

O quarto paradigma: a ciência da análise de dados

Bruno Ferrero¹ Gabriel Pato¹

¹Instituto de Matemática e Estatística

Universidade de São Paulo

23 de Novembro de 2011

Roteiro

O 4º Paradigma

Dados e as ciências climáticas

Dados e as ciências biológicas

Conclusão

O quarto paradigma: e-Science

blog *The Fourth Paradigm*

- ciência e quebra de paradigmas;
- teóricos;
- ciência empírica;
- simulações;
- e-Science;

O quarto paradigma: e-Science

blog *The Fourth Paradigm*

- ciência e quebra de paradigmas;
- teóricos;
- ciência empírica;
- simulações;
- e-Science;

O quarto paradigma: e-Science

blog *The Fourth Paradigm*

- ciência e quebra de paradigmas;
- teóricos;
- ciência empírica;
- simulações;
- e-Science;

O quarto paradigma: e-Science

blog *The Fourth Paradigm*

- ciência e quebra de paradigmas;
- teóricos;
- ciência empírica;
- simulações;
- e-Science;

O quarto paradigma: e-Science

e-Science

- uso intensivo de técnicas computacionais;
 - ex. novas técnicas de visualização de dados;
- grande volume de dados ;
- grupos de pesquisadores (diversos colaboradores);
- ferramenta X paradigma;

Dados e o clima

- estudar o clima → analisar dado;
- duas fontes de dados:
 - observacional;
 - dados numéricos;
- dado é caro!!;

Dados climáticos

- estudar o clima → analisar dado;
- duas fontes de dados:
 - observacional;
 - dados numéricos;
- dado climático observado é caro!!;
- avanço da computação → mais dados numéricos;

Dados climáticos

- estudar o clima → analisar dado;
- duas fontes de dados:
 - observacional;
 - dados numéricos;
- dado climático observado é caro!!;
- avanço da computação → mais dados numéricos;

Dados climáticos

- estudar o clima → analisar dado;
- duas fontes de dados:
 - observacional;
 - dados numéricos;
- dado climático observado é caro!!;
- avanço da computação → mais dados numéricos;

Dados climáticos

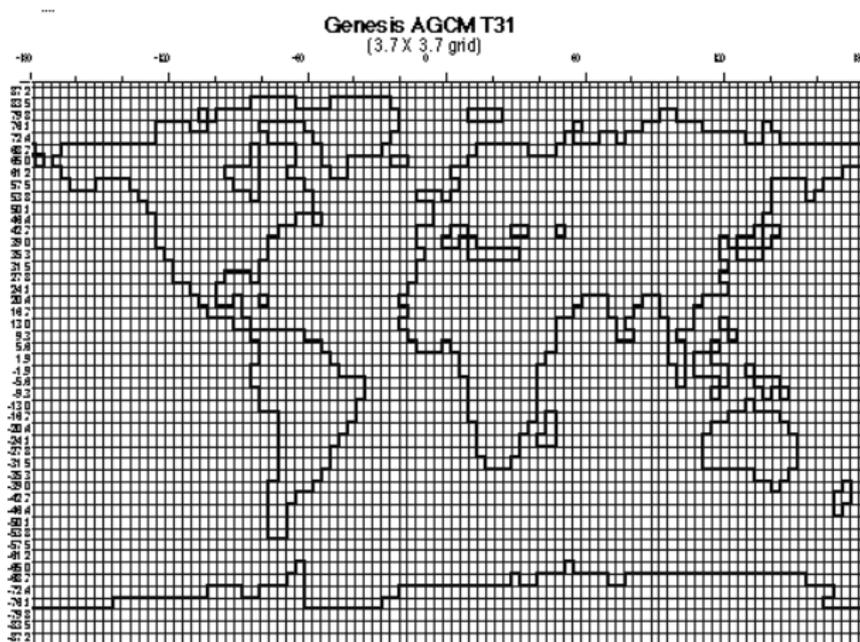
- estudar o clima → analisar dado;
- duas fontes de dados:
 - observacional;
 - dados numéricos;
- dado climático observado é caro!!;
- avanço da computação → mais dados numéricos;

