

A linear law for the speech sonority *and a stochastic model suggested by this law*

Pierre Collet, Didier Demolin, Antonio Galves
and Jesús García

*Some recent results of the **Tycho Brahe Project***

A linear law for the speech sonority

- A rough measure of **sonority** was introduced in Galves *et al.* (2002) as a tool to discriminate between rhythmic classes of languages.
- An empirical analysis of a multi-lingual corpus puts in evidence a **linear relationship** between the **mean sonority** and the **mean increment of the sonority** across sentences of the sample.
- There is a simple way to explain this linearity. This is the **goal** of this presentation.

A linear law for the speech sonority

- A rough measure of **sonority** was introduced in Galves *et al.* (2002) as a tool to discriminate between rhythmic classes of languages.
- An empirical analysis of a multi-lingual corpus puts in evidence a **linear relationship** between the **mean sonority** and the **mean increment of the sonority** across sentences of the sample.
- There is a simple way to explain this linearity. This is the **goal** of this presentation.

A linear law for the speech sonority

- A rough measure of **sonority** was introduced in Galves *et al.* (2002) as a tool to discriminate between rhythmic classes of languages.
- An empirical analysis of a multi-lingual corpus puts in evidence a **linear relationship** between the **mean sonority** and the **mean increment of the sonority** across sentences of the sample.
- There is a simple way to explain this linearity. This is the **goal** of this presentation.

Acknowledgments

Special thanks to Emmanuel Dupoux, Sharon Pepperkamp and Frank Ramus for making the data from the *Laboratoire de Sciences Cognitives et Psycholinguistique (EHESS/CNRS)* used in this paper available to us.

The rhythmic classes conjecture

It has been conjectured in the linguistic literature that languages are divided in two or maybe three **rhythmic classes** (Lloyd 1940, Pike 1945, Abercrombie 1967, ...).

- **Morse code** or **stress-timed** languages: English, European Portuguese, Dutch, Finish, Polish, ...
- **Machine-gun** or **Syllable-timed** languages: Brazilian Portuguese, Catalan, French, Italian, Spanish,...
- **Mora-timed** languages: Japanese, Fidji, ...

The rhythmic classes conjecture

It has been conjectured in the linguistic literature that languages are divided in two or maybe three **rhythmic classes** (Lloyd 1940, Pike 1945, Abercrombie 1967, ...).

- **Morse code** or **stress-timed** languages: English, European Portuguese, Dutch, Finish, Polish, ...
- **Machine-gun** or **Syllable-timed** languages: Brazilian Portuguese, Catalan, French, Italian, Spanish,...
- **Mora-timed** languages: Japanese, Fidji, ...

The rhythmic classes conjecture

It has been conjectured in the linguistic literature that languages are divided in two or maybe three **rhythmic classes** (Lloyd 1940, Pike 1945, Abercrombie 1967, ...).

- **Morse code** or **stress-timed** languages: English, European Portuguese, Dutch, Finish, Polish, ...
- **Machine-gun** or **Syllable-timed** languages: Brazilian Portuguese, Catalan, French, Italian, Spanish,...
- **Mora-timed** languages: Japanese, Fidji, ...

But for nearly half a century

- There was **no real definition of the rhythmic properties** characterizing a class.
- And no **empirical correlates** of the rhythmic properties in the **speech signal** was known.

But for nearly half a century

- There was **no real definition of the rhythmic properties** characterizing a class.
- And no **empirical correlates** of the rhythmic properties in the **speech signal** was known.

Correlates of linguistic rhythm in the speech signal

Ramus, Nespors and Mehler (1999) gave for the first time evidence that simple statistics of the speech signal could discriminate between different **rhythmic classes**.

- RNM analyzed the acoustic signal of 20 sentences produced by 4 speakers of each of the following languages: English, Polish, Dutch, Catalan, Spanish, Italian, French and Japanese.
- The chosen sentences were segmented into vocalic and consonantal intervals.
- For each language, the empirical standard deviation of the durations of the consonantal intervals (ΔC) and the proportion of time spent in vocalic intervals ($\%V$) were computed.

RNM's approach

- RNM analyzed the acoustic signal of 20 sentences produced by 4 speakers of each of the following languages: English, Polish, Dutch, Catalan, Spanish, Italian, French and Japanese.
- The chosen sentences were segmented into **vocalic** and **consonantal** intervals.
- For each language, the empirical **standard deviation of the durations of the consonantal intervals (ΔC)** and the **proportion of time spent in vocalic intervals ($\%V$)** were computed.

