

Aprendizado, minimização do arrependimento e equilíbrio

(Learning, Regret Minimization, and Equilibria)

Victor Alberto Romero

Instituto de Matemática e Estatística
Universidade de São Paulo

Teoria dos Jogos e Algorítmica

- O problema principal é *tomar decisões em ambientes desconhecidos*.
- A aproximação usada neste capítulo é desenhar algoritmos adaptativos que podem aprender entendendo a dinâmica do sistema.
- A estratégia usada é a análise do arrependimento (*Regret*), ou seja, o cálculo da diferença entre a melhor utilidade que poderia ter sido e a utilidade real obtida. Temos três análises possíveis do arrependimento

- O problema principal é *tomar decisões em ambientes desconhecidos*.
- A aproximação usada neste capítulo é desenhar algoritmos adaptativos que podem aprender entendendo a dinâmica do sistema.
- A estratégia usada é a análise do arrependimento (*Regret*), ou seja, o cálculo da diferença entre a melhor utilidade que poderia ter sido e a utilidade real obtida. Temos três análises possíveis do arrependimento

- O problema principal é *tomar decisões em ambientes desconhecidos*.
- A aproximação usada neste capítulo é desenhar algoritmos adaptativos que podem aprender entendendo a dinâmica do sistema.
- A estratégia usada é a análise do arrependimento (*Regret*), ou seja, o cálculo da diferença entre a melhor utilidade que poderia ter sido e a utilidade real obtida. Temos três análises possíveis do arrependimento

Introdução

Arrependimento Externo (External Regret)

- Também conhecido como *Combining expert advice*
- Compara com a melhor opção individual que poderia ter sido (a melhor em retrospectiva)
- Dá uma metodologia geral para o desenvolvimento de algoritmos online.

Introdução

Arrependimento Externo (External Regret)

- Também conhecido como *Combining expert advice*
- Compara com a melhor opção individual que poderia ter sido (a melhor em retrospectiva)
- Dá uma metodologia geral para o desenvolvimento de algoritmos online.

Introdução

Arrependimento Externo (External Regret)

- Também conhecido como *Combining expert advice*
- Compara com a melhor opção individual que poderia ter sido (a melhor em retrospectiva)
- Dá uma metodologia geral para o desenvolvimento de algoritmos online.

- No capítulo, o autor faz um desenvolvimento muito interessante: começa como um algoritmo guloso simples, e sobre ele vai construindo até chegar num algoritmo no qual a perda é, no máximo, $O\left(\sqrt{T\log N}\right)$ a mais do que a melhor escolha, onde T é o número de passos dados até o momento.
- Ou seja, o máximo afastamento do ótimo por lance é:

$$O\left(\sqrt{\frac{\log N}{T}}\right)$$

- No capítulo, o autor faz um desenvolvimento muito interessante: começa como um algoritmo guloso simples, e sobre ele vai construindo até chegar num algoritmo no qual a perda é, no máximo, $O\left(\sqrt{T\log N}\right)$ a mais do que a melhor escolha, onde T é o número de passos dados até o momento.
- Ou seja, o máximo afastamento do ótimo por lance é:

$$O\left(\sqrt{\frac{\log N}{T}}\right)$$

Introdução

Arrependimento interno e de troca (*Internal Regret* and *Swap Regret*)

- Arrependimento interno (*Internal Regret*)
 - Considera a historia
 - Permite que **uma** ação passada seja substituída por outra
- Arrependimento de troca (*Swap Regret*)
 - Considera a historia
 - Permite qualquer mapeamento de $\{1, \dots, N\}$ para $\{1, \dots, N\}$ e pode ser aumentado num fator de N

Introdução

Arrependimento interno e de troca (*Internal Regret* and *Swap Regret*)

- Arrependimento interno (*Internal Regret*)
 - Considera a historia
 - Permite que **uma** ação passada seja substituída por outra
- Arrependimento de troca (*Swap Regret*)
 - Considera a historia
 - Permite qualquer mapeamento de $\{1, \dots, N\}$ para $\{1, \dots, N\}$ e pode ser aumentado num fator de N

Introdução

Modelo de informação parcial (*partial information model*)

- Também conhecido como *Multiarmed bandit* - *MAB*
- Só pode ver as perdas de sua escolha atual, não as perdas das outras escolhas (e por isso, não pode saber diretamente se é a melhor escolha ou não!!)

