

Bieito Beceiro, Jorge González-Domínguez, Juan Touriño
CITIC & Grupo de Arquitectura de Computadores, UDC

*“Aceleración dun algoritmo de
selección de características usando
computación de altas prestacións”*

III Congreso **XoveTIC** Talento científico



1. Introducción

Contexto e impacto do problema

2. Desenvolvemento

Deseño e implementación da solución

3. Probas de rendemento

Estudo da eficacia da solución proposta

4. Conclusións

Impacto dos resultados

Contexto

Disciplinas en auge na actualidade



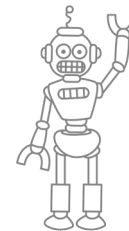
VIVENDAS
INTELIXENTES



VISIÓN
ARTIFICIAL



PROCESAMENTO DE
LINGUAXE



ROBÓTICA



CONDUCIÓN
AUTÓNOMA



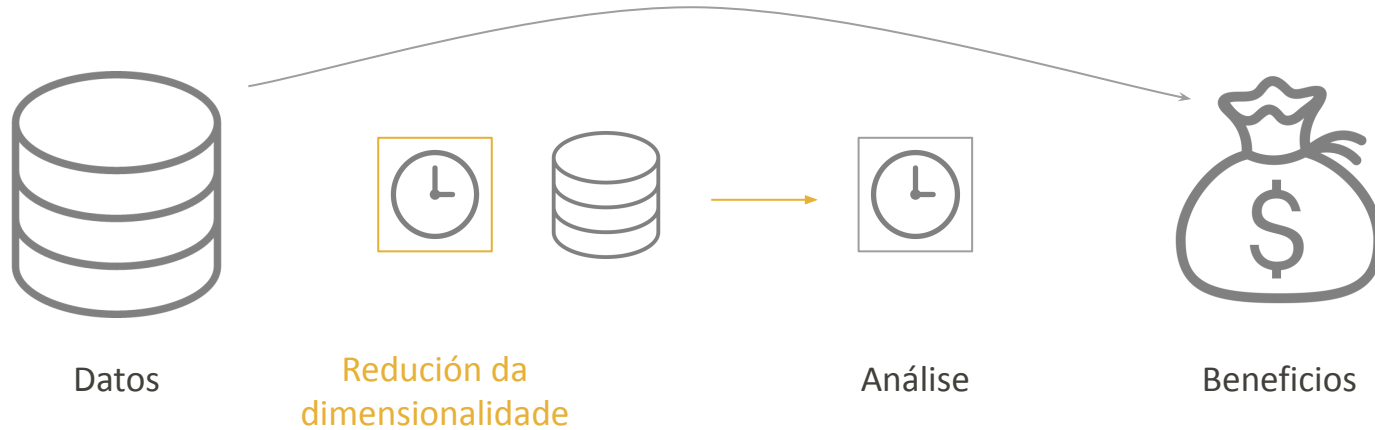
SAÚDE



MÁRKETING E
BUSINESS INTELLIGENCE

Contexto

Reducción da dimensionalidade



Contexto

FEAST: A Feature Selection Toolbox

- Selección de características mediante filtrado
 - Obtemos o subconxunto que mellor describe ao dataset orixinal como paso previo á análise
- Ganancia de información como métrica
 - Maximizar relevancia intentando reducir redundancia
- **Método JMI** (*Joint Mutual Information*)
 - Mellores resultados respecto a precisión, estabilidade e flexibilidade
 - Custo computacional moi alto
 - Potenciado por problemas “*big p, small n*”, cada vez máis habituais

