

# BÀI TOÁN GÁN NHÃN TỪ LOẠI

Lê Duy Thức ([thuc.leduy.int@gmail.com](mailto:thuc.leduy.int@gmail.com))

Phần 1 : Ví dụ HMM cho bài toán gán nhãn:

Cho câu -> **Tôi đi học**

Đây là câu trên gồm các từ đã dán nhãn:

**Tôi/S đi/V học/N**

Trong đó: S, V, N là các **nhãn** tương ứng

+ Ps: đầu câu

+ S : chủ từ (tôi, tao, mẹ, mày)

+ N: danh từ chung (cơm, cây, nước,...)

+ V động từ

(có cả động nhãn nhưng chỉ vd 3 nhãn cơ bản thôi cho dễ hiểu)

Phần 2:

- Cho tập dữ liệu dc huấn luyện (lấy vd 5 dòng trong 15k dòng)

*start/Ps Lại/R đi/V trên/E đường/N Nguyễn\_Ái\_Quốc/Np ./CH end/P\_s*

*start/Ps Anh/N ta/P tưởng/V bà/Nc mẹ/N “/CH bên/X ”/CH lắm/R ./CH end/P\_s*

*start/Ps Trái\_tìm/N ơi/T ./CH vồ/V ra/R và/C đừng/R bao\_giờ/P rung/V nữa/R ./CH end/P\_s*

*start/Ps N.Anh/X bé\_nhỏ/A yêu\_dấu/V đêm/N nay/P ở/V đâu/P .../X end/P\_s*

*start/Ps thể/T rồi/C không/R còn/R mơ/V được/V gì/P nữa/R ./CH cứ/R trượt/V theo/E những/L  
đường\_cong/N mềm\_mại/A .../X end/P\_s*

*start/Ps Mẹ/N mình/P chưa/R khóc/V bao\_giờ/P trong/E những/L lần/N tiễn/V con/N đi/V ./CH end/P\_s*

*start/Ps Thanh/Np bảo/V không/R dám/V kể/V tình\_hình/N thật/A của/E đơn\_vị/N cho/E gia\_đình/N vì/E  
nó/P sợ/V gia\_đình/N không/R yên\_tâm/A ./CH end/P\_s*

*start/Ps Càng/R nghĩ/V càng/R buồn/V ./CH nhất\_là/X trong/E những/L ngày/N này/P ./CH end/P\_s*

*start/Ps Thế/T đây/P ./CH mẹ/N mình/P chắc/V là/C không/R thế/P ./CH end/P\_s*

Mục tiêu của mô hình HMM cho bài toán này là từ **một tập dữ liệu đã được gán nhãn** cho sẵn ( cái đồng trên ), người dùng nhập vô 1 câu mới, **chương trình tự gán nhãn lại**

Ví dụ input: “*Tôi ăn cơm*” từ tập trên suy ra “*start/Ps Tôi/S ăn/V cơm/N*”

### ***Phần 3: mô tả chạy tay trên giấy***

**Input :** *Tôi ăn cơm*

Từ tập dữ liệu cho sẵn( đồng ở trên ), suy ra dc 2 bảng dưới đây, *làm sao suy ra – nói sau*

	S	V	N
<i>Start</i>	0.3	0.2	0.5
<i>N</i>	0.2	0.7	0.1
<i>V</i>	0.5	0.3	0.2
<i>S</i>	0.1	0.85	0.05

***Bảng Nhãn-nhãn***

	Tôi	ăn	cơm
<i>N</i>	0.4	0.2	0.5
<i>V</i>	0.1	0.6	0.05
<i>S</i>	0.5	0.1	0.05

***Bảng Từ-Nhãn***

***Bảng Từ-Nhãn:*** là xác xuất của từ đó là loại từ gì- nhãn gì.

VD: + nhìn vô thấy Tôi – S : 0.5 tức là “Tôi” thường xuất hiện trong câu với nhãn là S(chủ từ) chiếm 50%– như *Tôi ăn bánh* chẳng hạn. Còn “Tôi” – V : 0.1 thì tôi ở đây là V(động từ) --> rất hiếm có 10%,...

