Optimizing Subscription Services with Uplift Modeling

Dmitry Zolotukhin, Chief Data Officer, START

Plan

- 1. Background of our Uplift task
- 2. Data
- 3. Data preprocessing
- 4. Model training
- 5. Results

Company and Business Model

AVOD viewers watch ads, **SVOD** viewers buy a subscription



Task: Expanding Subscriber Base

How: offer discounts

Problem: don't offer a discount to someone who would subscribe anyway

Solution: Uplift modeling



Uplift Modeling

Uplift = difference in revenue from a user when targeted vs. not targeted

Problem: You can't both target and not target the same user!

Solution: Average Treatment Effect (ATE) — difference in revenue (or another metric) between test and control groups. This works because in A/B testing, users are randomly assigned to test or control

Requirements for the Uplift Model

- 1. Stability
- 2. Interpretability

These requirements stem from:

- 1. The high cost of acquiring users \rightarrow communication data is expensive
- 2. Movie/TV content is subject to elastic demand (frequent new releases)

Data

- Experiment MORE-2363 on users with a trial:
 - 17,000 users in the test group, offered a 30-day trial extension (handled by billing)
 - 17,000 users in the control group, no extension offered; auto-renewal disabled after users answered why they canceled
- We trained our first Uplift model on these data.



Нам есть что предложить

 К сожалению, иногда случаются технические неполадки, но мы их быстро отслеживаем и устраняем.
Не отключай подписку сейчас и получи

30 дней подписки бесплатно

пока не отключать подписку

ВСЁ-ТАКИ ОТКЛЮЧИТЬ ПОДПИСКУ

, hu

User Features

Number of active days

From the player:

- Number of watched projects
- Number of projects added to favorites
- Percentage of seasons completed
- Number of TV show views
- Number of movie views
- Number of successful auto-advances to the next episode
- Number of attempts at auto-advance to the next episode

From the project page:

- Number of times "Similar to this" was viewed
- Number of clicks on "Similar to this"

From search interaction:

- Average number of actions between a search and a watch
- Number of queries that yielded no results

Data Preprocessing

Step 1	Step 2	Step 3	
Split data into training and validation sets	Train the model	Check the validation score	

What is the Score?

Qini is the difference in the share of successful conversions in test versus control.

Model Prediction	From Test Group?	Converted into Subscription?	Qini
0.939249	1	1	1.000000
0.939237	1	1	2.000000
0.939125	0	1	0.875000
0.938707	1	1	2.173913
0.936231	1	0	2.173913
-0.925874	0	1	40.266077
-0.938097	0	1	40.019571



Data Preprocessing: Issues

No stability in quality: drastically different results with different random splits



Data Preprocessing: Issues

No stability in quality: drastically different results with different random splits



How to Fix It: Split by Registration Date

The oldest users (by registration date) go to train (80% of the total sample). The remaining 20% of more recently registered users go to validation

Train the model on the training set, get the score on the validation set

UpliftRandomForest

Train score: 0.338

Validation score: 0.330



Stratification Algorithm

Goal: ensure that the feature distributions in each subset are as close as possible to those in the overall population

Problem: many features, many values

Solution: choose the strongest feature (based on Gini coefficient between feature and subscription conversion), then binarize it

We used 3 features:

- Number of active days
- Whether the user converted
- Whether the user had the offer available

This produced 12 strata. We form clusters such that the proportions of each stratum remain consistent.



How to Fix It: Random Stratification

- 1. Take half of the sample for train; the rest is the pool from which we will draw data
- 2. From the pool, stratify and select 1,000 samples
- 3. Add these samples to train
- 4. Train the model
- 5. Evaluate on the remaining data
- 6. If Qini on the training and validation sets differs by no more than 5% and is better than random, augment the training set with those 1,000 samples
- 7. Otherwise, do not update the training set
- 8. Stop when the training set reaches at least 80% of the total data

("Stratified" means the sub-sample of 1,000 preserves the proportions of "active days," "converted or not," and "offer availability" as in the original dataset.)

							Metrics 3		
1 Start Time	Duration	Run Name		Source	Version	Models	Qini score diffe	Qini train score	Qini val scor
@ 11 days ago	6.0min		yst	🗆 ipykerne			0.01	0.313	0.323
@ 12 days ago	12.6min	-	yst	🗀 ipykerne	-	25	0.079	0.322	0.243
@ 12 days ago	14.0min	2	yst	🗆 ipykerne		40)	0.154	0.49	0.336
@ 12 days ago	3.0min	1	yst	🗅 ipykerne	-	÷.	0.009	0.316	0.325
@ 12 days ago	1,7min	10	yst	🗀 ipykerne		T	0.014	0.158	0.144
🙆 12 days ago	18.9min		yst	🗀 ipykerne		-	0.008	0.196	0.188
@ 12 days ago	50.4s	*	yst	🖂 ipykerne		÷2	0.849	0.871	0.023
@ 12 days ago	6.8min		yst	🗀 ipykerne		÷.	0.013	0.213	0.2
🗎 12 days ago	6.1min		yst	🗆 ipykerne		t.	0.019	0.22	0.201
@ 12 days ago	7.3min		yst	ipykerne_			0.011	0.276	0.287



How to Fix It: Thompson Sampling

Beta distribution for data clusters

Success = model score on train & validation is similar and exceeds random

Failure = not meeting one or both of these conditions

We incorporate the successes/failures of previous iterations. We account for changes in the external environment and training results.



