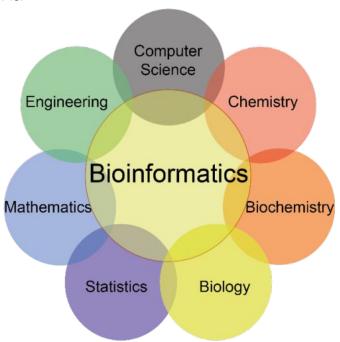


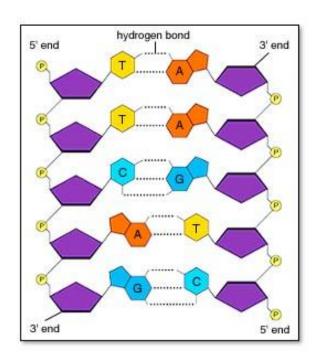
### ¿Qué es la bioinformática?

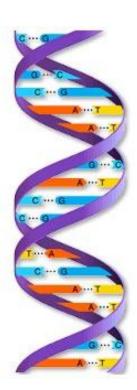
La bioinformática es la aplicación de tecnologías computacionales a la gestión y análisis de datos biológicos.

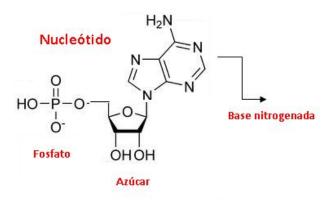
El término se acuñó en 1970 por Paulien Hogeweg y Ben Hesper



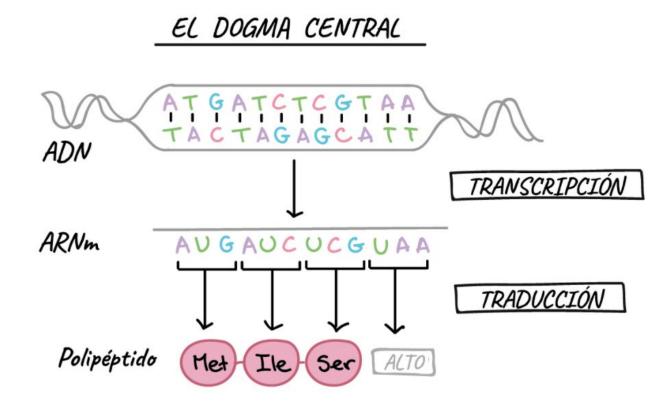
### ¿Qué es el ADN?