- RNM analyzed the acoustic signal of 20 sentences produced by 4 speakers of each of the following languages: English, Polish, Dutch, Catalan, Spanish, Italian, French and Japanese.
- The chosen sentences were segmented into **vocalic** and **consonantal** intervals.
- For each language, the empirical **standard deviation of the durations of the consonantal intervals** (ΔC) and the **proportion of time spent in vocalic intervals** ($\%V$) were computed.

Intuition behind the choice of these statistical parameters

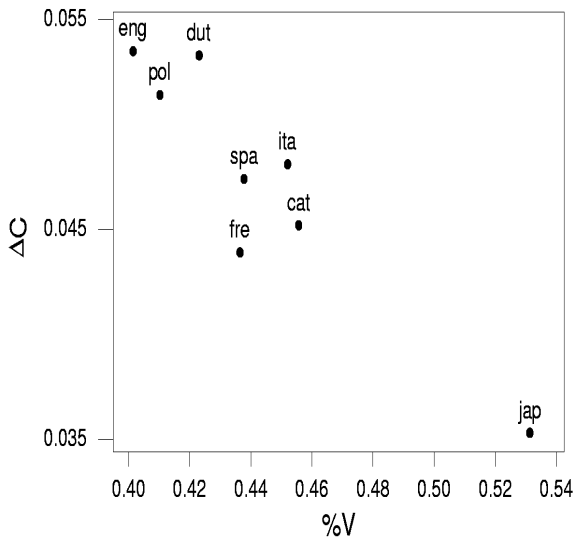
- Languages conjectured to be **syllable timed** (like *Italian*) spend a greater proportion of time in **vocalic** intervals than languages conjectured to be **stress-timed** (like *English* or *Dutch*). This justifies the choice of $%V$.
- Languages conjectured to be **stress-timed** display a greater variety and complexity of **consonantal** intervals than languages conjectured to be **syllable-timed**. This justifies the choice of δC .

It turns out that this was a good choice ...

Intuition behind the choice of these statistical parameters

- Languages conjectured to be **syllable timed** (like *Italian*) spend a greater proportion of time in **vocalic** intervals than languages conjectured to be **stress-timed** (like *English* or *Dutch*). This justifies the choice of $\%V$.
- Languages conjectured to be **stress-timed** display a greater variety and complexity of **consonantal** intervals than languages conjectured to be **syllable-timed**. This justifies the choice of δC .

It turns out that this was a good choice ...



Problems with RNM's approach

It is based on a hand label-ling of the speech signal.

- This is a **time-consuming** task.
- Moreover this hand label-ling is often based on decisions which are difficult to **reproduce in a homogeneous way**

This makes it difficult to reproduce RNM's approach on large samples.

Problems with RNM's approach

It is based on a hand label-ling of the speech signal.

- This is a **time-consuming** task.
- Moreover this hand label-ling is often based on decisions which are difficult to **reproduce in a homogeneous way**

This makes it difficult to reproduce RNM's approach on large samples.

A new approach to the problem

- Newborn babies are able to **discriminate rhythmic classes** with a signal filtered at 400Hz (Mehler et al. 1996).
- At this level, it is hard to distinguish nasals from vowels and glides from consonants.
- This strongly suggests that the discrimination of rhythmic classes by babies relies not on **fine-grained** distinctions between **vowels** and **consonants**, but on a **coarse-grained** perception of **sonority** in opposition to **obstruency**.

A new approach to the problem

- Newborn babies are able to **discriminate rhythmic classes** with a signal filtered at 400Hz (Mehler et al. 1996).
- At this level, it is hard to distinguish nasals from vowels and glides from consonants.
- This strongly suggests that the discrimination of rhythmic classes by babies relies not on **fine-grained** distinctions between **vowels** and **consonants**, but on a **coarse-grained** perception of **sonority** in opposition to **obstruency**.

A new approach to the problem

- Newborn babies are able to **discriminate rhythmic classes** with a signal filtered at 400Hz (Mehler et al. 1996).
- At this level, it is hard to distinguish nasals from vowels and glides from consonants.
- This strongly suggests that the discrimination of rhythmic classes by babies relies not on **fine-grained** distinctions between **vowels** and **consonants**, but on a **coarse-grained** perception of **sonority** in opposition to **obstruency**.