O autor deixa claro que os resultados obtidos nesse capítulo não são os melhores que se conhecem. Nas notas do capítulo tem uma lista de trabalhos que apresentam melhores algoritmos para este problema.

Introdução

Modelo de informação parcial (*partial information model*)

- Também conhecido como *Multiarmed bandit* - *MAB*
- Só pode ver as perdas de sua escolha atual, não as perdas das outras escolhas (e por isso, não pode saber diretamente se é a melhor escolha ou não!!)

O autor deixa claro que os resultados obtidos nesse capítulo não são os melhores que se conhecem. Nas notas do capítulo tem uma lista de trabalhos que apresentam melhores algoritmos para este problema.

Introdução

Modelo de informação parcial (*partial information model*)

- Também conhecido como *Multiarmed bandit* - *MAB*
- Só pode ver as perdas de sua escolha atual, não as perdas das outras escolhas (e por isso, não pode saber diretamente se é a melhor escolha ou não!!)

O autor deixa claro que os resultados obtidos nesse capítulo não são os melhores que se conhecem. Nas notas do capítulo tem uma lista de trabalhos que apresentam melhores algoritmos para este problema.

- Jogos de adversários **online** (processa só a entrada, sem memória de estados anteriores) como $X = \{1, \dots, N\}$ opções para cada jogador
- Em cada passo t o algoritmo H escolhe uma distribuição p^t sobre o conjunto de N ações.
- Depois disso o **adversário** escolhe um vetor de perdas $\ell^t \in [0, 1]^N$ (um valor de perda para cada escolha), onde $\ell_i^t \in [0, 1]$ é a perda se o jogador escolhe a opção i no momento t

- Jogos de adversários **online** (processa só a entrada, sem memória de estados anteriores) como $X = \{1, \dots, N\}$ opções para cada jogador
- Em cada passo t o algoritmo H escolhe uma distribuição p^t sobre o conjunto de N ações.
- Depois disso o adversário escolhe um vetor de perdas $\ell^t \in [0, 1]^N$ (um valor de perda para cada escolha), onde $\ell_i^t \in [0, 1]$ é a perda se o jogador escolhe a opção i no momento t

- Jogos de adversários **online** (processa só a entrada, sem memória de estados anteriores) como $X = \{1, \dots, N\}$ opções para cada jogador
- Em cada passo t o algoritmo H escolhe uma distribuição p^t sobre o conjunto de N ações.
- Depois disso o **adversário** escolhe um vetor de perdas $\ell^t \in [0, 1]^N$ (um valor de perda para cada escolha), onde $\ell_i^t \in [0, 1]$ é a perda se o jogador escolhe a opção i no momento t

Modelo e preliminares

Modelo com informação completa (*full information model*)

No modelo com informação completa (*full information model*) o algoritmo H recebe o vetor de perdas ℓ^t e tem um valor de perda:

$$\ell_H^t = \sum_{i=1}^N p_i^t \ell_i^t$$

que é a perda esperada quando o algoritmo escolhe a opção $i \in X$ com probabilidade p_i^t .

