

IA316

Youtube environment

Thomas Dahmen - Paul Chevalier

I. Environment presentation

II. Experiments & results

Real Youtube

Users = channels

Users can :

- post videos
- watch videos (watch time)
- like / dislike videos
- subscribe to channels
- create playlists and add videos to them

Our Youtube

Users \neq channels

Channels can post videos

Users can watch videos

Watch time $\in [0, 1]$ (\approx like)

Real Youtube

Videos :

- explicit content (3D tensor)
- implicit content (tags, channel, metadata...)

Our Youtube

Videos = feature vectors, **n dimensional, unit norm**.

Each feature = a content **keyword** (e.g. : humor, rock, sport, ...).

Sparse feature vectors = mix of a few keywords.

Same for users : tastes = sparse unit vectors of keywords.

Rewards

Reward = watch time $\in [0, 1]$.

→ maximize the watch time of recommended videos.

User \mathbf{u} , video \mathbf{v} .

Cosine similarity for tuple (\mathbf{u}, \mathbf{v}) : $s(\mathbf{u}, \mathbf{v})$.

Reward probability model for (\mathbf{u}, \mathbf{v}) :

- the higher the similarity, the higher the watch time in average
- extreme similarities (≈ 0 or 1) incur less variance

Rewards

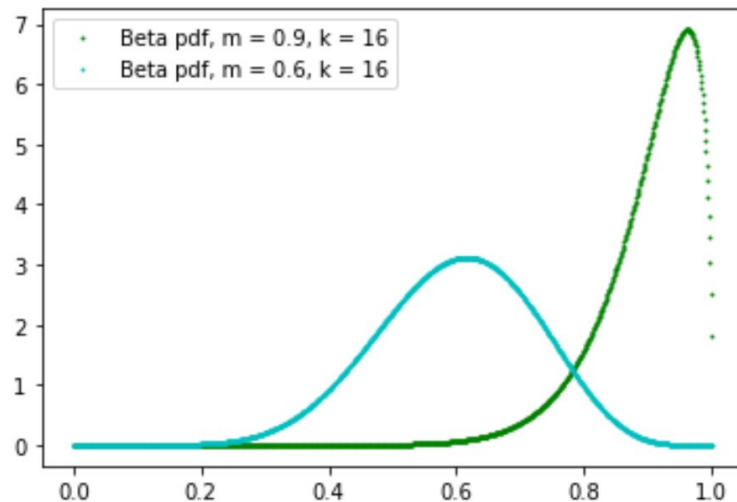
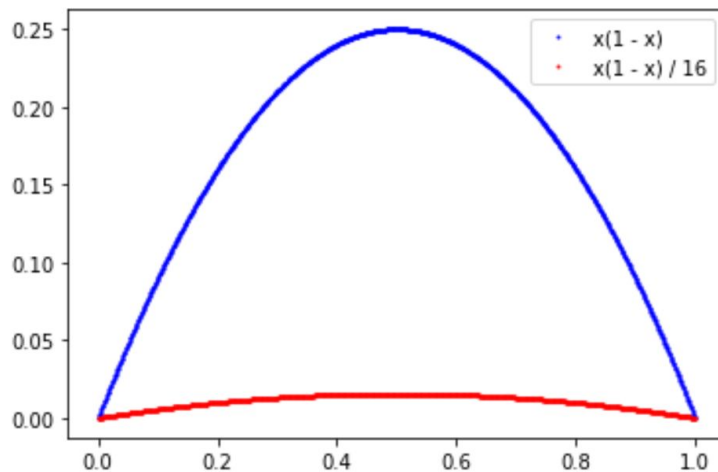
Mean = $m = s(\mathbf{u}, \mathbf{v})$

Var = $(m * (1 - m)) / k$

→ **Reward \sim Beta(a, b)**

$a = (k - 1) * s(\mathbf{u}, \mathbf{v})$

$b = (k - 1) * (1 - s(\mathbf{u}, \mathbf{v}))$



Evolution

Recommended videos can influence user tastes.

→ **the environment can be evolutive**

User \mathbf{u} (\mathbf{u}_0 at first), recommended video \mathbf{v} , watch time $\mathbf{r} \in [0, 1]$

