



Deloitte's table extraction tool: TableMiner

Saving Time for Deeper Analysis

Deloitte's table extraction tool "TableMiner" reproduces tables from unstructured (pdf) documents into spreadsheets, taking all-too-common dirty work out of daily life of the analyst.

The Need

Sound analysis is based on data... generally, the more, the better. Modern organizations increasingly rely on machines to analyze large volumes of data. This data must be structured, i.e. in the form of tables and databases that can be programmatically queried. The data age has ushered in widespread availability of structured data. Often, but not always. Some data remains "unstructured" – buried within narrative of

reports, or inserted as tables within published (digital) documents. The data may be available, yet it is not easily accessible for machine-enabled analysis.

The ubiquitous Portable Document Format (PDF) guarantees formatting consistency and in a generally compact filesize. It is also notoriously unhelpful to those seeking to extract tabular data from its contents. This difficulty lies in the fundamental design of PDFs to be easy on the eyes. Unlike other formats (MS or other Office formats), which store tabular data explicitly as embedded tables, PDFs store tables and text as vector graphics. Converting content to graphics preserves formatting at the cost of removing context: any formatting and structure is

lost when copying and pasting text out of a PDF document. Already a problem with e-documents (Office documents) saved as PDFs, scans saved as PDFs without embedded OCR (optical character recognition) are even more unwieldy.

The result: analysts are left with few options other than to manually transfer data to editable formats (spreadsheets) – a labor intensive and error-prone process. This binds qualified resources to menial tasks, representing a costly productivity drain, inviting fatigue-related manual errors, and leaving less time for value-added analytical work. 

Our Solution: TableMiner

Deloitte's table extraction tool "TableMiner" addresses this very issue, joining multiple Computer Vision and Natural Language Processing methods to provide an easy solution to an all too common problem.

TableMiner's neural networks scan each page for tabular data – irrespective of whether the document contains only a single or hundreds of tables in various formats and styles, even multiple per page. Once identified, tables are then automatically extracted and converted into a specified format, directly viewable in the TableMiner application or downloaded and viewed in a separate (MS or other) spreadsheet application.

It deftly handles so-called "dirty" scans without OCR – meaning: only a picture, no associated text meaning. TableMiner can automatically distinguish between e-documents saved as PDF, "clean" scans (with OCR) and "dirty" scans (without OCR). Finding a "dirty" scan, TableMiner first applies state-of-the-art OCR techniques: scanned tables are partitioned into smaller sub-boxes and characters are digitized. In other words, TableMiner "reads" the document and saves its meaning. TableMiner then summarily reconstructs the extracted information to form a text version of image.

TableMiner offers a convenient graphical user interface for the user to selectively search for and extract targeted tables. For larger jobs, TableMiner's batch processing feature saves valuable time, allowing the user to upload multiple documents, determine output format and let TableMiner get to work, automatically identifying and extracting all tables within the uploaded documents.

Fig. 1 – Proper recognition of columns

| Original PDF | Without TableMiner | With TableMiner | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--|--|--|---|---|---|-----|---|---|---|---|------|---|---|---|---|--------|---|---|---|---|---|------|---------|----------|------------|--|--|---|---|---|---|-----|---|---|---|---|------|---|---|---|---|--------|---|---|---|---|
| This is a sample unstructured document ... | Manual copy-pasting from the PDF results in all the columns being concatenated ... | The grid positioning of the values within the table remains intact ... | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <table border="1"> <thead> <tr> <th></th> <th>A</th> <th>B</th> <th>C</th> <th>D</th> </tr> </thead> <tbody> <tr> <td>Rev</td> <td>5</td> <td>6</td> <td>7</td> <td>8</td> </tr> <tr> <td>Cost</td> <td>2</td> <td>3</td> <td>4</td> <td>4</td> </tr> <tr> <td>Profit</td> <td>3</td> <td>3</td> <td>3</td> <td>4</td> </tr> </tbody> </table> | | A | B | C | D | Rev | 5 | 6 | 7 | 8 | Cost | 2 | 3 | 4 | 4 | Profit | 3 | 3 | 3 | 4 | <table border="1"> <tbody> <tr> <td>ABCD</td> </tr> <tr> <td>Rev5678</td> </tr> <tr> <td>Cost2344</td> </tr> <tr> <td>Profit3334</td> </tr> </tbody> </table> | ABCD | Rev5678 | Cost2344 | Profit3334 | <table border="1"> <thead> <tr> <th></th> <th>A</th> <th>B</th> <th>C</th> <th>D</th> </tr> </thead> <tbody> <tr> <td>Rev</td> <td>5</td> <td>6</td> <td>7</td> <td>8</td> </tr> <tr> <td>Cost</td> <td>2</td> <td>3</td> <td>4</td> <td>4</td> </tr> <tr> <td>Profit</td> <td>3</td> <td>3</td> <td>3</td> <td>4</td> </tr> </tbody> </table> | | A | B | C | D | Rev | 5 | 6 | 7 | 8 | Cost | 2 | 3 | 4 | 4 | Profit | 3 | 3 | 3 | 4 |
| | A | B | C | D | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Rev | 5 | 6 | 7 | 8 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Cost | 2 | 3 | 4 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Profit | 3 | 3 | 3 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ABCD | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Rev5678 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Cost2344 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Profit3334 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| | A | B | C | D | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Rev | 5 | 6 | 7 | 8 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Cost | 2 | 3 | 4 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Profit | 3 | 3 | 3 | 4 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ... with a table inserted within the text. | ... requiring as much, if not more work to clean than re-typing. | ... transferred to a spreadsheet for further analysis. | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Advantages/Benefits

- Shifts the analyst focus to what really counts: analysis vs data collection and aggregation.
- Reduced transmission error
- Automatically extracts tables from hundreds of documents via batch-processing.
- Reliably handles different table formats and types of PDF documents.
- Scanning throughout entire document
- Easy integration with existing applications and workflows via the TableMiner API.
- Can be hosted on the cloud for subscription service or implemented locally with client firewall.

Example Use Cases

- Facilitating balance sheet analysis (e.g. for underwriting SME/corporates).
- Various audit functions.
- Technical accounting/extraction of terms form contracts for input to systems.
- Extension of RPA capabilities.
- Exhaustive audit.
- Creating new and perfecting existing workflows: For example, a setup that directly forwards scanned documents to TableMiner via the API and stores a copy of the extracted tables.

Contacts

David Thogmartin

Leader aiStudio
dthogmartin@deloitte.de

This presentation contains general information only, and none of Deloitte GmbH Wirtschaftsprüfungsgesellschaft or Deloitte Touche Tohmatsu Limited ("DTTL"), any of DTTL's member firms, or any of the foregoing's affiliates (collectively, the "Deloitte Network") are, by means of this presentation, rendering professional advice or services. In particular this presentation cannot be used as a substitute for such professional advice. No entity in the Deloitte Network shall be responsible for any loss whatsoever sustained by any person who relies on this presentation. This presentation is to be treated confidential. Any disclosure to third parties – in whole or in part – is subject to our prior written consent.

Deloitte refers to one or more of Deloitte Touche Tohmatsu Limited, a UK private company limited by guarantee ("DTTL"), its network of member firms, and their related entities. DTTL and each of its member firms are legally separate and independent entities. DTTL (also referred to as "Deloitte Global") does not provide services to clients. Please see www.deloitte.com/de/UeberUns for a more detailed description of DTTL and its member firms.