# Chapter 1: Exercises & Labs

## Contents

# 1 Chapter 1 — Exercises & Labs (Application Mode)

Reward design is now backed both by the closed-form objective (Chapter 1, EQ-1.2) and by executable checks. The following labs keep theory and implementation coupled.

## 1.1 Lab 1.1 — Reward Aggregation in the Simulator

Goal: inspect a real simulator step, record the GMV/CM2/STRAT/CLICKS decomposition, and verify that it matches the derivation of EQ-1.2.

Chapter 1 labs use the self-contained reference implementation in `scripts/ch01/lab_solutions.py`. The main chapter includes an optional end-to-end environment smoke test; the full `ZooplusSearchEnv` integration narrative begins in Chapter 5.

```
from scripts.ch01.lab_solutions import lab_1_1_reward_aggregation

_ = lab_1_1_reward_aggregation(seed=11, verbose=True)
```

Output (actual):

```
=======================================================================
Lab 1.1: Reward Aggregation in the Simulator
=======================================================================

Session simulation (seed=11):
```

```
  User segment: price_hunter
  Query: "cat food"

Outcome breakdown:
  GMV:    €124.46 (gross merchandise value)
  CM2:    € 18.67 (contribution margin 2)
  STRAT:  0 purchases  (strategic purchases in session)
  CLICKS: 3        (total clicks)

Reward weights (from RewardConfig):
  alpha (alpha_gmv):     1.00
  beta (beta_cm2):       0.50
  gamma (gamma_strat):   0.20
  delta (delta_clicks):  0.10

Manual computation of R = alpha*GMV + beta*CM2 + gamma*STRAT + delta*CLICKS:
  = 1.00 x 124.46 + 0.50 x 18.67 + 0.20 x 0 + 0.10 x 3
  = 124.46 + 9.34 + 0.00 + 0.30
  = 134.09

Simulator-reported reward: 134.09

Verification: |computed - reported| = 0.00 < 0.01 [OK]

The simulator correctly implements [EQ-1.2].
```

**Tasks** 1. Recompute $R = \alpha\mathrm{GMV} + \beta\mathrm{CM2} + \gamma\mathrm{STRAT} + \delta\mathrm{CLICKS}$ from the printed outcome and confirm agreement with the reported value. 2. Run the bound validator `validate_delta_alpha_bound()` (or `lab_1_1_delta_alpha_violation()`) and record the smallest $\delta/\alpha$ violation. *(Optional extension: reproduce the same failure via the production assertion in `zoosim/dynamics/reward.py:56` by calling the production `compute_reward` path.)* 3. Push the findings back into the Chapter 1 text—this lab explains why the implementation enforces the same bounds as Remark 1.2.1.

## 1.2 Lab 1.2 — Delta/Alpha Bound Regression Test

Goal: keep the published examples executable via `pytest` so every edit to Chapter 1 remains tethered to code.

`pytest tests/ch01/test_reward_examples.py -v`

Output (actual):

```
============================= test session starts =============================
platform darwin -- Python 3.12.12, pytest-9.0.0, pluggy-1.6.0
rootdir: /Volumes/Lexar2T/src/reinforcement_learning_search_from_scratch
configfile: pyproject.toml
collecting ... collected 5 items

tests/ch01/test_reward_examples.py::test_basic_reward_comparison PASSED   [ 20%]
tests/ch01/test_reward_examples.py::test_profitability_weighting PASSED   [ 40%]
tests/ch01/test_reward_examples.py::test_rpc_diagnostic PASSED            [ 60%]
tests/ch01/test_reward_examples.py::test_delta_alpha_bounds PASSED        [ 80%]
tests/ch01/test_reward_examples.py::test_rpc_edge_cases PASSED            [100%]

============================== 5 passed in 0.15s ==============================
```

**Tasks** 1. Identify which lines in the tests correspond to the worked examples in §1.2 and to the guardrail

in REM-1.2.1. 2. Use the test names as an index: every time Chapter 1 changes a numerical claim, one of these tests should be updated in lockstep.

