

Una solución integrada con R para el Análisis de Interacciones entre genes con datos de Supervivencia en un estudio GWAS

Jesús Herranz, Antoni Picornell, María L. Calle, **Núria Malats**

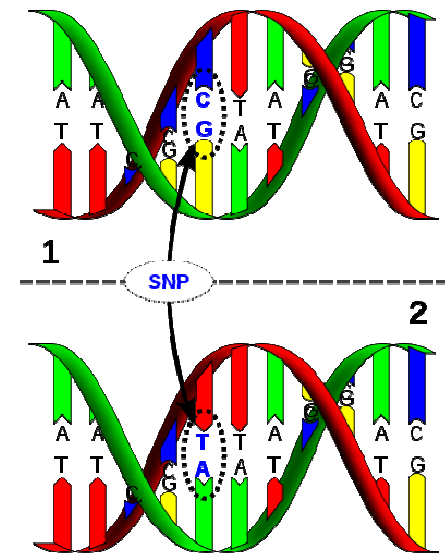
Epidemiología Genética y Molecular – CNIO

III Jornadas de Usuarios de R - 17 Noviembre 2011



Introducción

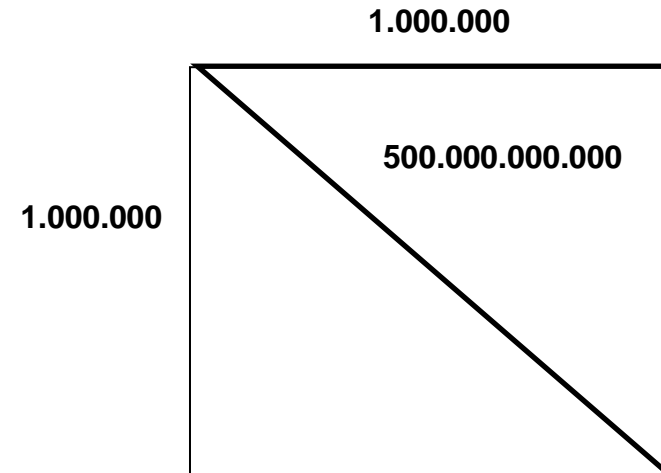
- Estudiar **variantes genéticas** implicadas en la **progresión** de enfermedades complejas
 - Factores pronóstico – Técnicas de **análisis de supervivencia**
- **Polimorfismos** – SNPs (Single Nucleotide Polymorphism)
 - Cambios en el genoma de 1 sola base
 - Genotipos:
 - Homocigotos comunes
 - Heterocigotos
 - Homocigotos variantes
- **GWAS** (Genome-wide Association Studies)
 - Estudios pangenómicos (genoma completo)
 - 100.000 - 1 Millón de SNPs



Introducción

- **Interacciones gen-gen (pares)**

- Alto coste computacional
- Miles de millones de interacciones



- **Número de Publicaciones - GWAS**

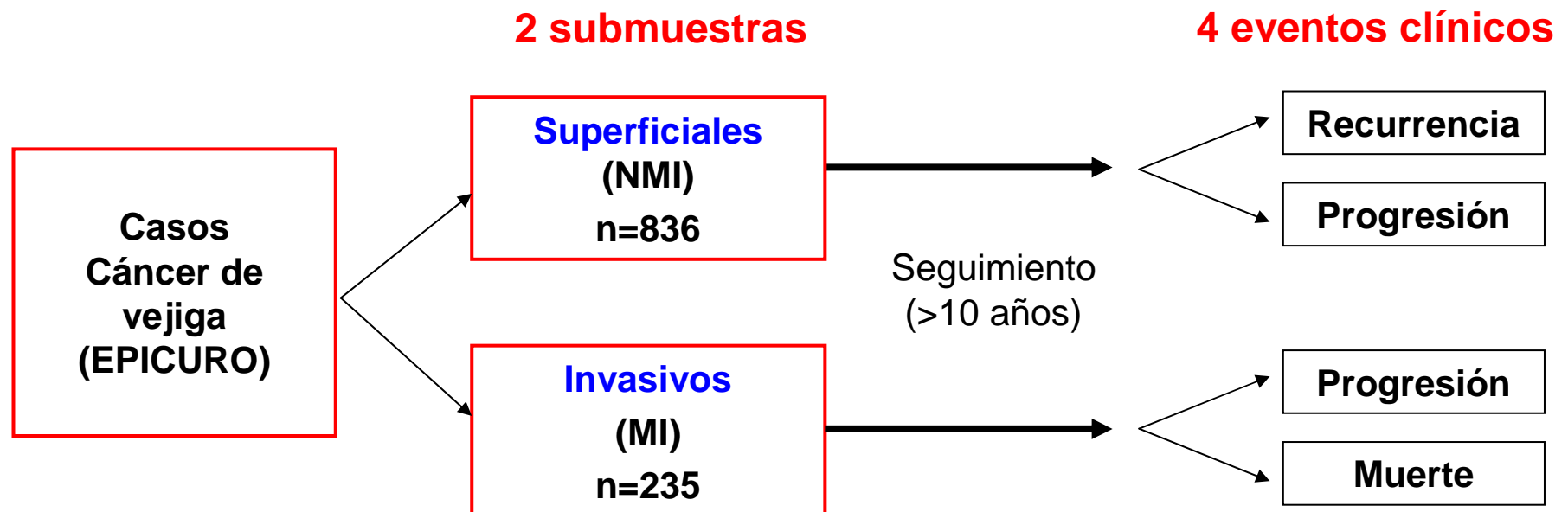
	Univariantes	Interacciones
Fac. Riesgo	> 1000	< 10
Fac. Pronóstico	< 10	0

