

Welcome

Kevin Jamieson (Instructor)

Lalit Jain (co-instructor)

Yuhao Wan (TA)

Interactive learning in ~~non-stochastic~~ environments
better modeled

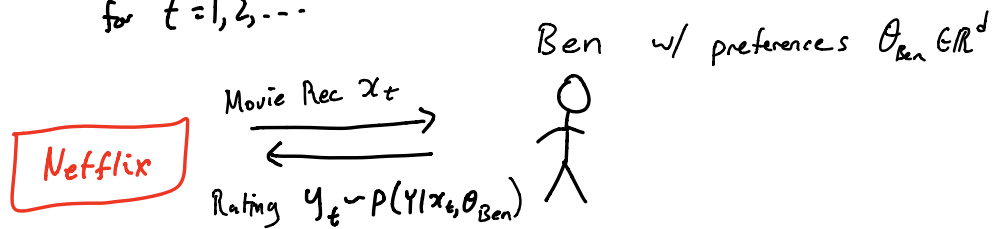
Machine learning relies on stochastic IID assumptions

$\{(x_i, y_i)\}_{i=1}^n$ where $(x_i, y_i) \stackrel{iid}{\sim} \mathcal{D}_{x,y}$

Learns predictor $\hat{h}: \mathcal{X} \rightarrow \mathcal{Y}$ using $\{(x_i, y_i)\}_{i=1}^n$

so that for new point $(x_a, y_a) \stackrel{iid}{\sim} \mathcal{D}_{x,y}$ $\hat{h}(x_a) \approx y_a$.

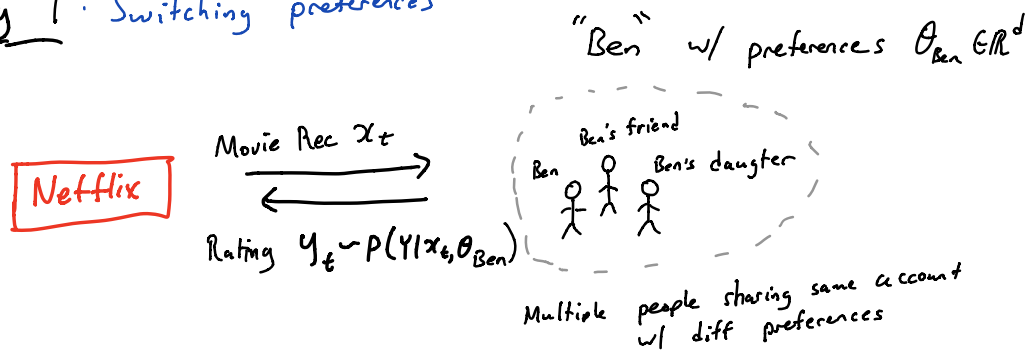
Model for $t=1, 2, \dots$



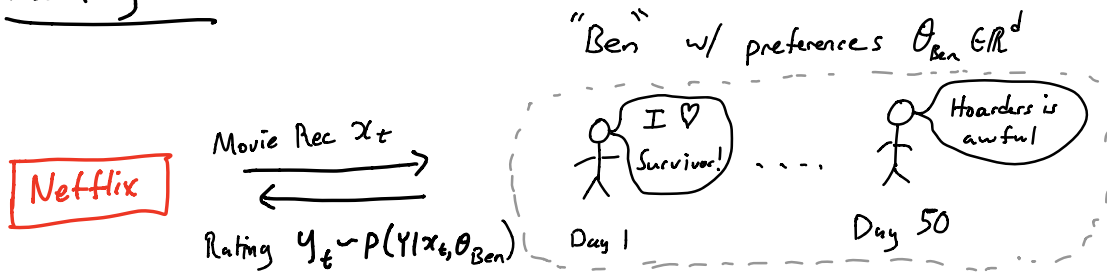
Things we worry about

- Trading off exploration versus exploitation (i.e. don't be too greedy)

Reality 1: Switching preferences



Reality 2: Time variation



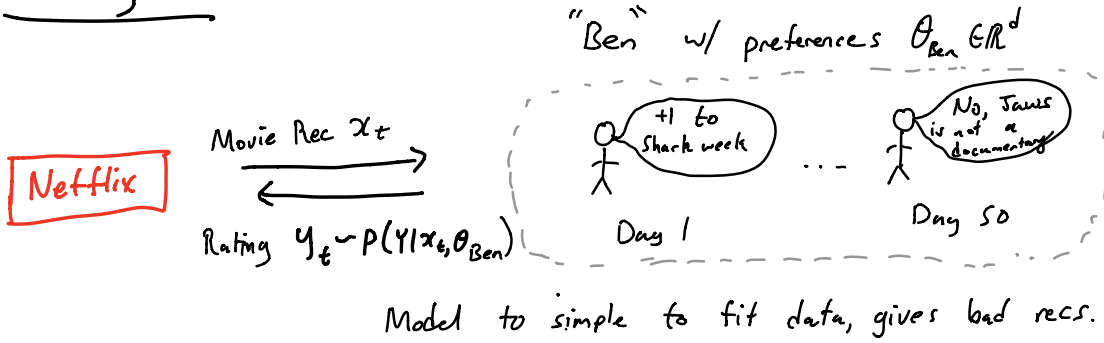
User preferences change over time...

