# "Ulpianus scripsit"? Using Artificial Intelligence for authorship attribution of ancient Roman law texts

**Thomas Rüfner**

Professor of Roman Law and Private Law

University of Trier

Judge in the Higher Regional Court

Koblenz

ruefner@uni-trier.de

─────────── ABSTRACT ───────────

The present article seeks to demonstrate the usefulness of artificial intelligence for the exploration of the sources of Roman law. In a small experiment, a so-called support vector machine was employed to determine to which Roman jurist a given source text can be attributed. While the results were not perfect, they are sufficient to show the potential of new technologies for future research in the area.

This paper has been subjected to double-blind peer review

# "Ulpianus scripsit"? Using Artificial Intelligence for authorship attribution of ancient Roman law texts[*]

SUMMARY: 1. Authorship attribution as a research topic in Roman law – 2. Authorship attribution and verification through machine learning techniques – 3. Preparation of the corpus – 4. Vectorization – 5. Training – 6. Results – 7. Conclusion.

## 1. Authorship attribution as a research topic in Roman law

Philologists and historians have long sought to determine the author of a text (or to exclude the authorship of a certain person) based on the text's linguistic and stylistic characteristics. Interestingly, an early example concerns a text purporting to be a Roman legal document, which had found its way into the *Corpus Iuris Canonici*[1]: In 1440 Lorenzo Valla demonstrated that the so-called *donatio Constantini* could not possibly be genuine due, in large part, to the document's language[2].

Valla's research shattered the authority of a document which constituted the basis of legal claims of the church at the time. The proof that the *donatio* was a fabrication potentially had concrete and far-reaching legal consequences. In later times, attempts to determine the authors of Roman legal texts were made for purely scholarly purposes. During the period of the 'hunt for interpolations', legal historians relied heavily on the assumption that the language and style of the classical jurists could be distinguished from that of the compilers of Justinian's Digest in the 6th century. This conviction was the basis for the attempt to detect post-classical changes and additions to the sources[3].

More recently, computers have been employed to analyse the linguistic structure and the style of texts in order to identify the author. As early as 1968, when digital technology was still in its early stages, Lothar Müller proposed to

---

[*] The contribution is intended for publication in the proceedings of the conference "*Dialogo transdisciplinare e identità del giurista*", organized by the Research Center "*Studi sulla Giustizia*" on 19th and 20th September 2022 at the University of Milan.

[1] Incorporated in the Decretum Gratiani as D. 96 c. 14.

[2] Cf. H. Craig, *Stylistic analysis and authorship studies*, in S. Schreibmann, R. Siemens, J. Unsworth (eds.), *A Companion to Digital Humanties*, Wiley, Hoboken, 2008, p. 282; H. Love, *Attributing Authorship*, CUP, Cambridge, 2002, p. 18 f.

[3] R. Repnow, *Überlegungen zur quantitativen Stilanalyse römischer Rechtstexte*, in *SDHI*, 2017, vol. 83, p. 101 f.

58

analyse Roman law texts of disputed authorship like *Ulpiani regulae* with the help of computers[4]. Sadly, his doctoral thesis on the subject was not accepted and remains unpublished[5]. From the 1970s onward, Tony Honoré sought to identify idiosyncratic features of individual Roman jurists' writings. Honoré even proposed a new mathematical method to measure the vocabulary richness of a given text[6]. On the basis of stylistic criteria, he tried to distinguish and identify the jurists who drafted rescripts of the imperial chancery as secretaries *a libellis* in the third century[7].

Despite the traditional importance of stylistic analysis and questions of authorship attribution in the field of Roman law[8], it seems that so far no attempt has been made to employ the techniques of artificial intelligence in order to determine the authors of Roman legal sources. This is somewhat surprising not only because of the general popularity of these techniques in recent years, but also because there is a considerable amount of literature demonstrating the potential of such methods for the resolution of questions regarding the authorship of historic sources, including Latin texts[9]. It seems worthwhile to conduct a few experiments and explore the potential of these new approaches.