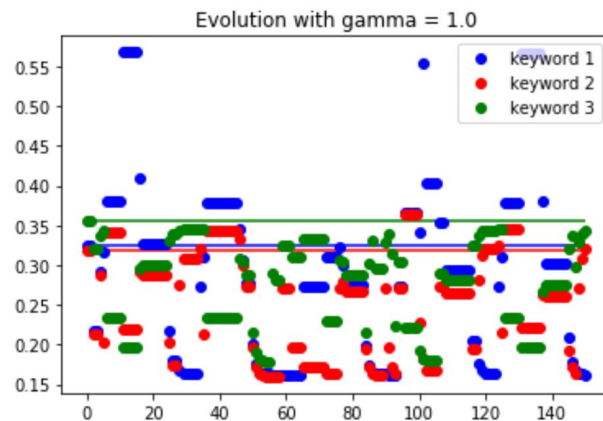
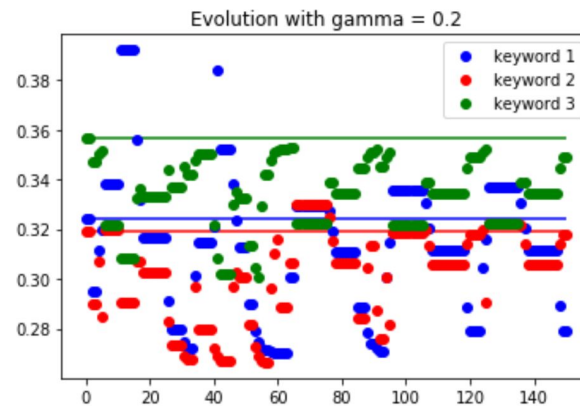
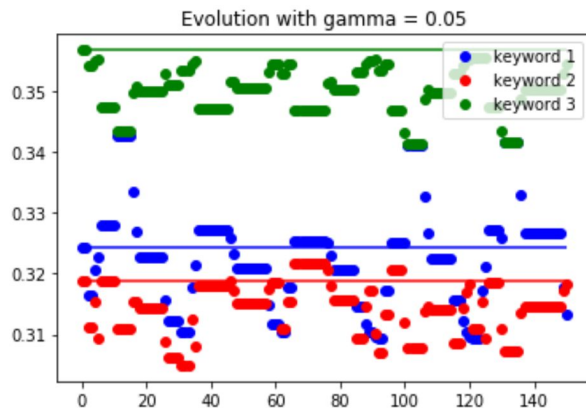
$\mathbf{r} \approx$ continuous version of “like”

With probability \mathbf{r} , let \mathbf{u} change :

$$\mathbf{u} \leftarrow \mathbf{u}_0 + \mathbf{u} + \gamma \times \mathbf{v}$$

γ : influence of the recommended videos

Evolution



Experiments

Agents used:

- Epsilon Greedy (epsilon = 0.1)
- Thompson Sampling
- Q-Learning

In two cases:

- Non evolutive environment
- Evolutive environment

Environment parameters

- 10 users
- 10 channels
- 3 videos per channel
- 100 keywords
- 3 keywords per user
- 3 keywords per video

Q Learning

Initialized

Q-Table		Actions					
		South (0)	North (1)	East (2)	West (3)	Pickup (4)	Dropoff (5)
States	0	0	0	0	0	0	0