Sonority as a basis for rhythmic class discrimination

- Galves, Garcia, Duarte and Galves (2002) suggests that it is possible to discriminate rhythmic classes of language, using a rough measure of the speech **sonority**.
- This measure is defined **directly from the spectrogram** of the signal, with no need of previous hand label-ling of the data
- Applied to the same linguistics samples considered in RNM, it produces the three conjectured clusters.

Sonority as a basis for rhythmic class discrimination

- Galves, Garcia, Duarte and Galves (2002) suggests that it is possible to discriminate rhythmic classes of language, using a rough measure of the speech **sonority**.
- This measure is defined **directly from the spectrogram** of the signal, with no need of previous hand label-ling of the data
- Applied to the same linguistics samples considered in RNM, it produces the three conjectured clusters.

Sonority as a basis for rhythmic class discrimination

- Galves, Garcia, Duarte and Galves (2002) suggests that it is possible to discriminate rhythmic classes of language, using a rough measure of the speech **sonority**.
- This measure is defined **directly from the spectrogram** of the signal, with no need of previous hand label-ling of the data
- Applied to the same linguistics samples considered in RNM, it produces the three conjectured clusters.

Defining the speech sonority

- We define a function which maps local windows of the acoustic signal on the interval $[0, 1]$.
- This function is close to 1 for spans displaying **regular** patterns, characteristic of **sonorant** portions of the signal.
- In contrast, regions in which the acoustic signal present a **chaotic** behavior, for instance regions corresponding to **stop consonants**, will correspond to intervals in which S_t will assume values close to 0, with important variations

Defining the speech sonority

- We define a function which maps local windows of the acoustic signal on the interval $[0, 1]$.
- This function is close to 1 for spans displaying **regular** patterns, characteristic of **sonorant** portions of the signal.
- In contrast, regions in which the acoustic signal present a **chaotic** behavior, for instance regions corresponding to **stop consonants**, will correspond to intervals in which S_t will assume values close to 0, with important variations

Defining the speech sonority

- We define a function which maps local windows of the acoustic signal on the interval $[0, 1]$.
- This function is close to 1 for spans displaying **regular** patterns, characteristic of **sonorant** portions of the signal.
- In contrast, regions in which the acoustic signal present a **chaotic** behavior, for instance regions corresponding to **stop consonants**, will correspond to intervals in which S_t will assume values close to 0, with important variations

- Denote by $c_t(f)$ the power spectral density at time t and frequency f .
- Time is discretized in steps of 2 milliseconds. The values of the spectrogram are estimated using a 25 milliseconds Gaussian window.
- We only consider frequencies from 80 Hz to 800 Hz, by steps of 20 Hz.

The normalized power spectral density is defined by

$$p_t(f) = \frac{c_t(f)}{\sum_{f'} c_t(f')} .$$

This defines a sequence of probability measures $\{p_t : t = 1, \dots, \}$.

- Denote by $c_t(f)$ the power spectral density at time t and frequency f .
- Time is discretized in steps of 2 milliseconds. The values of the spectrogram are estimated using a 25 milliseconds Gaussian window.
- We only consider frequencies from 80 Hz to 800 Hz, by steps of 20 Hz.

The normalized power spectral density is defined by

$$p_t(f) = \frac{c_t(f)}{\sum_{f'} c_t(f')} .$$

This defines a sequence of probability measures $\{p_t : t = 1, \dots, \}$.

- Denote by $c_t(f)$ the power spectral density at time t and frequency f .
- Time is discretized in steps of 2 milliseconds. The values of the spectrogram are estimated using a 25 milliseconds Gaussian window.
- We only consider frequencies from 80 Hz to 800 Hz, by steps of 20 Hz.

The normalized power spectral density is defined by

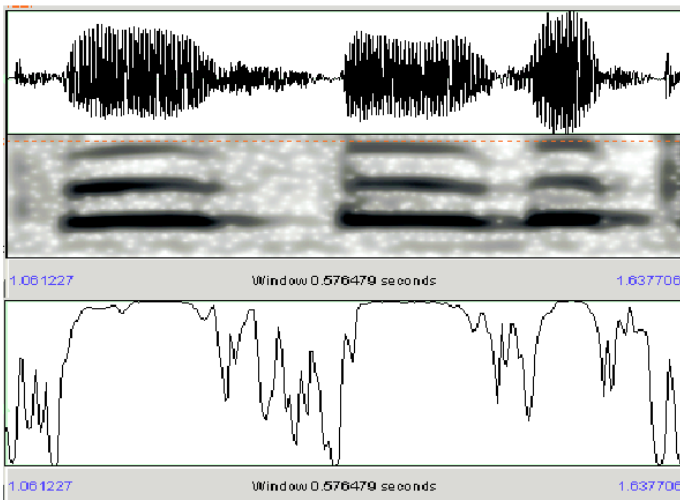
$$p_t(f) = \frac{c_t(f)}{\sum_{f'} c_t(f')} .$$

This defines a sequence of probability measures $\{p_t : t = 1, \dots, \}$.

