



Aprendizado por Reforço

AULA - 6

Off-Policy Learning

Anteriormente...

- Diferença Temporal

$$V(s) \leftarrow V(s) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s)]$$

- Deep Q-Network

$$L(\theta) = \left(r + \gamma \max_{a'} Q(s', a'; \phi) - Q(s, a; \theta) \right)^2$$



Anteriormente...

- Gradiente de Política

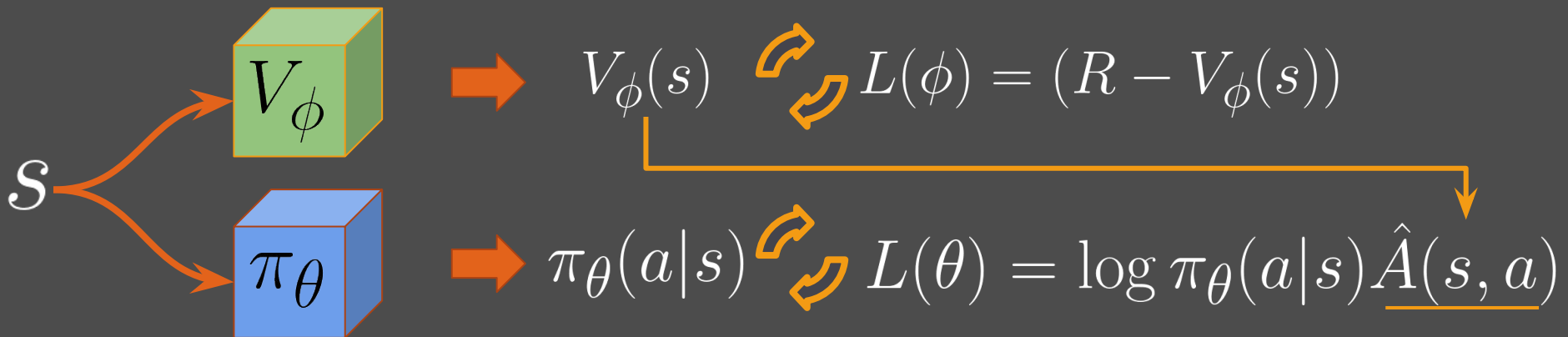
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) R_{i,t}$$

- Com Baseline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) [R_{i,t} - V_{\phi}(s_t)]$$

Anteriormente...