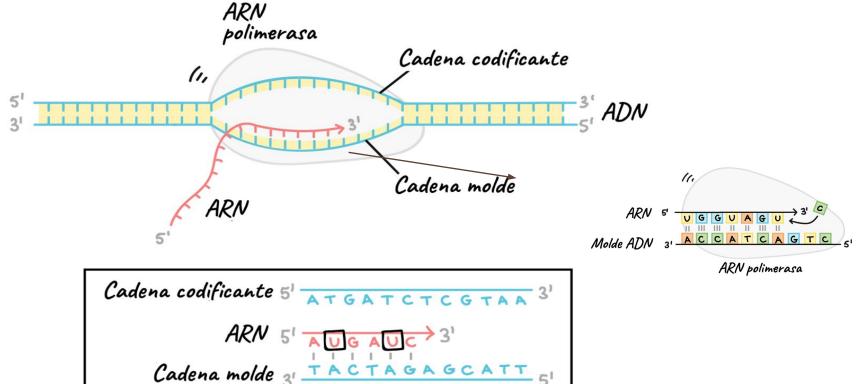


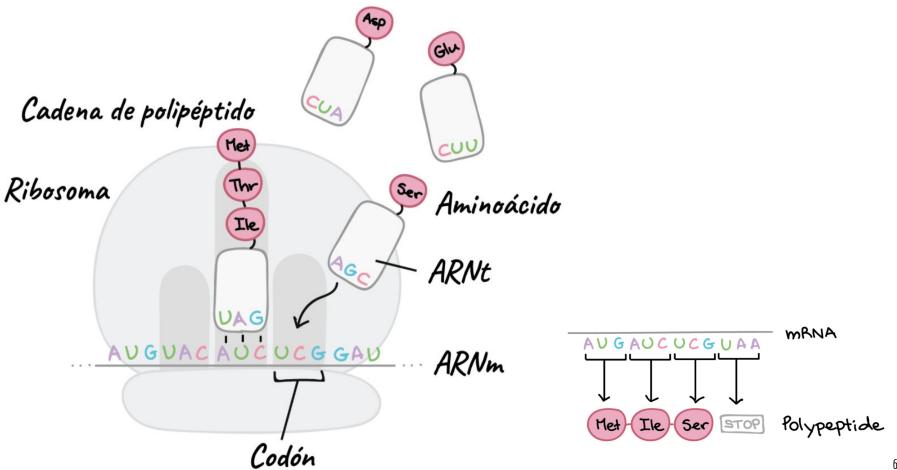


### De ADN a proteína



### De ADN a proteína

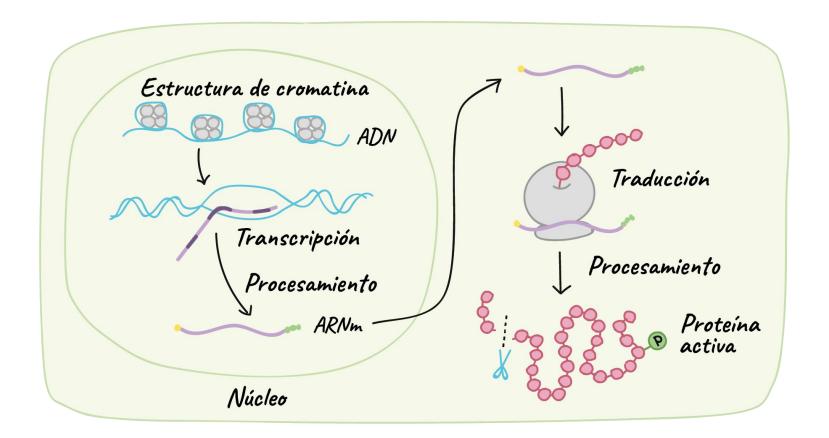




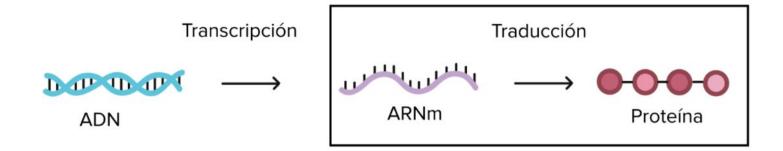
Second base of codon

					Occord ba	30 01	codon				
			U	С			Α	G			
		UUU	Phenylalanine	UCU		UAU	Tyrosine	UGU	Cysteine	U	
	U	UUC	phe Leucine leu	UCC	Serine ser	UAC	tyr	UGC	cys	С	
		UUA		UCA		UAA	STOP codon	UGA	STOP codon	Α	
		UUG		UCG		UAG	0101.000011	UGG	Tryptonphan trp	G	
		CUU	Leucine leu	CCU	Proline pro	CAU	Histidine his Glutamine	CGU		U	
6	С	CUC		CCC		CAC		CGC	Arginine arg	С	E
ot codon		CUA		CCA		CAA		CGA		Α	
		CUG	8	CCG		CAG	gin	CGG		G	Dase
First base		AUU		ACU	Threonine thr	AAU	Asparagine asn	AGU	Serine ser Arginine	U	2
st D	Α	AUC	Isoleucine ile	ACC		AAC		AGC		С	COGOI
Ì		AUA		ACA		AAA		AGA		Α	2
		AUG	Methionine met (start codon)	ACG		AAG	lys	AGG	arg	G	
		GUU		GCU	Alanine ala	GAU	Aspartic acid asp Glutamic acid glu	GGU		U	
	G	GUC	valle	GCC		GAC		GGC	Glycine gly	С	
		GUA		GCA		GAA		GGA		Α	
		GUG		GCG		GAG		GGG		G	