Definition of the **sonority**

$$S_t = e^{-\eta \sum_{i=1}^3 h(p_t | p_{t-i})},$$

where h is the **relative entropy** and η is a free parameter taking positive real values.



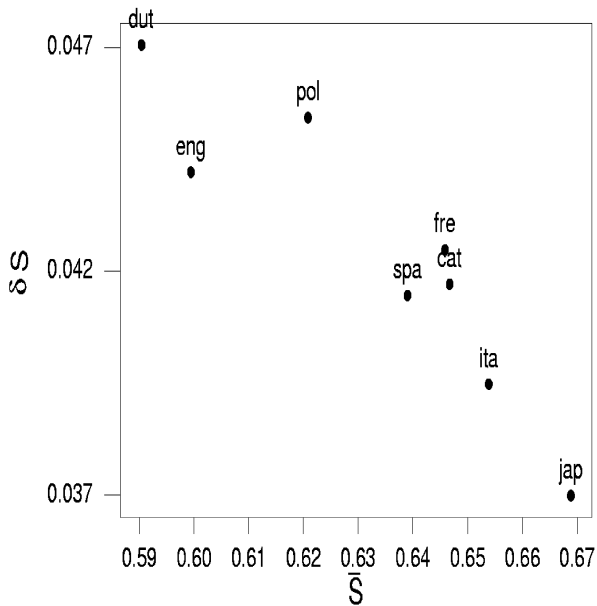
Define

- $\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$. (*This will play the role of %V.*)
- $\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t-1)|$. (*This will play the role of ΔC .*)

Reproducing RNM using the sonority

Define

- $\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$. (*This will play the role of %V.*)
- $\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t - 1)|$. (*This will play the role of ΔC .*)



A sound statistical basis to the clustering

- The pictures produced with both approaches suggest that the existence of three classes.
- Is this a real statistical fact ?
- The projected Kolmogorov-Smirnov test presented in Cuesta-Albertos, Fraiman and Ransford (2004) makes it possible to compare the laws of the stochastic processes producing the time evolutions of the sonority for the different sentences and languages.

A sound statistical basis to the clustering

- The pictures produced with both approaches suggest that the existence of three classes.
- **Is this a real statistical fact ?**
- The projected **Kolmogorov-Smirnov** test presented in **Cuesta-Albertos, Fraiman and Ransford (2004)** makes it possible to compare the laws of the stochastic processes producing the time evolutions of the **sonority** for the different sentences and languages.

A sound statistical basis to the clustering

- The pictures produced with both approaches suggest that the existence of three classes.
- **Is this a real statistical fact ?**
- The projected **Kolmogorov-Smirnov** test presented in **Cuesta-Albertos, Fraiman and Ransford (2004)** makes it possible to compare the laws of the stochastic processes producing the time evolutions of the **sonority** for the different sentences and languages.

The projected Kolmogorov-Smirnov test

- First choose a *direction* $W = (W(t))_{t \in [0, T]}$ at random. and then project the sonority trajectories in this direction.
- Then calculate the Kolmogorov-Smirnov statistic $D_W(S_I, S_{I'})$ for the two projected samples
- **Reject** the null hypothesis that the two samples belong to the same population if $D_W(S_I, S_{I'})$ is large enough. Otherwise **accept** it.

The projected Kolmogorov-Smirnov test

- First choose a *direction* $W = (W(t))_{t \in [0, T]}$ at random. and then project the sonority trajectories in this direction.
- Then calculate the Kolmogorov-Smirnov statistic $D_W(S_I, S_{I'})$ for the two projected samples
- **Reject** the null hypothesis that the two samples belong to the same population if $D_W(S_I, S_{I'})$ is large enough. Otherwise **accept** it.

The projected Kolmogorov-Smirnov test

- First choose a *direction* $W = (W(t))_{t \in [0, T]}$ at random. and then project the sonority trajectories in this direction.
- Then calculate the Kolmogorov-Smirnov statistic $D_W(S_I, S_{I'})$ for the two projected samples
- **Reject** the null hypothesis that the two samples belong to the same population if $D_W(S_I, S_{I'})$ is large enough. Otherwise **accept** it.