**Interacciones gen-gen
GWAS
Pronóstico - Supervivencia**

- **Objetivo: Estrategia para analizar todos los pares de interacciones de un estudio pangenómico que incluye 1 millón de SNPs con datos de supervivencia.**

Estudio Español Cáncer Vejiga / EPICURO

- Estudios multicéntrico realizado en 18 hospitales
 - 1219 pacientes de **cáncer de vejiga** (reclutamiento: 1998 – 2001)
 - 1271 controles – Estudio casos y controles – Factores riesgo
 - **1071 casos de cáncer de vejiga con información clínico-patológica, información genética y seguimiento**





Una solución integrada con R

- **Control de las 2 submuestras**

- Control de individuos y variables en cada análisis
- No duplicar información**

- **R scripts únicos**

- Análisis simultáneo de los 4 estudios**
- Parámetros de entrada
- Ficheros de salida con los resultados

- **Funciones con los análisis estadísticos**

- Incorporar nuevas técnicas estadísticas o modificaciones



Etapas del análisis

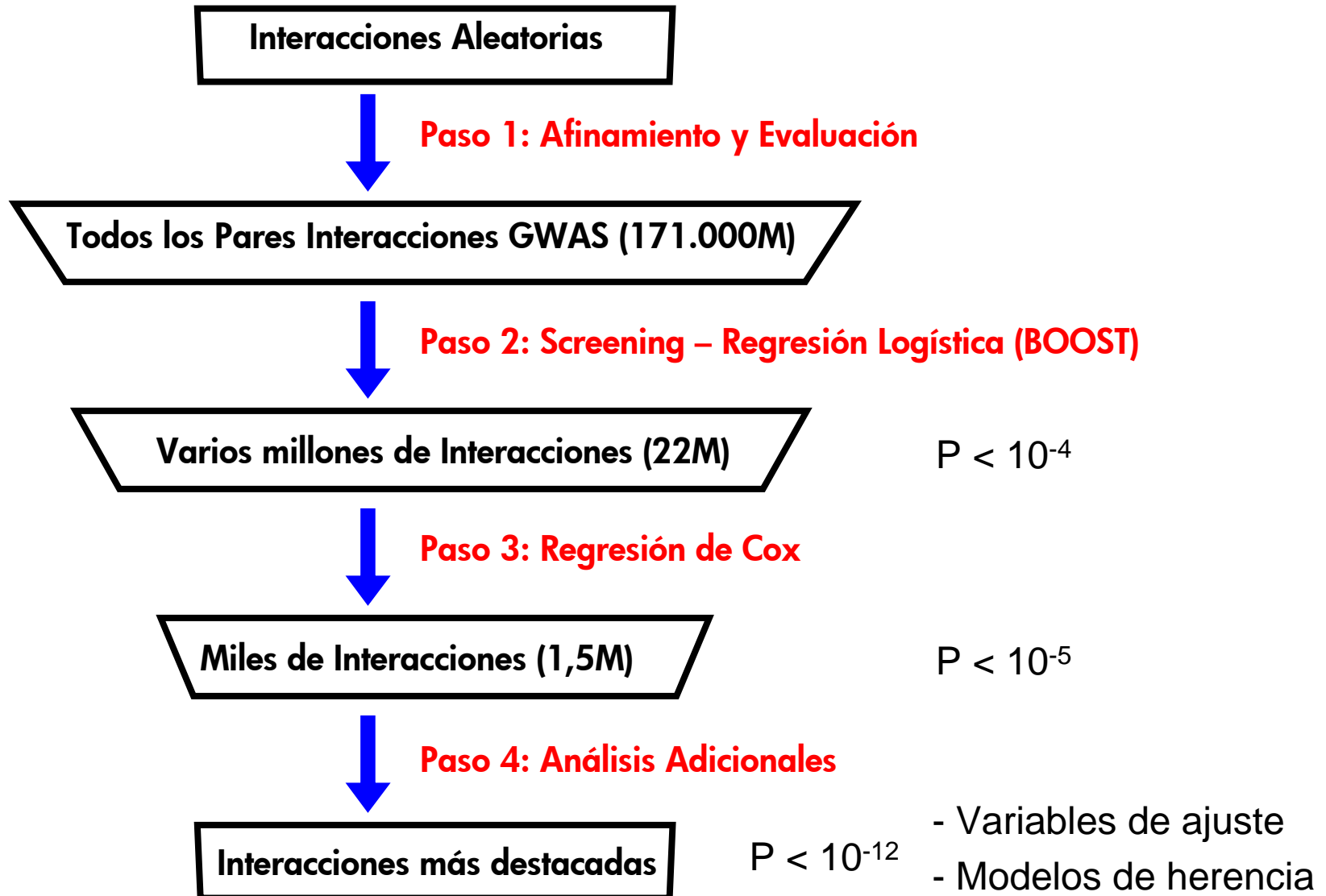
- **Preparación de los datos.** Etapas previas

1. Almacenamiento de los datos. Objetos ff
2. Imputación de missings
3. Criterios de inclusión de SNPs. Control de calidad
4. Reducción del número de variables para analizar

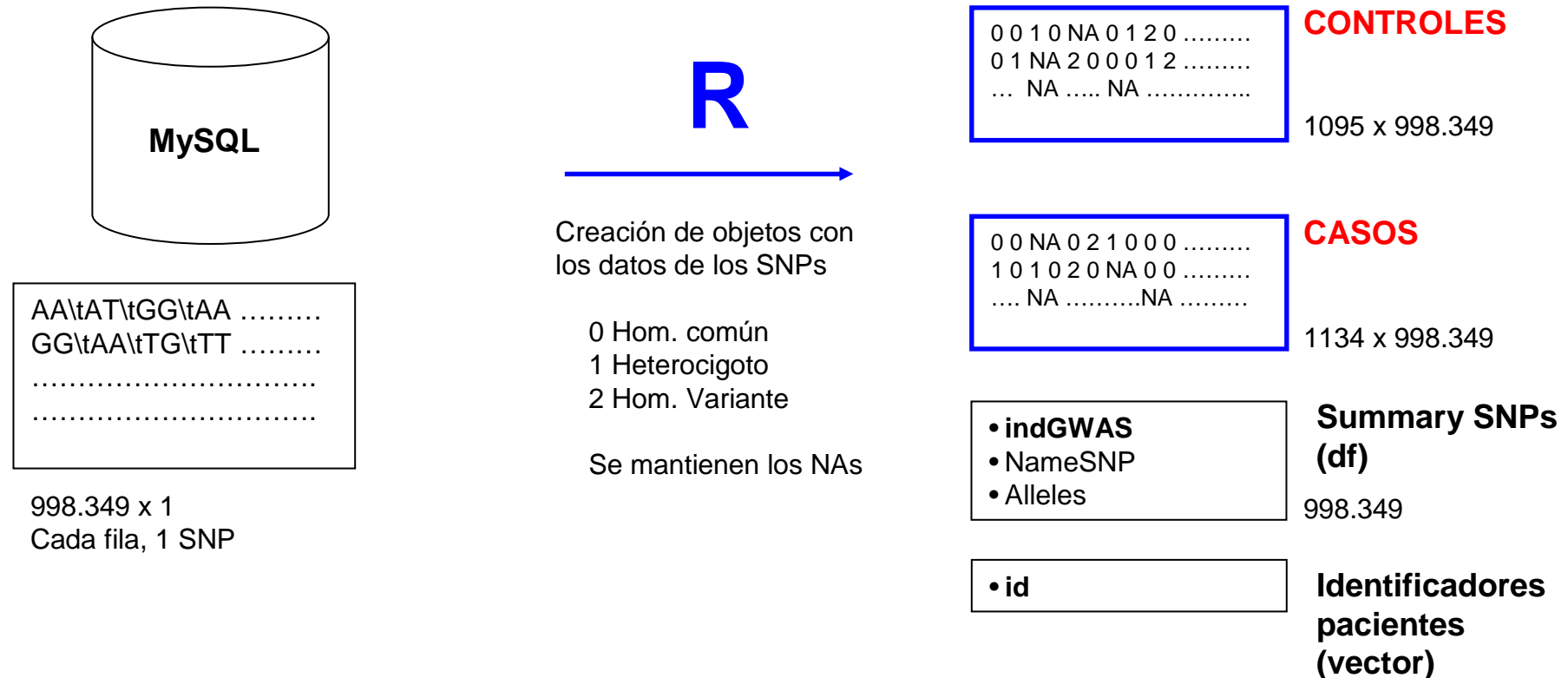
- **Estrategia analítica. Análisis Estadístico**

