AI/ML – Dangers of Feedback Loops

Group members:	Date:
Chosen domain (circle one): Music	• News • Shopping • Short-video • Jobs
1) What is a feedback loop? (2-	-3 sentences)
In your own words:	

2) Mini-Simulation: "RecoNarrow v1"

Day	Belief A (%)	Belief B (%)	#A shown (of 10)	#B shown (of 10)	Clicks A (20%)	Clicks B (15%)
0	60	40	6	4	1	1
1						
2						
3						

What happened? Why?

3) Where can harm occur? (tick all that apply and explain briefly)	
[] Loss of diversity (users see less variety) →	
[] Unfair exposure (Category B creators never get surfaced) \rightarrow	
[] Stereotype reinforcement (system doubles down on past patterns) \rightarrow	
[] User disengagement (boredom/echo chamber) →	

Stress test: cold start or trer	id shift
---	----------

Imagine a new Category C appears on Day 2 (22% click rate if shown, but the model doesn't know that).

If your system only shows A/B based on past clicks, how likely is C to be discovered?

What	simple	rule	could	ensure	exp	loration	so C	gets a	fair	chance	?
vviiat	Simple	I GIO	ooulu	CHOCK	CAP	oration	30 C	goto a	IUII	or idi ioc	•

5) Break the loop — Design fixes (choose at least three)
- Exploration quota (10–20% of slots reserved for under-shown items) → How/Trade-off:
- Diversity metric (no more than 4/10 from one category) → How/Trade-off:
- Debiasing the signal (normalise clicks by impressions) → How/Trade-off:
- Fair exposure constraints (per-creator/category floors) → How/Trade-off:
- User controls (toggle 'wider variety', reset) → How/Trade-off:
- Fresh-start windows (periodic re-seeding with broader data) → How/Trade-off:

6) Your group's recommendation (pitch)

Write a 5–7 sentence product note to your Head of Recommenders explaining the risks and your fixes.

7) Success metrics (pick at least two and define a target)

- Category diversity (entropy ↑ by ____% within 14 days)
- Creator exposure fairness (Gini ↓ to ≤ ____)
- User satisfaction (survey score ↑ to ≥ ____/5)
- Retention (7-day return rate ↑ to ____%)
- Novel discovery rate (share of clicks on 'new to user' ↑ to ____%)

8) Reflection (individual, 2–3 sentences)

What did you learn about feedback loops that surprised you?