Stable projected Kolmogorov-Smirnov test

Instead of taking only one random direction we will take many of them. For each pair of languages $l \neq l'$ we proceed as follows.

- Choose 100 independent direction $W_i : i = 1, \dots, 100$.
- Test if the samples corresponding to l and l' belong to the same population using the projected KS test using direction W_i .
- Build up the auxiliary random variable $Z_i(l, l')$ which takes the value 1 if the projected test in direction W_i rejects the null hypothesis, and takes the value 0 otherwise.
- **Reject** the null hypothesis if the average statistic

$$\bar{Z}(l, l') = \frac{1}{N} \sum_{i=1}^N Z_i(l, l') \geq c_\alpha.$$

Stable projected Kolmogorov-Smirnov test

Instead of taking only one random direction we will take many of them. For each pair of languages $l \neq l'$ we proceed as follows.

- Choose 100 independent direction $W_i : i = 1, \dots, 100$.
- Test if the samples corresponding to l and l' belong to the same population using the projected KS test using direction W_i .
- Build up the auxiliary random variable $Z_i(l, l')$ which takes the value 1 if the projected test in direction W_i rejects the null hypothesis, and takes the value 0 otherwise.
- **Reject** the null hypothesis if the average statistic

$$\bar{Z}(l, l') = \frac{1}{N} \sum_{i=1}^N Z_i(l, l') \geq c_\alpha.$$

Stable projected Kolmogorov-Smirnov test

Instead of taking only one random direction we will take many of them. For each pair of languages $l \neq l'$ we proceed as follows.

- Choose 100 independent direction $W_i : i = 1, \dots, 100$.
- Test if the samples corresponding to l and l' belong to the same population using the projected KS test using direction W_i .
- Build up the auxiliary random variable $Z_i(l, l')$ which takes the value 1 if the projected test in direction W_i rejects the null hypothesis, and takes the value 0 otherwise.
- **Reject** the null hypothesis if the average statistic

$$\bar{Z}(l, l') = \frac{1}{N} \sum_{i=1}^N Z_i(l, l') \geq c_\alpha.$$

Stable projected Kolmogorov-Smirnov test

Instead of taking only one random direction we will take many of them. For each pair of languages $l \neq l'$ we proceed as follows.

- Choose 100 independent direction $W_i : i = 1, \dots, 100$.
- Test if the samples corresponding to l and l' belong to the same population using the projected KS test using direction W_i .
- Build up the auxiliary random variable $Z_i(l, l')$ which takes the value 1 if the projected test in direction W_i rejects the null hypothesis, and takes the value 0 otherwise.
- **Reject** the null hypothesis if the average statistic

$$\bar{Z}(l, l') = \frac{1}{N} \sum_{i=1}^N Z_i(l, l') \geq c_\alpha.$$

Results using the sonority

language	pol	ital	fren	span	dut	eng	cat
jap	0.04	0.43	0.09	0.08	0.77	0.74	0.01
pol		0.03	0.0	0.0	0.60	0.21	0.03
ital			0.03	0.02	0.14	0.0	0.05
fren				0.0	0.50	0.19	0.06
span					0.41	0.08	0.03
dut						0.0	0.74
eng							0.58
$c_{0.05}^*(l)$	0.13	0.13	0.12	0.11	0.14	0.12	0.12
$c_{0.1}^*(l)$	0.05	0.05	0.06	0.05	0.05	0.04	0.05

Results with groups of languages

We performed a new test by grouping the sonority sample paths in three groups.

- In the first group we put together the 60 sonority sample paths of the conjectured syllable-timed languages, French, Italian and Spanish.
- The second group contains the 40 paths of the conjectured stress-timed languages, Dutch and English.
- Finally the 20 sonority paths of Japanese, which is conjectured to be a mora-timed language, remain in a third group.

Results with groups of languages

We performed a new test by grouping the sonority sample paths in three groups.

- In the first group we put together the 60 sonority sample paths of the conjectured syllable-timed languages, French, Italian and Spanish.
- The second group contains the 40 paths of the conjectured stress-timed languages, Dutch and English.
- Finally the 20 sonority paths of Japanese, which is conjectured to be a mora-timed language, remain in a third group.

Results with groups of languages

We performed a new test by grouping the sonority sample paths in three groups.

