

Hidden Markov Models

Rafael E. Banchs (rbanchs@gps.tsc.upc.edu)

Session 4: The EM algorithm & HMMs



Centre de Tecnologies i Aplicacions del Llenguatge i la Parla

UNIVERSITAT POLITÈCNICA DE CATALUNYA

Exercise 6: A data-driven inference experience

Consider the twelve pairs of sentences provided, which correspond to translations between two hypothetical languages: A and B.

Try to infer correspondences among tokens from both languages, in order to:

- Build a bilingual dictionary of tokens
- Compute translations for the provided test data



Review on HMMs

Few days ago, we defined HMMs, and recalled that they can be formally specified by a set of five elements:

$$\{ S, O, \pi, P, Q \}$$

- where, S is the set of states: $S = \{s_1, s_2, s_3 \dots s_m\}$;
- O is the observation alphabet: $O = \{o_1, o_2, o_3 \dots o_h\}$;
- and, π P and Q are the HMM associated probabilities.



Review on HMMs

- State sequence: $\mathbf{X} = X_1, X_2, X_3 \dots X_N$
- State alphabet (state set): $\mathbf{S} = \{s_1, s_2, s_3 \dots s_m\}$
- Observation sequence: $\mathbf{Y} = Y_1, Y_2, Y_3 \dots Y_N$
- Observation alphabet (symbols): $\mathbf{O} = \{o_1, o_2, o_3 \dots o_h\}$
- Initial state probabilities: $\boldsymbol{\pi} = \{\pi_j\}$ with $j \in \mathbf{S}$
- State transition probabilities: $\mathbf{P} = \{p_{ij}\}$ with $i, j \in \mathbf{S}$
- Symbol emission probabilities: $\mathbf{Q} = \{q_{ijk}\}$ with $i, j \in \mathbf{S}, k \in \mathbf{O}$



Review on HMMs

Last day we studied the problem of estimating the most probable state sequence $\mathbf{X} = X_1, X_2, X_3 \dots X_N$ that explains an observation sequence $\mathbf{Y} = Y_1, Y_2, Y_3 \dots Y_N$ for a given HMM $\lambda = \{\pi, \mathbf{P}, \mathbf{Q}\}$

$$\underset{X}{\operatorname{Argmax}} P(X / Y, \lambda) = \underset{X}{\operatorname{Argmax}} P(Y / X, \lambda) P(X / \lambda)$$

This problem constitutes actually an optimization problem and it can be efficiently solved by using the *Viterbi algorithm*.



Estimating HMM parameters: Model Training

The problem we are going to deal with today is the following:

- Consider a HMM defined as $\lambda = \{\pi, \mathbf{P}, \mathbf{Q}\}$, what would be the best model parameters: π , \mathbf{P} and \mathbf{Q} that explain an observation data set?
- Such data set is frequently called the training data set and it may contain information on observations only, or also on states.
- There are actually two different formulations for this problem depending on the available training data.



Mathematical formulation of the problem

The two different formulations are the following:

- The “easy” problem: $\hat{\lambda} = \underset{\lambda}{\text{Argmax}} P(X, Y | \lambda)$

when training data includes information about both state sequences and their corresponding observation sequences.

- The “difficult” problem: $\hat{\lambda} = \underset{\lambda}{\text{Argmax}} P(Y | \lambda)$

when training data only includes information about observation sequences.

Estimating parameters from states and observations

In the “easy” problem formulation, parameters are estimated just by counting from training data... So, for a given state sequence $X_0, X_1 \dots X_N$ and its corresponding observation sequence $Y_1, Y_2 \dots Y_n$

- $\hat{\pi} = \delta(X_0)$, i.e. only the element at index position X_0 is equal to 1
- $\hat{p}_{ij} = \frac{\text{total number of transitions from } X_i \text{ to } X_j}{\text{total number of transitions from } X_i}$
- $\hat{q}_{ijO_k} = \frac{\text{total number of transitions from } X_i \text{ to } X_j \text{ emitting } O_k}{\text{total number of transitions from } X_i \text{ to } X_j}$



Exercise 2 revisited

Consider the following sequences of states and observations for the HMM of exercise 2:

States

1233312333122311123123123122311222333232312231233322312231223122232231231233312231223122222312222232222233112312311231
22222223123222333222223312312312222333123223231222223312312223331222233312323311232222333332231223123233123112232222231
1222223122231231233122231122332231231223122231122233123122312311233223323233123112232222322331223231223123122223312
23233122222231223122222312223112323312222312222223112311223223123312322331123231223122231231233123333232222231123123222
22312312322233231122231123123122231231112312232331122223122222233122222311222232223122222231222222333331232

Observations

BBAABBAABAABAABBBBBBABAABBBABAABAABAABBAAAABAABAAAABABABABAABAABAABBBBAAAABAABAABAABA
ABAAAABABAAAAAABAAAAAABAAABBABBABAAAAAABBAAAABAABAABAABBBBAABAABAAAAAABABABBABAAAA
AAABAABBAAAAABBAAAAAABABAAAABBBAAAABAABBAABBABAABBAAAAAABAAAAABAABAABAABBBABA
BAABAABABAAAABAABBBAAABBAABAABAABBBAABAABAABAABABABBAABBABAABAABAABAABAABA
AABBABAAAAAABAABBBAAAAAABBAABAABAABAABAABAABAABAABAABAABAABAABAABAABAABAABAABA
AAAAAABBBAABBAABA
ABABABAABBAAAABA