If a method which enables us to have computers correctly attribute ancient Roman legal texts to their authors can be developed, various research questions can be answered. Ideally, Romanistic scholarship would be able definitely to answer the question whether the work known as *Ulpiani regulae* or *Tituli ex corpore Ulpiani* was actually authored by Ulpian, and to resolve similar issues regarding other works like Ulpian's *opiniones* or the *sententiae* attributed to Paul.

Even if no final answer can be given to questions of this order, it will be interesting to see to what extent computers can distinguish the styles of different jurists and if some jurists are more easy to identify than others. The experiments presented in this paper can thus be regarded as (further) tests of Savigny's famous assertion that the Roman jurists of the classical epoch were "fungible persons"

---

[4] L. Müller, *L'ordinateur et les textes de droit Romain*, in *Revue – Organisation internationale pour l'étude des langues anciennes par ordinateur*, 1968, p. 65–82, on *Ulpiani regulae*, in particular, cf. p. 66–69.

[5] M. Avenarius, *Der pseudo-ulpianische liber singularis regularum*, Wallstein, Göttingen, 2005, p. 56–58.

[6] T. Honoré, *Some simple measures of richness of vocabulary*, in *ALLC Bulletin*, 1979, vol. 7, p. 172–179.

[7] T. Honoré, *Emperors and Lawyers*, 2nd ed., OUP, Oxford, 1993.

[8] On the history of this area of research see the detailed explanation and bibliographic notes in R. Repnow, *Überlegungen zur quantitativen Stilanalyse römischer Rechtstexte*, cit., p. 101-104.

[9] See J. Kabala, *Computational authorship attribution in medieval Latin corpora: the case of the Monk of Lido (ca. 1101–08) and Gallus Anonymous (ca. 1113–17)*, in *Language resources & Evaluation*, 2020, vol. 54, p. 25–56 with further references on p. 28–30.

with few individual traits[10]. If their individual styles can be distinguished (even or at least) by a machine, then Savigny was obviously wrong.

## 2. Authorship attribution and verification through machine learning techniques

It is neither possible nor necessary for our purposes to give a definition of 'artificial intelligence' or to list all methods and techniques that are generally considered to fall into this category[11]. The techniques relevant for problems of authorship attribution and authorship verification[12] all belong to the field of machine learning and, more specifically, supervised learning. The machine (computer) is first trained with a set of input data for which the correct output is known. If the task is to recognize the author of a given text, the training data must consist in texts the author of which is known. The input data (which must be brought into a computer readable form, that is into the form of a series of numbers) are subjected to certain mathematical manipulations in order to produce an output. The mathematical operations applied to the data are then systematically changed until the output for the training data corresponds to the (known) solutions. If the machine is trained to attribute texts to authors, then the output should be the correct author for each text.

All attempts to determine the author of a text with the help of computers are based on the assumption that texts written by a given author contain certain characteristics that can be measured and quantified. Earlier methods of stylometry like those developed by Lothar Müller and Tony Honoré were based on certain measurable properties of the texts like the average length of sentences or the richness of the vocabulary of different jurists. By contrast, machine learning models are often described as 'black boxes'. It is impossible to tell precisely what characteristics of the text cause the computer to attribute it to one author or the other. This is an important difference between the modern techniques and earlier methods.

The existing approaches to the problem of determining the authors of Latin texts all appear to have been designed for large continuous texts[13]. It would be difficult to use the same approaches for research on the Roman legal sources,

---

[10] F.C. von Savigny, *Vom Beruf unserer Zeit für Gesetzgebung und Rechtswissenschaft*, Mohr und Zimmer, Heidelberg, 1814, p. 157.

[11] On the problem of definition see, for example, M. Barberis, *Giustizia predittiva: ausiliare e sostitutiva. Un approcio evolutivo*, in *Milan Law Review*, 2022, vol. 3, p. 3 f.; Th. Rüfner, *Juristische Herausforderungen der Künstlichen Intelligenz aus der Perspektive des Privatrechts*, in H.-G. Dederer, Y.-Ch. Shin (eds.), *Künstliche Intelligenz und juristische Herausforderungen*, Mohr Siebeck, Tübingen, 2021, p. 17–20.

