
django-data-migration Documentation

Release 0.2.1

Philipp Böhm

September 07, 2015

1	Contents:	3
1.1	Installing	3
1.2	Writing Migrations	3
1.3	Migrating your data	11
1.4	Troubleshooting	11
1.5	Changelog	11
	Python Module Index	13

`django-data-migration` is a reusable Django app that migrates your legacy data into your new django app. The only thing you have to supply is an appropriate SQL query that transforms your data from the old schema into your model structure. Dependencies between these migrations will be resolved automatically. Give it a try!

Warning: This documentation is a work in progress. Please open an issue on Github if any information are missing.

Contents:

1.1 Installing

Installation of `django-data-migration` is straight forward, as it only requires the following steps (assuming you have already set up `virtualenv` and `pip`).

1. Install using `pip`:

```
pip install django-data-migration
```

2. Add `django-data-migration` to your `requirements.txt`
3. Add to `INSTALLED_APPS`:

```
'data_migration',
```

4. Run `./manage.py migrate` or `./manage.py syncdb` to create the included models

1.2 Writing Migrations

1.2.1 What is a Migration?

A migration is a Python class that should be placed in a file called `data_migration_spec.py` in one of your app-directories. `django-data-migrations` searches in each app, included in `INSTALLED_APPS`, for this file and imports all from it automatically.

Your migration normally specifies the following things:

- A database connection to your legacy data (whereever this is)
- The model class, the migration should create instances for
- A corresponding SQL-Query, which maps the old DB-schema to the new Django-model-schema
 - You can specify what should be done with special columns, returned by the query (Many2Many-, ForeignKey-, One2One-Relations). With minimal configuration, these things can be migrated automatically.
- Dependencies to other models can be specified. This is used, to determine the order each migration can be applied. e.g. If a migration specifies a model as dependency, his migration will be executed before our migration will be processed.

- You can implement different hooks, where you normally manipulate the data returned by the query or do some things which are not possible by SQL itself.
- You can specify, if your migration should look for new instances on a second run. This is not the default case.

1.2.2 A complete Migration example

To give you an overview, how a common migration looks, the following listing shows a migration for a `Post` model. This is an excerpt from a `data_migration_spec.py` which can be found in a testing app, which is used by `django-data-migration` itself.

The complete app can be found [here](#) ...

```
class PostMigration(BaseMigration):
    query = """
    SELECT id,
           Title as title,
           Body as body,
           Posted as posted,
           Author as author,
           (
               SELECT
                   GROUP_CONCAT(id)
               FROM comments c
               WHERE c.Post = p.id
           ) as comments
    FROM posts p;
    """
    model = Post
    depends_on = [ Author, Comment ]
    column_description = {
        'author': is_a(Author, search_attr="id", fk=True, prefetch=False),
        'comments': is_a(Comment, search_attr="id", m2m=True,
                        delimiter=",", prefetch=False)
    }

    @classmethod
    def hook_after_save(self, instance, row):
        # because of the auto_now_add flag, we have to set it hard to this value
        instance.posted = row['posted']
        instance.save()
```

As you can see, `PostMigration` inherits from a class called `BaseMigration`. This is one of the classes which is listed here [Setup Database Connection](#).

1.2.3 Migration details

Setup Database Connection

`django-data-migration` should support as many databases as possible, so the connection part is not implemented directly for each database. You have to override the `open_db_connection` classmethod in your migration.

Tip: The connection handling should be implemented once in a `BaseMigration` where all other Migrations inherit from.

Important: `django-data-migration` requires that the database returns a `DictCursor`, where each row is a dict with column names as keys and the row as corresponding values.

SQLite

The following code implements an example database connection for SQLite:

```
import sqlite3

class BaseMigration(Migration):

    @classmethod
    def open_db_connection(self):
        conn = sqlite3.connect(.....)

    def dict_factory(cursor, row):
        d = {}
        for idx, col in enumerate(cursor.description):
            d[col[0]] = row[idx]
        return d

    conn.row_factory = dict_factory
    return conn
```

MySQL

You have to install the corresponding MySQL-Python-driver by executing:

```
pip install MySQL-python
```

The following code implements an example database connection for MySQL.

```
import MySQLdb

class BaseMigration(Migration):

    @classmethod
    def open_db_connection(self):
        return MySQLdb.connect(.....,
                               cursorclass=MySQLdb.cursors.DictCursor
        )
```

PostgreSQL

You have to install the corresponding PostgreSQL-Python-driver by executing:

```
pip install psycopg2
```

Important: a version of `psycopg2` \geq 2.5 is required as it allows the `cursor_factory` to be specified through `connect()` instead of `get_cursor()`.

