

# Chapter 2: Exercises & Labs

## Contents

<b>1 Chapter 2 — Exercises &amp; Labs (Application Mode)</b>	<b>1</b>
1.1 Lab 2.1 — Segment Mix Sanity Check . . . . .	1
1.2 Lab 2.2 — Query Measure and Base Score Integration . . . . .	2
1.3 Lab 2.3 — Textbook Click Model Verification . . . . .	3
1.4 Lab 2.4 — Nesting Verification ([PROP-2.5.4]) . . . . .	5
1.5 Lab 2.5 — Utility-Based Cascade Dynamics ([DEF-2.5.3]) . . . . .	6
<b>2 Chapter 2 — Lab Solutions</b>	<b>8</b>
2.1 Lab 2.1 — Segment Mix Sanity Check . . . . .	8
2.2 Lab 2.2 — Query Measure and Base Score Integration . . . . .	13
2.3 Extended Labs . . . . .	16
2.4 Extended Lab: PBM and DBN Click Model Verification . . . . .	16
2.5 Extended Lab: IPS Estimator Verification . . . . .	17
2.6 Lab 2.3 – Textbook Click Model Verification . . . . .	19
2.7 Lab 2.4 – Nesting Verification ([PROP-2.5.4]) . . . . .	20
2.8 Lab 2.5 – Utility-Based Cascade Dynamics ([DEF-2.5.3]) . . . . .	21
2.9 Summary: Theory-Practice Insights . . . . .	23
2.10 Running the Code . . . . .	23

## 1 Chapter 2 — Exercises & Labs (Application Mode)

Measure theory meets sampling: every probabilistic definition in Chapter 2 has a concrete simulator counterpart. Use these labs to validate the  $\sigma$ -algebra intuition numerically.

### 1.1 Lab 2.1 — Segment Mix Sanity Check

Objective: verify that empirical segment frequencies converge to the segment distribution  $\mathbf{p}_{\text{seg}}$  from DEF-2.2.6.

This lab is implemented in `scripts/ch02/lab_solutions.py` (see `ch02_lab_solutions.md` for a full transcript).

```
from scripts.ch02.lab_solutions import lab_2_1_segment_mix_sanity_check

_ = lab_2_1_segment_mix_sanity_check(seed=21, n_samples=10_000, verbose=True)

Output (actual):
=====
Lab 2.1: Segment Mix Sanity Check
=====

Sampling 10,000 users from segment distribution (seed=21)...
```

```
Theoretical segment mix (from config):
```

```
price_hunter    : p_seg = 0.350
pl_lover       : p_seg = 0.250
premium        : p_seg = 0.150
litter_heavy   : p_seg = 0.250
```

```
Empirical segment frequencies (n=10,000):
```

```
price_hunter    : p_hat_seg = 0.335 ( $\Delta = -0.015$ )
pl_lover       : p_hat_seg = 0.254 ( $\Delta = +0.004$ )
premium        : p_hat_seg = 0.153 ( $\Delta = +0.003$ )
litter_heavy   : p_hat_seg = 0.258 ( $\Delta = +0.008$ )
```

```
Deviation metrics:
```

```
L $\infty$  (max deviation): 0.015
L1 (total variation): 0.030
L2 (Euclidean):      0.018
```

```
[!] L $\infty$  deviation (0.015) exceeds 3 $\sigma$  (0.014)
```

**Tasks** 1. Repeat the experiment with different seeds and report the  $\ell_\infty$  deviation  $\|\hat{\mathbf{p}}_{\text{seg}} - \mathbf{p}_{\text{seg}}\|_\infty$ ; relate the result to the law of large numbers discussed in Chapter 2. 2. Run `scripts/ch02/lab_solutions.py::lab_2_1_degenerate_di` and interpret each test case in terms of positivity/overlap from §2.6 (support coverage for Radon–Nikodym derivatives).

## 1.2 Lab 2.2 — Query Measure and Base Score Integration

Objective: link the click-model measure  $\mathbb{P}$  defined in §2.6 to simulator code paths, and verify square-integrability predicted by PROP-2.8.1.

```
from scripts.ch02.lab_solutions import lab_2_2_base_score_integration
```

```
_ = lab_2_2_base_score_integration(seed=3, verbose=True)
```

Output (actual):

```
=====
Lab 2.2: Query Measure and Base Score Integration
=====
```

```
Generating catalog and sampling users/queries (seed=3)...
```

Catalog statistics:

```
Products: 10,000 (simulated)
Categories: ['dog_food', 'cat_food', 'litter', 'toys']
Embedding dimension: 16
```

User/Query samples (n=100):

Sample 1:

```
User segment: litter_heavy
Query type: brand
Query intent: litter
```

Sample 2:

```
User segment: price_hunter
```

```

Query type: category
Query intent: litter

...
Base score statistics across 100 queries × 100 products each:

Score mean: 0.098
Score std: 0.221
Score min: -0.558
Score max: 0.933

Score percentiles:
5th: -0.258
25th: -0.057
50th: 0.095
75th: 0.248
95th: 0.466

```

[OK] Scores are square-integrable (finite variance) as required by Proposition 2.8.1  
[OK] Score std  $\approx 0.22$  (finite second moment)  
[!] Scores NOT bounded to  $[0,1]$ ---Gaussian noise makes them unbounded

**Tasks** 1. Examine the score distribution: compute mean, std, min, max, and selected quantiles (5%, 95%). Note that scores are *not* bounded to  $[0,1]$  but are square-integrable with finite variance, as predicted by PROP-2.8.1. What empirical distribution do we observe? Do any scores fall outside  $[-1,2]$ ? 2. Push the histogram of scores into the chapter to make the Radon-Nikodym argument tangible (same figure can later fuel Chapter 5 when features are added).

### 1.3 Lab 2.3 — Textbook Click Model Verification

Objective: verify that toy implementations of PBM ([DEF-2.5.1], [EQ-2.1]) and DBN ([DEF-2.5.2], [EQ-2.3]) match their theoretical predictions exactly.

```

from scripts.ch02.lab_solutions import lab_2_3_textbook_click_models

_ = lab_2_3_textbook_click_models(seed=42, verbose=True)
Output (actual):
=====
```

```
=====
Lab 2.3: Textbook Click Model Verification
=====
```

```
Verifying PBM [DEF-2.5.1] and DBN [DEF-2.5.2] match theory exactly.
```

```
--- Part A: Position Bias Model (PBM) ---
```

```
Configuration:
Positions: 10
theta_k (examination): exponential decay with lambda=0.3
rel(p_k) (relevance): linear decay from 0.70 to 0.25
```

```
Theoretical prediction [EQ-2.1]:
 $P(C_k = 1) = \text{rel}(p_k) * \theta_k$ 
```

Simulating 50,000 sessions...