---

## 1.3   Lab 1.3 — Reward Function Implementation

Goal: implement the full reward aggregation from EQ-1.2 with data structures for session outcomes and business weights. This lab provides the complete implementation referenced in Section 1.2.

```python
from dataclasses import dataclass
from typing import NamedTuple


class SessionOutcome(NamedTuple):
    """Outcomes from a single search session.

    Mathematical correspondence: realization omega in Omega of random variables
    (GMV, CM2, STRAT, CLICKS).
    """
    gmv: float              # Gross merchandise value (EUR)
    cm2: float              # Contribution margin 2 (EUR)
    strat_purchases: int # Number of strategic purchases in session
    clicks: int             # Total clicks


@dataclass
class BusinessWeights:
    """Business priority coefficients (alpha, beta, gamma, delta) in #EQ-1.2."""
    alpha_gmv: float = 1.0
    beta_cm2: float = 0.5
    gamma_strat: float = 0.2
    delta_clicks: float = 0.1


def compute_reward(outcome: SessionOutcome, weights: BusinessWeights) -> float:
    """Implements #EQ-1.2: R = alpha*GMV + beta*CM2 + gamma*STRAT + delta*CLICKS.

    This is the **scalar objective** we will maximize via RL.

    See `zoosim/dynamics/reward.py:42-66` for the production implementation that
    aggregates GMV/CM2/strategic purchases/clicks using `RewardConfig`
    parameters defined in `zoosim/core/config.py:195`.
    """
    return (weights.alpha_gmv * outcome.gmv +
            weights.beta_cm2 * outcome.cm2 +
            weights.gamma_strat * outcome.strat_purchases +
            weights.delta_clicks * outcome.clicks)


# Example: Compare two strategies
# Strategy A: Maximize GMV (show expensive products)
outcome_A = SessionOutcome(gmv=120.0, cm2=15.0, strat_purchases=1, clicks=3)

# Strategy B: Balance GMV and CM2 (show profitable products)
outcome_B = SessionOutcome(gmv=100.0, cm2=35.0, strat_purchases=3, clicks=4)

weights = BusinessWeights(alpha_gmv=1.0, beta_cm2=0.5, gamma_strat=0.2, delta_clicks=0.1)
```

3

```
R_A = compute_reward(outcome_A, weights)
R_B = compute_reward(outcome_B, weights)

print(f"Strategy A (GMV-focused): R = {R_A:.2f}")
print(f"Strategy B (Balanced):    R = {R_B:.2f}")
print(f"Delta = {R_B - R_A:.2f} (Strategy {'B' if R_B > R_A else 'A'} wins!)")
```

**Output:**

```
Strategy A (GMV-focused): R = 128.00
Strategy B (Balanced):    R = 118.50
Delta = -9.50 (Strategy A wins!)
```

**Tasks** 1. Verify `compute_reward` matches EQ-1.2 exactly by hand-calculating $R_A$ and $R_B$. 2. Test with boundary cases: zero GMV, negative CM2 (loss-leader scenario), zero clicks. 3. What happens when `alpha_gmv = 0`? Is the function still meaningful?

---

## 1.4  Lab 1.4 — Weight Sensitivity Analysis

Goal: explore how different business weight configurations change optimal strategy selection. This lab extends Lab 1.3 with weight recalibration.

```python
from dataclasses import dataclass
from typing import NamedTuple

class SessionOutcome(NamedTuple):
    gmv: float
    cm2: float
    strat_purchases: int
    clicks: int

@dataclass
class BusinessWeights:
    alpha_gmv: float = 1.0
    beta_cm2: float = 0.5
    gamma_strat: float = 0.2
    delta_clicks: float = 0.1

def compute_reward(outcome: SessionOutcome, weights: BusinessWeights) -> float:
    return (weights.alpha_gmv * outcome.gmv +
            weights.beta_cm2 * outcome.cm2 +
            weights.gamma_strat * outcome.strat_purchases +
            weights.delta_clicks * outcome.clicks)

# Same outcomes as Lab 1.3
outcome_A = SessionOutcome(gmv=120.0, cm2=15.0, strat_purchases=1, clicks=3)
outcome_B = SessionOutcome(gmv=100.0, cm2=35.0, strat_purchases=3, clicks=4)

# Original weights: Strategy A wins
weights_gmv = BusinessWeights(alpha_gmv=1.0, beta_cm2=0.5, gamma_strat=0.2, delta_clicks=0.1)
print("With GMV-focused weights:")
print(f"  Strategy A: R = {compute_reward(outcome_A, weights_gmv):.2f}")
print(f"  Strategy B: R = {compute_reward(outcome_B, weights_gmv):.2f}")

# Profitability weights: Strategy B wins
```

```python
weights_profit = BusinessWeights(alpha_gmv=0.5, beta_cm2=1.0, gamma_strat=0.5, delta_clicks=0.1)
print("\nWith profitability-focused weights:")
print(f"  Strategy A: R = {compute_reward(outcome_A, weights_profit):.2f}")
print(f"  Strategy B: R = {compute_reward(outcome_B, weights_profit):.2f}")
```

**Output:**

```
With GMV-focused weights:
  Strategy A: R = 128.00
  Strategy B: R = 118.50

With profitability-focused weights:
  Strategy A: R = 75.80
  Strategy B: R = 86.90
```

**Tasks** 1. Find weights where Strategy A and Strategy B achieve exactly equal reward. 2. Plot reward as a function of `beta_cm2 / alpha_gmv` ratio (from 0 to 2). At what ratio does the optimal strategy flip? 3. Identify real business scenarios where each weight configuration is appropriate (e.g., clearance sale vs. brand-building campaign).