[12] On the distinction between authorship attribution and authorship verification see J. Kabala, *Computational authorship attribution in medieval Latin corpora*, cit., p. 29.

[13] See, as an example, J. Kabala, *Computational authorship attribution in medieval Latin corpora* cit., p. 33, who works with a corpus of texts with a minimum of 5000 words each.

because most of the writings of the Roman jurists only survive in Justinian's Digest. While the Digest preserves a large number of texts with a precise and mostly reliable attribution to an author in the *inscriptio* (and sometimes additional information on authors quoted within the text preserved in the Digest such as the *Ulpianus scripsit* in Macer 1 de apell. D. 2.8.15.1), the single fragments are relatively short.

For this reason, the method employed for the experiments presented here was borrowed from a blog post by Gareth Dwyer, which details how the authors of reviews of restaurants and other businesses on the Yelp platform can be determined through machine learning[14]. While the corpus of reviews used by Dwyer is quite different from the Roman sources as far as the language (English) and the subject matter (food and service quality of certain establishments) of the texts is concerned, the texts used are similar in length to the texts of the Digest.

Before the experiments conducted are explained in more detail, it should be noted that this paper is only a preliminary report on a research project that is far from complete. Neither the methodology employed nor the results reached have any claim to definiteness. They constitute no more than a first attempt to explore what is possible.

## 3. Preparation of the corpus

The texts for all experiments were taken from the Latin Library website which contains texts of the Digest and the Institutes of Gaius originating ultimately from Joseph Menner's Romtext database[15]. The text of Ulpian's *regulae* was downloaded from the website ancientrom.ru. The texts were divided up so that each paragraph was stored in a separate text file. All text files containing texts from the same work according to the inscription were stored in one folder. Thus, there was, e.g. one folder containing 5349 texts (paragraphs or fragments with no further subdivisions) from Ulpian's commentary on the edict and another folder containing 721 texts (paragraphs) from the institutes of Gaius.

The successful use of all machine learning methods depends on the availability of sufficient training data. The computer can only be expected to 'learn' the features which distinguish the style of a given author, if there are enough texts available on which the computer can be 'trained'. Therefore, as in Dwyer's experiment, all authors of whom less than 500 were contained in the collection of source texts were dropped. Additionally, for the training of the model all texts of doubtful attribution had to be left out. Thus, the texts from *Ulpiani regulae* (both

---

[14] G. Dwyer, *Yelp reviews: Authorship attribution with Python and scikit-learn*, in *Michael Kennedy on Technology*, https://blog.michaelckennedy.net/2017/06/21/yelp-reviews-authorship-attribution-with-python-and-scikit-learn/ (last visited on 9 January 2023).

[15] On the history of the Romtext project see G. Klingenberg, *Die Romtext-Datenbank*, in *Informatica e diritto*, 1995, vol. 4, p. 223–232.

those preserved outside the Digest and those ascribed to a *liber singularis regularum* or *regularum libri septem* in the Digest[16]) as well as the texts ascribed to Ulpian's *opiniones*, *pandectae* and *responsa* in the Digest were left out[17]. Likewise, the texts ascribed to *Pauli sententiae* in the Digest were not used. Finally, the texts ascribed to the *institutiones* of Gaius in the Digest were left out in order to avoid duplication since the *institutiones* as preserved outside the Digest were already included in the corpus.

These operations left a corpus with the texts by Ulpian (8762 texts), Paul (3653 texts), Gaius (1581 texts), Papinian (1156 texts), Pomponius (954 texts), Julian (816 texts), Scaevola (658 texts), Modestin (622), Marcianus (618 texts). These texts were loaded into the computer's memory in a randomized order.

## 4. Vectorization

In the next step, the texts in the corpus had to be transformed into a computer readable sequence of numbers. This was done using the TfidfVectorizer of the scikit-learn library for the Python programming language[18] with the parameter ngram_range=(1,2). This means that every text was transformed into a table which shows how frequently each word (unigram) and each combination of two words (bigram) present in the entire corpus occurs in the text. The 'table' contained more than 300,000 columns. Of course, of all the different word forms and combinations that are present in the entire collection of texts, only a few occur in a single given text. The table for each text was therefore a so-called sparse matrix. The columns for most words or bigrams contained the number zero[19].