- Actor-Critic



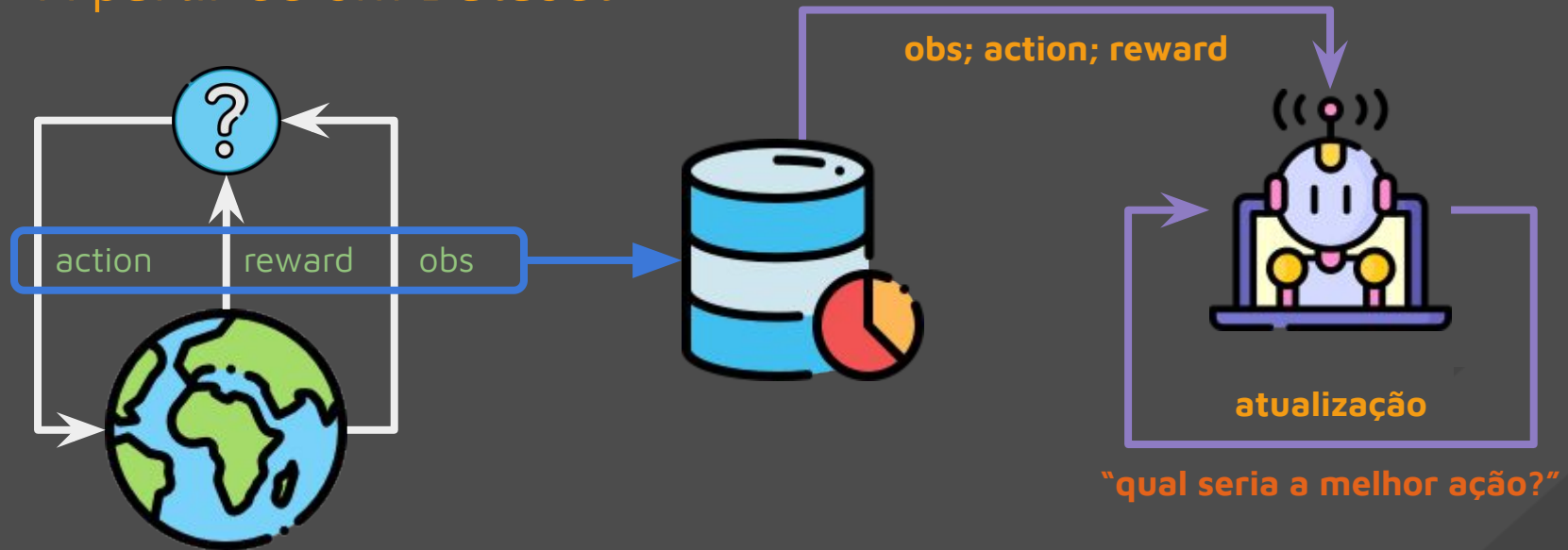
A3C - PPO - SAC



Off-Policy

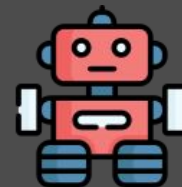
O que é aprender *Off-Policy*?

- Como o nome diz “fora da política”
- Com experiências de outras políticas
- A partir de um Dataset



O que é Offline e Online Learning?

- **Offline:** Não há acesso ao ambiente, agente aprende com experiências previamente coletadas
- **Online:** Há acesso ao ambiente, agente coleta novas experiências com uma certa frequência

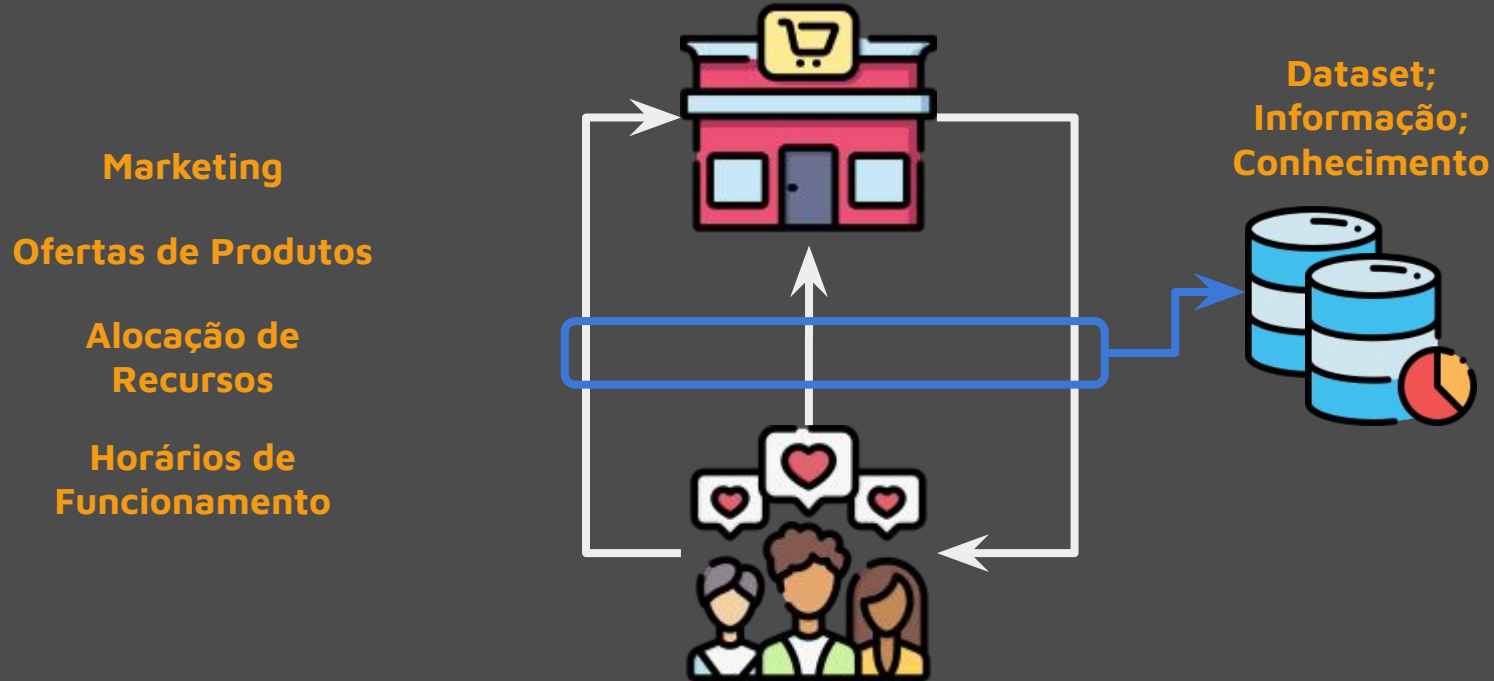


Por que Off-Policy/Offline RL é Importante?

