

Session 6 bis: Markov Chains and PageRank

Optimization and Computational Linear Algebra for Data Science

Be accurate !

Let x, v be vectors, S a subspace of \mathbb{R}^n and M an $n \times n$ matrix.

NO

- ❖ ~~$x = S$ or $x \subset S$~~
- ❖ ~~$S \in \mathbb{R}^n$~~
- ❖ ~~$\text{Span}(x, v) = \{ax + bv\}$~~
- ❖ ~~$\dim(M)$ or $\dim(x)$~~
- ❖ ~~$\text{Ker}(M) = 0$~~
- ❖ ~~$x + M$~~

YES

- $x \in S$
- $S \subset \mathbb{R}^n$
- $\{ax + bv \mid a, b \in \mathbb{R}\}$
- $\text{rank}(M)$
- $\text{Ker } M = \{0\}$

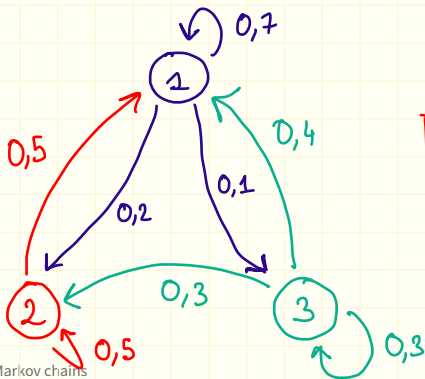
Contents

1. Markov chains
2. Perron-Frobenius Theorem
3. Application: PageRank
4. A first look at the Spectral theorem.

Markov chains

An example

- A cat has 3 states: ① Sleeping ② Eating ③ Playing
- We represent its evolution in time with a sequence:
 $X_0, X_1, X_2, \dots, X_t, \dots$ $\in \{1, 2, 3\}$ states at time t
"Markov chain"



Transition matrix:

$$P = \begin{pmatrix} 0,7 & 0,5 & 0,4 \\ 0,2 & 0,5 & 0,3 \\ 0,1 & 0 & 0,3 \end{pmatrix}$$

$$P(X_{t+1} = i \mid X_t = j) = P_{i,j}$$

Stochastic matrices

Definition

A matrix $P \in \mathbb{R}^{n \times n}$ is said to be *stochastic* if:

1. $P_{i,j} \geq 0$ for all $1 \leq i, j \leq n$. (non-negative entries.)
2. $\sum_{i=1}^n P_{i,j} = 1$, for all $1 \leq j \leq n$. (each column sum to 1)

Given a stochastic matrix $P \in \mathbb{R}^{n \times n}$, one can define a Markov chain over n states, and vice-versa.

Probability vectors

Question: after t steps, what is the probability of being in a given state $j \in \{1, \dots, n\}$?

• We have to compute:

$$x^{(t)} = \begin{pmatrix} \mathbb{P}(X_t = 1) \\ \vdots \\ \mathbb{P}(X_t = n) \end{pmatrix} \in \mathbb{R}^n$$

(definition)

• $x^{(t)} \in \Delta_n = \left\{ \underline{x} \in \mathbb{R}^n \mid \begin{array}{l} \underline{x}_i \geq 0 \text{ for all } i \\ \sum_{i=1}^n \underline{x}_i = 1 \end{array} \right\}$

set of all "probability vectors"

The key equation

Proposition

For all $t \geq 0$

$$x^{(t+1)} = Px^{(t)} \text{ and consequently, } x^{(t)} = P^t x^{(0)}.$$

Proof:

$$\begin{aligned} x_i^{(t+1)} &= \mathbb{P}(X_{t+1} = i) \\ &= \sum_{j=1}^n \underbrace{\mathbb{P}(X_{t+1} = i \mid X_t = j)}_{P_{ij}} \underbrace{\mathbb{P}(X_t = j)}_{x_j^{(t)}} \\ &= \sum_{j=1}^n P_{ij} x_j^{(t)} = (Px^{(t)})_i \end{aligned}$$

Long-term behavior

- Numerical simulations suggest that

$$x^{(t)} \xrightarrow{t \rightarrow +\infty} \nu \quad \text{for some } \nu \in \Delta_n.$$

- Since

$$\begin{array}{ccc} x^{(t+4)} = P x^{(t)} & & \\ \downarrow t \rightarrow +\infty & & \downarrow t \rightarrow +\infty \\ \boxed{\nu = P \nu} & & \end{array}$$

- ν has to be an eigenvector of P associated with the eigenvalue 1.