+ lưu ý 1 trong 1 cột vd Tôi xuất hiện là Tôi-N:0.4, Tôi-V:0.1, Tôi-S:0.5 . Tôi có thể là N V S nhưng cộng hết lại phải là 100% (0.4 + 0.1 + 0.5)

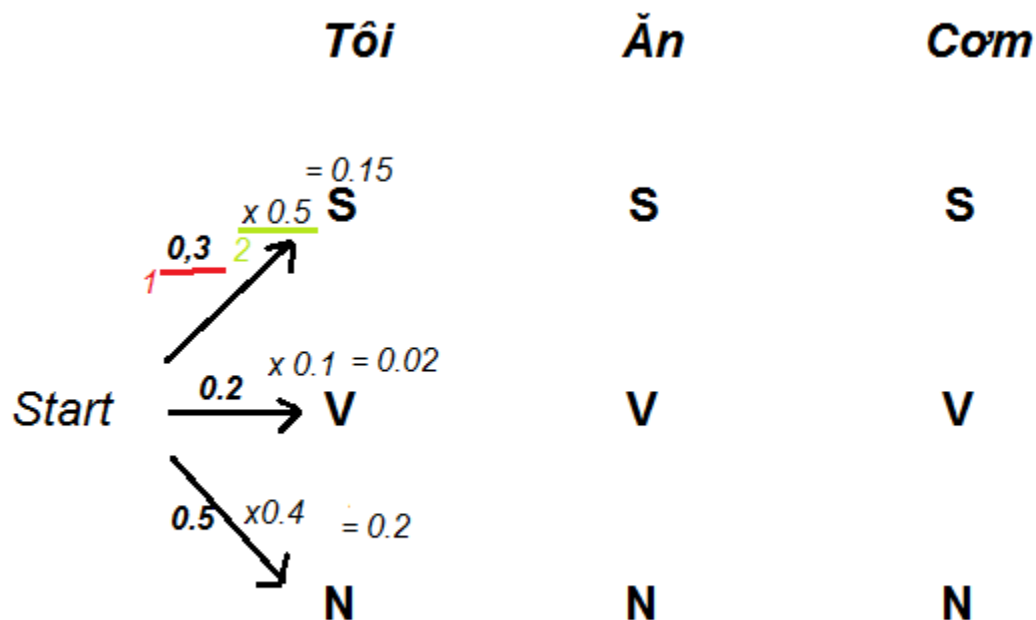
***Bảng Nhãn-Nhãn:***

VD + nhìn vô thấy S-V : 0.5 . Tức là: sau một S(chủ từ) thì khả năng xuất hiện V(danh từ) là 50% vd: “*Tao/S chạy/V*”, “*Mày/S ngủ/V*”. Hay V-N : 0.7 thì sau V(động từ) có tới 70% thường là danh từ(N) vd : “*đem/V cơm/N*”, “*Tao/S thích/V bánh/N*”

+ *Start* có nghĩa là bắt đầu câu: N-Start : 0.5 có nghĩa 1 câu thường bắt đầu với danh từ(N)