Modelo e preliminares

Modelo de informação parcial (*partial information model*)

No modelo de informação parcial (*partial information model*) o algoritmo H recebe $(\ell_{k^t}^t, k^t)$, onde k^t segue a distribuição p^t , e $\ell_H^t = \ell_{k^t}^t$ é sua perda. A perda da i -ésima ação durante os primeiros T passos é:

$$L_i^T = \sum_{t=1}^T \ell_i^t$$

e a perda de H é:

$$L_H^T = \sum_{t=1}^T \ell_H^t$$

Modelo e preliminares

Arrependimento externo (*external regret*)

A razão para calcular o arrependimento externo (*external regret*) é poder criar um algoritmo online que possa se aproximar ao melhor algoritmo de uma classe específica de algoritmos \mathcal{G} . Isto quer dizer, que queremos uma perda próxima a:

$$L_{\mathcal{G},\min}^T = \min_{g \in \mathcal{G}} L_g^T$$

Formalmente nós queremos minimizar o arrependimento externo:

$$R_{\mathcal{G}} = L_H^T - L_{\mathcal{G},\min}^T$$

Aqui \mathcal{G} é chamada a **classe de comparação**. Em geral, usa-se uma classe de comparação $\mathcal{G} = X$ (ou seja, todas as possíveis escolhas) e procura-se um algoritmo online com uma perda perto de

$$L_{\min}^T = \min_i L_i^T$$

sendo o arrependimento externo:

$$R = L_H^T - L_{\min}^T$$

- Normalmente o arrependimento externo usa uma classe de comparação fixa \mathcal{G} , mas poderia ter-se uma classe que dependa das ações escolhidas pelo algoritmo online. Neste caso podemos considerar regras de modificação que mudem as ações escolhidas pelo algoritmo online e produzam uma estratégia alternativa.
- Uma regra de modificação F recebe na entrada o histórico de ações e a ação atual escolhidas por H , e devolve uma ação (que **pode** ser diferente ou não). Chama-se F^t à função F no momento t

- Normalmente o arrependimento externo usa uma classe de comparação fixa \mathcal{G} , mas poderia ter-se uma classe que dependa das ações escolhidas pelo algoritmo online. Neste caso podemos considerar regras de modificação que mudem as ações escolhidas pelo algoritmo online e produzam uma estratégia alternativa.
- Uma **regra de modificação** F recebe na entrada o histórico de ações e a ação atual escolhidas por H , e devolve uma ação (que **pode** ser diferente ou não). Chama-se F^t à função F no momento t

Modelo e preliminares

Regras de modificação

Dada uma sequência de distribuições de probabilidade p^t usada pelo algoritmo H e uma regra de modificação F , pode-se definir uma nova sequência de distribuições de probabilidade $f^t = F^t(p^t)$, onde

$$f_i^t = \sum_{j:F^t(j)=i} p_j^t$$

A perda desta nova sequência é:

$$L_{H,F}^T = \sum_t \sum_i f_i^t \ell_i^t$$

no livro é: $L_{H,F} = \sum_t \sum_i f_i^t \ell_i^t$

Nota-se que no momento t , F muda a probabilidade de H associando a ação j a ação $F^t(j)$. Isto implica que a regra de modificação F gera uma distribuição diferente, que é função da distribuição de H , p^t .

Modelo e preliminares

Regras de modificação

Dada uma sequência de distribuições de probabilidade p^t usada pelo algoritmo H e uma regra de modificação F , pode-se definir uma nova sequência de distribuições de probabilidade $f^t = F^t(p^t)$, onde

$$f_i^t = \sum_{j:F^t(j)=i} p_j^t$$

A perda desta nova sequência é:

$$L_{H,F}^T = \sum_t \sum_i f_i^t \ell_i^t$$

no livro é: $L_{H,F} = \sum_t \sum_i f_i^t \ell_i^t$

Nota-se que no momento t , F muda a probabilidade de H associando a ação j a ação $F^t(j)$. Isto implica que a regra de modificação F gera uma distribuição diferente, que é função da distribuição de H , p^t .

Modelo e preliminares

Regras de modificação

Dada uma sequência de distribuições de probabilidade p^t usada pelo algoritmo H e uma regra de modificação F , pode-se definir uma nova sequência de distribuições de probabilidade $f^t = F^t(p^t)$, onde

$$f_i^t = \sum_{j:F^t(j)=i} p_j^t$$

A perda desta nova sequência é:

$$L_{H,F}^T = \sum_t \sum_i f_i^t \ell_i^t$$

no livro é: $L_{H,F} = \sum_t \sum_i f_i^t \ell_i^t$

Nota-se que no momento t , F muda a probabilidade de H associando a ação j a ação $F^t(j)$. Isto implica que a regra de modificação F gera uma distribuição diferente, que é função da distribuição de H , p^t .