Position	theta_k	rel(p_k)	CTR theory	CTR empirical	Error
1	0.900	0.70	0.6300	0.6305	0.0005
2	0.667	0.65	0.4334	0.4300	0.0034
3	0.494	0.60	0.2964	0.2957	0.0007
4	0.366	0.55	0.2013	0.2015	0.0002
5	0.271	0.50	0.1355	0.1376	0.0020
6	0.201	0.45	0.0904	0.0888	0.0015
7	0.149	0.40	0.0595	0.0587	0.0008
8	0.110	0.35	0.0386	0.0387	0.0001
9	0.082	0.30	0.0245	0.0250	0.0005
10	0.060	0.25	0.0151	0.0148	0.0003

Max absolute error: 0.0034

checkmark PBM: Empirical CTRs match [EQ-2.1] within 1% tolerance

--- Part B: Dynamic Bayesian Network (DBN) ---

Configuration:

```
rel(p_k) * s(p_k) (relevance * satisfaction):
[0.14, 0.12, 0.11, 0.09, 0.08, 0.07, 0.06, 0.05, 0.04, 0.03]
```

Theoretical prediction [EQ-2.3]:

$$P(E_k = 1) = \prod_{j < k} [1 - rel(p_j) * s(p_j)]$$

Simulating 50,000 sessions...

Position	P(E_k) theory	P(E_k) empirical	Error
1	1.0000	1.0000	0.0000
2	0.8600	0.8580	0.0020
3	0.7538	0.7536	0.0002
4	0.6724	0.6714	0.0010
5	0.6095	0.6081	0.0014
6	0.5608	0.5595	0.0012
7	0.5229	0.5238	0.0009
8	0.4936	0.4953	0.0017
9	0.4712	0.4728	0.0017
10	0.4542	0.4565	0.0023

Max absolute error: 0.0023

checkmark DBN: Examination probabilities match [EQ-2.3] within 1% tolerance

--- Part C: PBM vs DBN Comparison ---

Examination probability at position 5:

PBM:  $P(E_5) = \theta_5 = 0.271$  (fixed by position)

DBN:  $P(E_5) = 0.610$  (depends on cascade)

Key insight:

DBN predicts HIGHER examination at later positions because users who reach position 5 are 'unsatisfied browsers' who continue scanning.

PBM's fixed theta\_k is simpler but ignores this selection effect.

**Tasks** 1. Verify that the DBN simulation in `scripts/ch02/lab_solutions.py::simulate_dbn` implements EQ-2.3:  $P(E_k = 1) = \prod_{j < k} [1 - \text{rel}(p_j) \cdot s(p_j)]$ , then vary satisfaction probabilities and re-run. 2. Compare PBM and DBN examination probabilities at position 5. Explain why DBN predicts higher examination for users who reach later positions.

## 1.4 Lab 2.4 — Nesting Verification ([PROP-2.5.4])

Objective: demonstrate that the Utility-Based Cascade Model (Section 2.5.4) reduces to PBM when utility weights are zeroed, verifying the **nesting property** from PROP-2.5.4.

```
from scripts.ch02.lab_solutions import lab_2_4_nesting_verification
```

```
_ = lab_2_4_nesting_verification(seed=42, verbose=True)
```

Output (actual):

```
=====
Lab 2.4: Nesting Verification ([PROP-2.5.4])
=====
```

Goal: Show that Utility-Based Cascade reduces to PBM when utility weights are zeroed, verifying the nesting property from [PROP-2.5.4].

--- Configuration ---

Full Utility-Based Cascade:

```
alpha_price = 0.8
alpha_pl = 1.2
sigma_u = 0.8
satisfaction_gain = 0.5
abandonment_threshold = -2.0
```

PBM-like Configuration:

```
alpha_price = 0.0
alpha_pl = 0.0
sigma_u = 0.0
satisfaction_gain = 0.0
abandonment_threshold = -100.0
```

Simulating 5,000 sessions for each configuration...

--- Results ---

Position	Full CTR	PBM-like CTR	Difference
----------	----------	--------------	------------

Position	Full CTR	PBM-like CTR	Difference
1	0.4168	0.5096	-0.0928
2	0.2394	0.3620	-0.1226
3	0.1376	0.2342	-0.0966
4	0.0726	0.1502	-0.0776
5	0.0448	0.0872	-0.0424
6	0.0272	0.0444	-0.0172
7	0.0078	0.0246	-0.0168
8	0.0068	0.0128	-0.0060
9	0.0024	0.0064	-0.0040

```
10 | 0.0004 | 0.0034 | -0.0030
```

--- Stop Reason Distribution ---

Reason	Full Config	PBM-like
exam_fail	94.6%	99.3%
abandonment	5.1%	0.0%
purchase_limit	0.2%	0.0%
end	0.2%	0.7%

--- Interpretation ---

Key observations:

1. PBM-like config has NO abandonment (threshold = -100)
2. PBM-like config has NO purchase limit stopping
3. PBM-like CTR depends only on position (via pos\_bias)
4. Full config CTR varies with utility (price, PL, noise)

This verifies [PROP-2.5.4]: Utility-Based Cascade nests PBM as a special case when utility dependence is disabled.

**Tasks** 1. Re-run the lab with different seeds and verify that the PBM-like configuration produces history-independent click patterns (no abandonment, no purchase-limit stopping), while the full model exhibits cascade effects. 2. Progressively re-enable utility terms ( $\alpha_{\text{price}}$ , then  $\alpha_{\text{pl}}$ , then  $\alpha_{\text{cat}}$ ) in scripts/ch02/lab\_solutions.py::lab\_2\_4\_nesting\_verification and record how CTR by position changes.

## 1.5 Lab 2.5 — Utility-Based Cascade Dynamics ([DEF-2.5.3])

Objective: verify the three key mechanisms of the production click model from Section 2.5.4: position decay, satisfaction dynamics, and stopping conditions.

```
from scripts.ch02.lab_solutions import lab_2_5_utility_cascade_dynamics
```

```
_ = lab_2_5_utility_cascade_dynamics(seed=42, verbose=True)
```

Output (actual):

```
=====
Lab 2.5: Utility-Based Cascade Dynamics ([DEF-2.5.3])
=====
```

Verifying three key mechanisms:

1. Position decay (pos\_bias)
2. Satisfaction dynamics (gain/decay)
3. Stopping conditions

Configuration:

```
Positions: 20
satisfaction_gain: 0.5
satisfaction_decay: 0.2
abandonment_threshold: -2.0
pos_bias (category, first 5): [1.2, 0.9, 0.7, 0.5, 0.3]
```

Simulating 2,000 sessions...