### Phần 3 Cách chạy

#### B1:



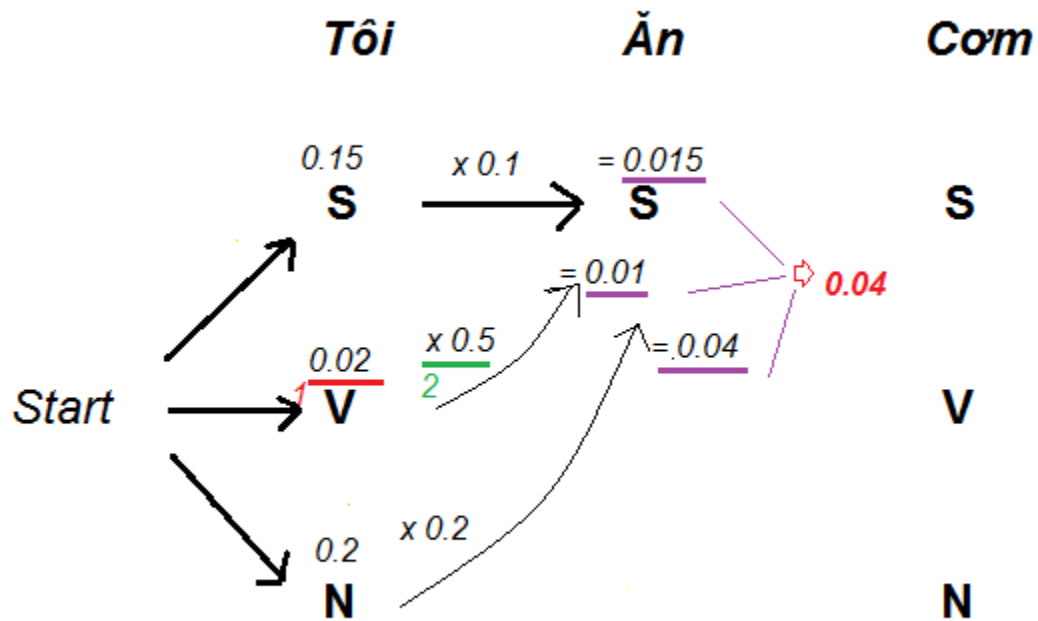
1: nhìn bảng Nhãn-nhãn: S-Start: 0.3

2 : Nhìn bảng Từ nhãn: Tôi – S: 0.5

⇒ Nhân ra vậy node S =  $0.3 \times 0.5 = 0.15$

Tương tự 2 nhãn V và N nhân ra 0.02 và 0.2

B2:



(Khác bước 1)

