Machine translation of mathematical text (using LATEX)

Tanya Schmah

Mathematics & Statistics, University of Ottawa, tschmah@uottawa.ca



u Ottawa.

Vitrine REL - OER Showcase, March 14th 2023

Outline

- Our specific project: translation of course notes for MAT 2122 and MAT 2125
- The General Problem
- Workarounds to Translating LaTeX documents
- A Solution: the PolyMath Translator
- The Future

Our specific project in 2022/23: Translation of course notes for MAT 2122 and MAT 2125

Original notes by Prof. Alistair Savage, University of Ottawa:

- MAT 2122 Multivariable Calculus (170 pages, in LaTeX). Translation done, pending final review by author.
- MAT 2125 Elementary Real Analysis (139 pages, in LaTeX) Translation half-done.

Extract from original "Multivariable Calculus" notes

Chapter 1

Differentiation

In this first chapter, we discuss differentiation in multiple dimensions. Some results will be review from previous courses, while other material will be new.

1.1 Open sets and boundaries

Many definitions in calculus require one to be able to get "near" a point. The precise formulation of this concept involves the notion of open sets, which we discuss here.

Definition 1.1.1 (Open ball). Let $\mathbf{a} \in \mathbb{R}^n$ and $\varepsilon > 0$. The *open ball* of radius ε centered at \mathbf{a} is

$$B_{\varepsilon}(\mathbf{a}) := \{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{a}\| < \varepsilon \}.$$



Definition 1.1.2 (Open set). A set $A \subseteq \mathbb{R}^n$ is open if

$$\forall \mathbf{a} \in A, \exists \varepsilon > 0 \text{ such that } B_{\varepsilon}(\mathbf{a}) \subseteq A.$$

In other words, A is open if every point of A is the center of an open ball contained in A.



Extract from "Calcul différentiel de plusieurs variables"

Dérivées partielles et différentielles

Dans ce premier chapitre, nous abordons la dérivation en dimensions supérieures. Certains résultats est une revu de cours préalables, tandis que d'autres seront nouveaux.

1.1 Ensembles ouverts et frontières

De nombreuses définitions en calcul différentiel nécessitent que l'on fasse l'étude «près» d'un point. La formulation précise de ces concepts nécessite la notion d'ensembles ouverts, que nous aborderons ici.

Définition 1.1.1 (Boule ouverte). Soient $\mathbf{a} \in \mathbb{R}^n$ et $\varepsilon > 0$. La boule ouverte de rayon ε et centre \mathbf{a} est

$$B_{\varepsilon}(\mathbf{a}) := \{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{a}\| < \varepsilon \}.$$

Définition 1.1.2 (Ensembles ouverts). Un ensemble $A \subseteq \mathbb{R}^n$ est dit ouvert si

 $\forall \mathbf{a} \in A, \exists \varepsilon > 0 \text{ tel que } B_{\varepsilon}(\mathbf{a}) \subseteq A.$

En d'autres termes, A est ouvert si chaque point de A est le centre d'une boule ouverte contenue dans A.



The General Problem

Translate a mathematical document (in LATEX, naturally) into a different natural language, e.g. from English to French.

written $J_x^k(M, N)$. The set $J^k(M, N)$ is the union of these sets, for all x. It is a smooth vector bundle over $M \times N$, called the k-jet bundle. The k-jet of f with source x is written $j^k f(x)$. In local coordinates, $j^k f(x)$ "is" the kth order Taylor expansion of f at z. The k-jet extension of $f : M \to N$ is the map

 $j^k f: M \longrightarrow J^k(M, N); \qquad x \longmapsto j^k f(x).$

If f is smooth, then $j^k f$ is as well, so there is a map

(1) $j^k : C^{\infty}(M, N) \longrightarrow C^{\infty}(M, J^k(M, N)),$

taking f to $j^k f$. This map is continuous, with respect to the strong topologies on domain and codomain [GG73].

Theorem 2.5 (Jet transversality). Let M and N be manifolds and let S be a submanifold of $J^k(M,N)$. Then the set of functions $f : M \to N$ such that $j^k f$ is transverse to S is residual in $C_k^{\infty}(M,N)$, and open dense if S is closed.

To apply jet transversality to vector fields, we need the modified version in Theorem 2.6. This result is known, but we are unaware of a proof in the literature. We will prove it from Theorem 2.5, using the globalisation technique in Lemma 2.8. We will re-use the same globalisation lemma in the proof of the new result in Theorem 2.9.



Translate a mathematical document.



► LATEX Google Translate? DeepL?

$ \overset{Google Transl.}{\underset{DeepL?}{\overset{Google Transl.}{\rightarrow}}} $	^{ate?} (<u>not</u> LATEX)
Original English LATEX (compiled)	Definition 2.1. Let $x \in \mathbb{F}^n$. The closed ball of radius r centered at x is $S_r(x) = \{y \in \mathbb{F}^n \mid d(x,y) \le r\}.$
Google Translate (results uncompilable)	<pre>\ begin {defn} Soit \$ x \ dans \ F ^ n\$.} La \ define {boule ferme de rayon \$ r \$ centre sur \$ x \$} est \$\$ \$\$ \$_r (x) = \ {y \ in \ F ^ n \ mid d (x, y) \ leq r \}. \$\$ \ end {defn}</pre>







A Solution: the PolyMath Translator

Developed together with Aditya Ohri



Summary of PolyMath Translator Results from 2021

Successful initial implementation using pandoc:

- Excellent translation quality (BLEU 53.5 on small test corpus).
- ► Output is LATEX that *usually* compiles without hand-correction.
- Moderate ease-of-use.