As a result of this process, the list of Latin texts was now present as a series of sparse matrices in the computer's memory. The names of the authors of the texts were stored in a separate list.

## 5. Training

After all these preparations, the training itself was started. In a first experiment, only 500 texts of the nine jurists with a least as many texts in the corpus

---

[16] On these different textual traditions see D. Liebs, *Ulpiani Regulae – Zwei Pseudoepigrafa*, in *Romanitas – Christianitas. Untersuchungen zur Geschichte und Literatur der römischen Kaiserzeit. Johennes Straub zum 70. Geburtstag*, Walter De Gruyter, Berlin – New York, 1982, p. 283 and 287.

[17] On the spuriousness of these works see T. Honoré, *Ulpian: Pioneer of Human Rights*, 2nd ed., OUP, Oxford, 2002, p. 212–215 and 217–226.

[18] On this library see F. Pedregrosa et alii, *Scikit-learn: machine learning in Python*, in *Journal of Machine Learning Research*, 2011, vol. 12, p. 2825–2830.

[19] Th. Joachims, *Text categorization with support vector machines: learning with many relevant features*, in C. Nédellec, C. Rouveirol (eds.), *Machine Learning: ECML-98*, Springer, Heidelberg, 1998, p. 140.

were used. Of these 4500 texts, 3600 were used for training (the training corpus) and 900 (=20%) were used for testing the results (the test corpus). The training was done using a so-called Linear Support Vector Classification. This method of classification is a variant of the support vector machine method, which is perhaps less widely known as a technique of machine learning than the use of neural networks, but particularly well suited for tasks of text classification[20]. It can be explained (with some degree of oversimplification) as a process designed to find a function which separates to classes of objects from each other:
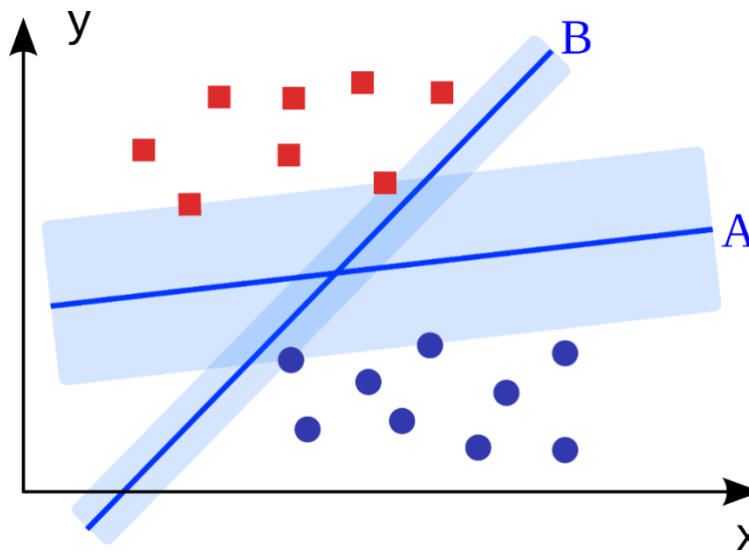


*Image 1: Two classes given as vectors with two possible separation lines and their corresponding margin areas between the class areas. Line A has a larger empty margin area than line B.*

Image by Fabian Bürger (https://commons.wikimedia.org/wiki/User:Ennepetaler86); original at https://commons.wikimedia.org/wiki/File:Svm_intro.svg. License: CC BY 3.0 (https://creativecommons.org/licenses/by/3.0/deed.en).

The image illustrates a simple example: The task is to find a function which separates the red squares from the blue discs. The lines A and B (which can be regarded as graphs of mathematical functions) both separate the two categories. However, line B comes very close to some red and blue items while line A leaves more room on both sides. It separates the two categories more clearly[21].