1 : Số vừa tính ở node V khi này 0.02

2 mũi tên từ V -> S thì nhìn bảng Nhân nhân: **Cột S hàng V** S-V:0.5

$$\Rightarrow 0.02 * 0.5 = 0.01$$

Tính tương tự được S->S=0.015

$$V->S=0.01$$

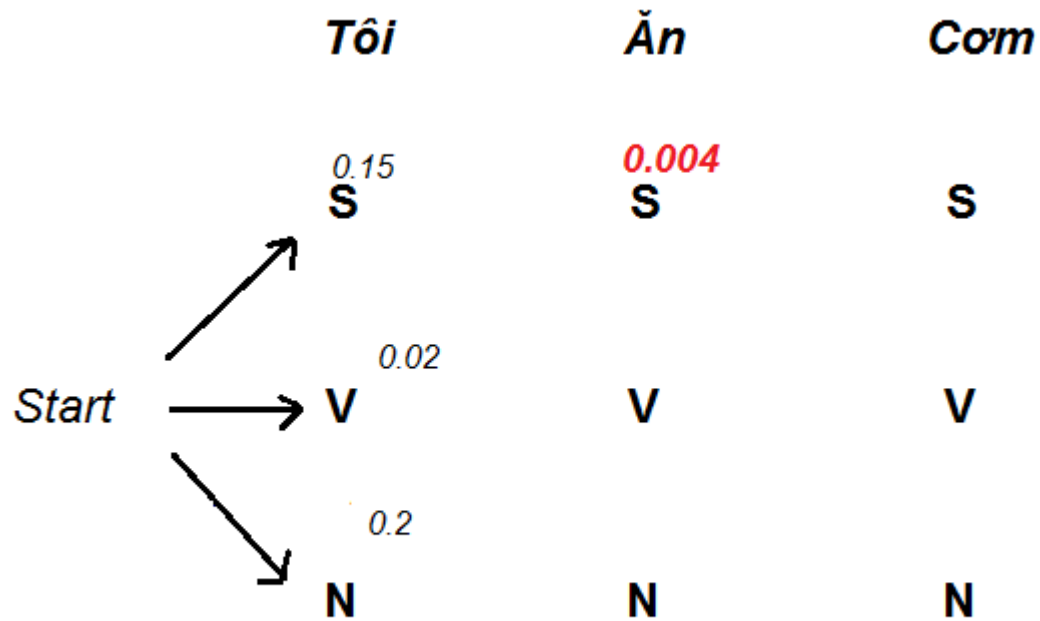
$$N \rightarrow S = 0,04$$

⇒ 0.04 lớn nhất. Lấy 0.04

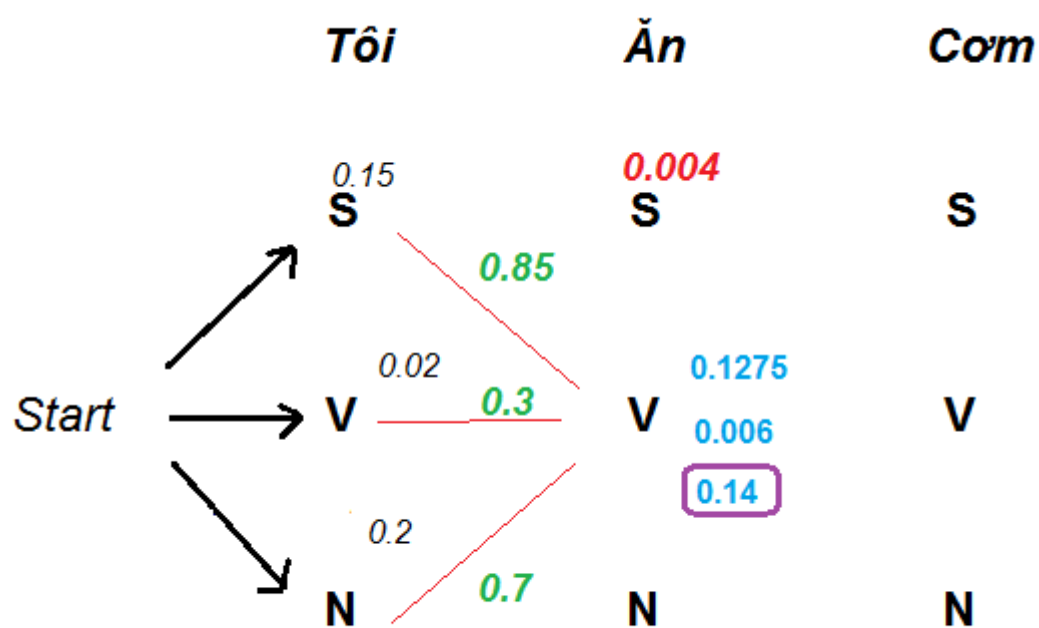
⇒ Ta đang tính cột “Ăn” ⇒ nhìn bảng Ăn-S : 0.1