- In the first group we put together the 60 sonority sample paths of the conjectured syllable-timed languages, French, Italian and Spanish.
- The second group contains the 40 paths of the conjectured stress-timed languages, Dutch and English.
- Finally the 20 sonority paths of Japanese, which is conjectured to be a mora-timed language, remain in a third group.

Results with groups of languages

category	mora-timed	stress-timed	$c_{0.05}^*(i)$
syllable-timed	0.32	0.70	0.24
stress-timed	0.82		0.09
$c_{0.05}^*(j)$	0.06	0.04	

Table: Values of $\bar{Z}(i, j)$ and bootstrap critical values for the three groups, with $N = 100$, $B = 1000$ and $\eta = 0.05$

The test found **significant** all the differences between groups. This **reinforces** the linguistic conjecture of existence of three different rhythmic classes.

These results are presented in a paper by Cuesta, Fraiman, Galves, Garcia and Svarc to appear in the Journal of Applied Statistics (2006).

2863 sentences from 15 different languages

- Plot \bar{S} vs. δS
- for 2863 sentences from 15 languages
- from a corpus belonging to the *Laboratoire de Sciences Cognitives et Psycholinguistique (EHESS/CNRS)*.
- Recall that

$$\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$$

and

$$\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t-1)| .$$

2863 sentences from 15 different languages

- Plot \bar{S} vs. δS
- for **2863 sentences** from **15 languages**
- from a corpus belonging to the *Laboratoire de Sciences Cognitives et Psycholinguistique (EHESS/CNRS)*.
- Recall that

$$\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$$

and

$$\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t-1)|.$$

2863 sentences from 15 different languages

- Plot \bar{S} vs. δS
- for **2863 sentences** from **15 languages**
- from a corpus belonging to the *Laboratoire de Sciences Cognitives et Psycholinguistique (EHESS/CNRS)*.
- Recall that

$$\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$$

and

$$\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t-1)|.$$

2863 sentences from 15 different languages

- Plot \bar{S} vs. δS
- for **2863 sentences** from **15 languages**
- from a corpus belonging to the *Laboratoire de Sciences Cognitives et Psycholinguistique (EHESS/CNRS)*.
- Recall that

$$\bar{S} = \frac{1}{T} \sum_{t=1}^T s(t)$$

and

$$\delta S = \frac{1}{T} \sum_{t=1}^T |s(t) - s(t-1)| .$$

A linear law for the speech sonority

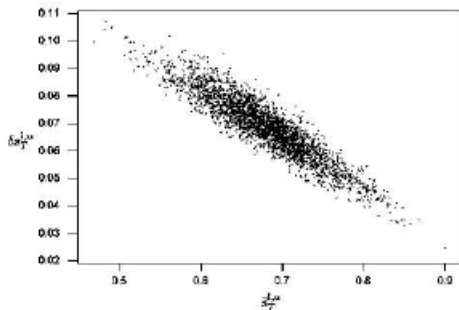


Figure: Mean sonority (horizontal axis) and mean increment of the sonority in absolute value (vertical axis) for 2863 sentences from 15 different languages

A model for the speech sonority

- A cross-linguistic exploratory analysis of the data shows that the sonority is quite **regular** in high level regions and displays **strong variations** below a certain level.
- This suggests modeling the sonority time evolutions for different languages by a family of **tied** quantized chains.

A model for the speech sonority

- A cross-linguistic exploratory analysis of the data shows that the sonority is quite **regular** in high level regions and displays **strong variations** below a certain level.
- This suggests modeling the sonority time evolutions for different languages by a family of **tied** quantized chains.

Tied quantized chains

- The chains are tied together by the assumption that the distribution of the sonority, conditioned on the fact that it belongs to a given region, is *universal*, i.e. language independent.
- In particular the partition in regions of sonority is assumed to be *language* independent.
- In this model the *specific features* characterizing each language are expressed by the *symbolic chain* indicating in which region of sonority the process is at each time step.

Tied quantized chains

- The chains are tied together by the assumption that the distribution of the sonority, conditioned on the fact that it belongs to a given region, is *universal*, i.e. language independent.
- In particular the partition in regions of sonority is assumed to be *language* independent.
- In this model the *specific features* characterizing each language are expressed by the *symbolic chain* indicating in which region of sonority the process is at each time step.

Tied quantized chains

- The chains are tied together by the assumption that the distribution of the sonority, conditioned on the fact that it belongs to a given region, is *universal*, i.e. language independent.
- In particular the partition in regions of sonority is assumed to be *language* independent.
- In this model the *specific features* characterizing each language are expressed by the *symbolic chain* indicating in which region of sonority the process is at each time step.