## El código genético es degenerado



### De ADN a proteína



### Campos de aplicación de la Bioinformática

- Medicina y medicina preventiva
- Microbiología
- Desarrollo de fármacos
- Resistencia antibiótica
- Biotecnología
- Agricultura
- Veterinaria
- Análisis forenses
- Muchas otros

### Formas de aplicación de la Bioinformática

Genómica

Proteómica

Transcriptómica

Metabolómica

### Ejemplo

#### Molde de ADN:

5'-ATGAAGTTTGGCACTTAA-3' 3'-TACTTCAAACCGTGAATT-5'



### Ejemplo: Mioglobina

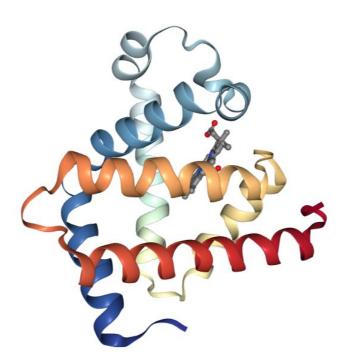
#### Secuencia:

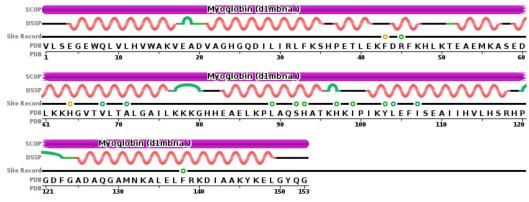
GAGCATGTTGGCCTGGTCCTTTG CTAGGTACTGTAGAGCAGGTGAG AGAGTGAGGGGGGAAGGACTCCAA ATTAGACCAGTTCTTAGCCATGAA GCAGAGACTCTGAAGCCAGACTA CCTGGGTCCCAATCTTGGGCTTG GTATTTCCTCGCTGTGTGACTCTG GGTAAGTTACTTAACTTCTCTGTG CCTCAGTTCTCTCAAGTGTAAA...

#### Aminoácidos:

MGLSDGEWQLVLNVWGKVEADIPG HGQEVLIRLFKGHPETLEKFDKFKHL

### Ejemplo: Mioglobina





### Matemáticas

#BioinformaticsGRX

Alguien tendrá que traducir los requerimientos biológicos en un lenguaje entendible para la máquina



# Multitud de herramientas matemáticas a nuestra disposición.

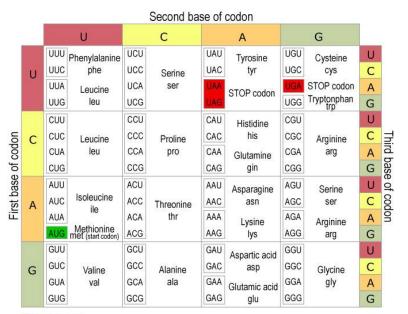
- Geometría diferencial
- Análisis
- Teoría de nudos
- Probabilidad y estadística
- Álgebra y álgebra abstracta

### 1. Empecemos por lo básico.

### ¿Cómo representamos el código genético?

#### Claves

- Simetrías
- Matrices
- Código degenerado
- Inmunidad al ruido



Clinical Tools, Inc.