---

[20] Th. Joachims, *Text categorization with support vector machines: learning with many relevant features*, cit., p. 137–142.

[21] For a more detailed explanation see K.P. Bennet, C. Campbell, *Support vector machines: Hype or hallelujah?*, in *ACM Special Interest Group on Knowledge Discovery in Data Newsletter*, 2000, vol. 2, p. 1–13.

The same technique which is here demonstrated for the separation of squares and discs in a two-dimensional space can be used for multi-dimensional objects like the sparse matrices representing the source texts.

## 6. Results

After the training, the model was applied to the 900 remaining texts. It succeeded in correctly attributing around 53% of the texts. This may not sound too impressive, but it should be noted that the computer did much better than guessing blindly: Since there were nine possible authors, the probability of correctly attributing a text with a random guess was only 11,11%! On the other hand, it should also be kept in mind that the task was made considerably easier by the fact that the test data, like the training data, were taken from a sample which contained an equal number of texts from each jurist.

In a second experiment, 500 texts from each author were used as training data and all remaining texts from the nine authors were used for testing. This means that all authors were equally represented in the training data whereas there were many more texts by Ulpian (8262) than by Marcianus (118) left for testing. Not unsurprisingly, this made the computer's task more arduous. The model was only able to classify correctly less than 45% of the texts in the test corpus. This was still better than guessing blindly, but not satisfactory.

In a third attempt, only texts by the six jurists with the largest number of texts were used: Ulpian, Paul, Gaius, Papinian, Pomponius, and Julian. In this setting, the computer was again able to attribute correctly the majority of the texts (around 52%), while the probability of a blind guess being correct was still only 16,66 %.

Admittedly, the attribution of the texts to certain jurists was far from reliable in all three experiments. There were, however, interesting differences. 51% of the texts that were not attributed correctly by the model were authored by Ulpian. This is hardly surprising because Ulpian is overrepresented in the test corpus. In fact, 59% of all texts were by Ulpian. On the other hand, only 22.7% of all texts in the test corpus were authored by Paul, but 35.2% of the texts that were attributed incorrectly were by Paul. This would seem to indicate the style of Paul is less recognizable than the style of Ulpian. Texts by Papinian make up 4.7% of the corpus but only 2.1% of the misattributed texts. Papinian's style seems to be even more recognizable for the computer than Ulpian's.

Finally, the attempt was made to guess the author of two disputed works ascribed to Ulpian. For these experiment, 618 texts from each of the nine jurists with the largest number of texts were used (since there are only 618 texts by Marcianus this is the largest number of texts which can be used if the training corpus must contain an equal number of texts from each jurist).

First, the model was used to attribute the 307 texts contained in the *Ulpiani regulae* or *Tituli ex corpore Ulpiani* to one of the nine jurists. The model ascribed