The following code implements an example database connection for PostgreSQL.

```
import psycopg2
import psycopg2.extras

class BaseMigration(Migration):

    @classmethod
    def open_db_connection(self):
        return psycopg2.connect(.....,
                                cursor_factory=psycopg2.extras.RealDictCursor
                                )
```

MS-SQL

@aidanlister contributed a sample DB connection for MS-SQL using pyodbc, which has to be installed first:

```
pip install pyodbc
```

The following code implements an example database connection for MS-SQL.

```
import pyodbc

class ConnectionWrapper(object):
    def __init__(self, cnxn):
        self.cnxn = cnxn

    def __getattr__(self, attr):
        return getattr(self.cnxn, attr)

    def cursor(self):
        return CursorWrapper(self.cnxn.cursor())

class CursorWrapper(object):
    def __init__(self, cursor):
        self.cursor = cursor

    def __getattr__(self, attr):
        return getattr(self.cursor, attr)

    def fetchone(self):
        row = self.cursor.fetchone()
        if not row:
            return None
        return dict((t[0], value) for t, value in zip(self.cursor.description, row))

    def fetchall(self):
        rows = self.cursor.fetchall()
        if not rows:
            return None

        dictrows = []
        for row in rows:
            row = dict((t[0], value) for t, value in zip(self.cursor.description, row))
            dictrows.append(row)
        return dictrows
```

```
class BaseMigration(Migration):
    @classmethod
    def open_db_connection(self):
        dsn = "DRIVER={SQL Server Native Client 11.0};SERVER=X;DATABASE=X;UID=X;PWD=X"

        cnxn = pyodbc.connect(dsn)
        wrapped_connection = ConnectionWrapper(cnxn)
        return wrapped_connection
```

What can be configured in every migration

In your migration classes you have several configuration options, which are listed below with a short description. For an in-depth explanation you can consult the paragraphs below.

`Migration.skip = False`

If *True*, this migration will be skipped and not processed.

`Migration.query = None`

An SQL-SELECT-query which returns the data that is processed and passed to the model-class-constructor. Please consult the documentation for an in-depth description for this attribute.

`Migration.model = None`

The Django model class where the query creates instances for. There could be only one migration for each model.

`Migration.depends_on = []`

A list of model classes the model requires to be migrated before itself. This includes normally all model classes which are listed in *column_description*.

`Migration.column_description = {}`

This dict contains information about special columns returned by the query (ForeignKey-, OneToOne- and Many2Many-Relations) or if it should be excluded from automatic processing.

`Migration.allow_updates = False`

If the following is set to *False*, the migration will be executed only once. Otherwise it will search for missing elements and creates them.

`Migration.search_attr = None`

This is a unique model field, which is used to search for existing model instances.

Example: for Django's User model it can be *username* or *id*

Important this attribute is required if *allow_updates* is set to *True*

Writing effective Migration-queries

Important: TODO

Define dependencies

Important: TODO

Describe special columns

Your `query` can include special columns, that are represented as special Django-relations (ForeignKey-, Many2Many- or One2One-Relations). Or you can exclude specific columns from automatic processing. You will normally define these settings with an invocation of the `is_a`-function, which does some tests and returns the required settings. This will then be used by `django-data-migration` in different places.

```
data_migration.migration.is_a(klass=None, search_attr=None, fk=False, m2m=False,
                              o2o=False, exclude=False, delimiter=';', skip_missing=False,
                              prefetch=True, assign_by_id=False)
```

Generates a uniform set of information out of the supplied data and does some validations. This function is used to build the `column_description` attribute in your migration. The information is used to translate your data from the old DB schema into django instances.

Parameters

- **klass** – A model class which has a reference to the current column
- **search_attr** – A model attribute which is used to filter for the right model instance
- **fk** – The specified column in query includes a ForeignKey-Reference
- **m2m** – The specified column in query includes (possible multiple) elements which will be represented by a ManyToMany-Reference
- **delimiter** – The character which separates multiple elements in a `m2m`-column
- **o2o** – The specified column in query includes a OneToOne-Reference
- **exclude** – The specified column should not be processed automatically, but can be accessed in any hook which includes the `row`-parameter
- **skip_missing** – defines the behaviour if any element in a defined relation (`fk`, `m2m` or `o2o`) can not be found. If set to `True`, missing elements are ignored. Otherwise an exception is raised.
- **prefetch** – If set to `True`, this will prefetch all existing instances from the related model into a lookup cache, so that less database queries will be made.
- **assign_by_id** – This will assign the related objects as Primary Key values (int) instead of whole objects. This decreases the memory usage but has a drawback, because the related object is not available until `save()` has been called on the model.

Some examples for `is_a` can be found here: [A complete Migration example](#).

Using Migration Hooks

`data_migration.migration.Migration` defines a number of different hook-functions which will be called at different places allowing you to customize the migration work at different levels.

```
class data_migration.migration.Migration
```

Baseclass for each data migration

```
classmethod hook_before_transformation(row)
```

Is called right before row is passed to the model constructor.