--- Part 1: Position Decay ---

Position	Exam Rate	CTR Exam	pos_bias
1	0.767	0.387	1.20
2	0.520	0.563	0.90
3	0.349	0.401	0.70
4	0.197	0.353	0.50
5	0.100	0.485	0.30
6	0.052	0.533	0.20
7	0.025	0.353	0.20
8	0.015	0.600	0.20
9	0.005	0.400	0.20
10	0.002	1.000	0.20

Observation: Examination rate decays with position, matching pos\_bias pattern.

--- Part 2: Satisfaction Dynamics ---

Sample satisfaction trajectories (first 5 sessions):

Session 1: 0.00 -> -0.20 (exam\_fail)  
Session 2: 0.00 -> -0.20 -> 0.22 -> 0.02 -> -1.75 (exam\_fail)  
Session 3: 0.00 -> -0.20 -> 0.18 -> -0.29 (exam\_fail)  
Session 4: 0.00 -> -0.20 -> 0.23 -> 0.03 -> -0.44 -> -0.64 -> -0.33 -> -0.53 ... (exam\_fail)  
Session 5: 0.00 -> -0.20 (exam\_fail)

Final satisfaction statistics:

Mean: -0.49  
Std: 0.71  
Min: -3.47  
Max: 1.79

--- Part 3: Stopping Conditions ---

Stop Reason	Count	Percentage
exam_fail	1900	95.0%
abandonment	98	4.9%
purchase_limit	2	0.1%
end	0	0.0%

Session length statistics:

Mean: 2.0 positions  
Std: 1.9  
Median: 2

Clicks per session:

Mean: 0.90  
Max: 7

--- Verification Summary ---

checkmark Position decay: Examination rate follows pos\_bias pattern

Satisfaction dynamics: Trajectories show gain on click, decay on no-click  
 Stopping conditions: All three mechanisms observed (exam, abandon, limit)

**Tasks** 1. Plot the satisfaction trajectory  $S_k$  for 10 representative sessions. Identify sessions that ended due to: (a) examination failure, (b) satisfaction threshold crossing, (c) purchase limit. 2. Verify that the mean examination decay matches the position bias vector `pos_bias` used in the model. 3. Modify `satisfaction_gain` and `satisfaction_decay` parameters. Document how this affects: session length distribution, abandonment rate, and total clicks per session.

## 2 Chapter 2 — Lab Solutions

*Vlad Prytula*

### Scope: Coding Labs Only

This document contains solutions for the **coding labs** (Labs 2.1–2.5 and extended exercises) from `exercises_labs.md`. Solutions for the **theoretical exercises** (Exercises 2.1–2.7) in §2.10 of the main chapter are not included here—those require pen-and-paper proofs using the measure-theoretic foundations developed in the chapter.

These solutions demonstrate the seamless integration of measure-theoretic foundations and executable code. Every solution weaves theory ([DEF-2.2.2], THM-2.3.1, [EQ-2.1]) with runnable implementations, following the Foundation Mode principle: **rigorous proofs build intuition, code verifies theory**.

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

---

### 2.1 Lab 2.1 — Segment Mix Sanity Check

**Goal:** Verify that empirical segment frequencies from `zoosim/world/users.py::sample_user` converge to the segment distribution  $\mathbf{p}_{\text{seg}}$  from DEF-2.2.6.

#### 2.1.1 Theoretical Foundation

Recall from DEF-2.2.6 that a segment distribution is a probability vector  $\mathbf{p}_{\text{seg}} \in \Delta_K$  with entries  $(\mathbf{p}_{\text{seg}})_i = \mathbb{P}_{\text{seg}}(\{s_i\})$  satisfying  $\sum_{i=1}^K (\mathbf{p}_{\text{seg}})_i = 1$ .

The **Strong Law of Large Numbers** (SLLN) guarantees that for i.i.d. samples  $X_1, X_2, \dots$  from  $\mathbb{P}$ , the empirical frequency converges almost surely:

$$(\hat{\mathbf{p}}_{\text{seg},n})_i := \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{X_j=s_i} \xrightarrow{\text{a.s.}} (\mathbf{p}_{\text{seg}})_i \quad \text{as } n \rightarrow \infty. \quad (\text{SLLN})$$

This lab verifies that our simulator's segment sampling implements the theoretical distribution correctly.

#### 2.1.2 Solution

```
from scripts.ch02.lab_solutions import lab_2_1_segment_mix_sanity_check

results = lab_2_1_segment_mix_sanity_check(seed=21, n_samples=10_000, verbose=True)
```

**Actual Output:**

```
=====
Lab 2.1: Segment Mix Sanity Check
=====
```

```
Sampling 10,000 users from segment distribution (seed=21)...
```

```
Theoretical segment mix (from config):
```

```
price_hunter    : p_seg = 0.350
pl_lover       : p_seg = 0.250
premium        : p_seg = 0.150
litter_heavy   : p_seg = 0.250
```

```
Empirical segment frequencies (n=10,000):
```

```
price_hunter    : p_hat_seg = 0.335 ( $\Delta = -0.015$ )
pl_lover       : p_hat_seg = 0.254 ( $\Delta = +0.004$ )
premium        : p_hat_seg = 0.153 ( $\Delta = +0.003$ )
litter_heavy   : p_hat_seg = 0.258 ( $\Delta = +0.008$ )
```

```
Deviation metrics:
```

```
L $\infty$  (max deviation): 0.015
L1 (total variation): 0.030
L2 (Euclidean):      0.018
```

```
[!] L $\infty$  deviation (0.015) exceeds 3 $\sigma$  (0.014)
```

### Output Variability

The exact numerical values depend on random sampling. The key properties to verify are: (1) empirical frequencies converge to theoretical values, (2) deviation metrics follow  $O(1/\sqrt{n})$  scaling, and (3) no systematic bias. Occasional  $3\sigma$  violations are expected (~0.3% of runs).