Simulação numérica do clima

- reproduzir a dinâmica do clima por meio de modelos matemáticos;
- alto poder de processamento;
- alta capacidade de armazenamento;
- super computadores, ex. *Earth Simulator*;
- mais processamento → modelos mais complexos → simulação mais próxima da realidade;

Simulação numérica do clima

gera um enorme volume de dados



- 110 variáveis climáticas;
- 360 x 180 x 30 resolução espacial (x, y, z);
- 6 horas frequência temporal

Simulando um período de 100 anos:

$$110 \times (4 \times 365 \times 100) \times (360 \times 180 \times 30) = \text{Muito dado!!}$$

- 110 variáveis climáticas;
- 360 x 180 x 30 resolução espacial (x, y, z);
- 6 horas frequência temporal

Simulando um período de 100 anos:

$$110 \times (4 \times 365 \times 100) \times (360 \times 180 \times 30) = \text{Muito dado!!}$$

- 110 variáveis climáticas;
- 360 x 180 x 30 resolução espacial (x, y, z);
- 6 horas frequência temporal

Simulando um período de 100 anos:

$$110 \times (4 \times 365 \times 100) \times (360 \times 180 \times 30) = \text{Muito dado!!}$$

- 110 variáveis climáticas;
- 360 x 180 x 30 resolução espacial (x, y, z);
- 6 horas frequência temporal

Simulando um período de 100 anos:

$$110 \times (4 \times 365 \times 100) \times (360 \times 180 \times 30) = \text{Muito dado!!}$$

- 110 variáveis climáticas;
- 360 x 180 x 30 resolução espacial (x, y, z);
- 6 horas frequência temporal

Simulando um período de 100 anos:

$$110 \times (4 \times 365 \times 100) \times (360 \times 180 \times 30) = \text{Muito dado!!}$$

Organização de dados

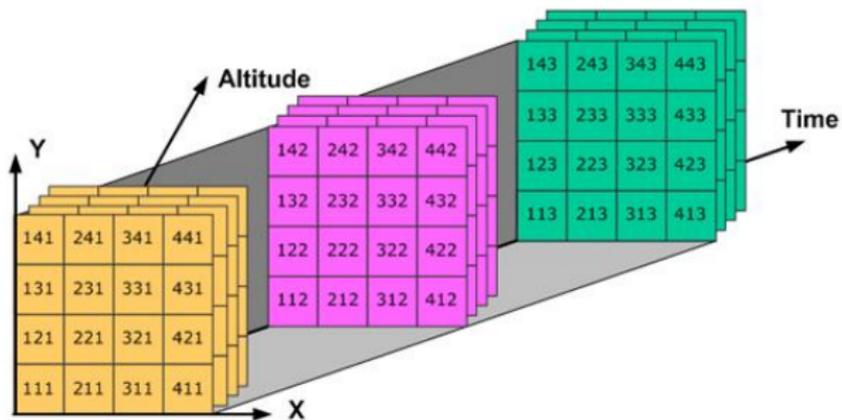
Características

- portabilidade;
- auto descritivo;
- padronização;

netCDF

- *software livre*
- auto descritivo;
- portabilidade
- binário com um cabeçalho ASCII;
- convenção;
- *gridded data*
- trabalha com múltiplas dimensões, $(x, y, z, time)$;

Gridded data



Exemplo

Panoply File Edit View Bookmarks Plot Window Help

Sources

23 nov 12:05 bruno

Remove Remove All Help Docs

Name	Long Name	Type
ccsm3.nc	ccsm3.nc	Remote File
height	height	—
huss	specific_humidity	[lon][lat][time]
lat	latitude	—
lat_bnds	lat_bnds	—
lon	longitude	—
lon_bnds	lon_bnds	—
plec	pressure	—
psl	air_pressure_at_sea_level	[lon][lat][time]
tas	air_temperature	[lon][lat][time]
tauv	surface_downward_eastward...	[lon][lat][time]
tauw	surface_downward_northward...	[lon][lat][time]
time	time	—
time_bnds	time_bnds	—
ua	eastward_wind	[lon][lat][event][time]
va	northward_wind	[lon][lat][event][time]

File "ccsm3.nc"