Exercise 2 revisited

- Since the initial state X_0 is state S_1 , then: $\hat{\pi} = [1 \ 0 \ 0]$

- Considering the state sequence, the total counts for state transitions are given by:

$$\mathbf{P}_{counts} = \begin{pmatrix} 24 & 96 & 0 \\ 0 & 167 & 130 \\ 95 & 35 & 53 \end{pmatrix}$$

- Normalizing \mathbf{P}_{counts} so that its row sums are equal to 1:

$$\hat{\mathbf{P}} = \begin{pmatrix} 0.2000 & 0.8000 & 0 \\ 0 & 0.5623 & 0.4377 \\ 0.5191 & 0.1913 & 0.2896 \end{pmatrix}$$



Exercise 2 revisited

- Considering both the state and observation sequences, the total state-transition counts for each observation symbol are:

$$\mathbf{Q}_{\mathbf{A} \text{ counts}} = \begin{pmatrix} 24 & 41 & 0 \\ 0 & 167 & 68 \\ 48 & 13 & 53 \end{pmatrix} \quad \mathbf{Q}_{\mathbf{B} \text{ counts}} = \begin{pmatrix} 0 & 55 & 0 \\ 0 & 0 & 62 \\ 47 & 22 & 0 \end{pmatrix}$$

- Normalizing each of these counts with respect to the total number of transitions for each pair of states (i.e. $\mathbf{P}_{\text{counts}}$)

$$\hat{\mathbf{Q}}_{\mathbf{A}} = \begin{pmatrix} 1.0000 & 0.4271 & \text{NaN} \\ \text{NaN} & 1.0000 & 0.5231 \\ 0.5053 & 0.3714 & 1.0000 \end{pmatrix} \quad \hat{\mathbf{Q}}_{\mathbf{B}} = \begin{pmatrix} 0 & 0.5729 & \text{NaN} \\ \text{NaN} & 0 & 0.4769 \\ 0.4947 & 0.6286 & 0 \end{pmatrix}$$

Exercise 2 revisited

So, the resulting estimates for $\lambda = \{\pi, \mathbf{P}, \mathbf{Q}\}$, are:

$$\hat{\pi} = [1 \ 0 \ 0]$$

$$\hat{\mathbf{P}} = \begin{pmatrix} 0.2000 & 0.8000 & 0 \\ 0 & 0.5623 & 0.4377 \\ 0.5191 & 0.1913 & 0.2896 \end{pmatrix}$$

$$\hat{\mathbf{Q}}_A = \begin{pmatrix} 1.0000 & 0.4271 & 0.5000 \\ 0.5000 & 1.0000 & 0.5231 \\ 0.5053 & 0.3714 & 1.0000 \end{pmatrix}$$

$$\hat{\mathbf{Q}}_B = \begin{pmatrix} 0 & 0.5729 & 0.5000 \\ 0.5000 & 0 & 0.4769 \\ 0.4947 & 0.6286 & 0 \end{pmatrix}$$

Which happen to match very well the actual HMM parameters...



Estimating parameters from observations only

In the “difficult” problem formulation, the state sequences that generate the observations are unknown...

- For this problem, we are interested in maximizing $P(Y / \lambda)$ over the model space λ , without caring about specific state sequences.
- This problem cannot be solved analytically, but suboptimal solutions can be computed in an iterative fashion.
- Local optimizations can be performed by using the ***Baum-Welch*** algorithm, that is a special case of ***Expectation Maximization***.



The Baum-Welch algorithm

- Let us define the vector \mathbf{g}_n which contains the probabilities of being at each state at step n for a given observation sequence \mathbf{Y} (It is a column vector of M elements)

$$g_n(i) = P(X_n = i \mid \mathbf{Y}, \lambda) \quad \text{for } i=1, 2 \dots M$$

*forward algorithm
probabilities*

backward algorithm probabilities

$$g_n(i) = \underbrace{f_n(i) * b_n(i)}_{P(X_n=i, Y/\lambda)} / \underbrace{f_n * b_n}_{P(Y/\lambda)}$$



The Baum-Welch algorithm

- Let us define the matrix \mathbf{W}_n which contains the probabilities of moving from actual state i to next state j at time step n for a given observation sequence \mathbf{Y} . (It is an $M \times M$ matrix)

$$w_n(i,j) = P(X_n=i, X_{n+1}=j \mid \mathbf{Y}, \lambda) \quad \text{for } i, j = 1, 2 \dots M$$

$$w_n(i) = \underbrace{f_n(i) * p_{ij} * q_{ijY_n} * b_{n+1}(j)}_{P(X_n=i, X_{n+1}=j, Y|\lambda)} / \underbrace{f_n * b_n}_{P(Y|\lambda)}$$