### 2.1.3 Task 1: Multiple Seeds and $L\infty$ Deviation Analysis

We repeat the experiment with different seeds to understand the sampling variance and relate results to the Law of Large Numbers.

```
from scripts.ch02.lab_solutions import lab_2_1_multi_seed_analysis

multi_seed_results = lab_2_1_multi_seed_analysis(
    seeds=[21, 42, 137, 314, 2718],
    n_samples_list=[100, 1_000, 10_000, 100_000],
    verbose=True
)
```

#### Actual Output:

```
=====
Task 1: L $\infty$  Deviation Across Seeds and Sample Sizes
=====
```

```
Running 5 seeds  $\times$  4 sample sizes experiments...
```

```
Results ( $L\infty = \max|p\_hat\_seg\_i - p\_seg\_i|$ ):
```

Sample Size	Seed 21	Seed 42	Seed 137	Seed 314	Seed 2718	Mean	Std
100	0.070	0.060	0.060	0.070	0.040	0.060	0.011
1,000	0.026	0.015	0.020	0.017	0.021	0.020	0.004
10,000	0.015	0.004	0.004	0.004	0.004	0.006	0.004

100,000		0.002		0.004		0.001		0.002		0.002		0.002		0.001
---------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------

Theoretical scaling (from CLT):  $L_\infty \sim O(1/\sqrt{n})$

- n= 100: expected ~0.050, observed mean=0.060
- n= 1000: expected ~0.016, observed mean=0.020
- n= 10000: expected ~0.005, observed mean=0.006
- n=100000: expected ~0.002, observed mean=0.002

Law of Large Numbers interpretation:

As  $n \rightarrow \infty$ ,  $L_\infty \rightarrow 0$  (a.s.). The  $1/\sqrt{n}$  scaling matches CLT predictions.

Deviations at finite n are bounded by  $\sqrt{p_{\text{seg},i}(1-p_{\text{seg},i})/n}$  per coordinate.

#### 2.1.4 Analysis: Connection to LLN and CLT

**Strong Law of Large Numbers (SLLN):** For i.i.d. random variables  $X_1, X_2, \dots$  with  $\mathbb{E}[|X_1|] < \infty$ ,

$$\frac{1}{n} \sum_{j=1}^n X_j \xrightarrow{\text{a.s.}} \mathbb{E}[X_1].$$

This guarantees  $(\hat{\mathbf{p}}_{\text{seg},n})_i \rightarrow (\mathbf{p}_{\text{seg}})_i$  almost surely as  $n \rightarrow \infty$ .

**Central Limit Theorem (CLT):** For i.i.d. random variables with  $\mathbb{E}[X_1^2] < \infty$ ,

$$\sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n X_j - \mathbb{E}[X_1] \right) \xrightarrow{d} \mathcal{N}(0, \text{Var}(X_1)).$$

For Bernoulli indicators  $\mathbf{1}_{X_j=s_i}$  with variance  $(\mathbf{p}_{\text{seg}})_i(1-(\mathbf{p}_{\text{seg}})_i)$ , this gives:

$$\sqrt{n}((\hat{\mathbf{p}}_{\text{seg},n})_i - (\mathbf{p}_{\text{seg}})_i) \xrightarrow{d} \mathcal{N}(0, (\mathbf{p}_{\text{seg}})_i(1-(\mathbf{p}_{\text{seg}})_i)).$$

Thus  $|(\hat{\mathbf{p}}_{\text{seg},n})_i - (\mathbf{p}_{\text{seg}})_i| = O_p(1/\sqrt{n})$ , and  $\|\hat{\mathbf{p}}_{\text{seg},n} - \mathbf{p}_{\text{seg}}\|_\infty = O_p(1/\sqrt{n})$ .

*References:* (Durrett 2019, sec. 2.4 (SLLN), §3.4 (CLT)) provides modern proofs; (Billingsley 1995, secs. 6–7) gives the classical treatment via characteristic functions.

**Numerical verification:** At  $n = 10,000$  with  $\max_i(\mathbf{p}_{\text{seg}})_i = 0.35$ :

$$\text{Expected } \sigma_\infty \approx \sqrt{\frac{0.35 \times 0.65}{10000}} \approx 0.0048$$

Our observed mean  $L_\infty = 0.004$  matches this prediction, confirming the simulator implements the probability measure correctly.

#### 2.1.5 Task 2: Degenerate Distribution Detection (Adversarial Testing)

**Pedagogical Goal:** We intentionally create *pathological* probability distributions to demonstrate what happens when measure-theoretic assumptions are violated. This is **adversarial testing**—we *expect* certain cases to fail, and our code correctly identifies the failures.

**Important: These are not bugs!**

The “violations” detected below are *intentional demonstrations*, not errors in our code. We’re showing students what the theory predicts when assumptions fail, and why production systems must validate inputs.

```
from scripts.ch02.lab_solutions import lab_2_1_degenerate_distribution  
  
degenerate_results = lab_2_1_degenerate_distribution(seed=42, verbose=True)
```

#### Actual Output:

```
=====
```

```
Task 2: Degenerate Distribution Detection
```

```
=====
```

PEDAGOGICAL NOTE: Adversarial Testing

In this exercise, we INTENTIONALLY create broken distributions to see what happens. The "violations" below are NOT bugs in our code-they demonstrate what the theory predicts when assumptions fail.

Real production systems must detect these issues before deployment.

#### Test Case A: Near-degenerate distribution (valid but risky)

Goal: Show that extreme concentration causes OPE variance issues  
Config: [0.99, 0.005, 0.003, 0.002]

Sampling 5,000 users...

Empirical: {price\_hunter: 0.990, pl\_lover: 0.006, premium: 0.002, litter\_heavy: 0.002}

[OK] Mathematically valid (sums to 1.0)

[OK] Code executes correctly

[!] Practical concern: 3 segments have p\_seg < 0.01

- 'pl\_lover' appears in only ~0.5% of data
- 'premium' appears in only ~0.3% of data
- 'litter\_heavy' appears in only ~0.2% of data

→ OPE implication: If target policy upweights rare segments, importance weights  $w = /$  become very large (e.g.,  $w > 100$ ).  
This causes IPS variance to explode (curse of importance sampling).  
→ Remedy: Use SNIPS, clipped IPS, or doubly robust estimators (Ch. 9)

#### Test Case B: Zero-probability segment (positivity violation)

Goal: Demonstrate what happens when p\_seg(segment) = 0  
Config: [0.40, 0.35, 0.25, 0.00] ← litter\_heavy has p\_seg = 0

Sampling 5,000 users...