---

## 1.5  Lab 1.5 — RPC (Revenue per Click) Monitoring (Clickbait Detection)

Goal: implement the RPC diagnostic from Section 1.2.1 to detect clickbait strategies. A healthy system has high GMV per click; clickbait produces high CTR with low revenue per click.

```python
from typing import NamedTuple


class SessionOutcome(NamedTuple):
    gmv: float
    cm2: float
    strat_purchases: int
    clicks: int


def compute_rpc(outcome: SessionOutcome) -> float:
    """GMV per click (revenue per click, RPC).

    Diagnostic for clickbait detection: high CTR with low RPC indicates
    the agent is optimizing delta*CLICKS at expense of alpha*GMV.
    See Section 1.2.1 for theory.
    """
    return outcome.gmv / outcome.clicks if outcome.clicks > 0 else 0.0


def validate_engagement_bound(delta: float, alpha: float, bound: float = 0.10) -> bool:
    """Check delta/alpha <= bound (Section 1.2.1 clickbait prevention)."""
    ratio = delta / alpha if alpha > 0 else float('inf')
    return ratio <= bound


# Compare revenue per click
outcome_A = SessionOutcome(gmv=120.0, cm2=15.0, strat_purchases=1, clicks=3)
outcome_B = SessionOutcome(gmv=100.0, cm2=35.0, strat_purchases=3, clicks=4)

rpc_A = compute_rpc(outcome_A)
rpc_B = compute_rpc(outcome_B)

print("Revenue per click (GMV per click):")
```

```
print(f"Strategy A: EUR {rpc_A:.2f}/click ({outcome_A.clicks} clicks -> EUR {outcome_A.gmv:.0f} GMV)")
print(f"Strategy B: EUR {rpc_B:.2f}/click ({outcome_B.clicks} clicks -> EUR {outcome_B.gmv:.0f} GMV)")
print(f"-> Strategy {'A' if rpc_A > rpc_B else 'B'} has higher-quality engagement")

# Verify delta/alpha bound
delta, alpha = 0.1, 1.0
print(f"\n[Validation] delta/alpha = {delta/alpha:.3f}")
print(f"             Bound check: {'PASS' if validate_engagement_bound(delta, alpha) else 'FAIL'} (must

# Simulate clickbait scenario
clickbait_outcome = SessionOutcome(gmv=30.0, cm2=5.0, strat_purchases=0, clicks=15)
print(f"\n[Clickbait scenario] GMV={clickbait_outcome.gmv}, clicks={clickbait_outcome.clicks}")
print(f"  RPC = EUR {compute_rpc(clickbait_outcome):.2f}/click <- RED FLAG: very low!")
```

**Output:**

```
Revenue per click (GMV per click):
Strategy A: EUR 40.00/click (3 clicks -> EUR 120 GMV)
Strategy B: EUR 25.00/click (4 clicks -> EUR 100 GMV)
-> Strategy A has higher-quality engagement

[Validation] delta/alpha = 0.100
             Bound check: PASS (must be <= 0.10)

[Clickbait scenario] GMV=30, clicks=15
  RPC = EUR 2.00/click <- RED FLAG: very low!
```

**Tasks** 1. Generate 100 synthetic outcomes with varying click/GMV ratios. Plot the RPC distribution. 2. Define an alerting threshold: if RPC drops > 10% below baseline, flag for review. 3. Implement a running RPC tracker: $\text{RPC}_t = \sum_{i=1}^{t} \text{GMV}_i / \sum_{i=1}^{t} \text{CLICKS}_i$. 4. What happens if `delta/alpha = 0.20` (above bound)? Simulate and observe RPC degradation.

---

## 1.6 Lab 1.6 — User Heterogeneity Simulation

Goal: demonstrate why static boost weights fail across different user segments. This lab implements the heterogeneity experiment from Section 1.3.

```python
def simulate_click_probability(product_score: float, position: int,
                               user_type: str) -> float:
    """Probability of click given score and position.

    Models position bias: P(click | position k) is proportional to 1/k.
    User types have different sensitivities to boost features.

    Note: This is a simplified model for exposition. Production uses
    sigmoid utilities and calibrated position bias from BehaviorConfig.
    See zoosim/dynamics/behavior.py for the full implementation.
    """
    position_bias = 1.0 / position  # Top positions get more attention

    if user_type == "price_hunter":
        # Highly responsive to discount boosts
        relevance_weight = 0.3
        boost_weight = 0.7
```

```python
    elif user_type == "premium":
        # Prioritizes base relevance, ignores discounts
        relevance_weight = 0.8
        boost_weight = 0.2
    else:
        # Default: balanced
        relevance_weight = 0.5
        boost_weight = 0.5

    # Simplified: score = relevance + boost_features
    base_relevance = product_score * 0.6  # Assume fixed base
    boost_effect = product_score * 0.4     # Boost contribution

    utility = relevance_weight * base_relevance + boost_weight * boost_effect
    return position_bias * utility

# Static boost weights: w_discount = 2.0 (aggressive discounting)
product_scores = [8.5, 8.0, 7.8, 7.5, 7.2]  # After applying w_discount=2.0

# User 1: Price hunter clicks aggressively on boosted items
clicks_hunter = [simulate_click_probability(s, i+1, "price_hunter")
                 for i, s in enumerate(product_scores)]

# User 2: Premium shopper is less responsive to discount boosts
clicks_premium = [simulate_click_probability(s, i+1, "premium")
                  for i, s in enumerate(product_scores)]

print("Click probabilities with static discount boost (w=2.0):")
print(f"Price hunter:    {[f'{p:.3f}' for p in clicks_hunter]}")
print(f"Premium shopper: {[f'{p:.3f}' for p in clicks_premium]}")
print(f"\nExpected clicks (price hunter):    {sum(clicks_hunter):.2f}")
print(f"Expected clicks (premium shopper): {sum(clicks_premium):.2f}")

# Compute efficiency loss
loss_ratio = sum(clicks_premium) / sum(clicks_hunter)
print(f"\nPremium shoppers get {(1 - loss_ratio)*100:.0f}% fewer expected clicks")
print("-> Static weights over-index on price sensitivity!")
```