The Baum-Welch algorithm

- So, starting from an initial model $\lambda = \{\pi, \mathbf{P}, \mathbf{Q}\}$, an updated model can be computed as follows:

$$\hat{\pi} = \mathbf{g}_0$$

$$\hat{P}_{ij} = \frac{\text{expected transitions from } X_i \text{ to } X_j}{\text{expected transitions from } X_i} = \frac{\sum_n w_n(i,j)}{\sum_n g_n(i)}$$

transition “counts”
normalization

$$\hat{Q}_{ijO_k} = \frac{\text{expected transitions from } X_i \text{ to } X_j \text{ emitting } O_k}{\text{expected transitions from } X_i \text{ to } X_j} = \frac{\sum_{nO_k} w_n(i,j)}{\sum_n w_n(i,j)}$$

transition “counts”
for a given observation

The Baum-Welch algorithm

- By using the updated model $\hat{\lambda}$, new values for f_n , b_n , g_n and W_n can be computed, from which a new model update can be performed, and so on...
- This re-estimation process is guaranteed to improve $P(Y/\lambda)$ so that at each new re-estimation $P(Y/\lambda_{new}) \geq P(Y/\lambda_{previous})$
- However, the procedure gets easily trapped in suboptimal solutions which strongly depend on the initial model parameters.



Exercise 2 revisited again

Now let us consider the previous observation sequence only:

Observations

BBAABBAABAABAABBABAAAABBABAABAABAABBAAAAABAAAABABABABAABAABAABBBAABAABAABAABAA
 ABAAAABABAAAAAABAAAAAABAABBABBABAAAAAABBAAAABAABAABAABBBAABAABAABAABAABAABAABAABA
 AAABAABBAAAAAABBAAAAAABABAAAABBAAAAAABAABAABBAAAAAABAABAABAABAABAABAABAABAABAABAABA
 BAABAABABAAAABAABAABBBAABBAAABAABAABBBAABAABAABAABAABAABAABAABAABAABAABAABAABAABAABA
 AABBBABAAAAAABAABBBAAAAABBAAAAABAAAABABAAAAAABAABAABAABAABAABAABAABAABAABAABAABAABA
 AAAAAAABBBAABA
 ABABABAABBAAAAAABAAAAAABA

And we also need a first guess for the model λ :

- Let us consider the following ones:

$$\boldsymbol{\pi} = [1 \ 0 \ 0], \quad \text{and} \quad p_{ij} = 1/3, \quad q_{ijA} = q_{ijB} = 1/2 \quad \text{for} \quad i, j = 1, 2, 3$$



Exercise 2 revisited again

After a hundred iterations of the *Baum-Welch* algorithm we reach the following model values:

$$\hat{\boldsymbol{\pi}} = [1.0000 \ 0 \ 0]$$

$$\hat{\mathbf{P}} = \begin{pmatrix} 0.3990 & 0.3005 & 0.3005 \\ 0.3088 & 0.3456 & 0.3456 \\ 0.3088 & 0.3456 & 0.3456 \end{pmatrix}$$

$$\hat{\mathbf{Q}}_{\mathbf{A}} = \begin{pmatrix} 0.9385 & 0.8095 & 0.8095 \\ 0.7516 & 0.6042 & 0.6042 \\ 0.7516 & 0.6042 & 0.6042 \end{pmatrix}$$

$$\hat{\mathbf{Q}}_{\mathbf{B}} = \begin{pmatrix} 0.0615 & 0.1905 & 0.1905 \\ 0.2484 & 0.3958 & 0.3958 \\ 0.2484 & 0.3958 & 0.3958 \end{pmatrix}$$

Note that these values have nothing to do with the actual model parameters... ***However they do represent an equivalent HMM !!!***

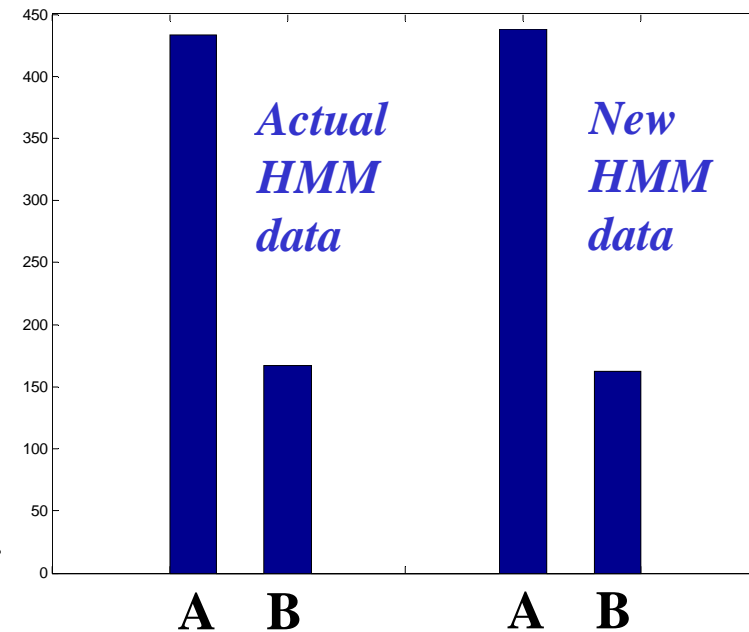


Exercise 2 revisited again

Consider now the state sequence that was used with the original HMM for generating the training observation sequence.

If we use the same state sequence to generate observations with the new estimated HMM parameters, we should reproduce (statistically) the original observation sequence.

These are the resulting histograms..