Dada uma sequência de distribuições de probabilidade p^t usada pelo algoritmo H e uma regra de modificação F , pode-se definir uma nova sequência de distribuições de probabilidade $f^t = F^t(p^t)$, onde

$$f_i^t = \sum_{j:F^t(j)=i} p_j^t$$

A perda desta nova sequência é:

$$L_{H,F}^T = \sum_t \sum_i f_i^t \ell_i^t$$

no livro é: $L_{H,F} = \sum_t \sum_i f_i^t \ell_i^t$

Nota-se que no momento t , F muda a probabilidade de H associando a ação j a ação $F^t(j)$. Isto implica que a regra de modificação F gera uma distribuição diferente, que é função da distribuição de H , p^t .

- O artigo foca-se num conjunto finito \mathcal{F} de regras de modificação que não requer memória.
- Dada uma sequência de vetores de perda, o arrependimento de usar o algoritmo online H com relação às regras de modificação \mathcal{F} é:

$$R_{\mathcal{F}} = \max_{F \in \mathcal{F}} \{L_H^T - L_{H,F}^T\}$$

- O artigo foca-se num conjunto finito \mathcal{F} de regras de modificação que não requer memória.
- Dada uma sequência de vetores de perda, o arrependimento de usar o algoritmo online H com relação às regras de modificação \mathcal{F} é:

$$R_{\mathcal{F}} = \max_{F \in \mathcal{F}} \{L_H^T - L_{H,F}^T\}$$

- Pode-se ver que o arrependimento externo é equivalente a ter um conjunto de \mathcal{F}^{ex} de N regras de modificação F_i , onde cada F_i tem como saída a ação i . Para o **arrependimento interno** (*internal regret*), o conjunto \mathcal{F}^{in} consiste em $N(N-1)$ regras de modificação $F_{i,j}$, onde $F_{i,j}(i) = j$ e $F_{i,j}(i') = i'$, para $i' \neq i$. Então o arrependimento interno de H é:

$$\max_{F \in \mathcal{F}^{in}} \{L_H^T - L_{H,F}^T\} = \max_{i,j \in X} \left\{ \sum_{t=1}^T p_i^t (\ell_i^t - \ell_j^t) \right\}$$

Modelo e preliminares

Arrependimento de troca (*swap regret*)

- Um caso mais geral das regras de modificação sem memória é o arrependimento de troca (*swap regret*), definido pela classe \mathcal{F}^{sw} . \mathcal{F}^{sw} inclui as N^N funções $F : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$, onde a função F troca a ação atual de H , de i para $F(i)$. O arrependimento de troca para H é:

$$\max_{F \in \mathcal{F}^{sw}} \{L_H^T - L_{H,F}^T\} = \sum_{i=1}^N \max_{j \in X} \left\{ \sum_{t=1}^T p_i^t (\ell_i^t - \ell_j^t) \right\}$$

- É importante notar que:

$$\mathcal{F}^{ex} \subseteq \mathcal{F}^{sw}$$

$$\mathcal{F}^{in} \subseteq \mathcal{F}^{sw}$$

Modelo e preliminares

Arrependimento de troca (*swap regret*)

- Um caso mais geral das regras de modificação sem memória é o arrependimento de troca (*swap regret*), definido pela classe \mathcal{F}^{sw} . \mathcal{F}^{sw} inclui as N^N funções $F : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$, onde a função F troca a ação atual de H , de i para $F(i)$. O arrependimento de troca para H é:

$$\max_{F \in \mathcal{F}^{sw}} \{L_H^T - L_{H,F}^T\} = \sum_{i=1}^N \max_{j \in X} \left\{ \sum_{t=1}^T p_i^t (\ell_i^t - \ell_j^t) \right\}$$