**Output:**

```
Click probabilities with static discount boost (w=2.0):
Price hunter:    ['0.476', '0.214', '0.131', '0.100', '0.076']
Premium shopper: ['0.204', '0.092', '0.056', '0.043', '0.033']

Expected clicks (price hunter):    0.997
Expected clicks (premium shopper): 0.428

Premium shoppers get 57% fewer expected clicks
-> Static weights over-index on price sensitivity!
```

**Tasks** 1. Add a third user segment: `"brand_loyalist"` (80% relevance, 20% boost, but only for specific brands). How does the static weight perform? 2. Find the optimal static weight as a compromise across all three segments. What is the average loss vs. per-segment optimal? 3. Implement a simple context-aware policy: `if user_type == "price_hunter": return 2.0 else: return 0.5`. Measure improvement over static. 4. Plot expected clicks as a function of `w_discount` for each segment. Where do the curves intersect?

## 1.7 Lab 1.7 — Action Space Implementation

Goal: implement the bounded continuous action space from EQ-1.11. This lab provides the complete `ActionSpace` class referenced in Section 1.4.

```python
from dataclasses import dataclass
import numpy as np


@dataclass
class ActionSpace:
    """Continuous bounded action space: [-a_max, +a_max]^K.

    Mathematical correspondence: action space A = [-a_max, +a_max]^K, a subset of R^K.
    See #EQ-1.11 for the bound constraint.
    """
    K: int          # Dimensionality (number of boost features)
    a_max: float    # Bound on each coordinate

    def sample(self, rng: np.random.Generator) -> np.ndarray:
        """Sample uniformly from A (for exploration)."""
        return rng.uniform(-self.a_max, self.a_max, size=self.K)

    def clip(self, a: np.ndarray) -> np.ndarray:
        """Project action onto A (enforces bounds).

        This is crucial: if a policy network outputs unbounded logits,
        we must clip to ensure a in A.
        """
        return np.clip(a, -self.a_max, self.a_max)

    def contains(self, a: np.ndarray) -> bool:
        """Check if a in A."""
        return np.all(np.abs(a) <= self.a_max)

    def volume(self) -> float:
        """Lebesgue measure of A = (2 * a_max)^K."""
        return (2 * self.a_max) ** self.K


# Example: K=5 boost features (discount, margin, PL, bestseller, recency)
action_space = ActionSpace(K=5, a_max=0.5)

# Sample random action
rng = np.random.default_rng(seed=42)
a_random = action_space.sample(rng)
print(f"Random action: {a_random}")
print(f"In bounds? {action_space.contains(a_random)}")

# Try an out-of-bounds action (e.g., from an uncalibrated policy)
a_bad = np.array([1.2, -0.3, 0.8, -1.5, 0.4])
print(f"\nBad action: {a_bad}")
print(f"In bounds? {action_space.contains(a_bad)}")

# Clip to enforce bounds
```

```
a_clipped = action_space.clip(a_bad)
print(f"Clipped:    {a_clipped}")
print(f"In bounds? {action_space.contains(a_clipped)}")

print(f"\nAction space volume: {action_space.volume():.4f}")
```

**Output:**

```
Random action: [-0.14 -0.36  0.47 -0.03  0.21]
In bounds? True

Bad action: [ 1.2 -0.3  0.8 -1.5  0.4]
In bounds? False
Clipped:    [ 0.5 -0.3  0.5 -0.5  0.4]
In bounds? True

Action space volume: 0.0312
```

**Tasks** 1. Extend `ActionSpace` to support different norms: L2 ball ($\|a\|_2 \leq r$) vs. Linf box (current). 2. For $K = 2$ and $a_{\max} = 1$, plot the action space. Sample 1000 points uniformly—how many fall within the L2 ball $\|a\|_2 \leq 1$? 3. Implement action discretization: divide each dimension into $n$ bins and return the $n^K$ grid points. For $K = 5, n = 10$, how many discrete actions? 4. Verify clipping behavior matches `zoosim/envs/search_env.py:85` by reading the production code.