Perron-Frobenius Theorem

Invariant measure

Definition

A vector $\mu \in \Delta_n$ is called an invariant measure for the transition matrix P if

$$\mu = P\mu,$$

i.e. if μ is an eigenvector of P associated with the eigenvalue 1.

"invariant": if X_t is distributed according to μ . ($x^{(t)} = \mu$) then

$$x^{(t+1)} = P x^{(t)} = P \mu = \mu$$

Therefore X_{t+1} is also distributed according to μ .

Perron-Frobenius Theorem

Theorem

Let P be a stochastic matrix such that there exists $k > 1$ such that all the entries of P^k are strictly positive. Then the following holds:

- 1 is an eigenvalue of P and there exists an eigenvector $\mu \in \Delta_n$ associated to 1. $(P\mu = \mu)$
- The eigenvalue 1 has multiplicity 1: $\text{Ker}(P - \text{Id}) = \text{Span}(\mu)$.
- For all $x \in \Delta_n$, $P^t x \xrightarrow{t \rightarrow \infty} \mu$.

Consequence

Corollary

Let P be a stochastic matrix such that there exists $k > 1$ such that all the entries of P^k are strictly positive.

Then there exists a unique invariant measure μ and for all initial condition $x^{(0)} \in \Delta_n$,

$$x^{(t)} = P^t x^{(0)} \xrightarrow{t \rightarrow \infty} \mu.$$

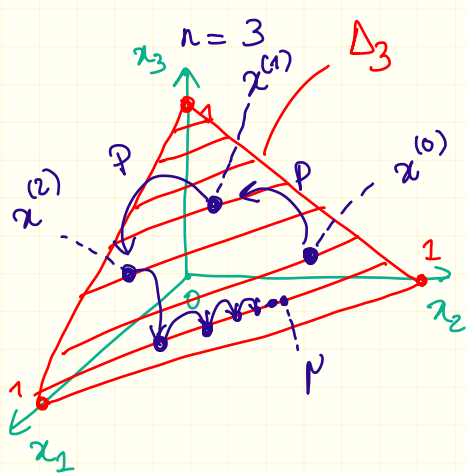
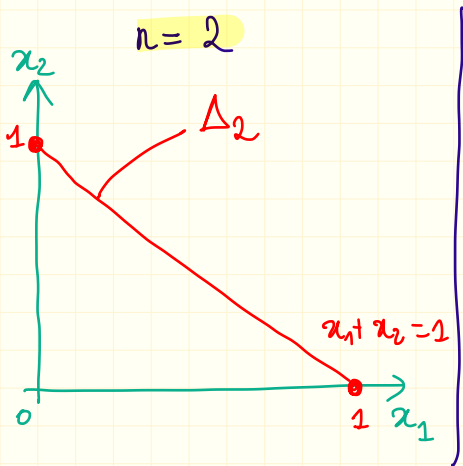
Point 1

pt. 2

pt. 3.

Proof: Geometrical observations

$$\Delta_n = \left\{ x \in \mathbb{R}^n \mid \begin{array}{l} x_i \geq 0 \text{ for all } i \\ x_1 + \dots + x_n = 1 \end{array} \right\}$$



Proof: contraction

We will prove the theorem in the case where $P_{i,j} > 0$ for all i, j .

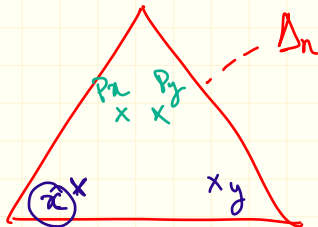
Lemma

The mapping

$$\begin{aligned}\varphi: \Delta_n &\rightarrow \Delta_n \\ x &\mapsto Px\end{aligned}$$

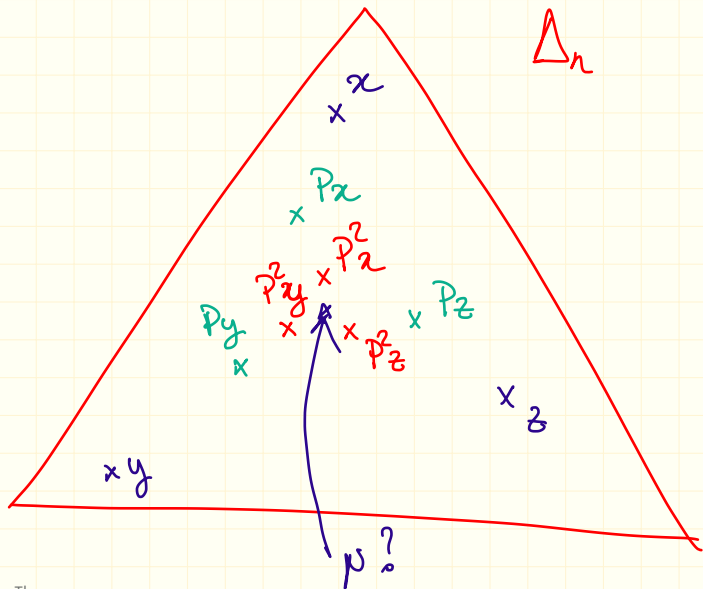
is a contraction mapping for the ℓ_1 -norm: there exists $c \in (0, 1)$ such that for all $x, y \in \Delta_n$:

$$\|Px - Py\|_1 \leq c \|x - y\|_1.$$



~~≠~~ \approx

Geometric picture



Proof of Perron-Frobenius

- Let $p \in \Delta_n$ be a minimizer of $x \mapsto \|Px - x\|_1$ on Δ_n . (we admit that \exists^a minimizer exists)

- Then, $Pp = p$. Indeed, if $\|Pp - p\|_1 > 0$,

$$\|P(Pp) - Pp\|_1 \leq c \|Pp - p\|_1 < \|Pp - p\|_1$$

"contraction"

this contradicts the optimality of p .

- Let $x \in \Delta_n$. $\|P^t x - p\|_1 = \|P^t x - P^t p\|_1$

$$0 < c < 1 \quad \leq c^t \|x - p\|_1 \xrightarrow{t \rightarrow \infty} 0$$

Proof of Perron-Frobenius

• let $x \in \mathbb{R}^n$ such that $\boxed{Px = x}$

$$x = P^t x$$

$$= P^t (x_1 e_1 + \dots + x_n e_n)$$

$$= x_1 \underbrace{P^t e_1} + \dots + x_n \underbrace{P^t e_n}$$

because
 $e_1, \dots, e_n \in \Delta_n$

$\xrightarrow{t \rightarrow \infty} N$

$\xrightarrow{t \rightarrow \infty} N$

$$\boxed{x = (x_1 + \dots + x_n) N}$$

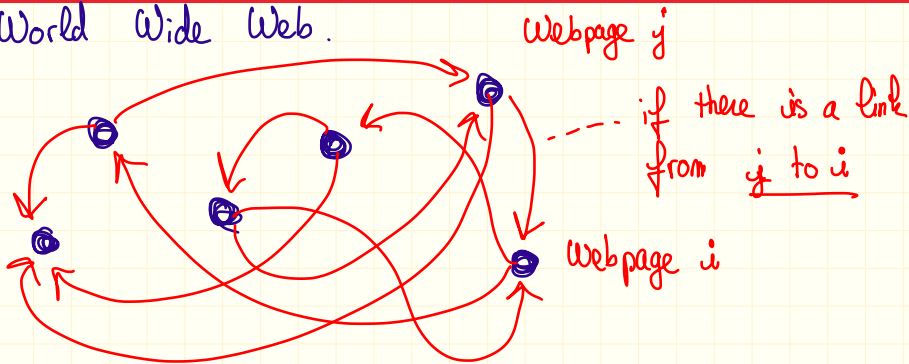
$$\xrightarrow{t \rightarrow \infty} (x_1 + \dots + x_n) N \in \text{Span}(\mu)$$

$$\boxed{x \in \text{Span}(\mu)}$$

PageRank

Ordering the Web

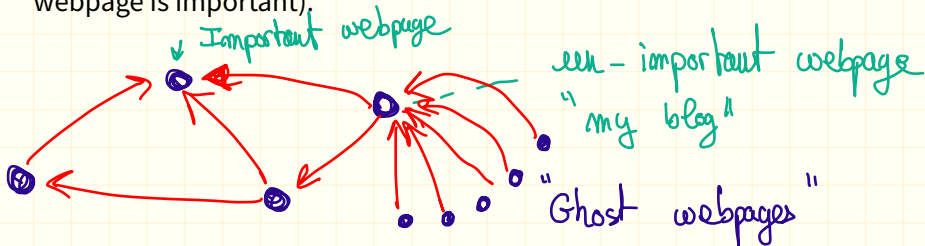
World Wide Web.



Goal: rank the webpages from the most "important" one to the less important one.

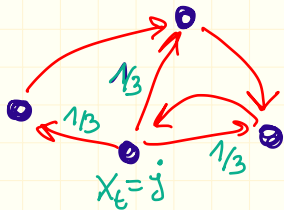
Naive attempt

First idea: rank the webpages according to their number of *incomming links*. (The more incomming links, the more the webpage is important).



That does not work : one can create thousand of "ghost webpages" pointing to a single webpage.