- É importante notar que:

$$\mathcal{F}^{ex} \subseteq \mathcal{F}^{sw}$$

$$\mathcal{F}^{in} \subseteq \mathcal{F}^{sw}$$

Minimização do arrependimento externo

- Más notícias: Não é possível garantir um baixo arrependimento em relação à sequência ótima de decisões em retrospectiva. (Borodin e El-Yaniv, 1998; Sleator e Tarjan, 1985)
- Como alternativa, o autor propõe focar em classes de comparação mais restritas. Vamos usar a classe \mathcal{G}_{all} , que é o conjunto de funções que fazem o mapeamento de instantes de tempo $\{1, \dots, T\}$ para escolhas $X = \{1, \dots, N\}$

Minimização do arrependimento externo

- Más notícias: Não é possível garantir um baixo arrependimento em relação à sequência ótima de decisões em retrospectiva. (Borodin e El-Yaniv, 1998; Sleator e Tarjan, 1985)
- Como alternativa, o autor propõe focar em classes de comparação mais restritas. Vamos usar a classe \mathcal{G}_{all} , que é o conjunto de funções que fazem o mapeamento de instantes de tempo $\{1, \dots, T\}$ para escolhas $X = \{1, \dots, N\}$

Minimização do arrependimento externo

Teorema

Teorema

Para todo algoritmo online H existe uma sequência de vetores de perda T como a qual o arrependimento $R_{g_{all}}$ é pelo menos $T(1 - \frac{1}{N})$.

Demonstração.

Para cada instante t , a escolha i^t de mais baixa probabilidade p_i^t toma um valor de perda de 0 e todas as outras escolhas tomam o valor de 1. Como:

$$\min_i \{p_i^t\} \leq \frac{1}{N}$$

Então, a perda de H ao longo de T passos é pelo menos $T(1 - \frac{1}{N})$. Também pode-se definir um $g \in \mathcal{G}_{all}$, tal que $g(t) = i_t$ e tem um perda total de 0. □

Minimização do arrependimento externo

Teorema

Teorema

Para todo algoritmo online H existe uma sequência de vetores de perda T como a qual o arrependimento $R_{\mathcal{G}_{all}}$ é pelo menos $T(1 - \frac{1}{N})$.

Demonstração.

Para cada instante t , a escolha i^t de mais baixa probabilidade p_i^t toma um valor de perda de 0 e todas as outras escolhas tomam o valor de 1. Como:

$$\min_i \{p_i^t\} \leq \frac{1}{N}$$

Então, a perda de H ao longo de T passos é pelo menos $T(1 - \frac{1}{N})$. Também pode-se definir um $g \in \mathcal{G}_{all}$, tal que $g(t) = i_t$ e tem um perda total de 0. □

Minimização do arrependimento externo

- Para simplificar, vamos supôr que todas as perdas são ou 0 ou 1.
- Vamos chamar $L_i^t = \sum_{\tau=1}^t \ell_i^\tau$ a perda acumulada até o momento t da ação i
- Um algoritmo guloso vai escolher para cada momento t uma ação:

$$x^t = \arg \min_{i \in X} L_i^{t-1}$$

Minimização do arrependimento externo

- Para simplificar, vamos supôr que todas as perdas são ou 0 ou 1.
- Vamos chamar $L_i^t = \sum_{\tau=1}^t \ell_i^\tau$ a perda acumulada até o momento t da ação i
- Um algoritmo guloso vai escolher para cada momento t uma ação:

$$x^t = \arg \min_{i \in X} L_i^{t-1}$$

Minimização do arrependimento externo

- Para simplificar, vamos supôr que todas as perdas são ou 0 ou 1.
- Vamos chamar $L_i^t = \sum_{\tau=1}^t \ell_i^\tau$ a perda acumulada até o momento t da ação i
- Um algoritmo guloso vai escolher para cada momento t uma ação:

$$x^t = \arg \min_{i \in X} L_i^{t-1}$$

Algoritmo Guloso

Inicializa: $x^1 = 1$

Em cada momento t :

