
CrowData Documentation

Release 0.1

Gabriela Rodriguez & Manuel Aristaran

June 02, 2014

1	When to use Crowdata?	3
1.1	Contents	3
2	Similar projects	7
3	Credits	9
4	Contributions	11
5	Indices and tables	13

CrowData is a tool to collaborate on the verification or release of data that otherwise would be hard or impossible to get via automatic tools.

When to use Crowdata?

The screenshot shows the 'Gastos del Senado' page on lanacion.com. The page features a header with the site logo and navigation links. The main content area includes a large image of the Argentine Senate chamber, a text block inviting users to collaborate on a database of public spending documents from 2010-2012, and a progress bar showing 6657 documents reviewed. There is also a 'Liberá un documento' button and social sharing options.

- VozData is a website from La Nacion in Argentina to convert scanned PDF documents from senate spendings into an usable dataset. Collaborating to free data from PDFs.

1.1 Contents

1.1.1 Technical Introduction

The basic features for 'Crowdata' are

- Store a set of documents (PDF or other formats supported by Document Cloud)
- Define a form, via the admin, for the information that wants to be extracted from the documents.
- Present users with a document and the form to allow anybody, that is registered to the website, to send us the information they see in the document.

- Have access to download the CSV of all the information extracted from the PDFs by the users.

1.1.2 How To Install It Locally

1. Python 2.7.5
2. We recommend the use of *virtualenv* <<http://virtualenv.org>> — Install it.
3. Create a virtual environment and activate it:

```
$ virtualenv ~/.python-envs/crowdata
$ . ~/.python-envs/crowdata/bin/activate
```

4. Get the source code:

```
$ git clone https://github.com/crowdata/crowdata.git crowdata
$ cd crowdata
```

5. Install dependencies:

```
$ pip install -r requirements.txt
```

(If you are using Ubuntu, you may need to install *python-dev* before dependencies.)

6. Create PostgreSQL database:

```
$ createuser -s -h localhost crow_user
$ createdb -O crow_user -h localhost crowdata_development
```

7. Create extensions for doing *trigram matching* <<http://www.postgresql.org/docs/9.2/static/pgtrgm.html>> and *removing accents* <<http://www.postgresql.org/docs/9.1/static/unaccent.html>> in PostgreSQL:

```
$ psql -ucrow_user
crow_user=# \c crowdata_development
crowdata_development=# CREATE EXTENSION pg_trgm;
crowdata_development=# CREATE EXTENSION unaccent;
```

*Note: In Debian/Ubuntu you need to install *postgresql-contrib-9.1* and *geospatial libraries*.*

8. We keep local settings out of GIT. You will need to copy *local_settings.py.example* to *local_settings.py*. You will need to edit the database settings there.:

```
DATABASES = {
    'default': {
        'ENGINE': 'django.db.backends.postgresql_psycopg2', # Add 'postgresql_psycopg2', 'postgr
        'NAME': 'crowdata_development', # Or path to database file if using
        'USER': 'crow_user',
        'PASSWORD': '',
        'HOST': '',
        'PORT': '',
    }
}
```

9. Initialize the database:

```
$ python manage.py syncdb
$ python manage.py migrate --all
```

10. Start the development server:

Similar projects

Crowdata was inspired in the project from ‘ProPublica <<http://www.propublica.org>>’_ called ‘Free the Files <[It was born from a need that La Nacion had to transform scanned image PDFs into a comprehensible and structured dataset, and also to ask for community’s help to catalog those spendings that call their attention.](https://projects.p</p></div><div data-bbox=)

Here some of the projects that do the same for some specific cases.

- [Free the Files](#)
- [Yanukovych Leaks](#)
- [How to crowdsource MPs’ expenses](#)

Credits

‘Crowdata’ is an open source project that was born when Manuel Aristaran was an Open News fellow at La Nacion in 2013. It was finally released as free software when Gabriela Rodriguez continued it for VozData in 2014. Thanks to Cristian Bertelegni and La Nacion for contributing to the code.

Now it relies on contributions from people and organizations. Please, use it, comment on it and make improvements by pull requests in [GitHub](#).

Contributions

- Fork the repo
- Clone your fork
- Make a branch of your changes
- Make a pull request through GitHub, and clearly describe your changes

Indices and tables

- *genindex*
- *modindex*
- *search*