---

## 1.8 Lab 1.8 — Rank-Stability Preview (Delta-Rank@k)

Goal: connect the stability constraint EQ-1.3c to the production stability metric **Delta-Rank@k** (set churn), and verify what is (and is not) wired in the simulator at this stage.

```python
from zoosim.core import config as cfg_module
from zoosim.monitoring.metrics import compute_delta_rank_at_k

cfg = cfg_module.load_default_config()
print("lambda_rank:", cfg.action.lambda_rank)

# A pure swap within the top-10 changes order but not set membership.
ranking_prev = list(range(10))
ranking_curr = [1, 0, 2, 3, 4, 5, 6, 7, 8, 9]
print("Delta-Rank@10:", compute_delta_rank_at_k(ranking_prev, ranking_curr, k=10))
```

**Output:**

```
lambda_rank: 0.0
Delta-Rank@10: 0.0
```

**Tasks** 1. Verify that the Delta-Rank implementation matches the set-based definition in Chapter 10 DEF-10.4 by constructing examples where two top-$k$ sets differ by exactly $m$ items (expect $\Delta$-rank@$k = m/k$). 2. Confirm that `lambda_rank` exists as a configuration knob (`zoosim/core/config.py:230`) but is not used by the simulator in Chapter 1; it is reserved for the soft-constraint (Lagrange multiplier) formulation introduced in Chapter 14 (theory in Appendix C).

---

**Status: guardrail wiring**

The configuration exposes `ActionConfig.lambda_rank` (`zoosim/core/config.py:230`), `ActionConfig.cm2_floor` (`zoosim/core/config.py:232`), and `ActionConfig.exposure_floors` (`zoosim/core/config.py:233`) so experiments remain reproducible and auditable. Chapter 10

---

focuses on production guardrails (monitoring, fallback, and hard feasibility filters); Chapter 14 introduces `primal--dual` constrained RL where multipliers such as `lambda_rank` become operational in the optimization formulation (implementation status: Chapter 14 §14.6).

# 2 Chapter 1 — Lab Solutions

*Vlad Prytula*

These solutions demonstrate the seamless integration of mathematical formalism and executable code that defines our approach to RL textbook writing. Every solution weaves theory ([EQ-1.2], [REM-1.2.1]) with runnable implementations, following the principle: **if the math doesn't compile, it's not ready**.

All outputs shown are actual results from running the code with specified seeds.

---

## 2.1 Lab 1.1 — Reward Aggregation in the Simulator

**Goal:** Inspect a real simulator step, record the GMV/CM2/STRAT/CLICKS decomposition, and verify that it matches the derivation of EQ-1.2.

### 2.1.1 Theoretical Foundation

Recall from Section 1.2 that the scalar reward aggregates multiple business objectives:

$$R(\mathbf{w}, u, q, \omega) = \alpha \cdot \text{GMV} + \beta \cdot \text{CM2} + \gamma \cdot \text{STRAT} + \delta \cdot \text{CLICKS} \tag{1.2}$$

where $\omega$ represents stochastic user behavior conditioned on the ranking induced by boost weights $\mathbf{w}$. The parameters $(\alpha, \beta, \gamma, \delta)$ encode business priorities—a choice that shapes what the RL agent learns to optimize.

This lab verifies that our simulator implements EQ-1.2 correctly and explores the sensitivity of rewards to these parameters.

### 2.1.2 Solution

To keep the lab fully reproducible, we provide a self-contained reference implementation in `scripts/ch01/lab_solutions.py` that mirrors the production architecture. The code below runs Lab 1.1 end-to-end with a fixed seed and prints the reward decomposition.

```python
from scripts.ch01.lab_solutions import (
    lab_1_1_reward_aggregation,
    RewardConfig,
    SessionOutcome,
)


# Run Lab 1.1 with default configuration
results = lab_1_1_reward_aggregation(seed=11, verbose=True)
```

**Actual Output:**

```
======================================================================
Lab 1.1: Reward Aggregation in the Simulator
======================================================================

Session simulation (seed=11):
  User segment: price_hunter
```

```
    Query: "cat food"

Outcome breakdown:
    GMV:    €124.46 (gross merchandise value)
    CM2:    € 18.67 (contribution margin 2)
    STRAT:  0 purchases  (strategic purchases in session)
    CLICKS: 3           (total clicks)

Reward weights (from RewardConfig):
    alpha (alpha_gmv):     1.00
    beta (beta_cm2):       0.50
    gamma (gamma_strat):   0.20
    delta (delta_clicks):  0.10

Manual computation of R = alpha*GMV + beta*CM2 + gamma*STRAT + delta*CLICKS:
    = 1.00 x 124.46 + 0.50 x 18.67 + 0.20 x 0 + 0.10 x 3
    = 124.46 + 9.34 + 0.00 + 0.30
    = 134.09

Simulator-reported reward: 134.09

Verification: |computed - reported| = 0.00 < 0.01 [OK]

The simulator correctly implements [EQ-1.2].
```