- $L_{\min}^{t-1} = \min_{i \in X} L_i^{t-1}$
- $S^{t-1} = \left\{ i : L_i^{t-1} = L_{\min}^{t-1} \right\}$
- $x^t = \min S^{t-1}$

Teorema

O algoritmo guloso tem uma perda acumulada

$L_G^T \leq N \cdot L_{\min}^T + (N - 1)$ para qualquer sequência de perdas

Minimização do arrependimento externo

Algoritmo Guloso

Algoritmo Guloso

Inicializa: $x^1 = 1$

Em cada momento t :

- $L_{\min}^{t-1} = \min_{i \in X} L_i^{t-1}$
- $S^{t-1} = \left\{ i : L_i^{t-1} = L_{\min}^{t-1} \right\}$
- $x^t = \min S^{t-1}$

Teorema

O algoritmo guloso tem uma perda acumulada

$$L_G^T \leq N \cdot L_{\min}^T + (N - 1) \text{ para qualquer sequência de perdas}$$

Demonstração.

Em cada momento t no qual o algoritmo perda 1 e o L_{\min}^t não mude, pelo menos uma ação é removida do S^t . Isto ocorre no máximo N vezes até que L_{\min}^t incremente em 1. Ou seja, o algoritmo tem uma perda de máximo N entre incrementos de L_{\min}^t :

$$L_G^t \leq N - |S^t| + N \cdot L_{\min}^t$$



Pode-se ver que o algoritmo é $O(N)$ vezes pior que o ótimo. Vamos ver que qualquer algoritmo determinístico tem pelo menos o mesmo fator.

Demonstração.

Em cada momento t no qual o algoritmo perda 1 e o L_{\min}^t não mude, pelo menos uma ação é removida do S^t . Isto ocorre no máximo N vezes até que L_{\min}^t incremente em 1. Ou seja, o algoritmo tem uma perda de máximo N entre incrementos de L_{\min}^t :

$$L_G^t \leq N - |S^t| + N \cdot L_{\min}^t$$



Pode-se ver que o algoritmo é $O(N)$ vezes pior que o ótimo. Vamos ver que qualquer algoritmo determinístico tem pelo menos o mesmo fator.

Teorema

Para qualquer algoritmo determinístico D existe uma sequência de perdas para a qual $L_D^T = T$ e $L_{min}^T = \lfloor \frac{T}{N} \rfloor$

Nota-se que este teorema implica que $L_D^T \geq N \cdot L_{min}^T + (T \bmod N)$ o que é um limite superior para o algoritmo guloso.

Teorema

Para qualquer algoritmo determinístico D existe uma sequência de perdas para a qual $L_D^T = T$ e $L_{min}^T = \lfloor \frac{T}{N} \rfloor$

Nota-se que este teorema implica que $L_D^T \geq N \cdot L_{min}^T + (T \bmod N)$ o que é um limite superior para o algoritmo guloso.

Demonstração.

Para a prova, vamos fixar um algoritmo online D e vamos dizer que x^t é sua saída no momento t . Vamos então a gerar uma sequência assim: No momento t vamos definir a perda para x^t de 1 e para qualquer outra ação de 0. Isto garante que o algoritmo D tem uma perda de 1 em cada passo, então $L_D^T = T$.

Como existem N ações diferentes, tem que existir uma ação que o algoritmo D tenha escolhido pelo menos $\lfloor \frac{T}{N} \rfloor$ vezes. Então, por construção, como só as ações escolhidas por D tem perdas, implica que $L_{\min}^T \leq \lfloor \frac{T}{N} \rfloor$ □

Minimização do arrependimento externo

Algoritmo Guloso Aleatório

Dá para pensar então que pode acontecer se o algoritmo guloso tem como saída uma distribuição de probabilidade. Vamos chamar esse algoritmo *Guloso Aleatório*

Algoritmo Guloso Aleatório

Inicializa: $p_i^1 = 1/N$ para $i \in X$

Em cada momento t :