### Una posible idea:

Discernir subconjuntos dentro del código genético y utilizar con ellos lenguaje binario:

TABLE 2.2 Three Binary Subalphabets According to Three Types of Binary-Opposite Attributes in a Set of Nitrogenous Bases C, A, G, U<sup>a</sup>

N	Symbol of a Genetic Letter	$\mathbf{C}$	Α	G	U/T
1	0 <sub>1</sub> , pyrimidines (one ring in a molecule) 1 <sub>1</sub> , purines (two rings in a molecule)	01	11	11	01
2	0 <sub>2</sub> , a letter with amino-mutating property (amino) 1 <sub>2</sub> , a letter without it (keto)	02	02	12	12
3	0 <sub>3</sub> , a letter with three hydrogen bonds 1 <sub>3</sub> , a letter with two hydrogen bonds	03	13	03	13

### Aprovechar los algoritmos existentes:

Matrices de Hadamard y producto de Kronecker:

Si denominamos:

$$H_{n+1} = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}^n$$

Como el Kernel de la familia de matrices, inmediatamente podemos establecer una conexión con nuestro código genético:

$$P(n) = \begin{pmatrix} C & A \\ U & G \end{pmatrix}^n$$

						00(	0)	01(1)	10(2)	11(3)
					00		CC	CA	AC	AA
	0		1			(0)	0000	0001	0010	0011
							(0)	(1)	(2)	(3)
$P^{(1)}=$		С	A		$P^{(2)}=$	01	CU	CG	AU	AG
	0	00	01	;		(1)	0100	0101	0110	0111
		(0)	(1)				(4)	(5)	(6)	(7)
		U	G		10		UC	UA	GC	GA
	1	10	11			(2)	1000	1001	1010	1011
		(2)	(3)				(8)	(9)	(10)	(11)
				•		11	UU	UG	GU	GG
						(3)	1100	1101	1110	1111
							(12)	(13)	(14)	(15)

	000 (0)	001(1)	010 (2)	011 (3)	100 (4)	101 (5)	110 (6)	111 (7)
<u>000</u> <u>(0)</u>	CCC 000000 (0)	CCA 000001 (1)	CAC 000010 (2)	CAA 000011 (3)	ACC 000100 (4)	ACA 000101 (5)	AAC 000110 (6)	AAA 000111 (7)
<u>001</u> (1)	CCU 001000 (8)	CCG 001001 (9)	CAU 001010 (10)	CAG 001011 (11)	ACU 001100 (12)	ACG 001101 (13)	AAU 001110 (14)	AAG 001111 (15)
<u>010</u> (2)	CUC 010000 (16)	CUA 010001 (17)	CGC 010010 (18)	CGA 010011 (19)	AUC 010100 (20)	AUA 010101 (21)	AGC 010110 (22)	AGA 010111 (23)
<u>011</u> <u>(3)</u>	CUU 011000 (24)	CUG 011001 (25)	CGU 011010 (26)	CGG 011011 (27)	AUU 011100 (28)	AUG 011101 (29)	AGU 011110 (30)	AGG 011111 (31)
100 (4)	UCC 100000 (32)	UCA 100001 (33)	UAC 100010 (34)	UAA 100011 (35)	GCC 100100 (36)	GCA 100101 (37)	GAC 100110 (38)	GAA 100111 (39)
101 (5)	UCU 101000 (40)	UCG 101001 (41)	UAU 101010 (42)	UAG 101011 (43)	GCU 101100 (44)	GCG 101101 (45)	GAU 101110 (46)	GAG 101111 (47)
110 (6)	UUC 110000 (48)	UUA 110001 (49)	UGC 110010 (50)	UGA 110011 (51)	GUC 001100 (52)	GUA 110101 (53)	GGC 110110 (54)	GGA 110111 (55)
<u>111</u> <u>(7)</u>	UUU 111000 (56)	UUG 111001 (57)	UGU 111010 (58)	UGG 111011 (59)	GUU 111100 (60)	GUG 111101 (61)	GGU 111110 (62)	GGG 111111 (63)

### Z. Alineación

### Bases

Una alineación entre dos (o más) secuencias es una comparación entre los componentes de cada secuencia.

El análisis básico de secuencias consiste en preguntar si hay dos o más secuencias relacionadas.

Una verdadera alineación de las secuencias biológicas es aquella que refleja la relación evolutiva entre dos o más secuencias que comparten un antepasado común antepasado común.