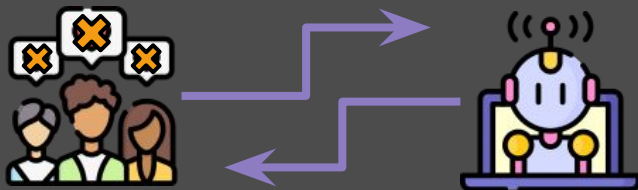
- O mundo possui uma quantidade enorme de dados
- Podemos usar esses dados para:
 - Minimizar Custos
 - Minimizar Riscos
 - Aproveitar Experiências



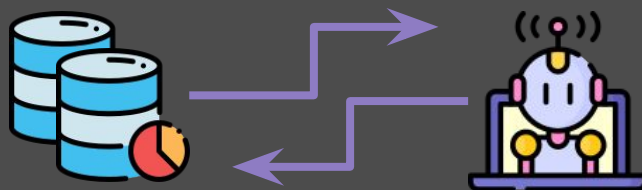
No Mundo Real...



Dados ou Ambiente?



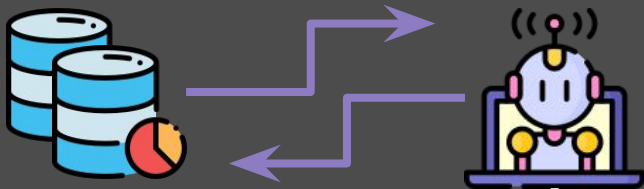
- Interagir com o ambiente pode ser caro e arriscado, pois erros precisam ser cometidos para explorar e aprender
- Permite exploração
- Otimização direta do ambiente



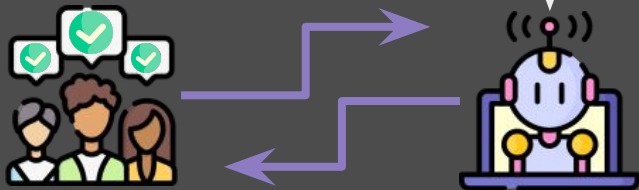
- Interagir com os dados é seguro, e, muitas vezes, mais rápido.
- Aproveita o conhecimento de políticas de coleta (humano ou modelos)
- Não há espaço para exploração
- Otimização limitada pelos dados

Podemos ter Ambos?

1



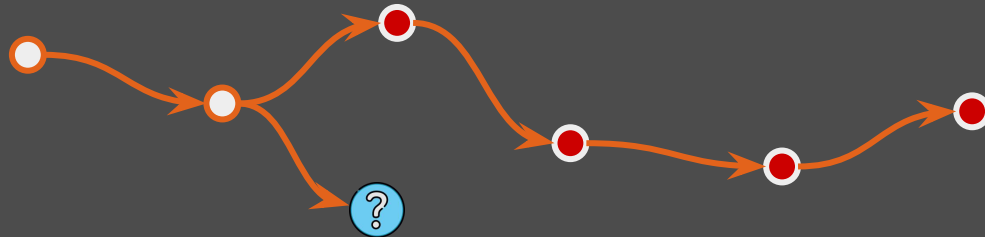
2



- Se a situação permitir
- Treino com dados primeiro
- Atuação no ambiente depois
- Aprender no ambiente
- Modelo menos aleatório
- Aplicação com menos riscos