- $L_{\min}^{t-1} = \min_{i \in X} L_i^{t-1}$
- $S^{t-1} = \left\{ i : L_i^{t-1} = L_{\min}^{t-1} \right\}$
- $p_i^t = \frac{1}{|S^{t-1}|}$ para $i \in S^{t-1}$ e $p_i^t = 0$ em outro contrario

Minimização do arrependimento externo

Algoritmo Guloso Aleatório

Dá para pensar então que pode acontecer se o algoritmo guloso tem como saída uma distribuição de probabilidade. Vamos chamar esse algoritmo *Guloso Aleatório*

Algoritmo Guloso Aleatório

Inicializa: $p_i^1 = 1/N$ para $i \in X$

Em cada momento t :

- $L_{\min}^{t-1} = \min_{i \in X} L_i^{t-1}$
- $S^{t-1} = \left\{ i : L_i^{t-1} = L_{\min}^{t-1} \right\}$
- $p_i^t = \frac{1}{|S^{t-1}|}$ para $i \in S^{t-1}$ e $p_i^t = 0$ em outro contrario

Minimização do arrependimento externo

Algoritmo Guloso Aleatório

Teorema

O algoritmo guloso aleatório tem, para qualquer sequência de perdas, uma perda de $L_{RG}^T \leq (\ln N) + (1 + \ln N)L_{min}^T$

Demonstração.

Como no caso anterior, vai se provar que entre incrementos sucessivos de L_{min}^T o algoritmo incorre numa perda pelo menos $1 + \ln N$. Seja t_j o momento onde L_{min}^t atinge por primeira vez uma perda j . Vamos ver que acontece entre os instantes t_j e t_{j+1} . Para qualquer instante de tempo t temos que $1 \leq |S^t| \leq N$.

(continua)



Minimização do arrependimento externo

Algoritmo Guloso Aleatório

Teorema

O algoritmo guloso aleatório tem, para qualquer sequência de perdas, uma perda de $L_{RG}^T \leq (\ln N) + (1 + \ln N)L_{min}^T$

Demonstração.

Como no caso anterior, vai se provar que entre incrementos sucessivos de L_{min}^T o algoritmo incorre numa perda pelo menos $1 + \ln N$. Seja t_j o momento onde L_{min}^T atinge por primeira vez uma perda j . Vamos ver que acontece entre os instantes t_j e t_{j+1} . Para qualquer instante de tempo t temos que $1 \leq |S^t| \leq N$.

(continua)



Demonstração.

(*continuação*) Agora, se num instante $t \in (t_j, t_{j+1}]$ o comprimento de S^t encolhe de n' para $n' - k$, a perda do algoritmo vai ser $\frac{k}{n'}$ (dado que toda ação tem peso de $\frac{1}{n'}$. Por último, pode-se ver que podemos fazer uma aproximação de $\frac{k}{n'}$ por $\frac{1}{n'} + \frac{1}{(n'-1)} + \dots + \frac{1}{(n'-k+1)}$. Isto gera uma perda para o intervalo $(t_j, t_{j+1}]$ de pelo menos:

$$\frac{1}{N} + \frac{1}{(N-1)} + \frac{1}{(N-2)} + \dots + \frac{1}{1} \leq 1 + \ln N$$

Por indução pode-se ver que

$$L_{RG}^t \leq \left(\frac{1}{N} + \frac{1}{(N-1)} + \frac{1}{(N-2)} + \dots + \frac{1}{(|S^t|+1)} \right) + (1 + \ln N)L_{\min}^t$$



Minimização do arrependimento externo

Algoritmo aleatório de peso maioritário

Uma ideia que surge naturalmente quando se pensa no algoritmo guloso aleatório é mudar as probabilidades para que não todas sejam iguais. Vamos associar para cada escolha i , como uma perda total L_i , um peso $w_i = (i - \eta)^{L_i}$, e uma probabilidade $p_i = \frac{w_i}{\sum_{j=1}^N w_j}$. Aqui η é uma constante pequena e ajustável. No capítulo usam $\eta = 0.01$