roughly one half of the texts to Gaius (49.05%) and only 3.6% to Ulpian. A significant portion of the texts was attributed to Modestinus.
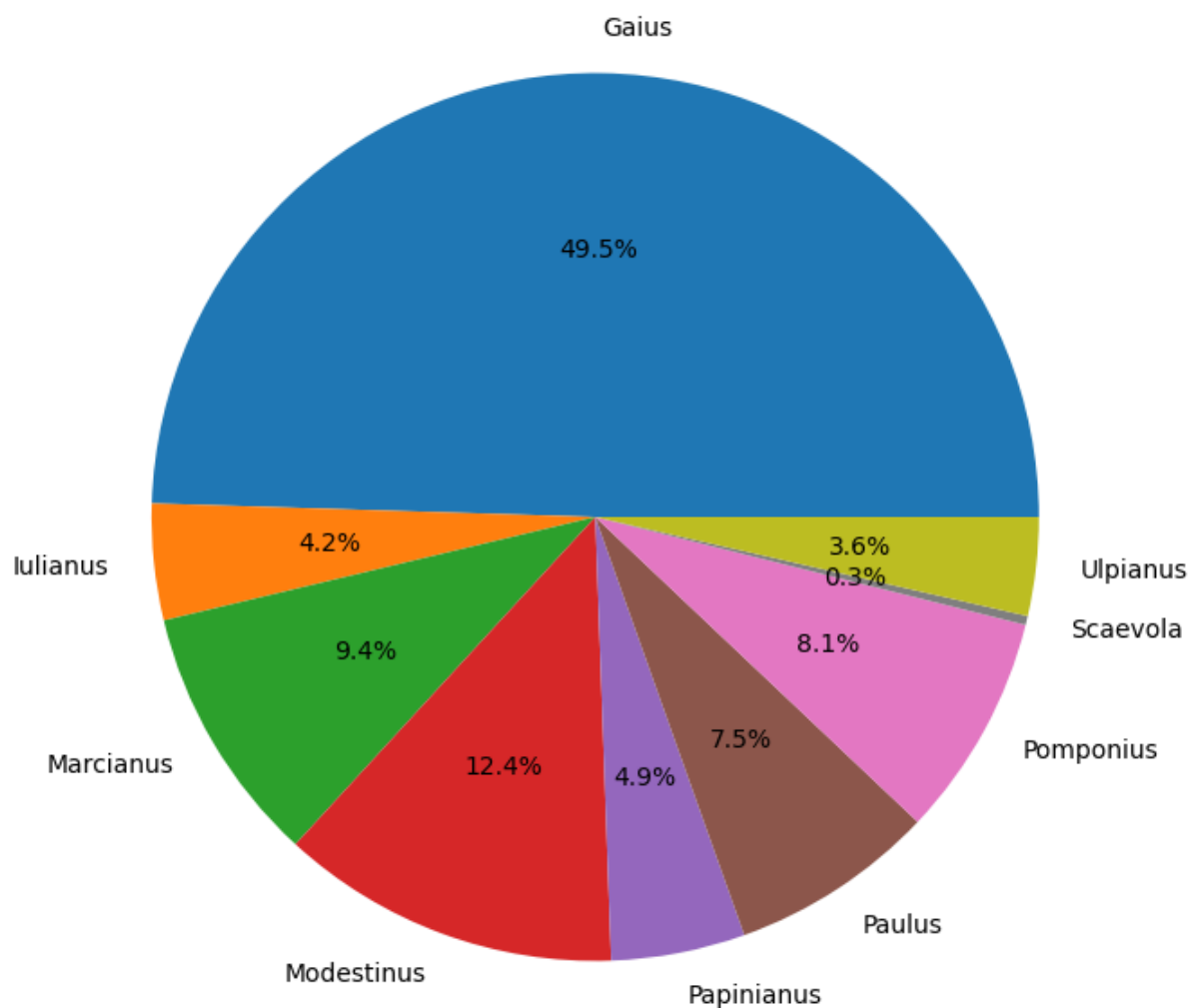


*Image 2: Attribution of the texts of Ulpiani regulae to the nine jurists most frequently represented in the Digest.*

While this result confirms the conclusion that the work was not authored by Ulpian[22], it should not be assumed that the *regulae* must have been written by Gaius. To conclude from the results that Gaius was the author would overlook the

---

[22] In this sense T. Honoré, *Ulpian,* cit., p. 211 f.; D. Liebs, *Ulpiani Regulae,* cit., p. 284; M. Avenarius, *Der pseudo-ulpianische liber singularis regularum*, cit., p. 531; for Ulpian as the author L. Müller, *L'ordinateur et les textes de droit Romain*, cit., p. 69; F. Mercogliano, *"Tituli ex corpore Ulpiani": storia di un testo*, Jovene, Napoli, 1997, p. 105.

fact that the author may well be a jurist who is not among the nine authors featured in the training corpus[23].

The fact that Gaius is assumed frequently by the model, while the late classical luminaries Ulpian, Paul and Papinian are all underrepresented would seem to lend support to the conclusion of Avenarius that the work was written in the high classical period[24]. However, the relative frequency of the attribution of texts to Modestinus contradicts that conclusion. It should be remembered that similarities in style and (especially) vocabulary may also be due to the genre to which texts belong rather than to the individual style of the author[25]. The frequent attribution of texts to Gaius is perhaps better explained by the fact that the *regulae* and the institutes of Gaius belong to similar, though not identical genres of legal literature.

