, s') in one line. Add a two-word reason (e.g., "low energy"). This is your experience replay buffer.

Weekly review checklist. Sample 20 past entries, ask "what surprised me?", adjust one rule, and refresh goals (target network). Lower ε slightly if returns are stabilizing.

Exploration budget. Decide experiments/week you can sustain (e.g., 3 tiny tries). Tie this to a decaying ε schedule so you explore less as the policy improves. We visualize a sample plan in Figure 13.

Guardrails. Define no-go actions in advance (e.g., no late-night shopping). Use defaults that make the good action easy (fruit bowl in sight) and the risky one inconvenient (sweets on a high shelf).

Visualization: Exploration Budget by Week

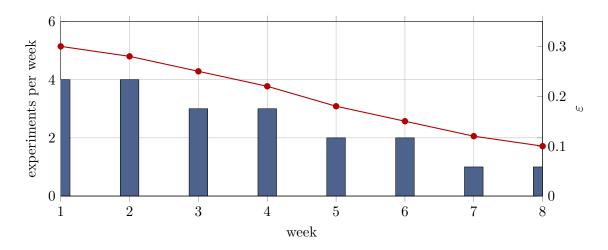


Figure 13: An exploration plan couples a shrinking ε with a realistic number of weekly experiments. Early weeks try more small changes; later weeks consolidate.

11 Limits, Ethics, and Human Factors

Real life is messy. The MDP lens is useful, but we should respect its limits and build in safeguards.

Partial observability. Some important variables are hidden (stress, hormones, social dynamics). Your state is an approximation, so value estimates have uncertainty. Treat confident rules differently from shaky ones. We illustrate uncertainty bands in Figure 14.

Reward hacking. If a proxy reward can be gamed (steps without sleep, inbox zero without outcomes), the policy may drift away from true goals. Add audits: weekly ask "did the proxy serve the purpose?" and shrink shaping if it distorts choices (Figure 3).

Fairness and consent. Optimize with care in relationships and teams. Use policies that respect boundaries and avoid nudging others without consent. Compassion outranks cleverness.

Cognitive load. Keep the model small. When in doubt, simplify states and default to a safe policy. Automate with habits so the system runs even on low-energy days.

Visualization: Uncertainty in Estimated Value

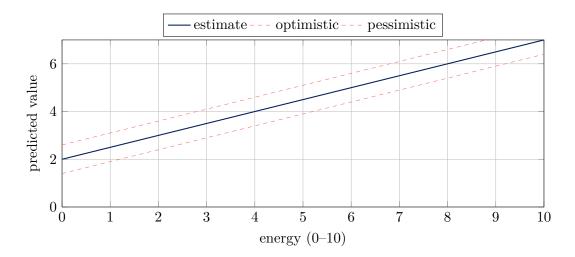


Figure 14: Uncertainty band for a value estimate: optimistic and pessimistic lines around a central prediction. Use caution near crossings or when the band is wide.

12 Related Work

This paper stands on three pillars: control theory and dynamic programming, reinforcement learning (RL), and the behavioral sciences on habits and choice.

Dynamic programming and MDPs. The Bellman principle traces to Bellman's original work on dynamic programming [1]. For formal MDP treatments and solution methods beyond the scope here (policy iteration, linear programming), see Puterman [2]. These texts ground the optimality equations we use for building habits.

Reinforcement learning. Sutton and Barto [3] synthesize prediction and control with value functions and temporal-difference learning. Deep Q-Learning (DQL) popularized function approximation for Q(s,a) using deep networks; the canonical reference is Mnih et al. (Atari) [4]. Our "human DQL" adapts the same ingredients—TD targets, experience replay, and target networks—into journaling, reflection, and goal stabilization.

Habits and behavior change. Wood [5] reviews how repetition under stable cues automates behavior, aligning with our policy-from-Q view. Clear [6] offers practitioner heuristics (cues, friction, identity) that fit naturally into the MDP mapping. On nudges and defaults as low-friction policy design, see Thaler and Sunstein [7]. Intertemporal choice and discounting are discussed widely; a concise overview appears in Loewenstein and Prelec [8].