How to Fix It: Thompson Sampling

- 1. Half of the data = train; the other half is a pool
- 2. We split the pool into 4 stratified clusters
- 3. Select a cluster from which to draw samples (via Beta distribution)
- 4. Take at least 1,000 samples from that cluster
- 5. Add them to train
- 6. Train the model
- 7. Evaluate on the remaining data
- 8. If Qini on train and validation differs by $\leq 5\%$ and is better than random, add the samples to train (α += 1)
- 9. Else do not update train (β += 1)
- 10. Stop when the training set reaches at least 80% of all data

It often fails to converge reliably.

	1 Start Time	Duration			Version	Models	Metrics (
			Run Name	rce			Qini score diffe	Qini train score	Qini val score	train_size
	@ Z days ago	3.1min		sykeme	80	14	0.01	0.153	0.163	29822
	🛞 2 days ago	13.0min		zykerne_	83	1.8	0.008	0.369	0.377	29822
	e 2 days ago	5.8min		zykerne_	X		0.014	0.195	0.181	29254
	@ 2 days ago	50.4min		zykerne_			0.018	0.191	0.173	29797
0	© Z days ago	8.7min		zykerne	÷.		0.008	0.137	0.129	29254
	@ 2 days ago	15.2min	Q	zykerne	4. C	4	5.895e-4	0.094	0.093	29822
	© Z days ago	6.4min		zykerne_	÷.		0.001	0.215	0.213	29254
	🛛 2 days ago	19.3min	3	zykerne		3. *	0.019	0.228	0.209	29797
	🛛 2 days ago	3.6h		zykerne_	6	1	0.588	0.568	-0.02	26254
	© 2 days ago	35.2min	2	zykerne	÷	22	0.006	0.184	0.19	29822
	© 2 days ago	3.1h	a,	zykerne_	2	8 4	0.67	0.789	0.119	18254
	Ø 2 døys ago	3.2h	28	zykerne	8		0.468	0.557	0.09	21254
	@ Z days ago	14.6min	8	zykerne	8	*	0.013	0.295	0.308	29254
	@ 2 days ago	2.9h		zykerne_	8		0.007	0.081	0.088	29797
	2 days ago	12.7min		zykerne_			0.012	0.16	0.148	29254



How to Fix It: Thompson Sampling + Registration Date

- 1. Sort data by registration date
- 2. Oldest $50\% \rightarrow$ train; remaining 50% is the pool
- 3. Split the pool into 4 stratified clusters
- 4. Select which cluster to pull samples from (Beta distribution)
- 5. Take at least 1,000 samples from that cluster
- 6. Add them to train
- 7. Train the model
- 8. Evaluate on the remaining data
- 9. If Qini on train and validation differs by $\leq 5\%$ and exceeds random, expand train ($\alpha += 1$)
- 10. Else, do not update train (β += 1)
- 11. Stop when train is $\geq 80\%$ of the total data

Works more reliably, yields higher scores-great!

								Metrics >		
	1 Start Time	Duration	Run Name	r .	Source	Version	Models	Qini score diffe	Qini train score	Qini val score
	© 2 days ago	11.9min		yst	🗆 ipykerne		-	0.011	0.417	0.407
	@ 2 days ago	1.1h	54.	yst	🖂 ipykerne	S	(+)	0.001	0.442	0.441
	© 2 days ago	1.8h	<u>.</u>	yst	ipykerne	*	(*):	0.007	0.395	0.388
	@ 2 days ago	44.8min		yst	🔲 ipykerne	15	2.52	0.014	0.364	0.378
	© 2 days ago	16.6min	-	yst	ipykerne	-		0.003	0.405	0.402
	⊘ 2 days ago	3.3h		yst	🗆 ipykerne	4		0.346	0.442	0.096
0	© 2 days ago	59.9min	12	yst	🛄 ipykerne	3.	(in)	0.019	0.47	0.489
0	😔 2 days ago	1.7h		yst	🗆 ipykerne	*		0.003	0.376	0.373
	@ 2 days ago	3.0h	15	yst	ipykerne		592	0.014	0.412	0.398
	@ 2 days ago	43.2min	-	yst	ipykerne			0.017	0.44	0.457
	@ 2 days ago	28.2min	S2	yst	🗆 ipykerne	Q	+	0.01	0.446	0.436
	@ 2 days ago	10.3min	54	yst	ipykerne_	×		0.008	0.326	0.335
	Ø 2 days ago	38.4min	10	yst	🗀 ipykerne		(14))	0.01	0.437	0.447
0	@ 2 days ago	10.1min		yst	ipykerne		(T)	7.332e-4	0.302	0.302
	@ 2 days ago	33.3min		yst	ipykerne		(*)	0.013	0.351	0.338



Design AB

- 1. Select new trial users
- 2. Select 5% of them to not receive an offer (control group), and 5% to receive an offer
- 3. For the remaining 90%:
 - a. Score the users with the model
 - b. Take the top 20% by model score
 - c. Send them the offer
 - d. Do not send the offer to the other 80%
- 4. Compare the test group's conversion to the control group (it should be higher)
- 5. Compare the conversion among the test group who did not receive the offer (80%) to the control group (they should be the same)



Next Steps

Experimenting with:

- Discount amount
- Offer duration (1–3–6 months)
- Timing of discount (1–2–3 weeks after start)
- Number of existing subscriptions a user has when the offer is made