Finally, the texts ascribed to Ulpian's *opiniones* were explored in the same way. Again, the results make it unlikely that the work was in fact written by Ulpian:

---

[23] The task is thus one of authorship verification rather than authorship attribution, see above, n. 12.

[24] M. Avenarius, *Der pseudo-ulpianische liber singularis regularum*, cit., p. 531.

[25] R. Repnow, *Überlegungen zur quantitativen Stilanalyse römischer Rechtstexte*, cit., p. 111.
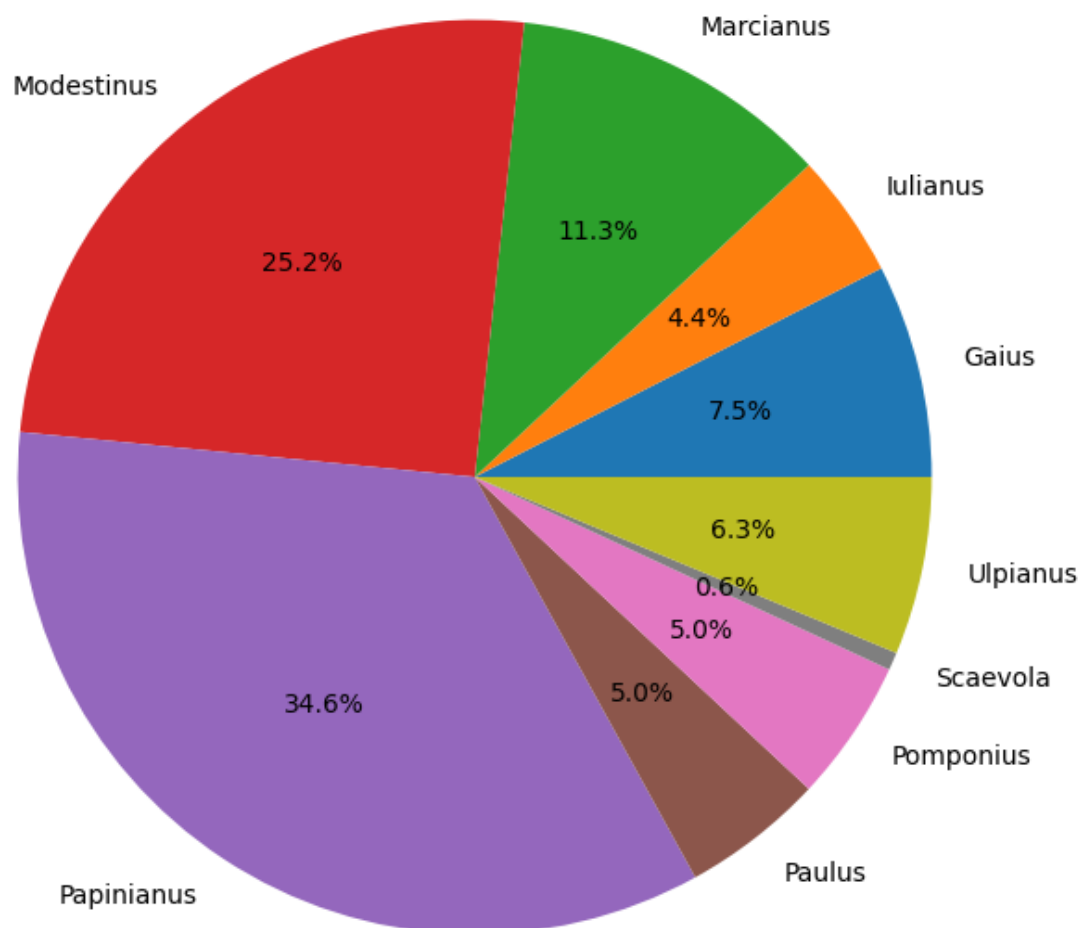
*Image 3: Attribution of the texts of Ulpian's opinions to the nine jurists most frequently represented in the Digest.*

Many texts are attributed to Modestinus. Papinian is even more prominent.