A model for the speech sonority

- Components of the model: two families of **stochastic chains**

$$\{(S_t^l)_{t \in \mathbf{Z}} : l \in \mathcal{L}\} \quad \text{and} \quad \{(X_t^l)_{t \in \mathbf{Z}} : l \in \mathcal{L}\}.$$

- \mathcal{L} is the *set of all languages*
- The chains $(S_t^l)_{t \in \mathbf{Z}}$ take values in the interval $[0, 1]$.
- The chains $(X_t^l)_{t \in \mathbf{Z}}$ take values in the binary alphabet $A = \{0, 1\}$.

We will assume that these chains are **stationary** and **ergodic**.
These chains are tied together by the following **assumption**.

A model for the speech sonority

- Components of the model: two families of **stochastic chains**

$$\{(S_t^l)_{t \in \mathbb{Z}} : l \in \mathcal{L}\} \quad \text{and} \quad \{(X_t^l)_{t \in \mathbb{Z}} : l \in \mathcal{L}\}.$$

- \mathcal{L} is the **set of all languages**
- The chains $(S_t^l)_{t \in \mathbb{Z}}$ take values in the interval $[0, 1]$.
- The chains $(X_t^l)_{t \in \mathbb{Z}}$ take values in the binary alphabet $A = \{0, 1\}$.

We will assume that these chains are **stationary** and **ergodic**.
These chains are tied together by the following **assumption**.

A model for the speech sonority

- Components of the model: two families of **stochastic chains**

$$\{(S_t^l)_{t \in \mathbb{Z}} : l \in \mathcal{L}\} \quad \text{and} \quad \{(X_t^l)_{t \in \mathbb{Z}} : l \in \mathcal{L}\}.$$

- \mathcal{L} is the **set of all languages**
- The chains $(S_t^l)_{t \in \mathbb{Z}}$ take values in the interval $[0, 1]$.
- The chains $(X_t^l)_{t \in \mathbb{Z}}$ take values in the binary alphabet $A = \{0, 1\}$.

We will assume that these chains are **stationary** and **ergodic**.
These chains are tied together by the following **assumption**.

A model for the speech sonority

- Components of the model: two families of **stochastic chains**

$$\{(S_t^l)_{t \in \mathbf{Z}} : l \in \mathcal{L}\} \quad \text{and} \quad \{(X_t^l)_{t \in \mathbf{Z}} : l \in \mathcal{L}\}.$$

- \mathcal{L} is the **set of all languages**
- The chains $(S_t^l)_{t \in \mathbf{Z}}$ take values in the interval $[0, 1]$.
- The chains $(X_t^l)_{t \in \mathbf{Z}}$ take values in the binary alphabet $A = \{0, 1\}$.

We will assume that these chains are **stationary** and **ergodic**.
These chains are tied together by the following **assumption**.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S_t^l \in B \mid X_t^l = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S_t^l \in B, S_{t+1}^l \in C \mid X_t^l = i, X_{t+1}^l = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S_t^l \in B \mid X_t^l = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S_t^l \in B, S_{t+1}^l \in C \mid X_t^l = i, X_{t+1}^l = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S'_t \in B \mid X'_t = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S'_t \in B, S'_{t+1} \in C \mid X'_t = i, X'_{t+1} = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S'_t \in B \mid X'_t = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S'_t \in B, S'_{t+1} \in C \mid X'_t = i, X'_{t+1} = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S_t^l \in B \mid X_t^l = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S_t^l \in B, S_{t+1}^l \in C \mid X_t^l = i, X_{t+1}^l = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

Assumption

- There exist **probability distributions**
 - π_i on $[0, 1]$
 - $\pi_{(i,j)}$ on $[0, 1]^2$
- indexed by symbols i and j in the alphabet $A = \{0, 1\}$
- which are **language independent**
- such that for any $l \in \mathcal{L}$

$$\mathbb{P} \left\{ S_t^l \in B \mid X_t^l = j \right\} = \pi_j(B), \quad (1)$$

and

$$\mathbb{P} \left\{ S_t^l \in B, S_{t+1}^l \in C \mid X_t^l = i, X_{t+1}^l = j \right\} = \pi_{i,j}(B \times C), \quad (2)$$

where B and C are Borel subsets of $[0, 1]$.