- No se pueden analizar 500.000 millones de interacciones con Regresión de Cox
- Etapa de screening con Regresión Logística (BOOST - C)
- Análisis posteriores con Regresión de Cox
- Alternativas: Survival - MDR (Multifactor Dimensionality Reduction)

Descripción de la Estrategia Analítica



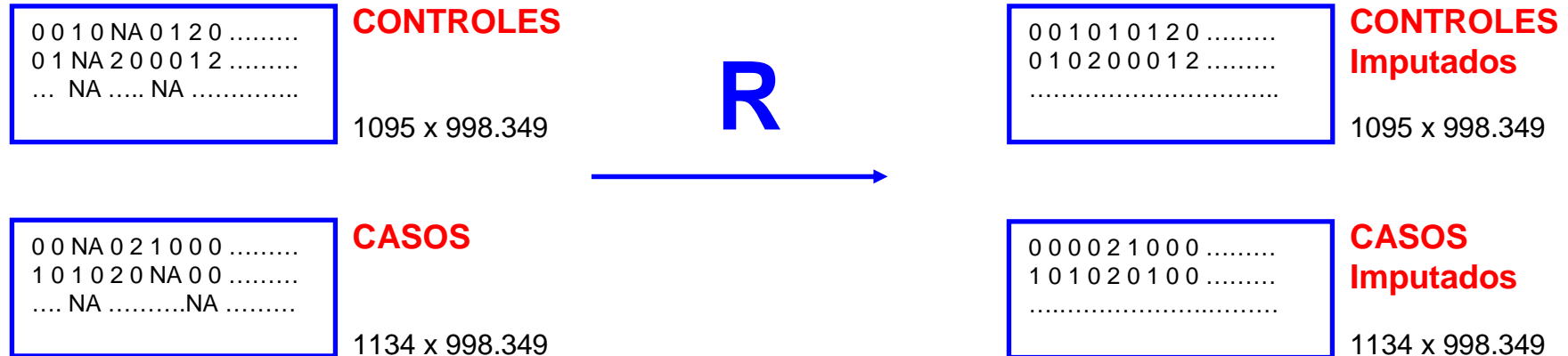
Etapa 1: Almacenar los datos del GWAS (ff)



■ La librería ff permite crear objetos

- Capacidad de almacenamiento, rápido acceso, no usa memoria de R
- Columnas: **nombres de los SNPs**
- Filas: **identificadores de los pacientes**

Etapa 2: Imputación de Missings



- Imputación por **Random Forests**
 - Ventana con los SNPs más próximos (correlacionados, LD)
- Otros métodos de imputación
 - **Mantener los datos originales sin imputar**
 - Crear otros ficheros de datos imputados

Etapa 3: Control de Calidad

```
0 0 1 0 NA 0 1 2 0 .....  
0 1 NA 2 0 0 0 1 2 .....  
... NA ..... NA .....
```