## 7. Conclusion

The results reached up to now are far from conclusive. The methodology needs to be refined. In particular, the stylistic characteristics of different genres of legal literature deserve more scholarly attention. Even so, it seems clear that the exploration of the Roman legal texts with the methods of machine learning has the potential of shedding new light on many long-standing issues of Romanistic schlolarship.

# Bibliography

M. Avenarius, *Der pseudo-ulpianische liber singularis regularum*, Wallstein, Göttingen, 2005.

M. Barberis, *Giustizia predittiva: ausiliare e sostitutiva. Un approcio evolutivo*, in *Milan Law Review*, 2022, vol. 3, p. 2–18.

K.P. Bennet, C. Campbell, *Support vector machines: Hype or hallelujah?*, in *ACM Special Interest Group on Knowledge Discovery in Data Newsletter*, 2000, vol. 2, p. 1–13.

H. Craig, *Stylistic analysis and authorship studies*, in S. Schreibmann, R. Siemens, J. Unsworth (eds.)*, A Companion to Digital Humanties,* Wiley, Hoboken, 2008, p. 273–288.

G. Dwyer, *Yelp reviews: Authorship attribution with Python and scikit-learn*, in *Michael Kennedy on Technology*, https://blog.michaelckennedy.net/2017/06/21/yelp-reviews-authorship-attribution-with-python-and-scikit-learn/ (last visited on 9 January 2023).

T. Honoré, *Emperors and Lawyers*, 2nd ed., OUP, Oxford, 1993.

T. Honoré, *Some simple measures of richness of vocabulary*, in *ALLC Bulletin*, 1979, vol. 7, p. 172–179.

T. Honoré, *Ulpian: Pioneer of Human Rights*, 2nd ed., OUP, Oxford, 2002.

Th. Joachims, *Text categorization with support vector machines: learning with many relevant features*, in C. Nédellec, C. Rouveirol (eds.), *Machine Learning: ECML-98*, Springer, Heidelberg, 1998, p. 137–142.

J. Kabala, *Computational authorship attribution in medieval Latin corpora: the case of the Monk of Lido (ca. 1101–08) and Gallus Anonymous (ca. 1113–17)*, in *Language resources & Evaluation*, 2020, vol. 54, p. 25–56.

G. Klingenberg, *Die Romtext-Datenbank*, in *Informatica e diritto*, 1995, vol. 4, p. 223–232.

D. Liebs, *Ulpiani Regulae – Zwei Pseudoepigrafa*, in *Romanitas – Christianitas. Untersuchungen zur Geschichte und Literatur der römischen Kaiserzeit. Johennes Straub zum 70. Geburtstag*, Walter De Gruyter, Berlin – New York, 1982, p. 282–292, revised version at https://freidok.uni-freiburg.de/dnb/download/8684 (last visited on 9 January 2023).

H. Love, *Attributing Authorship*, CUP, Cambridge, 2002.

L. Müller, *L'ordinateur et les textes de droit Romain*, in *Revue – Organisation internationale pour l'étude des langues anciennes par ordinateur,1100/1800*, 1968, p. 65–82.

F. Pedregrosa *et alii*, *Scikit-learn: machine learning in Python*, in *Journal of Machine Learning Research*, 2011, vol. 12, p. 2825–2830.

R. Repnow, *Überlegungen zur quantitativen Stilanalyse römischer Rechtstexte*, in *SDHI*, 2017, vol. 83, p. 101–129.

Th. Rüfner, *Juristische Herausforderungen der Künstlichen Intelligenz aus der Perspektive des Privatrechts*, in H.-G. Dederer, Y-Ch. Shin (eds.), *Künstliche Intelligenz und juristische Herausforderungen*, Mohr Siebeck, Tübingen, 2021, p. 15–42.

F.C. von Savigny, *Vom Beruf unserer Zeit für Gesetzgebung und Rechtswissenschaft*, Mohr und Zimmer, Heidelberg, 1814.