1. Introducción

Contexto e impacto do problema

2. Desenvolvemento

Deseño e implementación da solución

3. Probas de rendemento

Estudo da eficacia da solución proposta

4. Conclusións

Impacto dos resultados

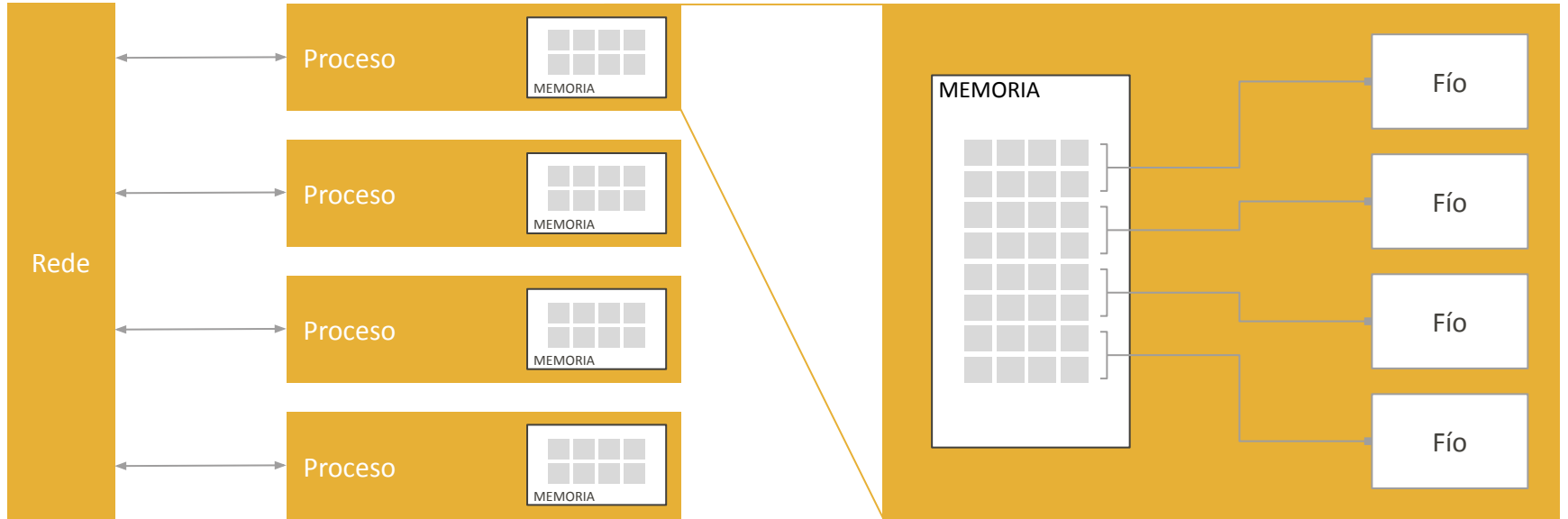
Desenvolvemento

Visión xeral

- Implementación paralela híbrida
 - Paralelismo explícito e implícito con descomposición de dominio
- Lectura semi-distribuída
- Compresión de rango

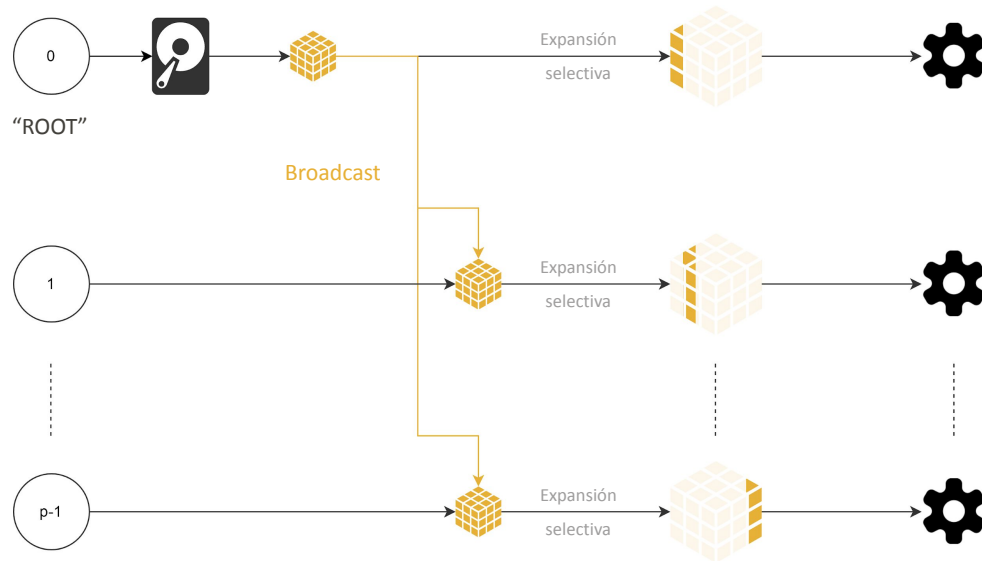
Desenvolvemento

Paralelismo explícito e implícito



Desenvolvemento

Lectura semi-distribuída



Desenvolvemento

Compresión de rango

Característica orixinal

2000	1000	1000	0	1000
------	------	------	---	------

Rango: $[0, 2000]$

2001 valores posibles
3 valores utilizados

2000	→	0
1000	→	1
0	→	2

Característica transformada

0	1	1	2	1
---	---	---	---	---

Rango: $[0, 2]$

3 valores posibles
3 valores utilizados

1. Introducción

Contexto e impacto do problema

2. Desenvolvemento

Deseño e implementación da solución

3. Probas de rendemento

Estudo da eficacia da solución proposta

4. Conclusións

Impacto dos resultados

Probas de rendemento

Entorno de probas

- Clúster privado de 16 nodos
 - CPU: 2 x octa-core con Hyperthreading (16 núcleos / 32 fíos lóxicos)
 - RAM: 64 GB
 - **En total:** 256 núcleos (ata 512 fíos lóxicos) e 1024GB de memoria RAM
- Uso de datasets representativos
 - Ratio de mostras-características
 - Bi-clase ou multiclase
 - E un dataset de gran tamaño (~512GB)

Probas de rendemento

Resultados

- Método *Joint Mutual Information*
 - Aceleracións de ata 198x para 256 elementos de procesado (eficiencia > 75%)
- Lectura semi-distribuída
 - Posibilidade de analizar un dataset que previamente era imposible
- Compresión de rango
 - Aforro de memoria e tempo de cómputo de JMI cun preprocesado mínimo

Probas de rendemento

Nun caso real...

Seleccionar 200 características de News20

FEAST-JMI

~30 minutos

~a duración de 3 exposicións

Parallel-JMI

~10 segundos

~o tempo que tardei en
dicir esta oración

Probas de rendemento

Nun caso real...

Seleccionar 200 características de E2006

FEAST-JMI

imposible

o dataset non cabe na memoria do nodo

Parallel-JMI

~13 minutos

1. Introducción

Contexto e impacto do problema

2. Desenvolvemento

Deseño e implementación da solución

3. Probas de rendemento

Estudo da eficacia da solución proposta

4. Conclusións

Impacto dos resultados

Conclusións

- O algoritmo forma parte dunha librería *open-source* que se publicará proximamente
- Boa aceleración e escalabilidade
- Posibilidade de analizar datasets de tamaño superior á memoria dun nodo
- Moi bos resultados para a optimización de compresión de rango

III Congreso **XoveTIC** Talento científico

Grazas!

Bieito Beceiro Fernández

Outubro 2020