Comments on exercise 6

The data in exercise 6 are actually Chinese and English sentences.

A common procedure used in machine translation and other NLP applications is the so called word-to-word alignment, which consists of inferring word correspondences between two languages.

This procedure can be implemented by means of HMMs as it is described in Vogel *et.al.* (1999).

However, in practice HMMs are combined with other models in order to achieve better results.



HMM-based word-to-word alignment model

The model can be represented as:

**Probability of the Chinese sentence
ZH given the English sentence EN**

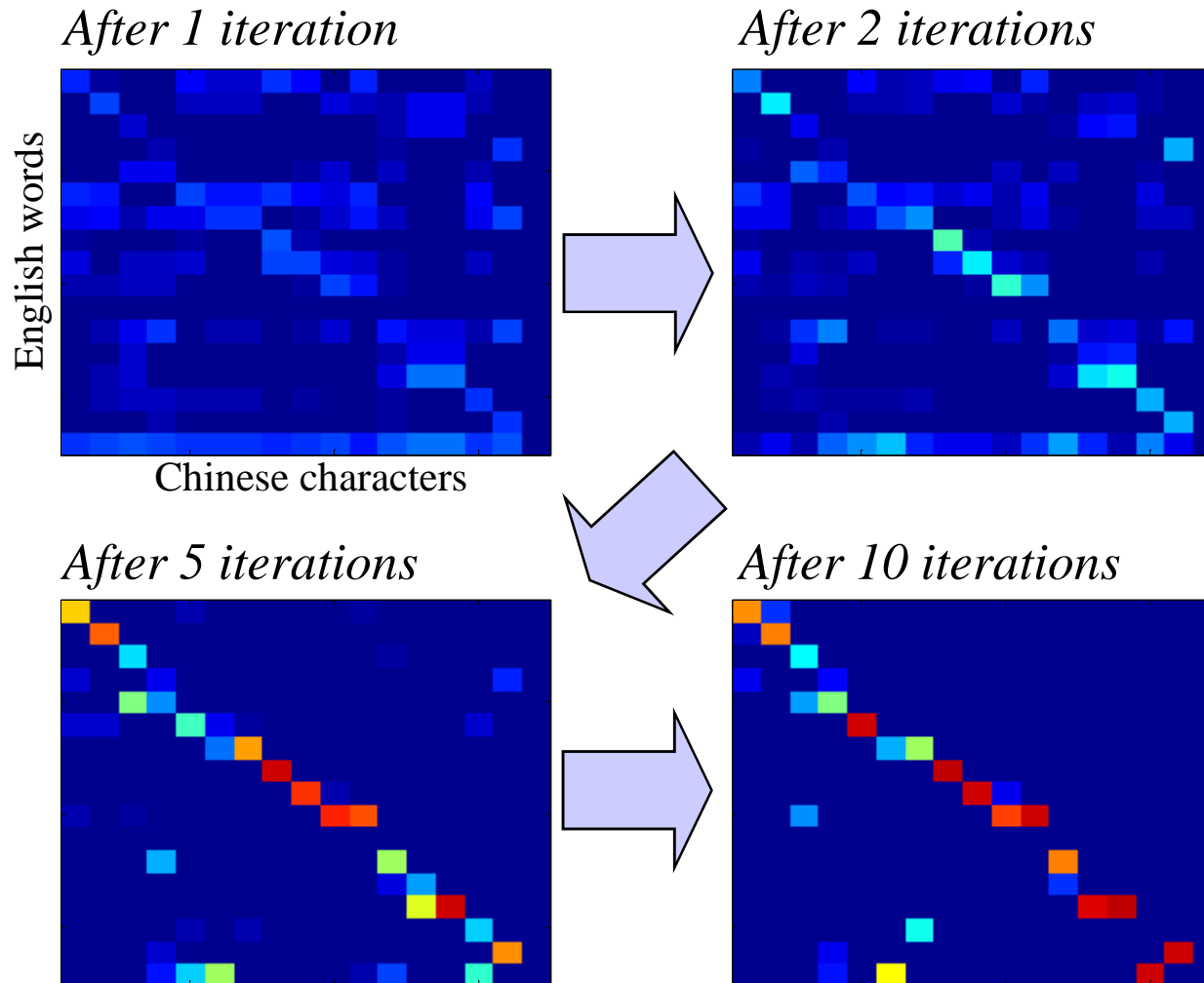
$$P(ZH | EN) = \sum_{a_j} \prod_j p(a_j | a_{j-1}) p(zh_j | en_{a_j})$$

Alignment probability:
Probability of the j^{th} Chinese word to be linked to the a_j^{th} English word, given that the $(j-1)^{\text{th}}$ Chinese word is linked to the a_{j-1}^{th} English word

Lexical Probability:
Probability of the j^{th} Chinese word given the a_j^{th} English word