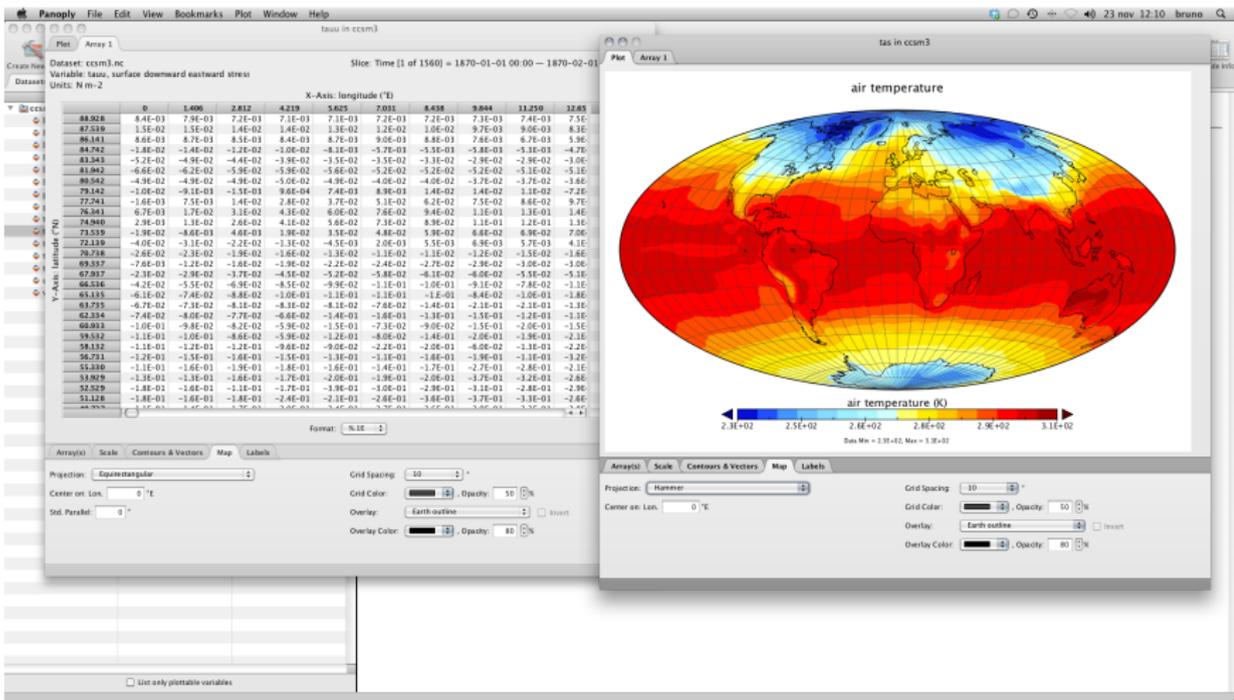
```

netcdf $@ds://1/vambd.Leg.usp.br:8986/thread/dodsC20c3m/ata/m0/ccsm3.nc {
dimensions:
    lon = 256;
    bnds = 2;
    lat = 128;
    time = 1548;
    plec = 17;
variables:
    double lon_bnds(lon=254, bnds=2);
    double lat_bnds(lat=128, bnds=2);
    double time_bnds(time=1548, bnds=2);
    float tauv(time=1548, lat=128, lon=256);
    {
        _CoordinateAxes = "time lat lon ";
        comment = "Created using NCL code CCM3_atm Jef.nc on machine mineral.epg.ucar.edu";
        missing_value = 1.0E20; // float
        _FillValue = 1.0E0f; // float
        cell_methods = "time; mean (interval: 1 month)";
        history = "No change";
        original_name = "tauv";
        standard_name = "surface_downward_eastward_stress";
        units = "N m-2";
        long_name = "surface downward eastward stress";
        cell_method = "time; mean";
        float tauw(time=1548, lat=128, lon=256);
    {
        _CoordinateAxes = "time lat lon ";
        comment = "Created using NCL code CCM3_atm Jef.nc on machine mineral.epg.ucar.edu";
        missing_value = 1.0E20; // float
        _FillValue = 1.0E0f; // float
        cell_methods = "time; mean (interval: 1 month)";
        history = "No change";
        original_name = "tauw";
        standard_name = "surface_downward_northward_stress";
        units = "N m-2";
        long_name = "surface downward northward stress";
        cell_method = "time; mean";
    float huss(time=1548, lat=128, lon=256);
    {
        _CoordinateAxes = "height time lat lon ";
        comment = "Created using NCL code CCM3_atm Jef.nc on machine mineral.epg.ucar.edu";
        missing_value = 1.0E20; // float
        _FillValue = 1.0E0f; // float
        cell_methods = "time; mean (interval: 1 month)";
        note = "Use 'huss' with caution. There are incorrect values over Antarctica in January. This data is from the CCM3 model coupler code, not the atmospheric model";
        history = "Added height coordinate";
        coordinates = "height";
        original_name = "kg/kg";
        standard_name = "q00000";
        units = "kg kg-1";
        long_name = "specific humidity";
        cell_method = "time; mean";
    float psl(time=1548, lat=128, lon=256);
    {
        _CoordinateAxes = "time lat lon ";
        comment = "Created using NCL code CCM3_atm Jef.nc on machine mineral.epg.ucar.edu";
        missing_value = 1.0E20; // float
        _FillValue = 1.0E0f; // float
        cell_methods = "time; mean (interval: 1 month)";
        history = "No change";
        original_name = "Pa";
        standard_name = "psl";
    }
}