The random surfer



Imagine a "drunk surfer"

- if at time t , the surfer is on a webpage $X_t = j$

he clicks on an outgoing link selected at random.

$$P(X_{t+1} = i \mid X_t = j) = \begin{cases} 0 & \text{if there is no link } j \rightarrow i \\ \frac{1}{\deg(j)} & \text{otherwise} \end{cases}$$

where $\deg(j)$ is the number of outgoing links from j

PageRank Algorithm

This defines a Markov chain of transition matrix:

$$P_{i,j} = \begin{cases} 1/\deg(j) & \text{if there is a link } j \rightarrow i \\ 0 & \text{otherwise,} \end{cases}$$

- After a long time, the surfer is more likely to be on an *important webpage*.
- If μ is the invariant measure of P (provided P verifies the hypotheses of Perron-Frobenius), we take

$$\mu_i = \text{« importance of webpage } i \text{ »}.$$

\approx prob. of being at i after a long time

PageRank Algorithm

Issue



Google considered the transition matrix:

$$G = \alpha P + \frac{1-\alpha}{N} \begin{pmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{pmatrix}$$

$\alpha \approx 0,85$

total number
of pages

This avoids being trapped !

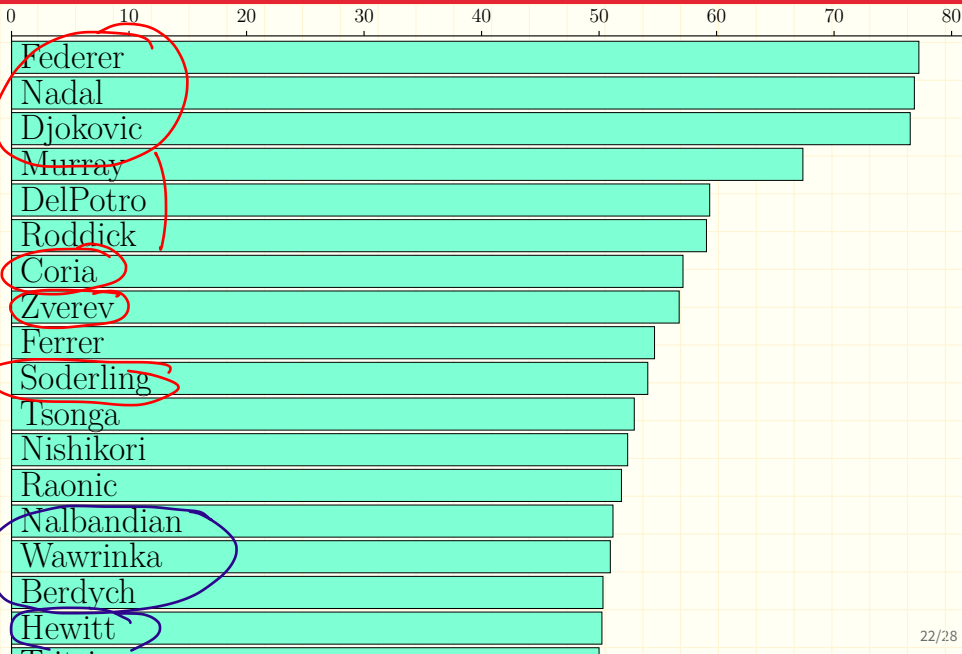
Application: ranking Tennis players

Goal: rank the following players:

Federer, Nadal, Djokovic, Murray, Del Potro, Roddick, Coria, Zverev, Ferrer, Soderling, Tsonga, Nishikori, Raonic, Nalbandian, Wawrinka, Berdych, Hewitt, Tsitsipas, Monfils, Gonzalez, Thiem, Ljubicic, Davydenko, Cilic, Pouille, Safin, Isner, Dimitrov, Medvedev, Ferrero, Goffin, Bautista Agut, Sock, Gasquet, Simon, Blake, Monaco, Coric, Stepanek, Khachanov, Almagro, Robredo, Verdasco, Anderson, Youzhny, Baghdatis, Dolgoplov, Kohlschreiber, Fognini, Melzer, Paire, Querrey, Tomic, Basilashvili.

Data: Head-to Head records (number of times that player x has defeated player y)

Ranking by % of victories



The random spectator

Imagine the following « random spectator »:

- ❖ At time t , the spectator believes that player j is the best:
 $X_t = j$.
- ❖ Then, he picks a game of player j uniformly at random:
 - ❖ if player j wins, then the spectator still believes that j is the best: $X_{t+1} = j$.
 - ❖ otherwise, the spectator changes his mind and now believes that player i who defeated j is the best: $X_{t+1} = i$.

