

What is the value of a wrong proof?

A case study

Frank Buschman

A thesis presented for the degree of
Bachelor of Mathematics



Thesis supervisor: Vincent Coumans

First corrector: Luca Consoli

Second corrector: Michael Mueger

Department of Mathematics, FNWI

Radboud University

September 2022

Contents

0	Abstract	2
1	Introduction	2
2	What is a wrong proof?	3
2.1	The Cartesian story	3
2.2	Gaps in proofs	5
2.3	Formalizability	7
2.4	Peer review process	9
2.5	Conclusion	11
3	Values of proofs	12
3.1	Didactics	12
3.1.1	Demonstrating proof techniques	13
3.1.2	Reconceive mathematical domains	13
3.1.3	Promote discovery	14
3.1.4	Explanatoriness	14
3.1.5	Motivatedness	15
3.1.6	Visual intuitiveness	16
3.2	Aesthetics	16
3.2.1	Platonist beauty	16
3.2.2	Kantian beauty	17
3.2.3	Representation	17
3.3	Miscellaneous	18
3.3.1	Purity	18
3.3.2	Fitting transparency	18
3.3.3	Fitting simplicity	19
3.3.4	Systematization	19
3.4	Potential values in wrong proof	19
3.4.1	Didactics	19
3.4.2	Aesthetics	20
3.4.3	Miscellaneous	21
4	The fundamental theorem of algebra	22
4.1	A proof and a gap	22
4.2	Gauss' first proof	23
4.2.1	Preparatory lemmas	23
4.2.2	The proof	25
4.3	Ostrowski's adaptation	27
4.3.1	Part one	28
4.3.2	Part two	29
4.3.3	Part three	30
4.3.4	Part four	33
4.3.5	Part five	35

5	Value analysis	35
5.1	Didactic values	36
5.1.1	Demonstrating proof techniques	36
5.1.2	Reconcieve mathematical domains	37
5.1.3	Promote discovery	38
5.1.4	Explanatoriness	38
5.1.5	Motivatedness	39
5.1.6	Visual intuitiveness	40
5.2	Aesthetic values	40
5.2.1	Kantian beauty	40
5.2.2	Representation	41
5.3	Miscellaneous	41
5.3.1	Purity	41
5.3.2	Systematization	42
6	Discussion	42

0 Abstract

In this thesis we investigate the value of wrong proofs. To this end, we compare Gauss' faulty first proof of the Fundamental theorem of Algebra, to Ostrowski's correct version of this proof. The theoretical background work for this case study is divided into two parts, first the notion of *wrong proofs* will be explored, after that some proof values found in literature will be presented. The results found in the value analysis are that the two proofs do not differ much for the values analysed. This suggests that when a wrong proof gets disregarded solely on invalidity, some mathematical value might be lost. Subsequent research is necessary for deeper statements, this thesis is the first looking into value of wrong proofs and only considers one case, thus can best be seen as a preliminary study into this topic. One can conclude on the basis of this thesis that an incorrect proof can uphold some proof values normally found in correct proofs.

1 Introduction

Mathematical proofs show the validity of a theorem. This is the standard conception of proofs. This is embodied by formal logic. In formal logic, mathematicians have formulated precise and check-able rules for the validity of a proof. Validation is thus an integral part of our conception of proving.

This is also visible in publication practices. When referees evaluate a mathematics article for publication, one of the criteria they look at is the correctness of the proofs [Geist et al., 2010]. This suggests that proofs are only valuable when they are correct. In this thesis, I assess whether this is the case, I try to answer the question 'What is the value of a wrong proof?'

There are two main reasons for suggesting that there might be some value in wrong proofs. The first is the study of values related to proofs, other than

validity. This research has not only focused on proofs of new mathematical theorems, but also on re-proving existing theorems. New proofs of existing theorems also get published. These new proofs have value beyond validity [Dawson, 2006]. Thus a published piece of research does not always have to prove a new theorem to be valuable. There is a wide array of values aside from validity as discussed in Chapter 3 that are found in proofs. Wrong proofs might also uphold some of these values just as proofs that re-prove theorem do.

Secondly there exists a database (viXra.org) where wrong proofs can get published. Such a database would have no reason to exist if there is no value in reading wrong proofs. This suggests that in practice some mathematicians might find wrong proofs valuable.

No previous research has been done on the value of wrong proofs. Thus this thesis is quite exploratory. To assess whether wrong proofs can have value, I will investigate a wrong proof to see if it has some value. To do this, I first need to determine what constitutes a wrong proof in mathematical practice. This is less clear than in the case of logic previously discussed. In Chapter 2 the considerations on when something is a wrong proof are laid out. Based on this reflection, I have chosen two proofs of the Fundamental theorem of Algebra.

The first proof is by Gauss [Struik, 2014]. This was deemed correct when it was published in 1799. Mathematicians later discovered a mistake in the proof. This mistake was quite fundamental to the proof and it could not easily be corrected. The second proof is by Ostrowski [1920]. This is an adaptation of Gauss' proof such that the mistake is solved. In this adaptation he tried to stay as true to Gauss' proof as possible. These two proofs are thus quite similar. Because of the similarity, it is easier to compare the values of these two proofs.