Content producers react to changing preferences...

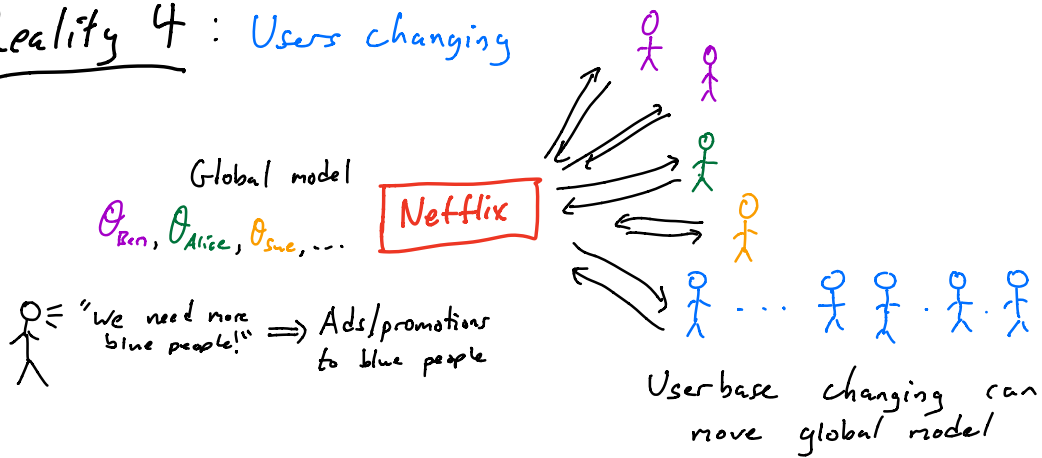
User preferences react to changing content...

⋮

Reality 3 : Model misspecification



Reality 4 : Users changing



This class is about

- Modeling
- Robustness
- Guarantees
- Practical algorithms (mostly)

Inspired by real-world scenarios we have encountered.

Class Logistics, evaluation, plan
zoom etiquette

Multi-armed Bandit

Each "arm" is a source
s.t. when "pulled" emits
some observation $X \in \mathbb{R}$

$$n \geq 2, [n] = \{1, 2, \dots, n\}$$

Input n arms

for $t=1, 2, \dots, T$ $\swarrow X_t \in \mathbb{R}^n$
Nature chooses $X_{t,i} \in \mathbb{R}$ for $i=1, 2, \dots, n$ (secretly)

Player chooses arm $I_t \in [n]$ and
observes X_{t,I_t}

The regret of the player is defined as

$$R(T) = \underbrace{\max_{i \in [n]} \sum_{t=1}^T X_{t,i}}_{\text{The reward of playing the single best action in hindsight}} - \underbrace{\sum_{t=1}^T X_{t,I_t}}_{\text{reward obtained by the player}}$$

Typically, goal is to define an algorithm

$$\text{s.t. } R(T) = o(T) \text{ so that } \frac{R(T)}{T} \rightarrow 0$$

If we assume nothing about rewards at all

then \exists choice of rewards $\{X_t\}_t \subset \mathbb{R}^n$ s.t.

$$R(T) = \Omega(T) \text{ w. const prob.}$$

So, we make assumptions.

Adversarial setting:

Bounded rewards $X_{t,i} \in [0,1] \quad \forall t, i \in [n]$

\exists algorithm s.t. $\mathbb{E}[R(T)] \leq O(\sqrt{nT})$

Stochastic setting

\exists distributions over \mathbb{R} s.t. ν_1, \dots, ν_n

$X_{t,i} \stackrel{iid}{\sim} \nu_i$ for all $t \in \mathbb{N}$.

Typically we assume each ν_i has bounded support, 1st/2nd moment, or is sub-Gaussian.

\exists algorithm s.t. $\mathbb{E}[R(T)] \leq O(\min\{C_V \log(T), \sqrt{nT}\})$

where C_V depends on "hardness" of ν_1, \dots, ν_n

When is stochastic justified? (Not on T)