Manipulate the row data if it is required. Here you can bring the data in a suitable form which is not possible in SQL itself.

Parameters `row` – the dict which represents one row of the SQL query

classmethod hook_before_save (*instance, row*)

Is called right before the migrated instance is saved.

Do the changes, that make the instance valid, in this hook.

If the instance should not be committed, e.g. due to a runtime check failing, you may return `False` which will prevent the model's `save` method and `after_save` hooks from being called.

Parameters

- **instance** – the migrating instance which could be altered
- **row** – the dict which represents one row of the SQL query

classmethod hook_after_save (*instance, row*)

Is called right after the migrated instance has been saved initially.

This is the place where you can set the data for a `DateTimeField` with `auto_now_add`, where the date from the SQL query is not used otherwise.

Parameters

- **instance** – the migrating instance which could be altered
- **row** – the dict which represents one row of the SQL query

Note when you make changes to the instance you have to call `save()` manually

classmethod hook_update_existing (*instance, row*)

Is called for each existing instance when `allow_updates` is `True`

Parameters

- **instance** – the existing instance which can be updated
- **row** – contains the raw result without any transformation

Note It is YOUR responsibility to make sure, that this method can be called MULTIPLE times.
DO SOME CHECKS

classmethod hook_before_all ()

Is called before the migration will be migrated

Here you can execute some special setup code, which should be executed only once.

classmethod hook_after_all ()

Is called after the migration has been processed succesfully

Here you can set a debugger breakpoint for testing the migration results.

classmethod hook_row_count (*connection, cursor*)

Is called for getting the number of elements returned by `query`

This is useful because several database engines (including `sqlite3`) doesn't provide a real rowcount value. They return `-1` instead. With this hook, you can implement your own method for getting the row count, for example by issuing a special count-SQL-query (use the `connection` parameter).

It should return a numeric value which is displayed when migrating.

classmethod hook_error_creating_instance (*exception, row*)

Is called in case of an error on creating instances from the query

It produces some debug output and reraises the exception

Error-Handling

In case of an exception when creating the instances, a default error handler will be called, to print the current row to stderr and then reraise the exception.

class `data_migration.migration.Migration`

Baseclass for each data migration

classmethod `hook_error_creating_instance` (*exception, row*)

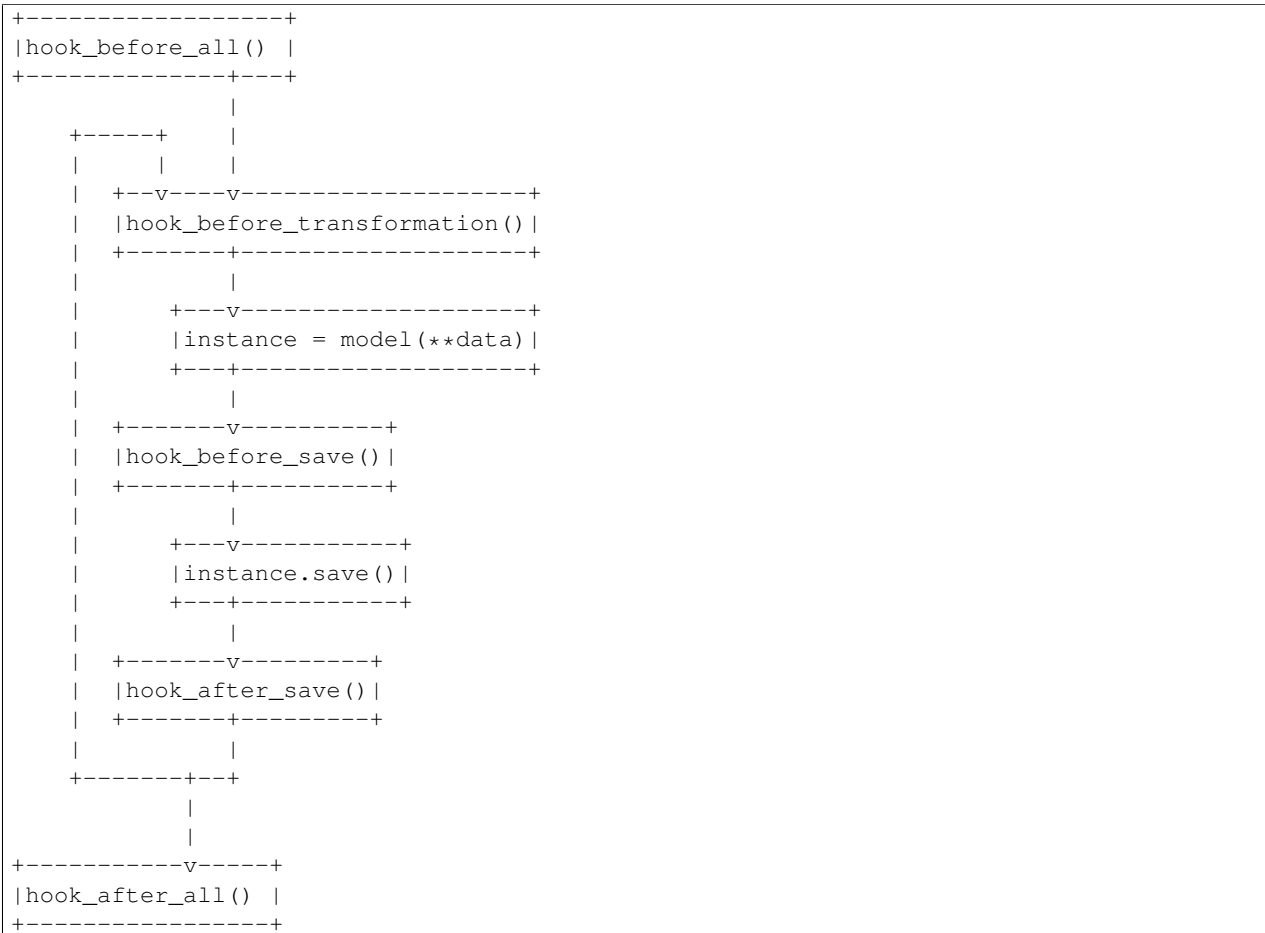
Is called in case of an error on creating instances from the query