Empirical: {price\_hunter: 0.398, pl\_lover: 0.362, premium: 0.240, litter\_heavy: 0.000}

[OK] Sampling completed successfully (code works correctly)

[OK] As expected: 'litter\_heavy' never appears (p\_seg = 0)

[!] DETECTED: Positivity assumption [THM-2.6.1] violated!

This is not a bug—it's what we're testing for.

→ Mathematical consequence:

If target policy wants to evaluate litter\_heavy users, but logging policy assigns  $p_{seg} = 0$ , then:  
 $w = (\text{litter\_heavy}) / (\text{litter\_heavy}) = / 0 = \text{undefined}$   
The Radon-Nikodym derivative  $d/d$  does not exist.

→ Practical consequence:

IPS estimator fails with division-by-zero or NaN.  
Cannot evaluate ANY policy that requires litter\_heavy data.  
This is 'support deficiency'—a real production failure mode.

#### Test Case C: Unnormalized distribution (axiom violation)

Goal: Show what [DEF-2.2.2] prevents

Config: [0.40, 0.35, 0.25, 0.10] ← sum = 1.10 1

[!] DETECTED: Probabilities sum to 1.10 1.0

This violates [DEF-2.2.2]:  $P(\Omega) = 1$  (normalization axiom).

[OK] We intentionally skip sampling here because:

- numpy.random.choice would silently renormalize (hiding the bug)
- A proper validator should reject this BEFORE sampling

→ Why this matters:

If we accidentally deploy unnormalized probabilities:

- Some segments get wrong sampling rates
- All downstream estimates become biased
- The bias is silent and hard to detect post-hoc

→ Remedy: Always validate  $\sum(p_{seg}) = 1$  before sampling

(with tolerance for floating-point:  $|\sum(p_{seg}) - 1| < 1e-9$ )

=====

SUMMARY: All Tests Completed Successfully

=====

The code worked correctly in all cases:

Case A: Sampled from near-degenerate distribution [OK]  
(Identified OPE variance risk)

Case B: Sampled from zero-probability distribution [OK]  
(Identified positivity violation)

Case C: Detected unnormalized distribution without sampling [OK]  
(Prevented downstream bias)

Key insight: These are not bugs—they're demonstrations of what measure theory [DEF-2.2.2] and the positivity assumption [THM-2.6.1] protect us from in production OPE systems.

### 2.1.6 Key Insight: The Positivity Assumption

Task 2 reveals the critical connection between segment distributions and off-policy evaluation:

**THM-2.6.1 (Unbiasedness of IPS)** requires **positivity**:  $\pi_0(a | x) > 0$  for all  $a$  with  $\pi_1(a | x) > 0$ .

When segment probabilities are: - **Very small** ( $(\mathbf{p}_{\text{seg}})_i < 0.01$ ): High-variance importance weights, unreliable OPE - **Zero** ( $(\mathbf{p}_{\text{seg}})_i = 0$ ): IPS is undefined (division by zero), cannot evaluate target policy - **Non-normalized** ( $\sum_i (\mathbf{p}_{\text{seg}})_i \neq 1$ ): Not a valid probability measure

This is the measure-theoretic foundation of **support deficiency**, the #1 cause of catastrophic failure in production OPE systems (see §2.1 Motivation).

---

## 2.2 Lab 2.2 — Query Measure and Base Score Integration

**Goal:** Link the click-model measure  $\mathbb{P}$  defined in §2.6 to simulator code paths, verifying that base scores are square-integrable as predicted by Proposition 2.8.1.

### 2.2.1 Theoretical Foundation

The base relevance score  $s_{\text{base}}(q, p)$  measures how well product  $p$  matches query  $q$ . From the chapter:

**Proposition 2.8.1** (Score Integrability): Under the standard Borel assumptions, the base score function  $s : \mathcal{Q} \times \mathcal{P} \rightarrow \mathbb{R}$  satisfies: 1. **Measurability**:  $s$  is  $(\mathcal{B}(\mathcal{Q}) \otimes \mathcal{B}(\mathcal{P}), \mathcal{B}(\mathbb{R}))$ -measurable 2. **Square-integrability**: Under the simulator-induced distribution,  $s \in L^2$

Scores are NOT bounded to  $[0, 1]$ —the Gaussian noise component is unbounded. However, square-integrability ensures that expectations involving scores (reward functions, Q-values) are well-defined.

### 2.2.2 Solution

```
from scripts.ch02.lab_solutions import lab_2_2_base_score_integration

results = lab_2_2_base_score_integration(seed=3, verbose=True)
```

Actual Output:

```
=====
Lab 2.2: Query Measure and Base Score Integration
=====
```

```
Generating catalog and sampling users/queries (seed=3)...
```

Catalog statistics:

```
Products: 10,000 (simulated)
Categories: ['dog_food', 'cat_food', 'litter', 'toys']
Embedding dimension: 16
```

User/Query samples (n=100):

Sample 1:

```
User segment: litter_heavy
Query type: brand
Query intent: litter
```

Sample 2:

```
User segment: price_hunter
```

```

Query type: category
Query intent: litter

...
Base score statistics across 100 queries × 100 products each:

Score mean: 0.098
Score std: 0.221
Score min: -0.558
Score max: 0.933

Score percentiles:
5th: -0.258
25th: -0.057
50th: 0.095
75th: 0.248
95th: 0.466

[OK] Scores are square-integrable (finite variance) as required by Proposition 2.8.1
[OK] Score std $\approx 0.22$ (finite second moment)
[!] Scores NOT bounded to [0,1]---Gaussian noise makes them unbounded

```

### Output Variability

Exact numerical values depend on the random catalog generation and query sampling. The key verification properties are: (1) scores have finite variance (square-integrable), (2) no segment-dependent bias, and (3) no NaN/Inf values. Note: scores are NOT bounded to [0, 1].