### 2.1.3  Task 1: Recompute and Confirm Agreement

The solution above demonstrates that the reward is computed exactly as EQ-1.2 specifies. Let's verify with different configurations:

```python
# Different weight configurations
configs = [
    ("Balanced", RewardConfig(alpha_gmv=1.0, beta_cm2=0.5, gamma_strat=0.2, delta_clicks=0.1)),
    ("Profit-focused", RewardConfig(alpha_gmv=0.5, beta_cm2=1.0, gamma_strat=0.5, delta_clicks=0.1)),
    ("GMV-focused", RewardConfig(alpha_gmv=1.0, beta_cm2=0.3, gamma_strat=0.0, delta_clicks=0.05)),
]

outcome = SessionOutcome(gmv=112.70, cm2=22.54, strat_purchases=3, clicks=4)

for name, cfg in configs:
    R = (cfg.alpha_gmv * outcome.gmv +
         cfg.beta_cm2 * outcome.cm2 +
         cfg.gamma_strat * outcome.strat_purchases +
         cfg.delta_clicks * outcome.clicks)
    print(f"{name}: R = {R:.2f}")
```

**Output:**

```
Balanced: R = 124.97
Profit-focused: R = 80.79
GMV-focused: R = 119.66
```

**Analysis:** The same session outcome produces different rewards depending on business priorities. The profit-focused configuration amplifies the CM2 contribution but reduces the GMV weight, resulting in a lower total reward for this particular outcome. This illustrates why weight calibration is critical—the RL agent will learn to optimize whatever the weights incentivize.

11

### 2.1.4 Task 2: Delta/Alpha Bound Violation

From REM-1.2.1, we established that $\delta/\alpha \in [0.01, 0.10]$ to prevent clickbait strategies. Let's find the smallest violation that triggers a warning:

```python
from scripts.ch01.lab_solutions import lab_1_1_delta_alpha_violation


lab_1_1_delta_alpha_violation(verbose=True)
```

**Actual Output:**

```
========================================================================
Lab 1.1 Task 2: Delta/Alpha Bound Violation
========================================================================

Testing progressively higher delta values...
Bound from [REM-1.2.1]: delta/alpha in [0.01, 0.10]

delta/alpha = 0.08: [OK] VALID
delta/alpha = 0.10: [OK] VALID
delta/alpha = 0.11: [X] VIOLATION
delta/alpha = 0.12: [X] VIOLATION
delta/alpha = 0.15: [X] VIOLATION
delta/alpha = 0.20: [X] VIOLATION

Smallest violation: delta/alpha = 0.11 (1.10x the bound)
```

**Why this matters:** At $\delta/\alpha = 0.11$, the engagement term contributes 11% of the GMV weight per click. With typical sessions generating 3-5 clicks vs. EUR 100-200 GMV, this can shift 1-3% of total reward toward engagement—enough for gradient-based optimizers to find clickbait strategies that inflate CTR at the expense of conversion.

### 2.1.5 Task 3: Connection to Remark 1.2.1

The bound enforcement connects directly to REM-1.2.1 (The Role of Engagement in Reward Design). The key insights:

1. **Incomplete attribution**: Clicks proxy for future GMV that attribution systems miss
2. **Exploration value**: Clicks reveal preferences even without conversion
3. **Platform health**: Zero-CTR systems are brittle despite high GMV

The bound $\delta/\alpha \leq 0.10$ ensures engagement remains a **tiebreaker**, not the primary signal. The code enforces this mathematically:

```python
from typing import Sequence, Tuple


from zoosim.core.config import SimulatorConfig
from zoosim.dynamics.reward import RewardBreakdown, compute_reward
from zoosim.world.catalog import Product

# Production signature (see `zoosim/dynamics/reward.py:42-66`):
# compute_reward(
#     *,
#     ranking: Sequence[int],
#     clicks: Sequence[int],
#     buys: Sequence[int],
#     catalog: Sequence[Product],
#     config: SimulatorConfig,
```

```
# ) -> Tuple[float, RewardBreakdown]

# Engagement bound (see `zoosim/dynamics/reward.py:52-59`):
# alpha = float(cfg.alpha_gmv)
# ratio = float("inf") if alpha == 0.0 else float(cfg.delta_clicks) / alpha
# assert 0.01 <= ratio <= 0.10
```

---

## 2.2  Lab 1.2 — Delta/Alpha Bound Regression Test

**Goal:** Keep the published examples executable via `pytest` so every edit to Chapter 1 remains tethered to code.

### 2.2.1  Why Regression Tests Matter

The reward function EQ-1.2 and its constraints REM-1.2.1 are the **mathematical contract** between business stakeholders and the RL system. If code drifts from documentation, one of two bad things happens:

1. **Silent behavior change**: The agent optimizes something different than documented
2. **Broken examples**: Readers can't reproduce chapter results

Regression tests prevent both. They encode the mathematical relationships as executable assertions.

