

MANUAL DE INSTALAÇÃO  
DO SIPAC –  
SISTEMA DE INFORMAÇÃO  
DA PRODUÇÃO  
AGROPECUÁRIA CAPIXABA

Elaboração: Cristian da Silva Antério

Bolsista PIBIT/FAPES/INCAPER

2016

---

# Table of Contents

Introdução	1.1
Instalação	1.2
Transformações	1.3

# Introdução

O Pentaho é uma suíte com várias ferramentas disponíveis para Inteligência de Negócios, incluindo:

- **Pentaho Data Integration:** Também conhecido como Kettle, é uma ferramenta de código livre para extração, transformação e carga de dados.
- **Pentaho Analysis Service:** Conhecido também como Mondrian OLAP server, é uma ferramenta de código livre para Análise Online e Processamento (OLAP).
- **Pentaho Data Mining:** Derivado do projeto Weka, um conjunto de ferramentas relacionadas com a mineração de dados.
- **Pentaho Dashboards:** Utilizado para mostrar, de uma maneira fácil, os dados trabalhados.

Neste tutorial trabalharemos com o Pentaho Data Integration(PDI), para carregar os arquivos dos arquivos csv, manipulando-os e os inserindo em uma base de dados. O PDI tem uma interface muito amigável e é muito simples de se manipular.

<sup>1</sup>. ETL é um processo para extrair dados de um sistema de arquivos, transformá-lo em um formato de banco de dados e colocá-lo no banco de dados. ↩

# Instalação

A Utilização das ferramentas da suíte Pentaho é bastante simples. Para este tutorial, precisaremos apenas de duas delas: Pentaho Data Integration (PDI) e Pentaho Business Intelligence Server (BI-Server). Os arquivos estão disponíveis em:

<http://sourceforge.net/projects/pentaho/files/>

Após efetuar o download, precisamos fazer algumas configurações, antes de começar a utilizar a ferramenta:

- Instalação Java (Linux):
  - Precisamos instalar o Oracle Java, pois com o JDK, padrão do linux não é compatível com o Pentaho. Sendo assim, abra o terminal e adicione o repositório Java com os seguintes comandos:
    - `sudo add-apt-repository ppa:webupd8team/java`
    - `sudo apt-get update`
    - `sudo apt-get install oracle-java7-installer`
  - Feito isso, vamos configurar as variáveis de ambiente. Vá até o caminho **/usr/lib/jvm** e veja se o arquivo **java-7-oracle** se encontra lá. Caso esteja, volte até a raiz e edite o arquivo **/etc/environment** com o seguinte comando:
    - `sudo nano /etc/environment`
    - E adicione as seguintes linhas:
      - `JAVA_HOME=/usr/lib/jvm/java-7-oracle`
      - `JRE_HOME=/usr/lib/jvm/java-7-oracle`
    - Salve o arquivo, reinicie o sistema e pronto!
- Para os usuários do Windows, sugiro este tutorial de como instalar o Java:
  - <https://www.youtube.com/watch?v=pZI7hdepiy0>
- Vá na pasta onde você baixou e extraiu o BI-Server, abra-a no terminal e execute o arquivo **start-pentaho.sh** da seguinte forma:
  - `./start-pentaho.sh`
  - Depois disso, abra o navegador no seguinte endereço: <http://localhost:8080> e seja feliz!
  - **OBS.:** Login: Admin, Senha: password.
  - **OBS<sup>2</sup>.:** Iniciar o BI-Server era apenas para testar se estava tudo funcionando  
... :p

# Transformações

Com tudo certo, vamos até a pasta do Pentaho Data Integration (PDI) e execute o arquivo

```
spoon.sh .
```

Nesta parte iremos mostrar como carregar arquivos CSV para o banco de dados. Utilizaremos aqui o postgresql.

-- Colocar link para configuração do pgAdmin e postgres.

Com o **spoon.sh** aberto, você precisa se conectar à sua base de dados

- No Painel do lado esquerdo, vá na aba **View** (ao lado de Design). Clique com o botão direito em **Conexões** e depois **Novo**. Na tela que se abrirá, escolha:

```
Tipo de conexão:PostgreSQL
Host Name: localhost
Database Name: [O nome do banco que você criou]
User Name: [Seu login]
Password: [Sua senha]
```

Clique em **Test**, se estiver tudo certo, clique em OK.

Sua base de dados estará listada em **Conexões**. Clique com o botão direito nela e então clique em **Share** (Compartilhar), para que ela fique disponível para todos os seus futuros trabalhos no PDI.