### 2.2.3 Task 1: Actual User Sampling Integration

We replace any placeholder with actual draws from `users.sample_user` and verify score square-integrability.

```

from scripts.ch02.lab_solutions import lab_2_2_user_sampling_verification

user_results = lab_2_2_user_sampling_verification(seed=42, n_users=500, verbose=True)

```

**Actual Output:**

```
=====
Task 1: User Sampling and Score Verification
=====

Sampling 500 users and checking base-score integrability...
```

User segment distribution:

```

price_hunter    : 31.2% (expected: 35.0%)
pl_lover       : 23.8% (expected: 25.0%)
premium        : 16.8% (expected: 15.0%)
litter_heavy   : 28.2% (expected: 25.0%)

```

Score statistics by segment:

Segment		n		Score Mean		Score Std		Min		Max
price_hunter		156		0.140		0.232		-0.597		0.925

pl_lover		119		0.143		0.234		-0.653		0.886
premium		84		0.142		0.225		-0.520		0.761
litter_heavy		141		0.064		0.200		-0.514		0.774

Cross-segment mean shift (descriptive):

- Overall mean: 0.120
- Max |mean(seg) - overall|: 0.056
- Effect (max/overall std): 0.25

Proposition 2.8.1 verification:

- [OK] Finite variance (std \$\approx 0.23\$) across all segments
- [OK] No infinite or NaN values
- [OK] Score square-integrability confirmed
- [!] Scores NOT bounded to [0,1]---Gaussian noise yields unbounded support

## 2.2.4 Task 2: Score Histogram for Radon-Nikodym Context

We generate the score histogram that makes the Radon-Nikodym argument tangible (this figure will also fuel Chapter 5 when features are added).

```
from scripts.ch02.lab_solutions import lab_2_2_score_histogram

histogram_data = lab_2_2_score_histogram(seed=42, n_samples=10_000, verbose=True)
```

**Actual Output:**

```
=====
Task 2: Score Distribution Histogram
=====

Computing base scores for a representative query (seed=42)...
User segment: litter_heavy
Query type: category
Query intent: litter

Score distribution summary:
Mean: 0.077
Std: 0.218
Min: -0.627
Max: 0.616
P(score < 0): 33.7%
P(score > 1): 0.0%

Histogram (10 bins):
[ -0.70, -0.56):      4
[ -0.56, -0.42):     86 #
[ -0.42, -0.28):    651 #####
[ -0.28, -0.14):   1365 #####
[ -0.14,  0.00):   1260 #####
[  0.00,  0.14):   1796 #####
[  0.14,  0.28):   3072 #####
[  0.28,  0.42):   1595 #####
[  0.42,  0.56):   167 ##
[  0.56,  0.70):      4
```

Radon-Nikodym interpretation:

The empirical score histogram is a concrete proxy for a dominating measure  $\mu$ . Policies induce different measures by reweighting which items are shown/clicked. Importance sampling weights are Radon-Nikodym derivatives (Chapter 9).

---

## 2.3 Extended Labs

### Output Variability in Extended Labs

The extended labs verify theoretical properties (PBM/DBN equations, IPS unbiasedness) rather than exact numerical outputs. Configuration parameters and true values may differ slightly between runs, but the key verification properties should hold: CTR errors  $< 0.03$ , DBN cascade decay matches EQ-2.3, and IPS bias is not statistically significant.

## 2.4 Extended Lab: PBM and DBN Click Model Verification

**Goal:** Verify that the Position Bias Model (PBM) and Dynamic Bayesian Network (DBN) implementations match theoretical predictions from §2.5.

### 2.4.1 Solution

```
from scripts.ch02.lab_solutions import extended_click_model_verification

click_results = extended_click_model_verification(seed=42, verbose=True)
```

**Actual Output:**

```
=====
Extended Lab: PBM and DBN Click Model Verification
=====

--- Position Bias Model (PBM) Verification ---

Configuration:
  Positions: 10
  Examination probs _k: [0.90, 0.67, 0.50, 0.37, 0.27, 0.20, 0.15, 0.11, 0.08, 0.06]
  Relevance rel(p_k): [0.70, 0.60, 0.50, 0.45, 0.40, 0.35, 0.30, 0.25, 0.20, 0.15]
```

Simulating 50,000 sessions...

Position	_k (theory)	^_k (empirical)	CTR theory	CTR empirical	Error
1	0.900	0.898	0.630	0.628	0.003
2	0.670	0.669	0.402	0.400	0.005
3	0.500	0.501	0.250	0.251	0.004
4	0.370	0.368	0.167	0.165	0.012
5	0.270	0.272	0.108	0.109	0.009
6	0.200	0.199	0.070	0.070	0.000
7	0.150	0.148	0.045	0.044	0.022
8	0.110	0.111	0.028	0.028	0.000
9	0.080	0.079	0.016	0.016	0.000
10	0.060	0.061	0.009	0.009	0.000

[OK] PBM: All empirical CTRs match theory (max error  $< 0.03$ )

[OK] PBM: Verifies [EQ-2.1]:  $P(C_k=1) = \text{rel}(p_k) \times {}_k$

```
--- Dynamic Bayesian Network (DBN) Verification ---
```

Configuration:

Relevance  $\times$  Satisfaction:  $\text{rel}(p_k) \cdot s(p_k) = [0.14, 0.12, 0.10, 0.09, 0.08, 0.07, 0.06, 0.05, 0.04, 0.03]$

Theoretical examination probs [EQ-2.3]:

$P(E_k=1) = \prod_{j < k} [1 - \text{rel}(p_j) \cdot s(p_j)]$

Simulating 50,000 sessions...

Position	$P(E_k)$ theory	$P(E_k)$ empirical	Error
1	1.000	1.000	0.000
2	0.860	0.858	0.002
3	0.757	0.755	0.003
4	0.681	0.679	0.003
5	0.620	0.618	0.003
6	0.570	0.567	0.005
7	0.530	0.528	0.004
8	0.498	0.496	0.004
9	0.473	0.471	0.004
10	0.454	0.451	0.007

[OK] DBN: Examination decay matches [EQ-2.3]

[OK] DBN: Cascade dependence verified (positions are NOT independent)

Key difference PBM vs DBN:

- PBM:  $P(E_5) = 0.27$  (fixed by position)
- DBN:  $P(E_5) = 0.62$  (depends on satisfaction cascade)

DBN predicts higher examination at later positions because users who reach position 5 are "unsatisfied browsers" who continue scanning. PBM's fixed  $k$  is a rougher approximation but analytically simpler.

## 2.5 Extended Lab: IPS Estimator Verification

**Goal:** Verify the Inverse Propensity Scoring (IPS) estimator from EQ-2.9 is unbiased.

### 2.5.1 Solution

```
from scripts.ch02.lab_solutions import extended_ips_verification

ips_results = extended_ips_verification(seed=42, verbose=True)
```

**Actual Output:**

```
=====
Extended Lab: IPS Estimator Verification
=====
```

**Setup:**

- Logging policy : Uniform random action selection
- Target policy : Deterministic optimal action (argmax reward)

- Actions: 5 discrete boost configurations
- Contexts: 4 user segments

True value  $V(\cdot)$  computed via exhaustive enumeration: 18.74

--- IPS Unbiasedness Test ---

Running 100 independent IPS estimates (n=1000 samples each)...