Applications. RL-style decision systems appear in operations, health, and finance; we adopt their framing but keep complexity low for personal use. The goal is not algorithmic novelty, but a faithful translation of robust ideas into day-to-day practice.

13 Conclusion

Life offers a steady stream of states and choices. The Markov lens helps us focus on what matters now; the Bellman principle aligns today's action with tomorrow's best continuation; Deep Q-Learning (DQL) supplies a practical learning loop to refine policies from experience.

With small, stable states and value-aligned rewards, you can turn this trio into habits that compound.

What to remember.

- Define a compact state and a simple reward that reflects values.
- Use the Bellman idea: good now plus sets you up for good later.
- Learn weekly: replay notes, update rules toward TD targets, stabilize goals.
- Keep exploration tiny but steady; protect yourself with guardrails.

One-Week Starter Plan

Pick one domain (health, money, focus). Define 4–6 state features and 3 actions. Write a one-line reward. For seven days, log 3–5 decisions. Each night, sample 5–10 entries and adjust one rule. On Sunday, freeze goals for next week, lower ε slightly, and keep the best rule visible. Repeat.

Appendix: Formal Bits

This appendix collects compact statements and figures that back the narrative with standard results. It can be read independently or used as a quick reference.

A.1 Bellman Operators and Contraction

Let (S, A, P, R, γ) be a finite Markov decision process (MDP) with $\gamma \in [0, 1)$. For a policy π , define the Bellman operator T^{π} on bounded functions $v: S \to \mathbb{R}$ as

$$(T^{\pi}v)(s) = \sum_{a} \pi(a \mid s) \Big[R(s,a) + \gamma \sum_{s'} P(s' \mid s,a) v(s') \Big].$$

Define the optimality operator T by

$$(Tv)(s) = \max_{a} \left[R(s, a) + \gamma \sum_{s'} P(s' \mid s, a) v(s') \right].$$

Theorem 1 (Contraction and Fixed Points). For the sup-norm $||v||_{\infty} = \max_{s} |v(s)|$ we have: (i) T^{π} and T are γ -contractions; (ii) each has a unique fixed point, namely V^{π} and V^{*} ; (iii) value iteration $v_{k+1} = Tv_k$ converges to V^{*} for any initial v_0 with error $||v_k - V^{*}||_{\infty} \leq \gamma^k ||v_0 - V^{*}||_{\infty}$.

Plain Words

Applying T repeatedly shrinks errors by a factor of γ . Start anywhere; keep applying T; you converge to the unique V^* .

A.2 Tabular Q-Learning

For visits (s_t, a_t, r_t, s_{t+1}) , the tabular update is

$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \alpha_t) Q_t(s_t, a_t) + \alpha_t \Big[r_t + \gamma \max_{a'} Q_t(s_{t+1}, a') \Big].$$

Convergence (tabular). If every state–action pair is visited infinitely often and the step sizes satisfy the Robbins–Monro conditions $\sum_t \alpha_t = \infty$ and $\sum_t \alpha_t^2 < \infty$, and if exploration is greedy in the limit (GLIE), then $Q_t \to Q^*$ with probability 1. With function approximation (e.g., neural nets), convergence is not guaranteed; target networks, experience replay, and small learning rates mitigate divergence in practice.

A.3 Example MDP: Dessert vs. Walk

We discretize a small slice of the running example. Let

$$S = \{(e, h) \mid e \in \{0, \dots, 6\}, h \in \{0, 1\}\}, \qquad A = \{\text{cake, walk}\}.$$

Rewards: $R((e, h), \text{cake}) = 31\{h = 1\} - 2$, and R((e, h), walk) = 1 + 0.6e. Transitions: energy next day drifts toward 4 with the action's effect, e.g., $e' = \min\{6, \max\{0, e + \eta_a + \xi\}\}$ with $\eta_\text{walk} = +1$, $\eta_\text{cake} = -1$, and small noise ξ . Hunger resets: h' = 0 after cake, otherwise h' flips to 1 with moderate probability. Discount $\gamma = 0.9$. This toy MDP supports value iteration and small Q-learning demos. See Figure 15 for synthetic convergence curves.