$$\Rightarrow 0.04 * 0.1 = 0.004$$

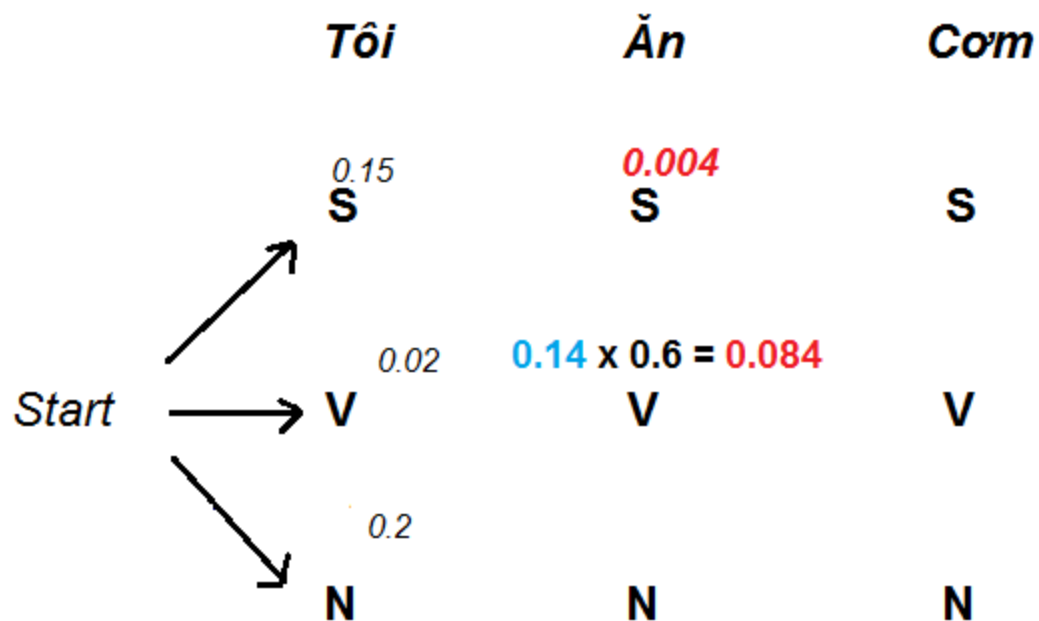
Nó sẽ ra như vậy



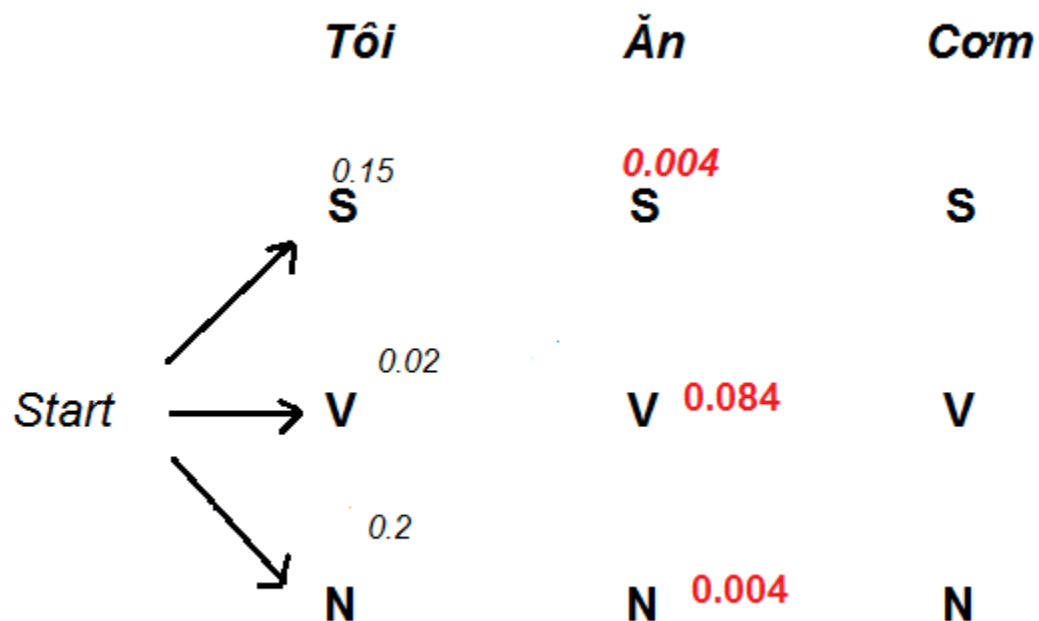
Tiếp tục y vậy làm với 2 node Ăn-V



0.14 lớn nhất lấy 0.14 ra và **nhân** với  $\text{Ăn-V} \cdot 0.6 = 0.084$

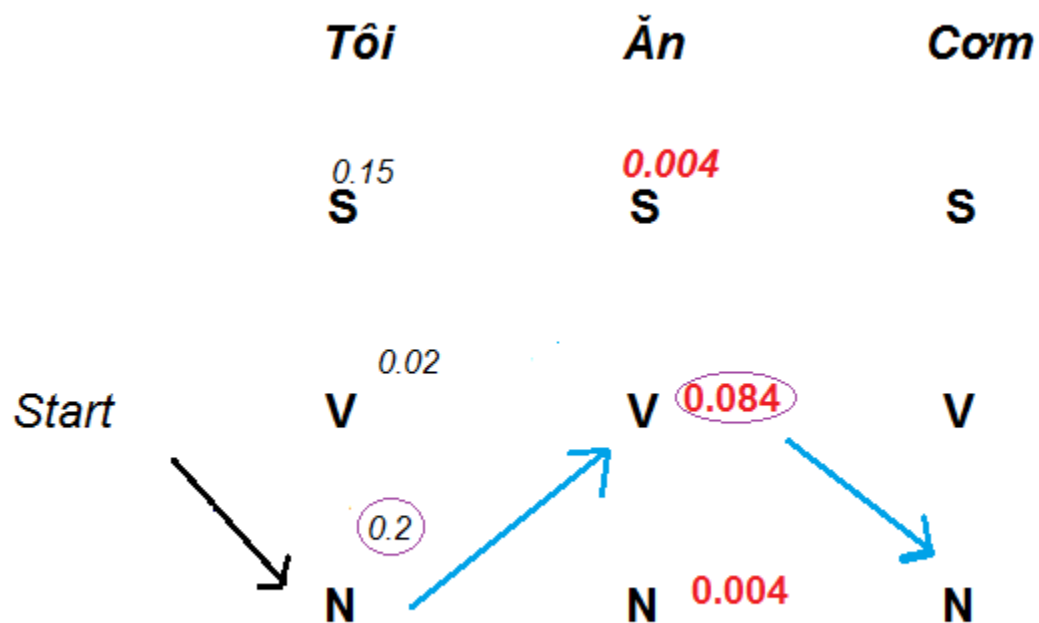


*Làm tiếp cho ra như vẬY*



*Thức mệt quá bạn làm tương tự cái cột “Cơm” luôn. Tương tự như trên*

*\*Sau đó cuối cùng vẽ đường đi từ Start qua các nhãn. Đi theo số lớn nhất như vậy nè*





- Rồi kết luận input “Tôi ăn cơm” đc gán nhãn Tôi/N ăn/V cơm/N dùng HMM như thế đấy
- Thức cho xác xuất ko chuẩn đáng lẽ kết quả ra **Tôi** phải là **S** chứ ko phải **N**