CONTROLES

1095 x 998.349

```
0 0 NA 0 2 1 0 0 0 .....  
1 0 1 0 2 0 NA 0 0 .....  
... NA ..... NA .....
```

CASOS

1134 x 998.349

R

P-HWE-ct > 0.00001
MAF > 0.02
Ht+HV > 10

Summary SNPs (df)

- indGWAS
- NameSNP
- Alleles
- Gen / Chr
- **Frec. Genotipos**
- **MAF**
- **Num. NAs**
- **HWE tests**
- **Included.NMI (880.000)**
- **Included.MI (840.000)**

998.349

- **SNPs eliminados del análisis**
 - **MAF (minor allele frequency):** poca variabilidad en la población
 - **HWE (Hardy-Weinberg equilibrium):** errores de genotipado
- **Otros criterios de inclusión**
 - No se recalcula MAF / HWE

Etapa 4: Linkage disequilibrium (LD)

```
0 0 1 0 NA 0 1 2 0 .....  
0 1 NA 2 0 0 0 1 2 .....  
... NA ..... NA .....
```

CONTROLES

1095 x 998.349

```
0 0 NA 0 2 1 0 0 0 .....  
1 0 1 0 2 0 NA 0 0 .....  
... NA ..... NA .....
```

CASOS

1134 x 998.349

R



LD < 0.9

Summary SNPs (df)

- indGWAS
- NameSNP
- Alleles
- Gen / Chr

- Frec. Genotipos
- MAF
- Num. NAs
- HWE tests

- **Included.NMI (880.000)**
- **Included.MI (840.000)**

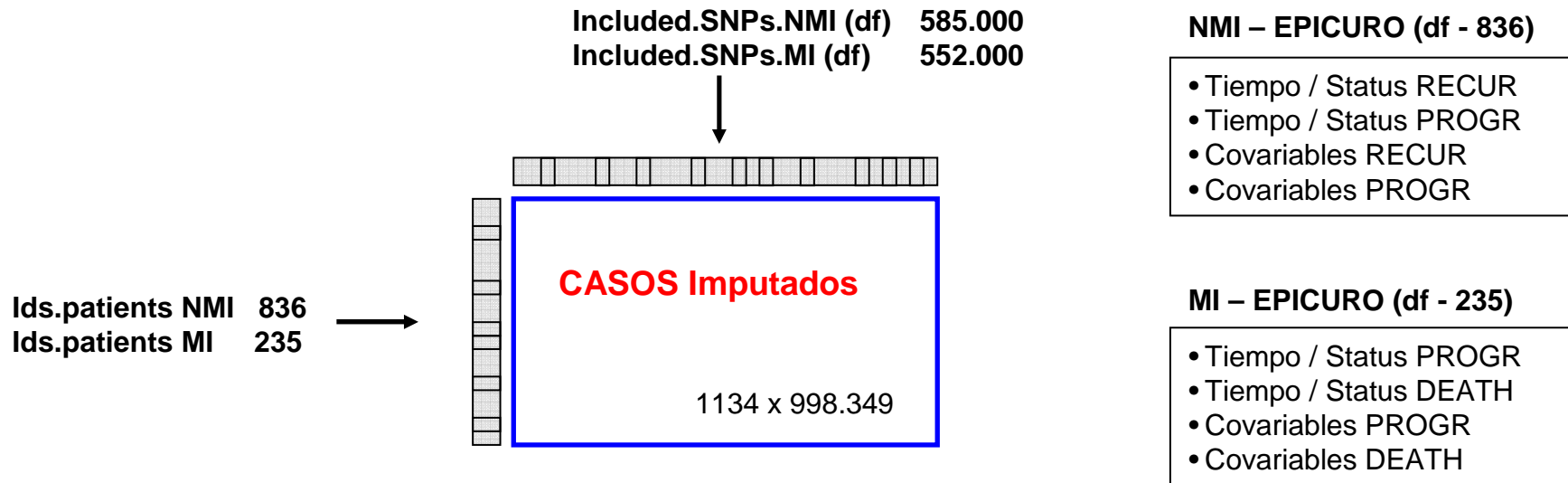
- **LD blocks**
- **LD represent (0/1)**

998.349

■ LD – SNPs correlacionados

- Se detectan **bloques de LD**
- Se detectan **singletons** (SNPs no correlacionados)
- Representantes del bloque:** entran en las fases de análisis de screening
- No representantes del bloque: entran en las fases finales (comp. múlt.)