Problema de Superestimação

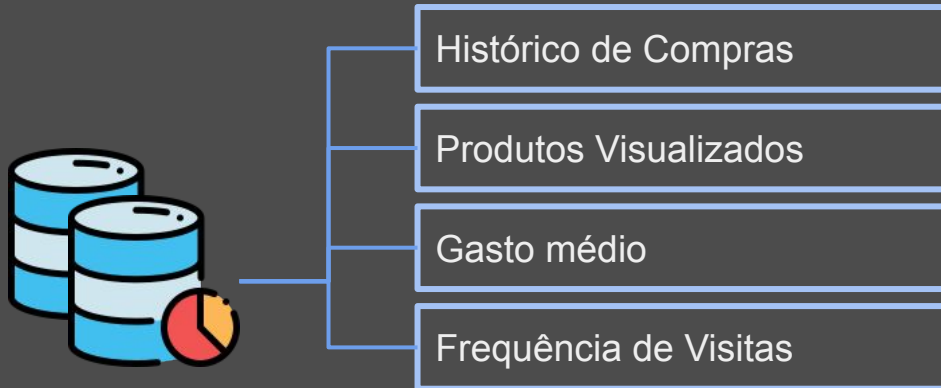
- Treinar com dados nem sempre é fácil
- Se os dados apresentam um caminho ruim, é fácil para o modelo superestimar o caminho não percorrido



- Quanto mais representativo, menor esse problema
- Dados muito representativos são incomuns, porque geralmente alguma política está sendo seguida para coletá-los

Dificuldades com os Dados

- Informações Disponíveis
- Há informações suficientes para tomar decisões?
- Construção da representação do estado
- Construção da Recompensa

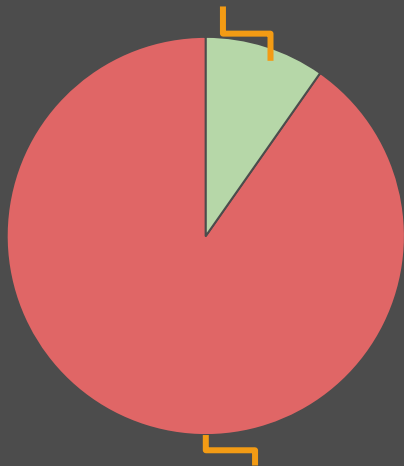


Dificuldades com os Dados

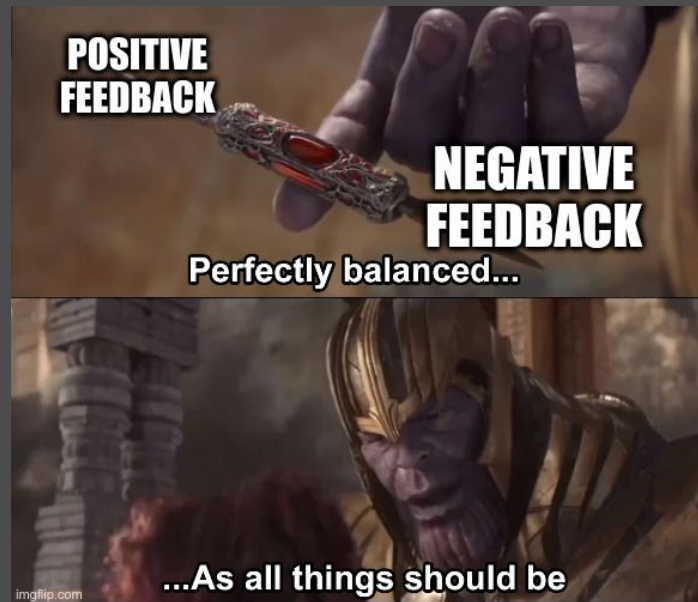


- Desproporcionalidade
- Talvez é necessário repensar as recompensas ou filtrar os dados

Sequências que geram Recompensas Positivas



Sequências que geram Recompensas Negativas/Zero



Off-Policy Estimators

- Tão bons quanto os dados disponíveis



Projeto com a Acordo Certo



acordocerto
Seu bem-estar financeiro.

- **Problema**
- Trazer usuários para a plataforma para quitar suas dívidas através de mensagens SMS
- Não se pode acionar usuários com menos de 3 dias do último contato
- Não se pode acionar todos os usuários (são muitos)

Dados Disponíveis

- Montar Recompensa

- cadastro +
- acordo ++
- acionamento -

- Montar Estados

- Estado Markoviano
- Histórico
- Último acionamento?
- Quantos acionamentos?