For this model we have

$$\mathbb{E} \left(|S'_t - S'_{t+1}| \right) = a + b\mathbb{E} \left(S'_t \right) + \epsilon',$$

where the constants a and b are language independent.

- $p^l(i) = \mathbb{P} \{X_t^l = i\}$
- $p^l(i, j) = \mathbb{P} \{X_t^l = i, X_{t+1}^l = j\}$
- $\theta(i) = \mathbb{E} (S_t^l | X_t^l = i)$ and
- $\theta(i, j) = \mathbb{E} (|S_t^l - S_{t+1}^l| | X_t^l = i, X_{t+1}^l = j)$.

- $p^l(i) = \mathbb{P} \{X_t^l = i\}$
- $p^l(i, j) = \mathbb{P} \{X_t^l = i, X_{t+1}^l = j\}$
- $\theta(i) = \mathbb{E} (S_t^l | X_t^l = i)$ and
- $\theta(i, j) = \mathbb{E} (|S_t^l - S_{t+1}^l| | X_t^l = i, X_{t+1}^l = j)$.

- $p'(i) = \mathbb{P} \{X'_t = i\}$
- $p'(i, j) = \mathbb{P} \{X'_t = i, X'_{t+1} = j\}$
- $\theta(i) = \mathbb{E} (S'_t | X'_t = i)$ and
- $\theta(i, j) = \mathbb{E} (|S'_t - S'_{t+1}| | X'_t = i, X'_{t+1} = j)$.

- $p^l(i) = \mathbb{P} \{X_t^l = i\}$
- $p^l(i, j) = \mathbb{P} \{X_t^l = i, X_{t+1}^l = j\}$
- $\theta(i) = \mathbb{E} (S_t^l | X_t^l = i)$ and
- $\theta(i, j) = \mathbb{E} (|S_t^l - S_{t+1}^l| | X_t^l = i, X_{t+1}^l = j)$.

Constants and correction term

$$a = \theta(0,0) - \theta(0) \frac{\theta(1,1) - \theta(0,0)}{\theta(1) - \theta(0)}, \quad b = \frac{\theta(1,1) - \theta(0,0)}{\theta(1) - \theta(0)},$$

and the correction ϵ^l is language dependent and defined as

$$\epsilon^l = p_{0,1}^l (\theta(1,0) + \theta(0,1) - \theta(0,0) - \theta(1,1)) .$$

Consequence

- The theorem together with the ergodicity of the chains (S_t') explains the linear relationship in the data.
- We observe that in the corpus

$$0.008 < \epsilon' < 0.009$$

and

$$0.065 < \mathbb{E} \left(|S_t' - S_{t+1}'| \right) < 0.075.$$

Consequence

- The theorem together with the ergodicity of the chains (S_t^l) explains the linear relationship in the data.
- We observe that in the corpus

$$0.008 < \epsilon^l < 0.009$$

and

$$0.065 < \mathbb{E} \left(|S_t^l - S_{t+1}^l| \right) < 0.075.$$

Further directions of research

- We still don't know what is a rhythmic feature!
- We still don't have a model for the rhythmic classes
- But the symbolic chains behind the sonority open new perspectives of research.
- An example of this is the discrimination between Brazilian and European Portuguese, modeling **stress** contours obtained by codifying **written** texts, using **Suffix Probabilistic Trees**.

Further directions of research

- We still don't know what is a rhythmic feature!
- We still don't have a model for the rhythmic classes
- But the symbolic chains behind the sonority open new perspectives of research.
- An example of this is the discrimination between Brazilian and European Portuguese, modeling **stress** contours obtained by codifying **written** texts, using **Suffix Probabilistic Trees**.

Further directions of research

- We still don't know what is a rhythmic feature!
- We still don't have a model for the rhythmic classes
- But the symbolic chains behind the sonority open new perspectives of research.
- An example of this is the discrimination between Brazilian and European Portuguese, modeling **stress** contours obtained by codifying **written** texts, using **Suffix Probabilistic Trees**.

Further directions of research

- We still don't know what is a rhythmic feature!
- We still don't have a model for the rhythmic classes
- But the symbolic chains behind the sonority open new perspectives of research.
- An example of this is the discrimination between Brazilian and European Portuguese, modeling **stress** contours obtained by codifying **written** texts, using **Suffix Probabilistic Trees**.

The results presented have been obtained by
Marzio Cassandro, Pierre Collet, Juan Cuesta, Denise Duarte,
Ricardo Fraiman, Charlotte Galves and Marcela Svarc.

www.ime.usp.br/~tycho