```

List only portable variables

Exemplo



Painel Intergovernamental para as Mudanças Climáticas - IPCC

- cenários de desenvolvimento para fazer projeções climáticas;
- uso intensivo de modelos climáticos;
- milhares de pesquisadores em todo o mundo;
- avaliação colaborativa das simulações;

Painel Intergovernamental para as Mudanças Climáticas - IPCC

- cenários de desenvolvimento para fazer projeções climáticas;
- uso intensivo de modelos climáticos;
- milhares de pesquisadores em todo o mundo;
- avaliação colaborativa das simulações;

Dados e as ciências biológicas

- Genoma;
- Proteoma;
- Metaboloma;

Genoma e Proteoma

Tipos de dados gerados

- sequência de nucleotídeos;
- sequência de aminoácidos;
- dobramento de proteínas;

Genoma

- Sequenciamento genético;
 - sequenciadores:
 - Pyrosequencing → 40 milhões de pares de bases por hora;
 - banco de dados:
 - GenBank
- necessidade de tratamento/análise;

GenBank

Características

- criação: 1982
- acesso livre
- atualmente, mais de 140 milhões de sequências

Formatos de arquivos

- ASN.1
- GenBank
- outros ...

FEATURES	Location/Qualifiers
source	1..428
	/organism="Macaca mulatta"
	/mol_type="mRNA"
	/strain="Indian"
	/db_xref="taxon:9544"
	/clone="IBIUW:32275"
	/sex="female"
	/dev_stage="adult"
	/lab_host="Electromax DH10B"
	/clone_lib="Katze_MMOV"
	/note="Organ: ovary; Vector: pDONR 222; Site_1: BsrG I; Site 2: BsrG I; Created from CloneMiner cDNA Library Construction kit (catalog #18249-029)"
ORIGIN	
1	ttggctcttc tacctgcaac cgaatgcttg atgaagccac cagtgcctcg acagaggagg
61	tggagaatga gctctatcgc atcgccagc agctggggat gacgttcac agtgtgggac
121	atcgccagag ccttgagaag ttctattcct tcgttctgaa actctgtgga ggaggaagat
181	gggagctgat gagaatcaaa gtggaatgaa gctccagctt ttagaaggag agccacactc
241	tggagggtcg gcagccctca ggagtgacca ggaggactgg cggggaagat cgagctcagg
301	ttcgccacat aggtcctgtg caggagcctt ggcggtggtg ggctgagccc gggctctggt
361	ttctgtgggg gacactgagt ctcccagtg tcagtctccc aggactctgc tgccctagcc
421	agagcctc

FEATURES describes biological features related to the sequence.

Conclusões

- importância da organização de dados;
- ambiente de compartilhamento;
- levando a um sistema de produção colaborativa.

Obrigado!