Lexical probabilities estimated in the ZH to EN direction

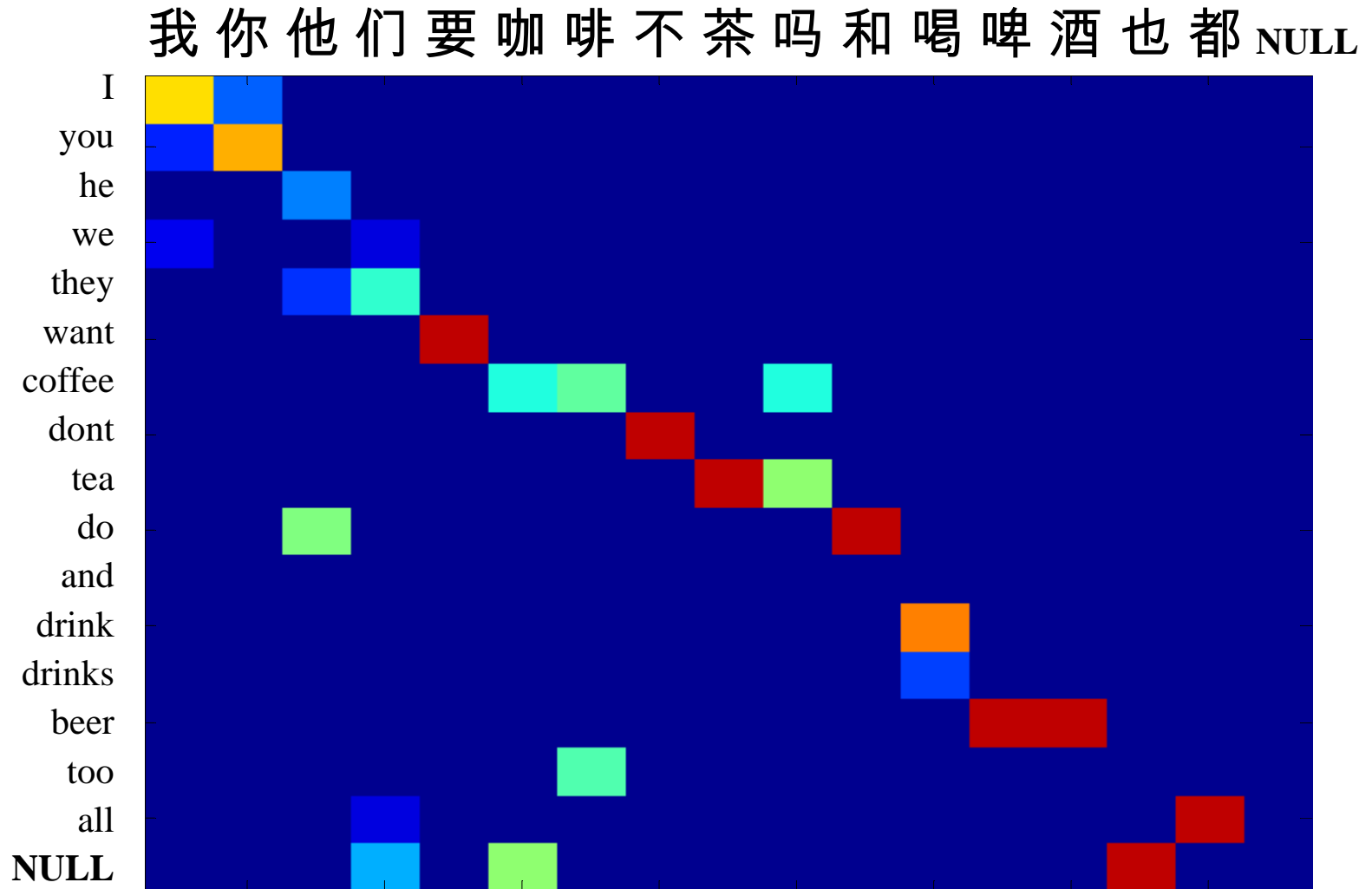


$P = 0$
 $P = 1$

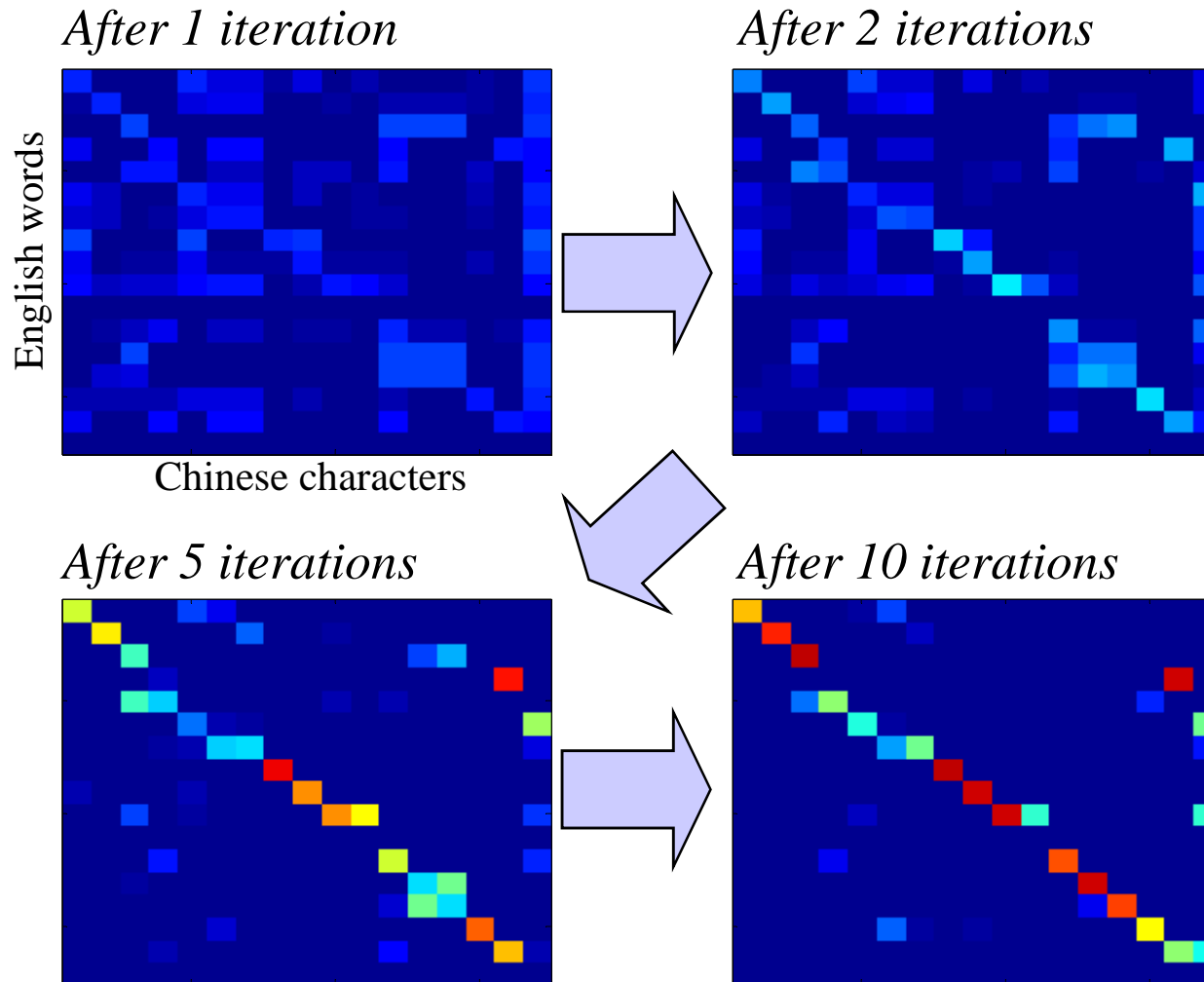


$P = 0$
 $P = 1$

Lexical probabilities after 100 iterations
 Computed in the ZH to EN direction



Lexical probabilities estimated in the EN to ZH direction

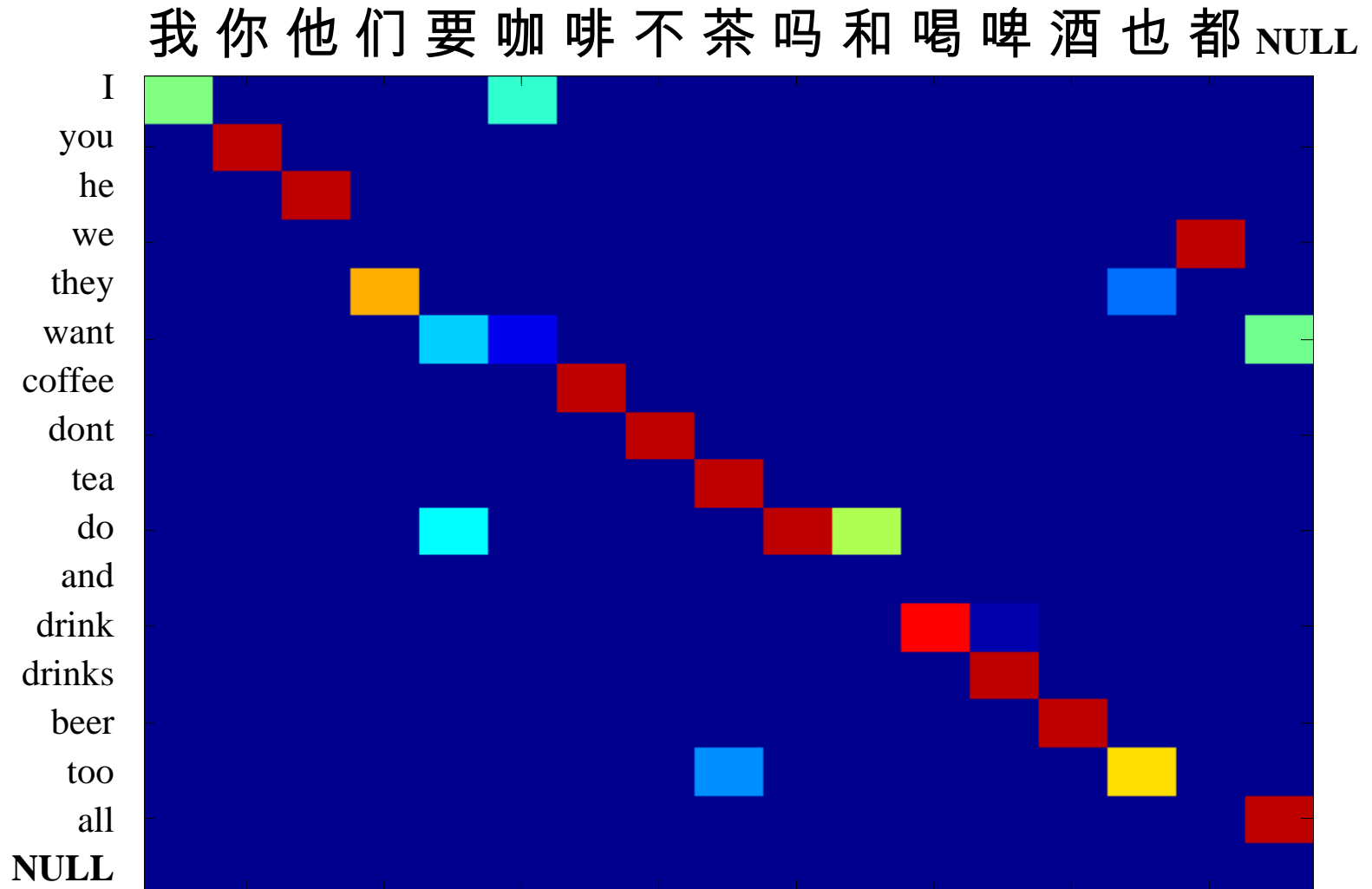


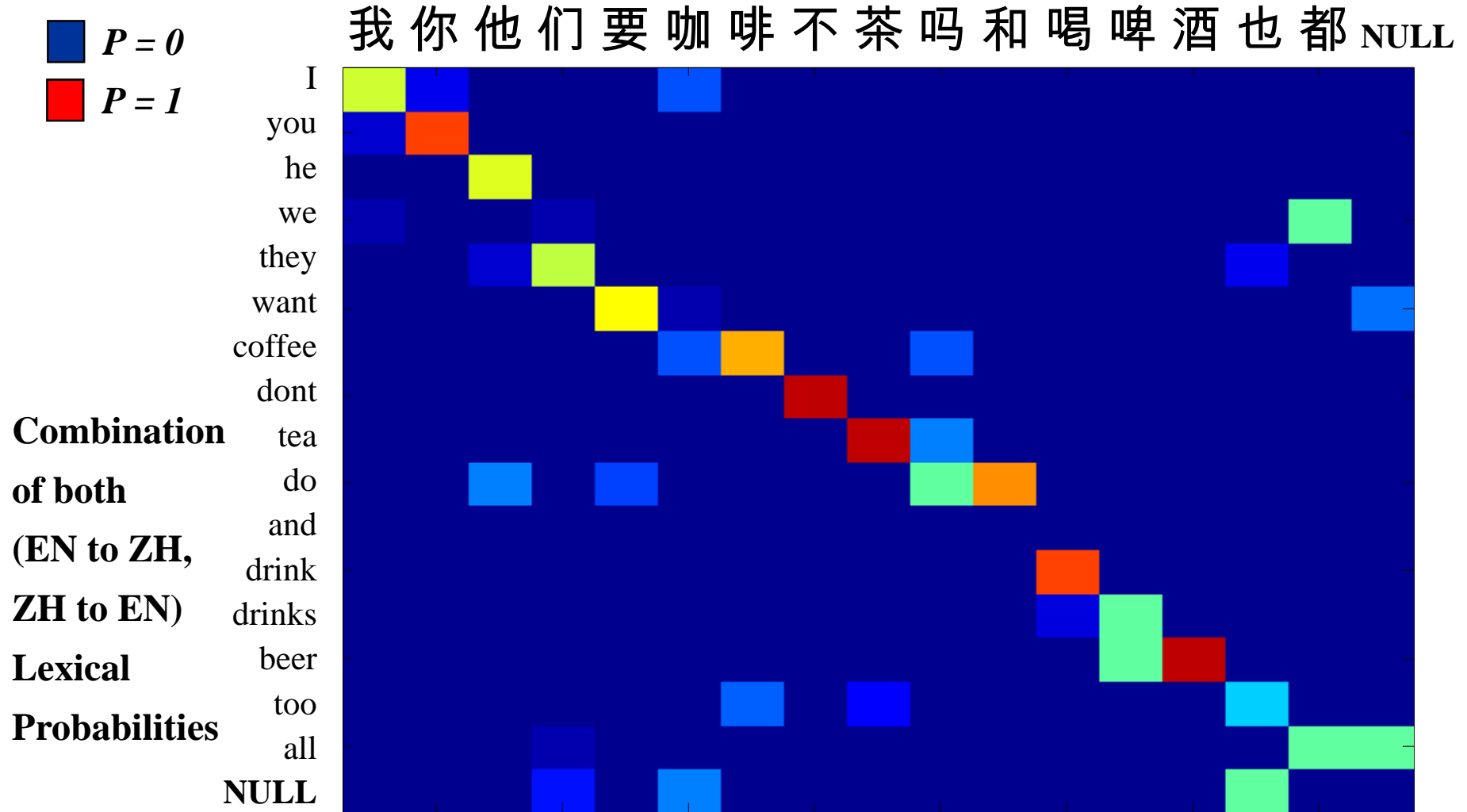
■ $P = 0$
■ $P = 1$



■ $P = 0$
 ■ $P = 1$

Lexical probabilities after 100 iterations