A Ohri, T Schmah (2021) Machine translation of mathematical text. IEEE Access 9, 38078-38086

Some limitations of original PolyMath Translator:

- Only supports English to French.
- Only translates single files.
- No user configuration or editable glossaries
- ► Doesn't use LATEX semantics to e.g. *not* translate verbatim environments, or comments, or certain arguments.

changes the LATEX commands

A Ohri, T Schmah (2021) Machine translation of mathematical text. IEEE Access 9, 38078-38086

A limitation of original PolyMath Translator :

It changes the LATEX commands.

```
For example, \textit{hello}
is translated to \emph{bonjour}
instead of \textit{bonjour}.
```

This is an inevitable consequence of using <code>pandoc</code> to translate $\[Mathebar{E}T_EX\]$ into an abstract internal representation and then back into $\[Mathebar{E}T_EX\]$.

This may be acceptable for some applications, e.g. browsing articles, but not for e.g. book authors.

A limitation of original PolyMath Translator :

It introduces errors into some LATEX commands.

For example, \includegraphics[scale=0.2]{file.png}
is translated to \includegraphics{file.png}

This is again a consequence of using pandoc to "translate" from LATEX to LATEX: Since pandoc is trying to *interpret* and *translate* every LATEX command, it will always fail on commands it doesn't know.

From the point of view of PolyMath, pandoc is trying to do <u>too much</u>: it's trying to understand the <u>semantics</u> of LateX when we mainly just need the <u>syntax</u>. What a translator needs to understand about LATEX

Mainly the syntax.

Plus, enough semantics to ...

- identify which arguments should be translated, e.g.:
 - Don't translate: math; label names;
 - Do translate: title, section names, text mode strings inside math environments;
- tokenize math expressions (inline and displayed environments);
- tokenize label references;
- translate other files referred to in \input and \include commands.

PolyMath Translator v0.2-dev, using TexSoup parser



TexSoup is a fault-tolerant, Python3 package for searching, navigating, and modifying LaTeX documents.

Created by Alvin Wan + contributors.

Inspired by ${\tt Beautiful Soup},$ a Python package for parsing HTML and XML documents.

https://github.com/alvinwan/TexSoup
https://texsoup.alvinwan.com

PolyMath Translator v0.2-dev, using TexSoup parser

This version of PolyMath leaves T_EX commands unchanged, and understands *just enough* semantics. It includes:

- Editable lists of which command and environment arguments to translate
- Tokenization of math expressions
- Tokenization of label references
- Translation of entire file trees using \input and \include.

Experience with the PolyMath Translator at uOttawa

- Automatic translation
 - A textbook for Intro to Math Models, which was automatically translated for a francophone student in an English-language class.
- Semi-automated translation (automatic + post-editing)
 - Course notes for Intro to ODEs and Multivariable Calculus

Example 2.1 Résoudre le problème de la problème à valeur initiale

$$y'(t) + 2y(t) = 10,$$
 $y(0) = 1.$

Cette équation se présente sous la forme (2.1) avec p(t) = 2 et g(t) = 10. On multiplie par une fonction $\mu(t)$ et on obtient

$$\mu(t)y'(t) + \underbrace{2\mu(t)}_{\mu'}y(t) = 10\mu(t).$$

Ainsi, l'équation définissant le facteur intégrant est $\mu'(t) = 2\mu(t)$, qui a une solution $\mu(t) = \exp(2t)$. En multipliant l'équation par le facteur intégrant,

Experience with the PolyMath Translator at uOttawa

Findings:

- PolyMath produces excellent translations that nonetheless need correction and polishing.
- The uncorrected automatic translation is already useful.
- Overall, a semi-automated professional translation process is about twice as fast as manual translation.
- PolyMath, and the TexSoup parser, still have issues and require an expert user.

In progress: the Ottawa Mathematical Term Bank

Machine translation makes heavy use of subject-specific *glossaries*, i.e. *dictionaries*.

A term bank is a like a glossary but more specific, containing information to disambiguate homonyms:

field	algebra	corps	Körper
field	database	champ	Datenfeld
finitely generated group		groupe de type fini	endlich erzeugte Gruppe

From a term bank, a glossary can be extracted and customised for each project.

"Translating mathematical text" vs. "Translating Later Accuments"

Most of the work in this project applies to any LATEX document.

Some of it is math-specific:

- Tokenization of math expressions (inline and displayed).
- Development of training corpora of mathematical sentence pairs (in English and French).
- Training of neural networks on math-heavy corpora.
- ► The Ottawa Mathematical Term Bank.

In general, domain-specific translators (e.g. in law and medicine) have given better results than general-purpose ones. Our work suggests that this will be true for mathematical text as well.

This is a big opportunity for the math community to improve communication, accessibility and inclusion.

A vision of the future

- An open source PolyMath Translator project, supported by:
- Open data: term banks, training corpora, and pre-trained deep learning models for many language pairs
- Improved math-specific translation, e.g. using content of math expressions, or using topic models to select appropriate translations for ambiguous terms.
- A web service (like "Google Translate for LaTeX")
- A "translate" button in preprint servers.



We have implemented a machine translation system, the PolyMath Translator, for LaTeX documents containi implementation translates English LaTeX to French LaTeX, attaining a BLEU score of 53.5 on a held-out test it produces LaTeX documents that can be compiled to PDF without further editing. The system first converst document into English sentences containing math tokens, using the pandoc universal document converter t

Contributions welcome!

Especially:

- Glossaries (or term banks) for many language pairs and many specialized subjects.
- Training data: pairs/sets of corresponding sentences in different languages.
- Technical advice:
 - setting up an open source project;
 - Hosting, naming, structuring term banks and training corpora.

Thank you!