### Bases

La base de la alineación y el análisis de la secuencia se basa en el hecho de que las secuencias biológicas se desarrollan a partir de secuencias preexistentes en lugar de ser inventado por la naturaleza desde el principio.

A esto debemos añadirle tres fenómenos principales:

- Mutaciones puntuales
- Inserciones
- Delecciones

### Bases

Además de los conceptos biológicos, debemos hacernos preguntas específicas con las que deberemos trabajar, y de las que dependerá nuestro método:

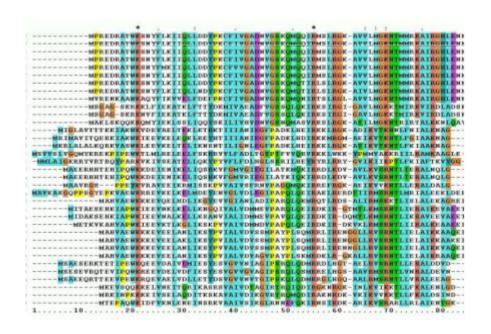
- ¿Qué tipo de alineamiento estamos buscando?
- ¿Qué puntuación aplicamos al alineamiento?
- ¿Qué algoritmo desarrollamos de forma óptima para llevar nuestro método a la máquina?
- ▶ ¿Cuál es la significación estadística de nuestro resultado?

### Ejemplos:

Alineamiento global: FASTA

Alineamiento local: BLAST

### Siguiente paso: Alineamiento múltiple.



## Análisis de componentes principales

### Definición:

"Principal component analysis (PCA) is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components (or sometimes, principal modes of variation)"

### Formas básicas:

#### Existen dos formas básicas de aplicar el ACP:

- Método basado en la matriz de correlación, cuando los datos no son dimensionalmente homogéneos o el orden de magnitud de las variables aleatorias medidas no es el mismo.
- 2. Método basado en la matriz de covarianzas, que se usa cuando los datos son dimensionalmente homogéneos y presentan valores medios similares.

### Caso particular: Correlación

$$\mathbf{R} = [r_{ij}] \in M_{m imes m} \qquad ext{donde} \qquad r_{ij} = rac{ ext{cov}(F_i, F_j)}{\sqrt{ ext{var}(F_i) ext{var}(F_j)}}$$

http://setosa.io/ev/principal-component-analysis/

# Multitud de herramientas matemáticas a nuestra disposición.

- Geometría diferencial
- Análisis
- Teoría de nudos
- Probabilidad y estadística
- Álgebra y álgebra abstracta

### Informática

Bioinformática

# Illumina MiSeq



#### ¡Ups!



## ¿Qué recursos tenemos?

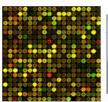






















Biological Data GUI usage CLI usage Ad hoc script

programming

Algorithm

Math and IT staff

Bioinformaticians

Biologists

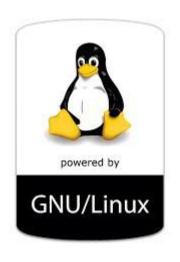
## ¿Qué podemos hacer?

- Programar
- Utilizar plataformas fiables y potentes

## ¿Qué podemos hacer?

- Aprender a dar órdenes al ordenador con total libertad
- Utilizar plataformas fiables y potentes

#### Otros sistemas operativos



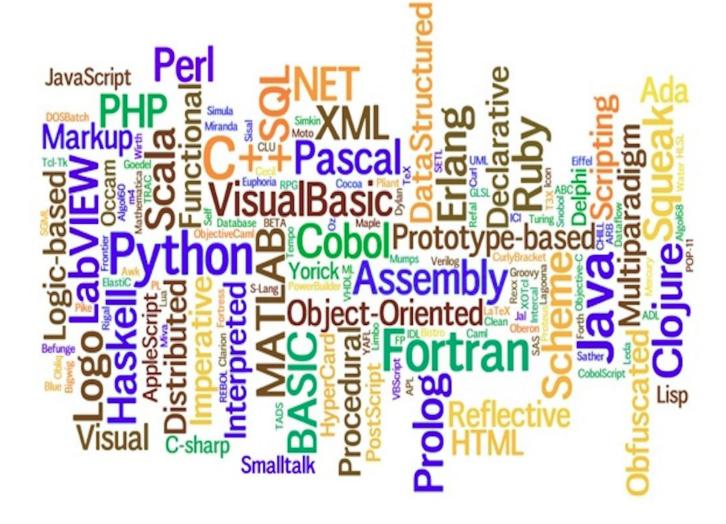


¿Programar?