Algoritmo aleatório de peso maioritário

Inicializa:

- $w_i^1 = 1$
- $p_i^1 = \frac{1}{N}$ para $i \in X$

Em cada momento t :

- se $\ell_i^{t-1} = 1$ então $w_i^t = w_i^{t-1}(1 - \eta)$; se não então $w_i^t = w_i^{t-1}$
- $p_i^t = \frac{w_i^t}{W^t}$, onde $W^t = \sum_{i \in X} w_i^t$

Minimização do arrependimento externo

Algoritmo aleatório de peso maioritário

Uma ideia que surge naturalmente quando se pensa no algoritmo guloso aleatório é mudar as probabilidades para que não todas sejam iguais. Vamos associar para cada escolha i , como uma perda total L_i , um peso $w_i = (i - \eta)^{L_i}$, e uma probabilidade $p_i = \frac{w_i}{\sum_{j=1}^N w_j}$. Aqui η é uma constante pequena e ajustável. No capítulo usam $\eta = 0.01$

Algoritmo aleatório de peso maioritário

Inicializa:

- $w_i^1 = 1$
- $p_i^1 = \frac{1}{N}$ para $i \in X$

Em cada momento t :

- se $\ell_i^{t-1} = 1$ então $w_i^t = w_i^{t-1}(1 - \eta)$; se não então $w_i^t = w_i^{t-1}$
- $p_i^t = \frac{w_i^t}{W^t}$, onde $W^t = \sum_{i \in X} w_i^t$

Teorema

Para qualquer $\eta \leq \frac{1}{2}$, a perda do algoritmo aleatório de peso maioritário como uma sequência binária $\{0, 1\}$ de perdas é:

$$L_{RWM}^T \leq (1 + \eta)L_{min}^T + \frac{\ln N}{\eta}$$

Definindo $\eta = \min\left\{\sqrt{\frac{(\ln N)}{T}}, \frac{1}{2}\right\}$, temos que

$$L_{RWM}^T \leq L_{min}^T + 2\sqrt{T \ln N}$$

- Mas, quanto podemos baixar?

Teorema

Vamos assumir $T < \log_2 N$. Então existe uma geração aleatória de perdas que, para qualquer algoritmo R_1 , gera uma perda esperada $E[L_{R_1}^T] = \frac{T}{2}$, e ainda assim $L_{min}^T = 0$

- Mas, quanto podemos baixar?

Teorema

Vamos assumir $T < \log_2 N$. Então existe uma geração aleatória de perdas que, para qualquer algoritmo $R1$, gera uma perda esperada $E[L_{R1}^T] = \frac{T}{2}$, e ainda assim $L_{min}^T = 0$

Demonstração.

Vamos considerar a seguinte sequência de perdas. No instante $t = 1$, vamos ter um subconjunto de $N/2$ ações como uma perda de 0, e o resto vai ter uma perda de 1. No instante $t = 2$, vamos escolher um subconjunto aleatório de $N/4$ ações (dos que no $t = 1$ tiveram uma perda de 0), que vão ter uma perda de 0 no instante $t = 2$ e ao resto damos uma perda de 1. Se repetimos o processo, em cada passo qualquer algoritmo online vai ter uma perda esperada de $1/2$. Mas ainda assim, pode-se ver que pelo menos para $T < \log_2 N$, sempre existe uma ação como uma perda total de 0. \square

- O limite imposto para o arrependimento que foi apresentado depende unicamente do número de passos dados T , e não depende do resultado da melhor ação. Isto é conhecido como limite de ordem zero. Pode-se encontrar outros limites:
 - *Primeiro Ordem*: Depende da perda da melhor escolha
 - *Segundo Ordem*: Depende da soma dos quadrados das perdas.
- Ter um arrependimento externo que seja proporcional à variância da melhor escolha.
- Problemas como um número de ações possíveis muito grande. Por exemplo, um problema onde as escolhas possíveis sejam todos os caminhos simples entre dois nós (combinatório)
- Analisar a dinâmica dos algoritmos de minimização de arrependimento.