	327	0	0	0	0	0	0

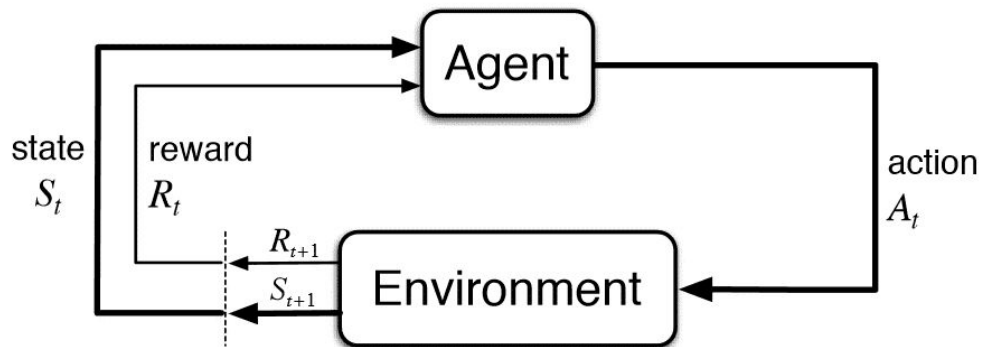
	499	0	0	0	0	0	0

Training

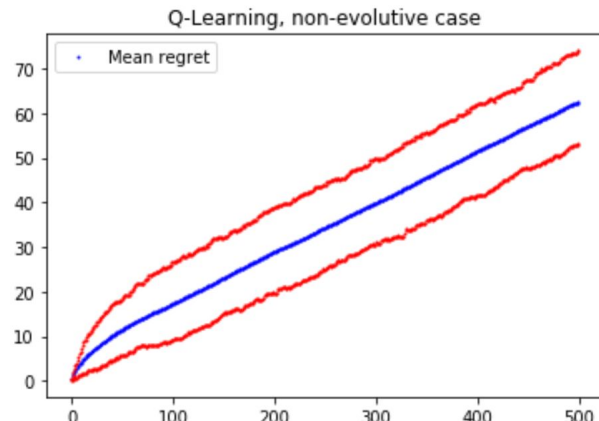
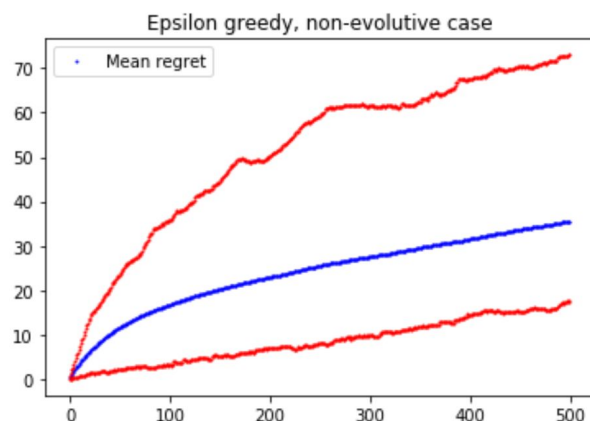
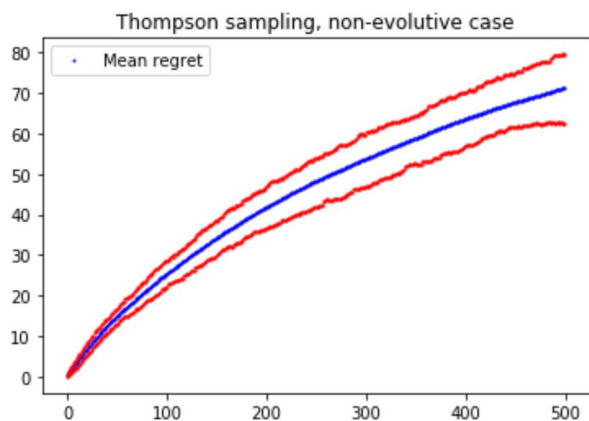
Q-Table		Actions					
		South (0)	North (1)	East (2)	West (3)	Pickup (4)	Dropoff (5)
States	0	0	0	0	0	0	0

	328	-2.30108105	-1.97092096	-2.30357004	-2.20591839	-10.3607344	-8.5583017

	499	9.96984239	4.02706992	12.96022777	29	3.32877873	3.38230603

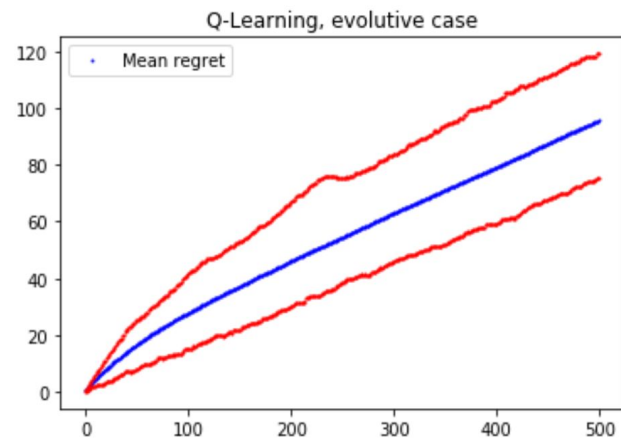
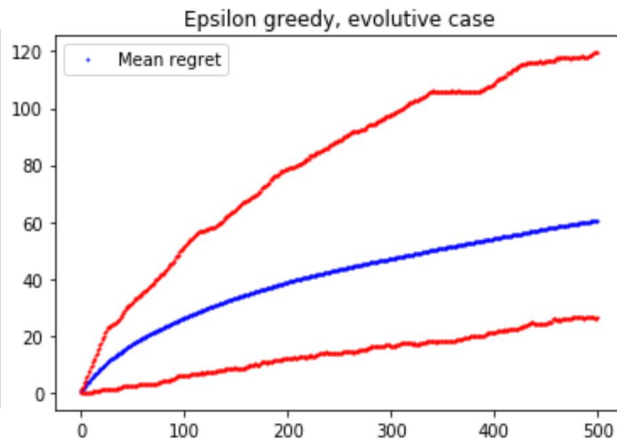
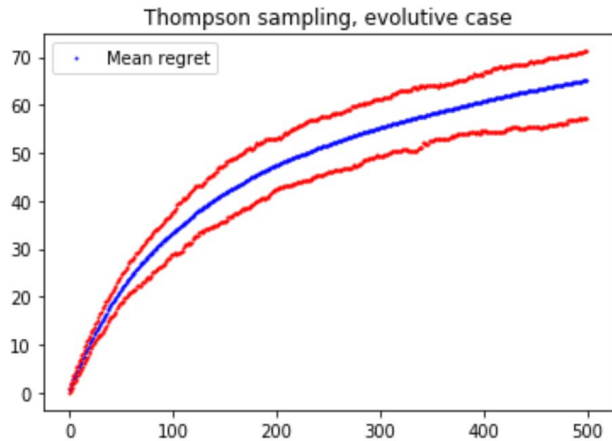


Experiments in non-evolutive case



- Time horizon : 500 steps
- Number of simulations : 100

Experiments in evolutive case



- Time horizon : 500 steps
- Number of simulations : 100

Conclusion

- Epsilon Greedy is better in average, but has a high variance
- Thompson Sampling is more stable
- Q Learning is the less satisfying agent
- In average, regret is higher for an evolutive environment

To go further...

- Experiment with a bigger environment (more users and videos)
- Change video/user ratio
- Use collaborative algorithms (such as top videos)
- Use Deep Q Learning