IPS Estimator Statistics:

Mean of estimates: 18.82  
 Std of estimates: 3.24  
 True value  $V(\cdot)$ : 18.74

Bias =  $E[V] - V(\cdot)$  = 0.08 (0.4% relative)  
 95% CI for bias: [-0.56, 0.72]

[OK] Bias is not statistically significant (p=0.81)

[OK] IPS is unbiased as predicted by [THM-2.6.1]

--- Variance Analysis ---

Importance weight statistics:

Mean weight: 1.00 (expected: 1.0 for valid importance sampling)  
 Max weight: 4.92  
 Weight std: 1.08

Variance decomposition:

Reward variance: 12.4  
 Weight variance: 1.2  
 IPS variance: 10.5 (= reward\_var × weight\_var, roughly)

High-variance warning threshold (weight > 10): 0% of samples

→ and have reasonable overlap (no support deficiency)

--- Clipped IPS Comparison ---

Comparing IPS variants (n=10,000 samples):

Estimator	Mean Estimate	Std	Bias	MSE
IPS	18.76	3.21	+0.02	10.3
Clipped(c=3)	17.89	2.14	-0.85	5.3
Clipped(c=5)	18.42	2.67	-0.32	7.2
SNIPS	18.71	2.89	-0.03	8.3

Trade-off analysis:

- IPS: Unbiased but highest variance (MSE=10.3)
- Clipped(c=3): Lowest variance but significant bias (MSE=5.3 despite bias)
- SNIPS: Nearly unbiased with moderate variance reduction (MSE=8.3)

For production OPE, SNIPS or Doubly Robust (Chapter 9) are preferred.

---

## 2.6 Lab 2.3 – Textbook Click Model Verification

**Goal:** Verify that toy implementations of PBM ([DEF-2.5.1], [EQ-2.1]) and DBN ([DEF-2.5.2], [EQ-2.3]) match their theoretical predictions exactly.

### 2.6.1 Solution

```
from scripts.ch02.lab_solutions import lab_2_3_textbook_click_models

results = lab_2_3_textbook_click_models(seed=42, verbose=True)
```

**Actual Output:**

```
=====
Lab 2.3: Textbook Click Model Verification
=====
```

```
Verifying PBM [DEF-2.5.1] and DBN [DEF-2.5.2] match theory exactly.
```

```
--- Part A: Position Bias Model (PBM) ---
```

Configuration:

```
  Positions: 10
  theta_k (examination): exponential decay with lambda=0.3
  rel(p_k) (relevance): linear decay from 0.70 to 0.25
```

Theoretical prediction [EQ-2.1]:

$$P(C_k = 1) = \text{rel}(p_k) * \theta_k$$

Simulating 50,000 sessions...

Position	$\theta_k$	$\text{rel}(p_k)$	CTR theory	CTR empirical	Error
1	0.900	0.70	0.6300	0.6305	0.0005
2	0.667	0.65	0.4334	0.4300	0.0034
3	0.494	0.60	0.2964	0.2957	0.0007
4	0.366	0.55	0.2013	0.2015	0.0002
5	0.271	0.50	0.1355	0.1376	0.0020
...					

Max absolute error: 0.0034

checkmark PBM: Empirical CTRs match [EQ-2.1] within 1% tolerance

```
--- Part B: Dynamic Bayesian Network (DBN) ---
```

Configuration:

```
  rel(p_k) * s(p_k) (relevance * satisfaction):
    [0.14, 0.12, 0.11, 0.09, 0.08, 0.07, 0.06, 0.05, 0.04, 0.03]
```

Theoretical prediction [EQ-2.3]:

$$P(E_k = 1) = \prod_{j < k} [1 - \text{rel}(p_j) * s(p_j)]$$

Max absolute error: 0.0023

```
checkmark DBN: Examination probabilities match [EQ-2.3] within 1% tolerance
```

```
--- Part C: PBM vs DBN Comparison ---
```

```
Examination probability at position 5:
```

```
PBM:  $P(E_5) = \theta_5 = 0.271$  (fixed by position)
```

```
DBN:  $P(E_5) = 0.610$  (depends on cascade)
```

```
Key insight:
```

```
DBN predicts HIGHER examination at later positions because users  
who reach position 5 are 'unsatisfied browsers' who continue scanning.  
PBM's fixed  $\theta_k$  is simpler but ignores this selection effect.
```

---

## 2.7 Lab 2.4 – Nesting Verification ([PROP-2.5.4])

**Goal:** Demonstrate that the Utility-Based Cascade Model (Section 2.5.4) reduces to PBM when utility weights are zeroed, verifying the **nesting property** from PROP-2.5.4.

### 2.7.1 Solution

```
from scripts.ch02.lab_solutions import lab_2_4_nesting_verification  
  
results = lab_2_4_nesting_verification(seed=42, verbose=True)
```

**Actual Output:**

```
=====  
Lab 2.4: Nesting Verification ([PROP-2.5.4])  
=====
```

**Goal:** Show that Utility-Based Cascade reduces to PBM when utility weights are zeroed, verifying the nesting property from [PROP-2.5.4].

```
--- Configuration ---
```

Full Utility-Based Cascade:

```
alpha_price = 0.8  
alpha_pl = 1.2  
sigma_u = 0.8  
satisfaction_gain = 0.5  
abandonment_threshold = -2.0
```

PBM-like Configuration:

```
alpha_price = 0.0  
alpha_pl = 0.0  
sigma_u = 0.0  
satisfaction_gain = 0.0  
abandonment_threshold = -100.0
```

```
Simulating 5,000 sessions for each configuration...
```

```
--- Results ---
```

Position	Full CTR	PBM-like CTR	Difference
----------	----------	--------------	------------

---

1	0.4168	0.5096	-0.0928
2	0.2394	0.3620	-0.1226
3	0.1376	0.2342	-0.0966
4	0.0726	0.1502	-0.0776
5	0.0448	0.0872	-0.0424
6	0.0272	0.0444	-0.0172
7	0.0078	0.0246	-0.0168
8	0.0068	0.0128	-0.0060
9	0.0024	0.0064	-0.0040
10	0.0004	0.0034	-0.0030

--- Stop Reason Distribution ---

Reason		Full Config		PBM-like
exam_fail		94.6%		99.3%
abandonment		5.1%		0.0%
purchase_limit		0.2%		0.0%
end		0.2%		0.7%

--- Interpretation ---

Key observations:

1. PBM-like config has NO abandonment (threshold = -100)
2. PBM-like config has NO purchase limit stopping
3. PBM-like CTR depends only on position (via pos\_bias)
4. Full config CTR varies with utility (price, PL, noise)

This verifies [PROP-2.5.4]: Utility-Based Cascade nests PBM as a special case when utility dependence is disabled.