Histórico de Acionamentos

Ações no Site (cadastro/acordo)

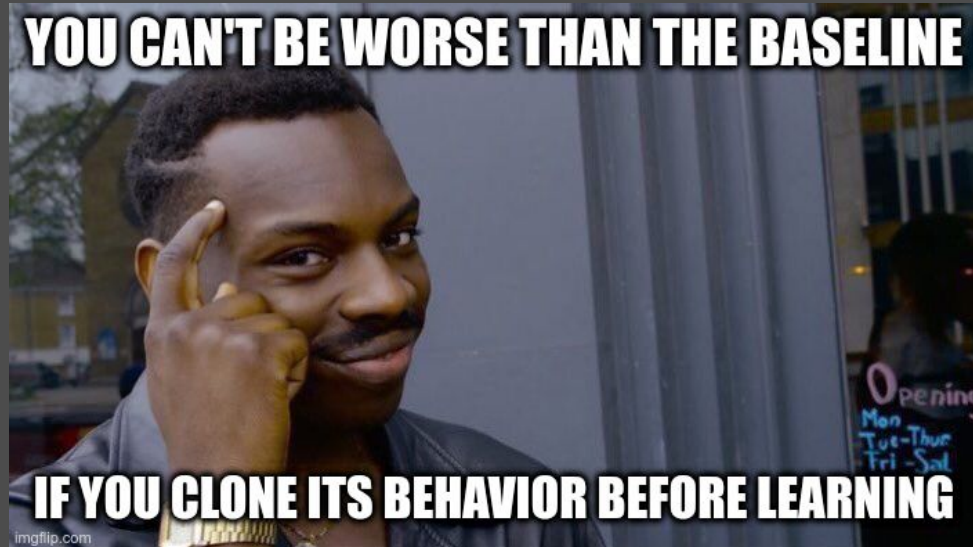
Profissão

Valor da Dívida

Renda média

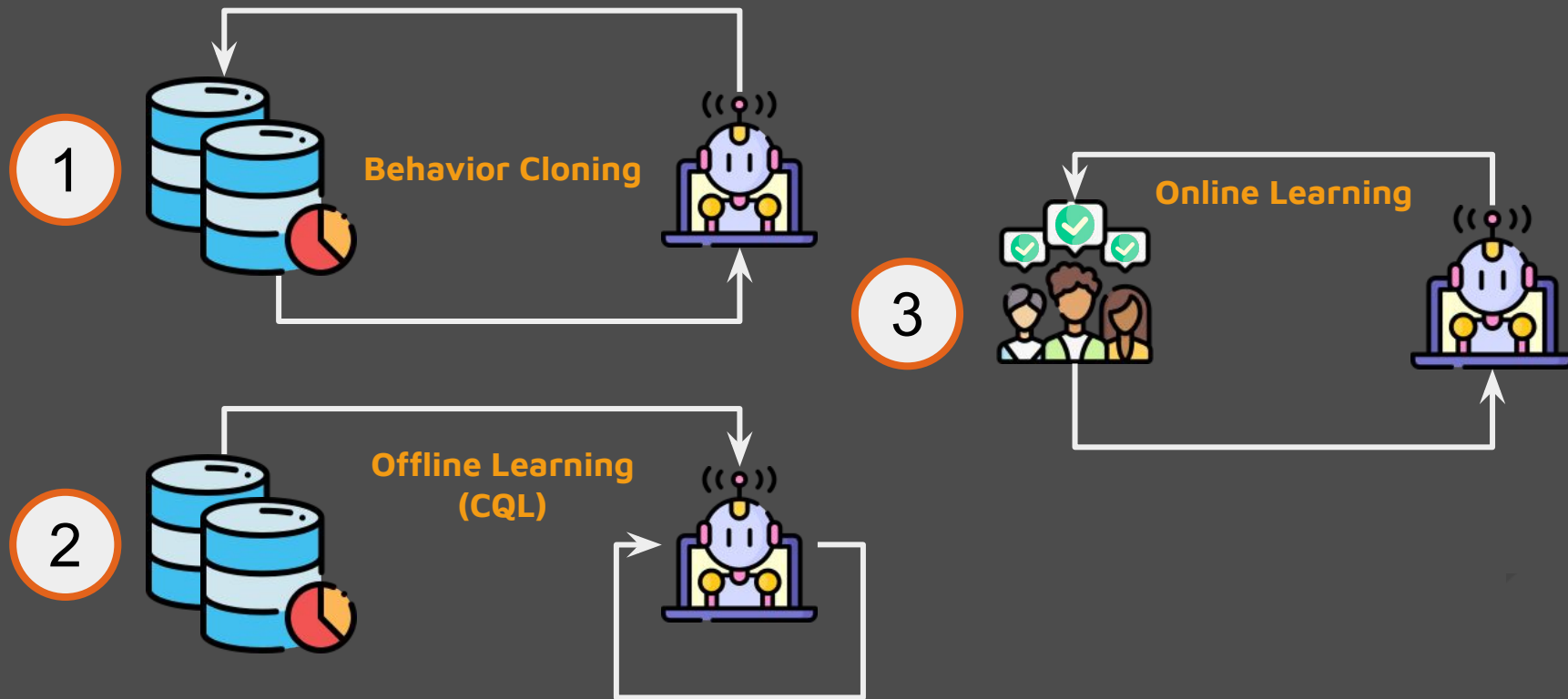
Quem Coletou os Dados?