Depois disso, podemos começar a brincar com alguns arquivos CSV's.

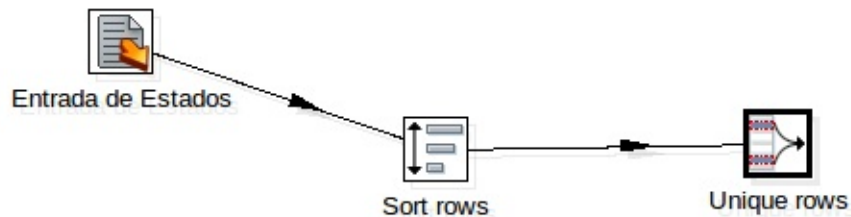
Neste primeiro exemplo, vamos utilizar o arquivo Cod\_Municipio.csv. O nosso objetivo é carregar todos os estados desse arquivo (sem repetição) e inserir em nossa base de dados. Para abri-lo, basta ir em **Arquivo > Novo > Transformações** (Ou apenas Ctrl + n).

Após isso, na aba **Design** no campo **Input**, procure por CSV input file e arraste-o para tela. Dê duplo clique neste ícone e em Step name, renomeie para Entrada de Estados (ou como você desejar). No campo abaixo, filename, navegue até a pasta que se encontra o arquivo Cod\_Municipio.csv.

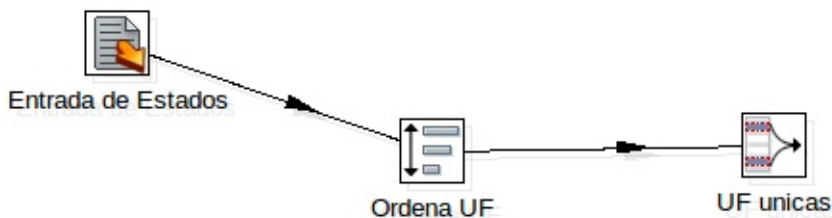
Em alguns arquivos CSV o delimitador é o ';' (ou um outro qualquer), neste caso, vamos mudar o campo delimiter para ';'. Logo após, clique em **Obtem Campos** e clique em OK. Você pode ver os dados carregados clicando em **Preview**. Após isso, vá em **Transform** e procure o ícone **Sort rows**. Nós iremos ordenar todos os campos do nosso arquivo, pois o passo seguinte o exige. Arraste-o para a tela, clique no passo Entrada de Estados

segurando a tecla shift e clique no passo atual (Sort Rows). Isso ligará os dois passos. Em seguida, dê duplo clique nele (redefine o nome, se quiser) e na opção **Fieldname**, seleciona o campo UF e OK.

Proximo passo é utilizar o Transform -> Unique Rows. Ligue-o com o passo anterior (Sort rows) e deverá ficar assim:



Para ficar mais organizado, vamos renomear 'Sort rows' para Ordena UF e Unique rows para UF únicas:



Duplo clique em UF únicas e no campo Fieldname escolha UF. Isso serve para que ele ignore as repetições que possivelmente encontrará.

O último passo é arrastar **Output -> Table output** para a tela. Table output será a nossa base de dados que conectamos no inicio desse tutorial. Renomeie esse passo para Tabela Estado. No campo Connection escolha a conexão que você criou e em Target table escolha a tabela **core\_estado**, que é onde iremos salvar nossos dados. Depois disso marque a opção **Specify database fields** e na aba **Database fields** faça o mapeamento como está abaixo:

The diagram shows a data pipeline with three components: 'Entrada de Estados' (Input States), 'Ordena UF' (Order States), and 'UF unicas' (Unique States). The 'UF unicas' component is connected to a 'Table output' step. Below this is a configuration window titled 'Saída a Tabela' (Output to Table).

**Saída a Tabela**

Nome do Step: *Table output*

Connection: *Incaper* [Edit... New... Wizard...]

Target schema: [ ] [Navega...]

Target table: *core\_estado* [Navega...]

Commit size: *1000*

Truncate table:

Ignore insert errors:

Specify database fields:

Main options Database fields

Colunas a inserir:

▲ #	Table field	Stream field
1	<i>codigo</i>	<i>UF</i>
2	<i>nome</i>	<i>Nome_UF</i>

[Get fields] [Enter field mapping]

Clique em **OK** e tudo pronto! Agora, para executar a transformação, basta apertar a tecla **F9**. Ou **Action > Run**