It produces some debug output and reraises the exception

You can override this hook in your migration if it requires special handling of errors. When this method returns without an exception, the next row from the query will be processed.

Hook-Flowchart

The following graphic shows each Hook-method and when it is called in contrast to the model handling which is done by `django-data-migration`.



Implement updateable Migrations

Important: TODO

1.3 Migrating your data

After you wrote all of your Migrations, you can put them in action by executing the `migrate_legacy_data` management command:

```
./manage.py migrate_legacy_data [--commit]
```

If you omit the `--commit`-flag, the data is not saved to the DB. This is useful when you develop your migrations and have some failing migrations, but the db is not cluttered with any incomplete data. When your migrations are successful you can add `--commit` and your data is saved when no errors occur.

Note: In older versions of this library, the management command is called `migrate_this_shit`. This has been deprecated, but it is still there. `migrate_legacy_data` should be more appropriate.

1.4 Troubleshooting

1.4.1 GROUP_CONCAT column size limit in Mysql is too low

Mysql has fairly low limit for rows that can be merged by the `GROUP_CONCAT` function. For large result sets, this has to be increased. This can be done with the following SQL statement, which can be executed in `open_db_connection`:

```
@classmethod
def open_db_connection(self):
    conn = MySQLdb.connect(...)
    cursor = conn.cursor()

    cursor.execute('SET SESSION group_concat_max_len = 60000000;')
    return conn
```

1.5 Changelog

1.5.1 Version 0.2.1

- `atomic()` is now used instead of `commit_on_success()` when it is available. This prevents deprecation warnings that are displayed with Django `>= 1.7`.

1.5.2 Version 0.2.0

- Introduced some performance improvements by implementing prefetching of related objects, that reduces the number of issued SQL-Queries dramatically. There is now the opportunity to assign related objects by their id instead of a full instance, which can reduce the memory usage.
- There are two new arguments in `is_a`: `prefetch=True` and `assign_by_id=False`. Because prefetching is enabled by default, it should bring a massive performance boost only by upgrading to this version
- Switching to a new minor release because of the changed behaviour in `get_object`

1.5.3 Version < 0.2.0

There is no explicit Changelog until 0.2.0. Use `git log` to get the information from git.

d

`data_migration.migration`, 7

A

allow_updates (data_migration.migration.Migration attribute), 7

C

column_description (data_migration.migration.Migration attribute), 7

D

data_migration.migration (module), 7

depends_on (data_migration.migration.Migration attribute), 7

H

hook_after_all() (data_migration.migration.Migration class method), 9

hook_after_save() (data_migration.migration.Migration class method), 9

hook_before_all() (data_migration.migration.Migration class method), 9

hook_before_save() (data_migration.migration.Migration class method), 8

hook_before_transformation()
(data_migration.migration.Migration class method), 8

hook_error_creating_instance()
(data_migration.migration.Migration class method), 9, 10

hook_row_count() (data_migration.migration.Migration class method), 9

hook_update_existing() (data_migration.migration.Migration class method), 9

I

is_a() (in module data_migration.migration), 8

M

Migration (class in data_migration.migration), 8, 10

model (data_migration.migration.Migration attribute), 7

Q

query (data_migration.migration.Migration attribute), 7

S

search_attr (data_migration.migration.Migration attribute), 7

skip (data_migration.migration.Migration attribute), 7