¿Cuál?



¿Python?



## ¿Por qué Python?

- Tiene una sintaxis sencilla
- Viene con bibliotecas por defecto con muchas tareas comunes
- Es libre
- Es interpretado
- Es muy utilizado (también en la comunidad científica)

#### Sintaxis sencilla

print("Hola mundo")

Código en Python

```
#include <stdio.h>
int main(void) {
   printf("Hello, World!\n");
   return 0;
```

Código en C

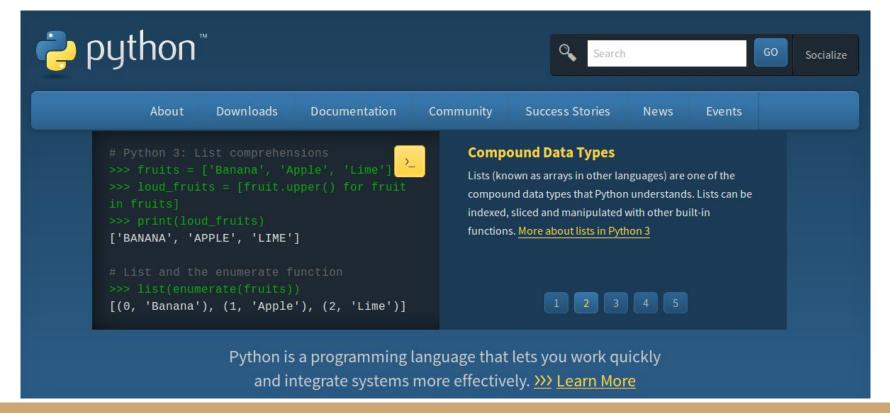
#### Consola o intérprete

```
ilia@ilia-Lenovo-G500:~$ python
Python 3.6.0 |Continuum Analytics, Inc.| (default, Dec 23 2016, 12:22:00)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-1)] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hola")
hola
>>> ■
```

#### https://repl.it/languages/python3

```
input ⊡
                                                         clear 🕶
                             Python 3.6.1 (default, Dec 2015, 13:05:11)
                             [GCC 4.8.2] on linux
print("hola")
                             hola
```

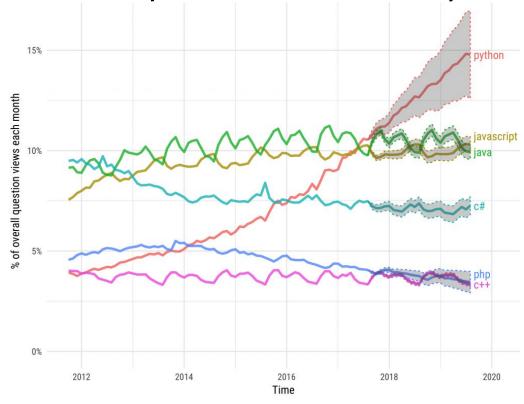
#### OpenSource - www.python.org



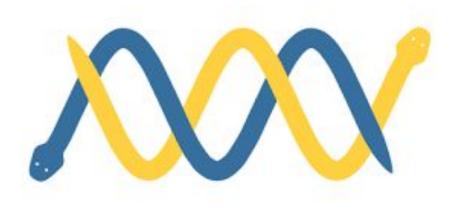
## Ranking de lenguajes según IEEE en 2017