Análisis simultáneos (4 estudios)



■ R Scripts únicos y flexibles

- Inclusión de otras **submuestras**
- Inclusión de otro **conjunto de SNPs** (criterios de inclusión, LD)
- Inclusión de otros **eventos clínicos** de interés (tiempo y status)
- Inclusión de otras **variables de ajuste**
- Parámetro de entrada:** estudio

Análisis estadístico (NMI – Progresión)

- NameSNP
- indGWAS
- indBOOST
- Modelo Herencia



BOOST NMI Progr

0	0011010120
0	1120000012
1	0000010210

836 x 585.000

Ids. NMI
836



Boost - C

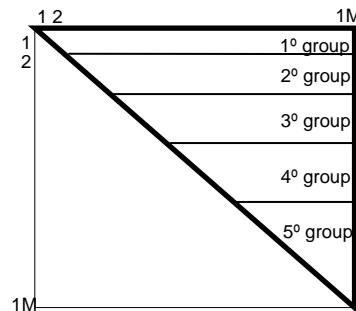


Regresión Logística
171.000 M interacciones
 $P < 10^{-4}$

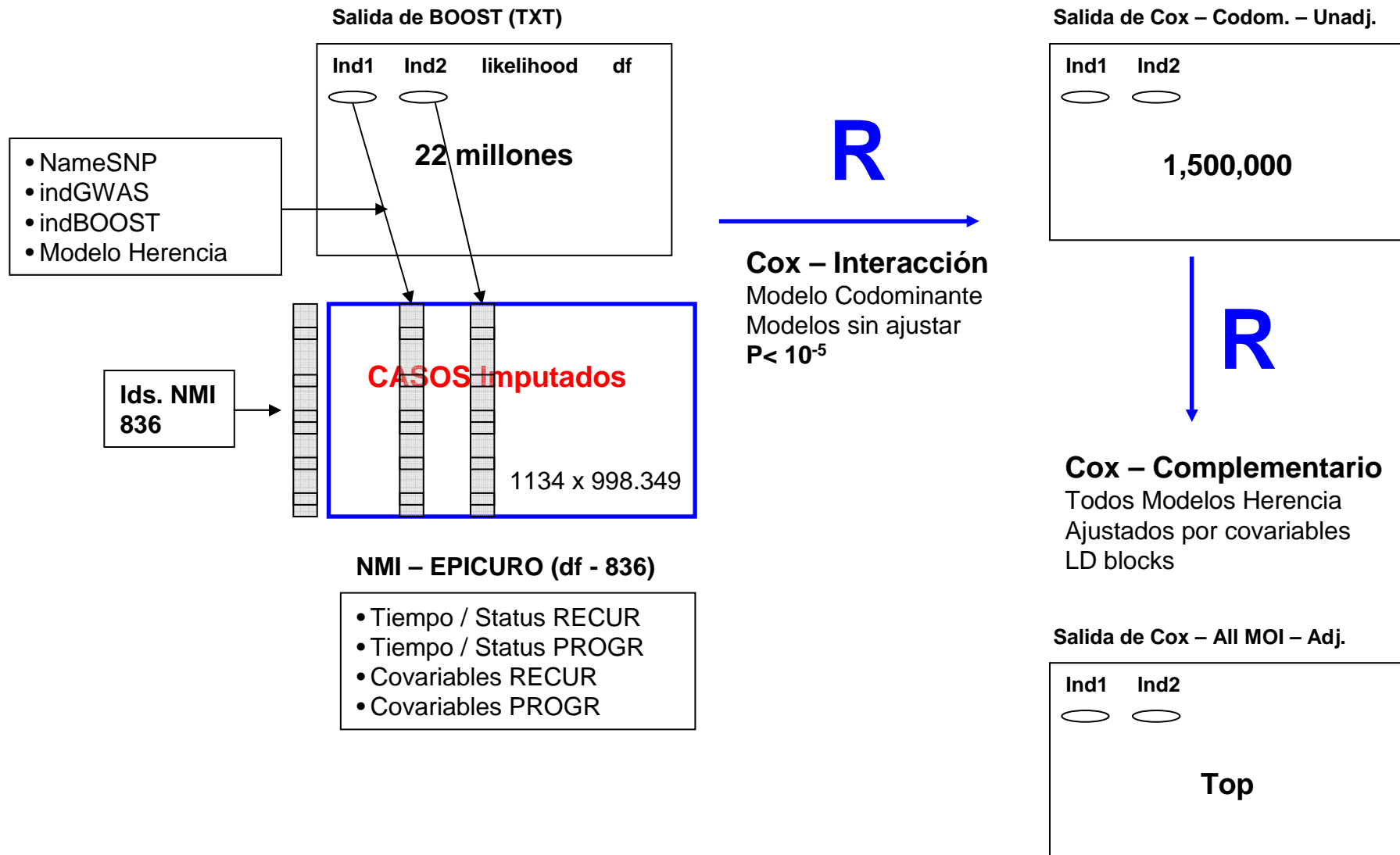
5 procesos (4 días)