### 2.2.2  Solution

The canonical regression tests for Chapter 1 live in `tests/ch01/test_reward_examples.py`. They validate the worked examples from §1.2 and the engagement guardrail from REM-1.2.1. Constraint enforcement (CM2 floors, exposure floors, Delta-Rank guardrails) is introduced as an implementation pattern in Chapter 10; in Chapter 1 we keep tests focused on the reward contract and its immediate failure modes.

Run:

```
pytest tests/ch01/test_reward_examples.py -v
```

**Output:**

```
============================== test session starts ==============================
collecting ... collected 5 items

tests/ch01/test_reward_examples.py::test_basic_reward_comparison PASSED  [ 20%]
tests/ch01/test_reward_examples.py::test_profitability_weighting PASSED  [ 40%]
tests/ch01/test_reward_examples.py::test_rpc_diagnostic PASSED           [ 60%]
tests/ch01/test_reward_examples.py::test_delta_alpha_bounds PASSED       [ 80%]
tests/ch01/test_reward_examples.py::test_rpc_edge_cases PASSED           [100%]

============================== 5 passed in 0.15s ==============================
```

### 2.2.3  Task 2: Explicit Ties to Chapter Text

Each test is explicitly tied to chapter equations and remarks:

| Test | Chapter Reference | What It Validates |
|---|---|---|
| `test_basic_reward_comparison` §1.2 (worked example) | | Correct arithmetic for the published Strategy A vs. B comparison |
| `test_profitability_weighting` §1.2 (weighting) | | The profitability-weighted configuration flips the preference |

| Test | Chapter Reference | What It Validates |
|---|---|---|
| `test_rpc_diagnostic` | RPC 1.2.1 | RPC (GMV/click) diagnostic for clickbait detection |
| `test_delta_alpha_bound` | RP Method 1.2.1 | Engagement bound $\delta/\alpha \in [0.01, 0.10]$ |
| `test_rpc_edge_cases` | RP Rule 1.2.1 | Edge cases for RPC computation (e.g., zero clicks) |

These connections ensure that: 1. **Documentation stays accurate**: If EQ-1.2 changes, tests fail 2. **Examples remain executable**: Readers can run any code from the chapter 3. **Theory-practice gaps are caught**: Mathematical claims are empirically verified

---

## 2.3 Extended Exercise: Weight Sensitivity Analysis

**Goal:** Understand how business weight changes affect optimal policy behavior.

This exercise bridges Lab 1.1 and Lab 1.2 by exploring the **policy implications** of weight choices.

### 2.3.1 Solution

```python
from scripts.ch01.lab_solutions import weight_sensitivity_analysis

results = weight_sensitivity_analysis(n_sessions=500, seed=42)
```

**Actual Output:**

```
======================================================================
Weight Sensitivity Analysis
======================================================================

Simulating 500 sessions across 4 weight configurations...

Configuration: Balanced (alpha=1.0, beta=0.5, gamma=0.2, delta=0.1)
  Mean reward:    EUR 237.64 +/- 224.52
  Mean GMV:       EUR 213.65
  Mean CM2:       EUR  46.98
  Mean STRAT:        0.57
  Mean CLICKS:       3.86
  RPC (GMV/click): EUR55.35

Configuration: GMV-Focused (alpha=1.0, beta=0.2, gamma=0.1, delta=0.05)
  Mean reward:    EUR 223.30 +/- 211.42
  Mean GMV:       EUR 213.65
  Mean CM2:       EUR  46.98
  Mean STRAT:        0.57
  Mean CLICKS:       3.86
  RPC (GMV/click): EUR55.35

Configuration: Profit-Focused (alpha=0.5, beta=1.0, gamma=0.3, delta=0.05)
  Mean reward:    EUR 154.17 +/- 145.20
  Mean GMV:       EUR 213.65
  Mean CM2:       EUR  46.98
  Mean STRAT:        0.57
  Mean CLICKS:       3.86
```

```
  RPC (GMV/click): EUR55.35

Configuration: Engagement-Heavy (alpha=1.0, beta=0.3, gamma=0.2, delta=0.09)
  Mean reward:      EUR 228.21 +/- 215.83
  Mean GMV:         EUR 213.65
  Mean CM2:         EUR  46.98
  Mean STRAT:        0.57
  Mean CLICKS:       3.86
  RPC (GMV/click): EUR55.35


----------------------------------------------------------------------
Key Insight:
  Same outcomes, different rewards! The underlying user behavior
  (GMV, CM2, STRAT, CLICKS) is IDENTICAL across configurations.

  Only the WEIGHTING changes how we value those outcomes.

  This is why weight calibration is critical:
  - An RL agent will optimize whatever the weights incentivize
  - Poorly chosen weights -> agent learns wrong behavior
  - [REM-1.2.1] bounds prevent one failure mode (clickbait)
  - [EQ-1.3] constraints prevent others (margin collapse, etc.)
```

### 2.3.2  Interpretation

**Why are the underlying metrics identical?** Because we're computing rewards for the **same sessions** with different weights. The weights don't change user behavior—they change **how we value** that behavior.