 Computed in the EN to ZH direction

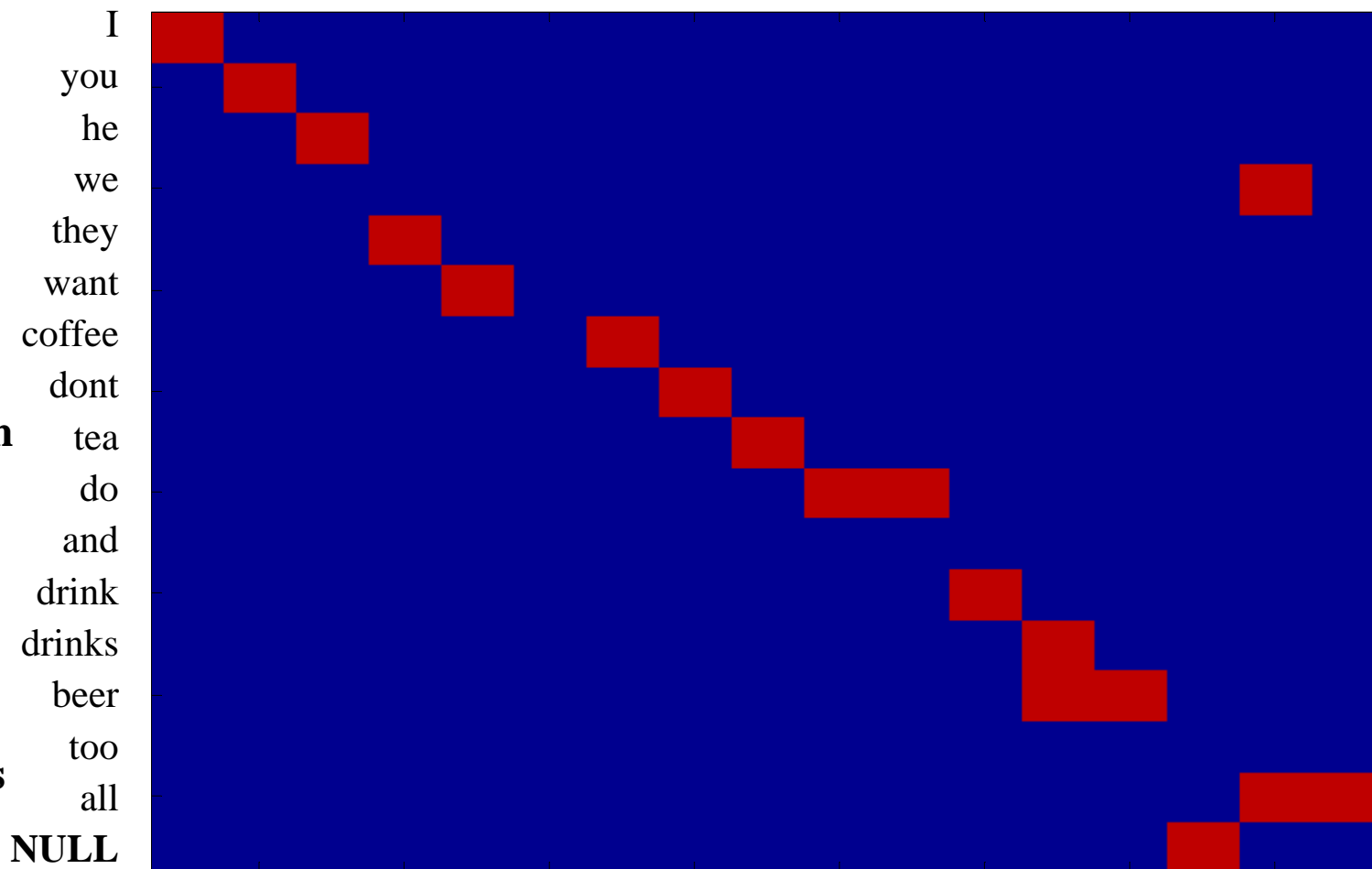




■ $P = 0$
■ $P = 1$

我 你 他们 要 咖啡 不 茶 吗 和 喝 啤酒 也 都 NULL

**Combination
of both
(EN to ZH,
ZH to EN)
Lexical
Probabilities**



Exercise 7

Consider again the HMM of exercise 2...

- 1.- Use MATLAB or any other language to implement the *Baum-Welch* algorithm.
- 2.- Generate a random sequence of 500 states and observations by using the HMM parameters of exercise 2
- 3.- Estimate some different models from the observation sequence generated in 2.- by considering different initial models λ_0 .



Important information about Exercise 7...

- If you already did the programming related to exercises 3: the forward and backward procedures...
- and you also did the programming related to exercise 5: Viterbi algorithm...

Programming exercise 7 should be very easy, and then...

YOU HAVE ALMOST READY THE FINAL PROJECT !!!

You only need an application. If interested we can talk now...



Important information about Exercise 7...

Some possible applications are:

- communication channel equalization,
- convolutional code decoding,
- part of speech tagging,
- automatic speech recognition,
- automatic word-to-word alignment,
- and any other idea you may have...



Hidden Markov Models

Rafael E. Banchs (rbanchs@gps.tsc.upc.edu)

Session 4: The EM algorithm & HMMs



Centre de Tecnologies i Aplicacions del Llenguatge i la Parla

UNIVERSITAT POLITÈCNICA DE CATALUNYA