Salida de BOOST (TXT)

Ind1	Ind2	likelihood	df
○	○		
22 millones			



Análisis estadístico (NMI – Progresión)



Resultados

Subphenotype	Outcome	Chr - Gen 1		Chr - Gen 2		MOI 1	MOI 2	P-Value	Threshold Informative SNF 8.60E-14	Threshold 0.4 - Factor 1.06E-13
NMI	Recurrence	16	<i>Gen1</i>	19	<i>Gen2</i>	A	A	5.53E-14	Sign.	Sign.
	Recurrence	14	<i>Gen3</i>	20	<i>Gen4</i>	R	A	7.31E-14	Sign.	Sign.
	Recurrence	11	<i>Gen5</i>	14	<i>Gen6</i>	D	C	1.91E-13	NS	NS
	Recurrence	10	<i>Gen7</i>	21	<i>Gen8</i>	D	D	1.94E-13	NS	NS
NMI	Progression	18	<i>Gen9</i>	X	<i>Gen10</i>	C	C	5.33E-15	Sign.	Sign.
	Progression	5	<i>Gen11</i>	14	<i>Gen12</i>	R	D	7.78E-14	Sign.	Sign.
	Progression	7	<i>Gen13</i>	10	<i>Gen14</i>	A	C	2.46E-13	NS	NS

Subphenotype	Outcome	Chr - Gen 1		Chr - Gen 2		MOI 1	MOI 2	P_Value	Threshold Informative SNF 9.46E-14	Threshold 0.4 - Factor 1.17E-13
MI	Progression	2	<i>Gen15</i>	8	<i>Gen16</i>	A	D	7.54E-14	Sign.	Sign.
	Progression	2	<i>Gen17</i>	16	<i>Gen18</i>	D	R	1.66E-13	NS	NS
MI	Death	8	<i>Gen19</i>	12	<i>Gen20</i>	D	D	4.11E-14	Sign.	Sign.
	Death	7	<i>Gen21</i>	10	<i>Gen22</i>	C	C	7.24E-14	Sign.	Sign.
	Death	7	<i>Gen23</i>	14	<i>Gen24</i>	R	D	9.87E-14	Sign.	Sign.
	Death	17	<i>Gen25</i>	22	<i>Gen26</i>	C	D	1.17E-13	NS	Sign.
	Death	1	<i>Gen27</i>	3	<i>Gen28</i>	D	D	1.56E-13	NS	NS

- Modelos de Cox **ajustados por covariables clínicas**
- Distintos **modelos de herencia** genética (dominante, recesivo, aditivo, codominante)
- Significación basada en **comparaciones múltiples**



Conclusiones

- **Estrategia analítica novedosa y viable**
 - Exhaustivamente todas los millones de Interacciones gen-gen
 - Análisis de supervivencia – Pronóstico – GWAS
 - Tiempo de computación aceptable (15-20 días)
- Hemos encontrado **varias interacciones gen-gen estadísticamente significativas**

Limitaciones

- **Interacciones perdidas** (con tamaños muestrales bajos)
- Los resultados deberían ser **replicados**

Agradecimientos



- Nuria Malats
- Toni Picornell
- Roger Milne
- Evangelina López
- Gaëlle Marene
- Mirari Márquez
- André Amaral
- Salman Tajuddin
- Matt Czachorowski



- Francisco X Real (CNIO)
- Stephen Chanock (NCI)
- Nathaniel Rothman (NCI)
- M. García-Closas (NCI)
- Debra Silverman (NCI)
- María L. Calle (Unv. Vic)
- Manolis Kogevinas (CREAL)