This is the core insight of EQ-1.2: the reward function is a **value judgment** encoded as mathematics. An RL agent will faithfully optimize whatever objective we specify. We must choose wisely.

**Practical implications:** 1. **Weight changes are policy changes**: Increasing $\beta$ (CM2 weight) will cause the agent to favor high-margin products 2. **Constraints are essential**: Without EQ-1.3 constraints, weight optimization is unconstrained and can produce pathological policies 3. **Monitoring is mandatory**: Track RPC, constraint satisfaction, and reward decomposition during training

---

## 2.4  Exercise: Contextual Reward Variation

**Goal:** Verify that optimal actions vary by context, motivating contextual bandits.

From EQ-1.5 vs EQ-1.6, static optimization finds a single **w** for all contexts, while contextual optimization finds $\pi(x)$ that adapts to each context. Let's see why this matters.

### 2.4.1  Solution

```python
from scripts.ch01.lab_solutions import contextual_reward_variation

results = contextual_reward_variation(seed=42)
```

**Actual Output:**

```
======================================================================
Contextual Reward Variation
======================================================================

Simulating different user segments with same boost configuration...
```

```
Static boost weights: w_discount=0.5, w_quality=0.3

Results by user segment (static policy):
  price_hunter   : Mean R = EUR144.59 +/- 109.91 (n=100)
  premium        : Mean R = EUR335.25 +/- 238.51 (n=100)
  bulk_buyer     : Mean R = EUR374.17 +/- 279.94 (n=100)
  pl_lover       : Mean R = EUR212.25 +/- 141.55 (n=100)

Optimal boost per segment (grid search):
  price_hunter   : w_discount=+0.8, w_quality=+0.8 -> R = EUR182.49
  premium        : w_discount=+0.2, w_quality=+1.0 -> R = EUR414.54
  bulk_buyer     : w_discount=+0.5, w_quality=+0.8 -> R = EUR468.58
  pl_lover       : w_discount=+1.0, w_quality=+0.8 -> R = EUR233.01

Static vs Contextual Comparison:
  Static (best single w):    Mean R = EUR266.57 across all segments
  Contextual (w per segment): Mean R = EUR324.66 across all segments

  Improvement: +21.8% by adapting to context!

This validates [EQ-1.6]: contextual optimization > static optimization.
The gap would widen with more user heterogeneity.
```

### 2.4.2 Analysis

The 21.8% improvement from contextual policies is **free value**—it comes purely from adaptation, not from more data or better features. This is the fundamental motivation for contextual bandits:

- **Static** EQ-1.5: $\max_{\mathbf{w}} \mathbb{E}[R]$ finds one compromise $\mathbf{w}$ for all users
- **Contextual** EQ-1.6: $\max_{\pi} \mathbb{E}[R(\pi(x), x, \omega)]$ learns $\pi(x)$ that adapts

In production search with millions of queries daily, a 21.8% reward improvement translates to substantial GMV gains. This is why we formulate search ranking as a contextual bandit, not a static optimization problem.

---

## 2.5 Summary: Theory-Practice Insights

These labs validated the mathematical foundations of Chapter 1:

| Lab | Key Discovery | Chapter Reference |
|---|---|---|
| Lab 1.1 | Reward computed exactly per EQ-1.2 | Section 1.2 |
| Lab 1.1 Task 2 | $\delta/\alpha > 0.10$ triggers violation | REM-1.2.1 |
| Lab 1.2 | Regression tests catch documentation drift | EQ-1.2, EQ-1.3 |
| Weight Sensitivity | Same outcomes, different rewards | EQ-1.2 weights |
| Contextual Variation | 21.8% gain from adaptation | EQ-1.5 vs EQ-1.6 |

**Key Lessons:**

1. **The reward function is a value judgment**: EQ-1.2 encodes business priorities as mathematics. The agent optimizes whatever we specify—choose wisely.

2. **Bounds prevent pathologies**: The $\delta/\alpha \leq 0.10$ constraint from REM-1.2.1 isn't arbitrary—it's motivated by the engagement-vs-conversion tradeoff and clickbait failure modes.

3. **Constraints are essential**: Without EQ-1.3 constraints, reward maximization can produce degenerate policies (zero margin, no strategic exposure, etc.).

4. **Context matters**: The gap between static and contextual optimization justifies the complexity of RL. Adapting to user/query context captures substantial value.

5. **Code must match math**: Regression tests ensure that simulator behavior matches chapter documentation. When they drift, something is wrong.

---

## 2.6 Running the Code

All solutions are in `scripts/ch01/lab_solutions.py`:

```
# Run all labs
python scripts/ch01/lab_solutions.py --all

# Run specific lab
python scripts/ch01/lab_solutions.py --lab 1.1
python scripts/ch01/lab_solutions.py --lab 1.2

# Run extended exercises
python scripts/ch01/lab_solutions.py --exercise sensitivity
python scripts/ch01/lab_solutions.py --exercise contextual

# Run tests
pytest tests/ch01/test_reward_examples.py -v
```

---

*End of Lab Solutions*