- Modelos previamente desenvolvidos pela empresa
- Há inteligência nos dados / nas trajetórias



*quase verdade

Fluxo de Treino



Resultados

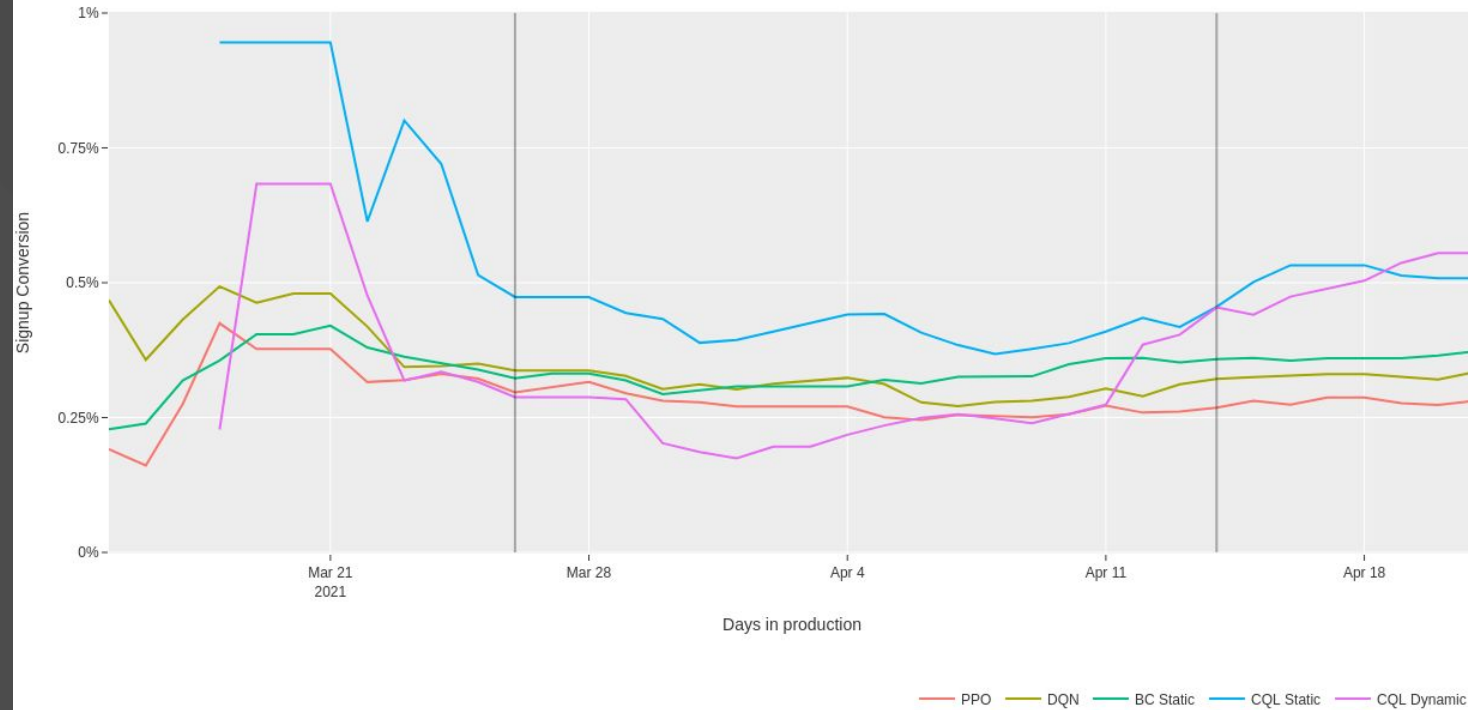
% Usuários cadastrados



Acionamentos via SMS



Signup Conversion



Conversão em Cadastros



É isso aí...

Próxima Aula: Tópicos Avançados