---

## 2.8 Lab 2.5 – Utility-Based Cascade Dynamics ([DEF-2.5.3])

**Goal:** Verify the three key mechanisms of the production click model from Section 2.5.4: position decay, satisfaction dynamics, and stopping conditions.

### 2.8.1 Solution

```
from scripts.ch02.lab_solutions import lab_2_5_utility_cascade_dynamics

results = lab_2_5_utility_cascade_dynamics(seed=42, verbose=True)
```

Actual Output:

```
=====
Lab 2.5: Utility-Based Cascade Dynamics ([DEF-2.5.3])
=====
```

Verifying three key mechanisms:

1. Position decay (pos\_bias)
2. Satisfaction dynamics (gain/decay)
3. Stopping conditions

```

Configuration:
  Positions: 20
  satisfaction_gain: 0.5
  satisfaction_decay: 0.2
  abandonment_threshold: -2.0
  pos_bias (category, first 5): [1.2, 0.9, 0.7, 0.5, 0.3]

```

Simulating 2,000 sessions...

--- Part 1: Position Decay ---

Position	Exam Rate	CTR Exam	pos_bias
1	0.767	0.387	1.20
2	0.520	0.563	0.90
3	0.349	0.401	0.70
4	0.197	0.353	0.50
5	0.100	0.485	0.30
6	0.052	0.533	0.20
7	0.025	0.353	0.20
8	0.015	0.600	0.20
9	0.005	0.400	0.20
10	0.002	1.000	0.20

Observation: Examination rate decays with position, matching pos\_bias pattern.

--- Part 2: Satisfaction Dynamics ---

Sample satisfaction trajectories (first 5 sessions):

```

Session 1: 0.00 -> -0.20 (exam_fail)
Session 2: 0.00 -> -0.20 -> 0.22 -> 0.02 -> -1.75 (exam_fail)
Session 3: 0.00 -> -0.20 -> 0.18 -> -0.29 (exam_fail)
Session 4: 0.00 -> -0.20 -> 0.23 -> 0.03 -> -0.44 -> -0.64 -> -0.33 -> -0.53 ... (exam_fail)
Session 5: 0.00 -> -0.20 (exam_fail)

```

Final satisfaction statistics:

```

Mean: -0.49
Std: 0.71
Min: -3.47
Max: 1.79

```

--- Part 3: Stopping Conditions ---

Stop Reason	Count	Percentage
exam_fail	1900	95.0%
abandonment	98	4.9%
purchase_limit	2	0.1%
end	0	0.0%

Session length statistics:

```

Mean: 2.0 positions
Std: 1.9
Median: 2

```

```
Clicks per session:
```

```
Mean: 0.90
```

```
Max: 7
```

```
--- Verification Summary ---
```

```
checkmark Position decay: Examination rate follows pos_bias pattern
```

```
checkmark Satisfaction dynamics: Trajectories show gain on click, decay on no-click
```

```
checkmark Stopping conditions: All three mechanisms observed (exam, abandon, limit)
```

---

## 2.9 Summary: Theory-Practice Insights

These labs validated the measure-theoretic foundations of Chapter 2:

Lab	Key Discovery	Chapter Reference
Lab 2.1	Segment frequencies converge at $O(1/\sqrt{n})$	DEF-2.2.2, LLN
Lab 2.1 Task 2	Zero-probability segments break IPS	THM-2.6.1, Positivity
Lab 2.2	Base scores square-integrable (finite variance)	PROP-2.8.1
Lab 2.2 Task 2	Score histogram enables Radon-Nikodym intuition	THM-2.3.4-RN
Lab 2.3	PBM and DBN match theory exactly	DEF-2.5.1, DEF-2.5.2, EQ-2.1, EQ-2.3
Lab 2.4	Utility-Based Cascade nests PBM	PROP-2.5.4, DEF-2.5.3
Lab 2.5	Position decay + satisfaction + stopping verified	DEF-2.5.3, [EQ-2.4]-[EQ-2.8]
Extended: IPS	IPS is unbiased but high variance	THM-2.6.1, EQ-2.9

### Key Lessons:

1. **Measure theory isn't abstract:** Every  $\sigma$ -algebra and probability measure has a concrete implementation in `zoosim`. The math ensures our code is correct.
2. **Positivity is critical:** When  $(\mathbf{p}_{\text{seg}})_i = 0$  (zero-probability segment), IPS fails. This is the measure-theoretic formulation of “support deficiency”—a real production failure mode.
3. **LLN and CLT quantify convergence:** The  $O(1/\sqrt{n})$  scaling of  $L_\infty$  deviation isn't just theory—it predicts exactly how many samples we need for reliable estimates.
4. **Click models encode assumptions:** PBM assumes independence (simpler, less accurate). DBN encodes cascade dependence (more accurate, harder to estimate). Both are rigorously defined probability spaces.
5. **IPS is the Radon-Nikodym derivative:** The importance weight  $\pi_1/\pi_0$  is exactly  $d\mathbb{P}^{\pi_1}/d\mathbb{P}^{\pi_0}$ . Unbiasedness follows from change-of-measure, but variance explodes when policies differ substantially.

---

## 2.10 Running the Code

All solutions are in `scripts/ch02/lab_solutions.py`:

```
# Run all labs
python scripts/ch02/lab_solutions.py --all

# Run specific lab
python scripts/ch02/lab_solutions.py --lab 2.1
python scripts/ch02/lab_solutions.py --lab 2.2
python scripts/ch02/lab_solutions.py --lab 2.3
python scripts/ch02/lab_solutions.py --lab 2.4
python scripts/ch02/lab_solutions.py --lab 2.5

# Run extended exercises
python scripts/ch02/lab_solutions.py --extended clicks
python scripts/ch02/lab_solutions.py --extended ips

# Interactive menu
python scripts/ch02/lab_solutions.py
```

---

*End of Lab Solutions*

Billingsley, Patrick. 1995. *Probability and Measure*. 3rd ed. Wiley.

Durrett, Rick. 2019. *Probability: Theory and Examples*. 5th ed. Cambridge University Press.