1. Python	⊕ 🖵	100.0
<b>2.</b> C	□무:	99.7
3. Java		99.4
<b>4.</b> C++	□ 🖵 🛊	97.2
<b>5.</b> C#	$\oplus$ $\Box$ $\Box$	88.6
<b>6</b> . R	$\Box$	88.1
7. JavaScript		85.5
8. PHP		81.4
<b>9.</b> Go	⊕ 🖵	76.1
10. Swift		75.3

#### Predicción de tráfico en stack overflow



## http://biopython.org/



#### Un ejemplito

```
dna = 'AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAA'

nt = ['A', 'C', 'G', 'T']

for n in nt:
   nt_total = dna.count(n)
   print('{}: {}'.format(n, nt_total))
```

#### ¿Puedo tener más interactividad?



#### http://jupyter.org/

- Interactividad
- Más lenguajes (+40): R, Perl, Octave...
- Compartir notebooks
- Incrustar imágenes, vídeos y componentes para manipular los datos

# https://try.jupyter.org/



#### Compartir y colaborar

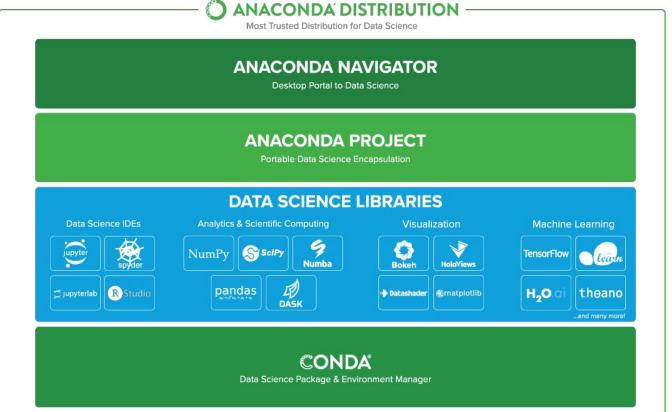




# Aplicación

- Libros
- Cursos
- Talleres

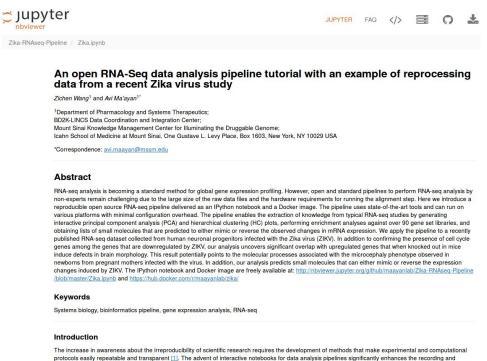
## https://www.anaconda.com/distribution/



## Compartir y Colaborar

Ciencia abierta

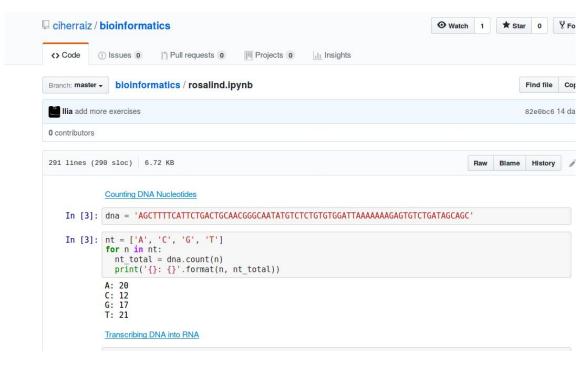
#### Artículo abierto sobre RNA



http://nbviewer.jupyter.org/github/maayanlab/Zika-RNAseg-Pipeline/blob/master/Zika.jpynb

sharing of data, source code, and figures [2]. In a subset of recent publications, an interactive notebook was published alongside customary manuscripts [3].

#### Ejercicios resueltos de Rosalind



#### ¡Muchas gracias!

Déjanos tu opinión: https://tinyurl.com/yc58omy7

#### Dónde encontrarnos



t.me/bioinformaticsGRX



https://github.com/bioinformaticsGRX



@bioinfGRX



bioinformaticsgrx.slack.com *(con invitación)*