The random spectator

Imagine the following « random spectator »:

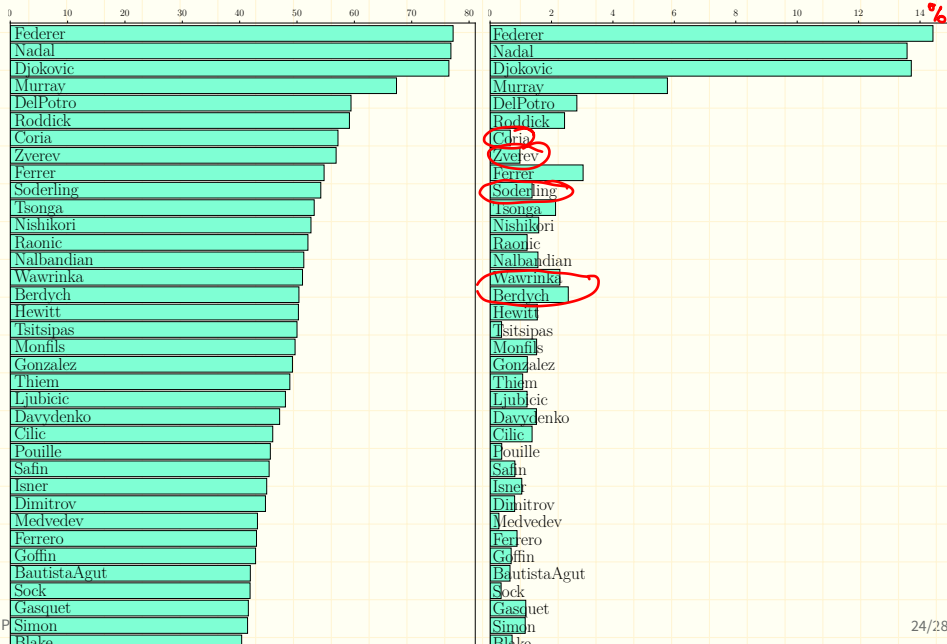
- At time t , the spectator believes that player j is the best:
 $X_t = j$.
- Then, he picks a game of player j uniformly at random:
 - if player j wins, then the spectator still believes that j is the best: $X_{t+1} = j$.
 - otherwise, the spectator changes his mind and now believes that player i who defeated j is the best: $X_{t+1} = i$.

This defines a transition matrix P . We rank the players according to the stationary distribution μ of

$$M = \alpha P + \frac{1 - \alpha}{N} J$$

$$\begin{pmatrix} 1 & - & 1 \\ 1 & - & 1 \\ 1 & - & 1 \end{pmatrix}$$

Naive ranking vs PageRank



The Spectral Theorem

The spectral theorem

Theorem

Let $A \in \mathbb{R}^{n \times n}$ be a **symmetric** matrix. Then there is a orthonormal basis of \mathbb{R}^n composed of eigenvectors of A .

There exists $v_1 \dots v_n \in \mathbb{R}^n$ such that

- $(v_1 \dots v_n)$ is an **orthonormal basis** of \mathbb{R}^n
- $A v_i = \lambda_i v_i$ for all i

If we consider $x \in \mathbb{R}^n$

$$\underline{x = \langle v_1, x \rangle v_1 + \dots + \langle v_n, x \rangle v_n}$$

The spectral theorem

$$Ax = \langle v_1, x \rangle A v_1 + \dots + \langle v_n, x \rangle A v_n$$

$$= \langle v_1, x \rangle \lambda_1 v_1 + \dots + \langle v_n, x \rangle \lambda_n v_n$$

$$P = \begin{pmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{pmatrix}$$

(this is an orthogonal matrix)

$$P^T x = \begin{pmatrix} \langle v_1, x \rangle \\ \vdots \\ \langle v_n, x \rangle \end{pmatrix}$$

$$\text{let } D = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \lambda_n \end{pmatrix}$$

$$D P^T x = \begin{pmatrix} \lambda_1 \langle v_1, x \rangle \\ \vdots \\ \lambda_n \langle v_n, x \rangle \end{pmatrix}$$

$$Ax = P D P^T x$$

$$A = P \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} P^T$$

Matrix formulation

Theorem

Let $A \in \mathbb{R}^{n \times n}$ be a **symmetric** matrix. Then there exists an orthogonal matrix P and a diagonal matrix D of sizes $n \times n$ such that

$$A = PDP^T.$$

- the columns of P are eigenvectors of A
- the entries on the diagonal of D are associated eigenvalues.

Questions?

Questions?