A.4 Pseudocode

The following two callout boxes summarize the core procedures we rely on in the main text. The first presents value iteration for finite MDPs; the second gives tabular Q-learning with GLIE exploration.

Value Iteration (finite MDP)

Input: T operator, tolerance ε ; Init: $v_0(s) = 0$. Repeat $k = 0, 1, 2, \ldots$ until $||v_k| + 1 - v_k||_{\infty} < \varepsilon$:

- 1. $v_k + 1 \leftarrow Tv_k$ (apply optimality operator statewise)
- 2. Optionally extract $\pi_k + 1(s) \in \arg\max_a [R(s, a) + \gamma \sum_s' P(s' \mid s, a)v_k + 1(s')]$

Return: $v_k + 1 \approx V^*$ and a greedy policy.

Tabular Q-Learning (GLIE)

Loop over episodes: for t = 0, 1, ... with ε -greedy behavior policy

- 1. Observe s, choose a by ε -greedy from Q.
- 2. Execute a, observe r, s', set $\delta \leftarrow r + \gamma \max_{a} a'Q(s', a') Q(s, a)$.
- 3. Update $Q(s, a) \leftarrow Q(s, a) + \alpha \delta$.
- 4. Decay ε slowly (GLIE) and choose $s \leftarrow s'$.

A.5 Notation

Table 1 collects symbols used throughout and their meanings.

Table 1: Core symbols and their meanings used across the appendix and main text.

Symbol	Meaning
\mathcal{S}, \mathcal{A}	state space, action space
$P(s' \mid s, a)$	transition kernel
R(s,a)	expected one-step reward
γ	discount factor in $[0,1)$
$\pi(a \mid s)$	policy (probability of action)
V, Q	state- and action-value functions
T^{π}, T	Bellman expectation and optimality operators

A.6 Convergence Illustration

Figure 15 illustrates, on synthetic data, the typical geometric contraction of value iteration and a representative decay of average TD errors for tabular Q-learning under mild conditions.

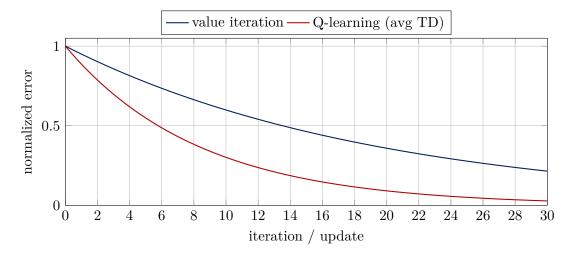


Figure 15: Synthetic convergence curves: value iteration contracts at roughly γ^k ; average temporal-difference (TD) errors in tabular Q-learning typically decay with small steps and sufficient exploration.

References

- [1] R. Bellman. Dynamic Programming. Princeton University Press, 1957.
- [2] M. L. Puterman. Markov Decision Processes: Discrete Stochastic Dynamic Programming. Wiley, 2014.
- [3] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. 2nd ed., MIT Press, 2018.
- [4] V. Mnih, K. Kavukcuoglu, D. Silver, et al. Human-level control through deep reinforcement learning. *Nature* 518, 2015.
- [5] W. Wood. Good Habits, Bad Habits. Farrar, Straus and Giroux, 2019.
- [6] J. Clear. Atomic Habits. Avery, 2018.
- [7] R. H. Thaler and C. R. Sunstein. Nudge. Yale University Press, 2008.
- [8] G. Loewenstein and D. Prelec. Anomalies in intertemporal choice: Evidence and an interpretation. *The Quarterly Journal of Economics*, 107(2), 1992.