This thesis is outlined as follows. In Chapter 2 I will try to better understand the notion of a 'wrong proof'. In Chapter 3 I will discuss several salient values other than validity of proofs found in literature. For each of these values I will try to understand whether a wrong proof could uphold it. In Chapter 4 I will present the proofs of Gauss and Ostrowski. In Chapter 5, I will compare these proofs on the basis of the values discussed in Chapter 3. In the last chapter I will discuss these results.

2 What is a wrong proof?

Before one can start to study the value of wrong proofs, it first has to be clear what a wrong proof is. This chapter explores the notion of wrong proofs. In the later chapters the proof values used will be presented and the case study will be executed.

2.1 The Cartesian story

Logic is a subject of mathematics. The proofs written in logic are precise and unyielding. These could be called *formal proofs*. The logical steps one can take are defined precisely, room for interpretation is squeezed out of these proofs. In

contrast there are informal, or everyday proofs. These are proofs produced and published in mathematical practice.

In logic one can check the correctness of proofs by logical soundness, e.g. check if every step is a predefined step of logic. In informal proofs this is not so easily done as “proofs are written in a way to make them easily understood by mathematicians. Routine logical steps are omitted. An enormous amount of context is assumed on the part of the reader” [Hales, 2008, p1371].

In order to understand when a proof is a wrong proof, we thus need a criterion that explains when an informal proof is proper¹. The “Cartesian story” [Fallis, 2003, p47] is a theory that tries to provide this criterion.

According to the Cartesian story a mathematician is justified in believing that a proof is proper when she recognizes it is a deduction “from true and known principles by the continuous and uninterrupted action of a mind that has a clear vision of each step in the process” [Descartes, 1927, p47].² This implies that in order to know that a proof is proper one has to have checked every step and held this to a background of knowledge to check whether that step is correct. This background of knowledge is often a wealth of proper proofs and proof techniques built by education and practice. Only after checking every step in such a way is she justified in believing it is a proper proof.

When one believes a proof to be proper there are three ways it can break the Cartesian story. The first is by accepting principles that are not “true and known” [Descartes, 1927, p47]. For example, this can be done by basing the proof on an unknown assumption. The second way is by not having a “clear vision of each step” [Descartes, 1927, p47]. For example, believing a proof step with an unseen caveat while accepting the proof is not in compliance with the Cartesian story. Lastly the proof can be discontinuous or interrupted. This implies a leap of faith, one has to believe a proof step could be made without checking whether it can.

One should note that the Cartesian story is written from the perspective of a mathematician, not mathematics as a whole. This implies that not only the author has to comply to the Cartesian story. When one reads a proof one still needs to fulfill the Cartesian story to be rightly convicted that the proof is proper without disregarding the Cartesian story [Fallis, 2003]. One could also soften the Cartesian story to having the viewpoint of mathematics as a whole. This would make sense as some mathematicians “accept results from the published literature as black boxes in their own research” [Geist et al, 2010, p159-160].³

This Cartesian story is a fine test for proofs to be proper, but only a basis. On the one hand the Cartesian story appeals to many a Mathematician as this is “the only method by which one truly justifies that one knows that A follows from B” [Azzouni, 1994, p169]. Thus one could say this is *the* test. In practice

¹Here proper is used in the sense that a thing or object is deemed suitable or appropriate for use. In contrast with right or correct which carry a stronger, more stern and more precise meaning.

²Original work published in 1701.

³More about the implications of published research and peer review in Section 2.4.

there is a caveat; there are proofs considered proper without appealing to the Cartesian story [Fallis, 2003]. These proper proofs leave some gaps that do not fulfill the Cartesian story [Fallis, 2003]. Thus the Cartesian story is not a complete account of when a proof is proper in practice. In order to understand this we will briefly discuss four different gaps that mathematicians leave in proofs in the next section and see their effect on the Cartesian story.

2.2 Gaps in proofs

In this section we will discuss the types of gaps one can find based upon the work of [Fallis, 2003], namely the inferential, enthymematic and untraversed gaps. The study of these gaps will shed light on what is accepted in a proper proof in practice, and where the Cartesian story is incomplete. This in turn refines our understanding of when a proof is wrong, which is the goal of this chapter.

One of the types of gaps left is an *inferential gap*. This is a gap that uses a proof step that is, in fact, not a proof step [Fallis, 2003] i.e. there is a step in the proof that is inconsistent with logic or accepted mathematical reasoning. These sometimes arise as small mistakes as “mathematicians ...are not infallible.” [Poincaré, 1952, p47]

Inferential gaps are consistent with the Cartesian story [Fallis, 2003]. When one recognizes such a gap one does not believe the proof to be proper. At least not until the gap has been consolidated. Thus the Cartesian story is not diffused by this type of gap as proofs with this type of gap are not deemed proper.

Enthymematic gaps are gaps where the author has traversed the steps of the proof but left them out nonetheless [Fallis, 2003]. This can be done for clarity, time consumption or simply out of laziness. This type of gap is often left in such a way that that a reader could easily fill the gaps of the proof. Mathematicians often leave them such that “a large majority of the experts in the specialized field should be able to validate the proof within a reasonable amount of time” [Andersen, 2020, p241].

As mentioned before, the Cartesian story is written from an individual perspective. When an author leaves an enthymematic gap the reader of the proof has to fill this gap in order to comply with the Cartesian story. Else she does not have a “a clear vision of each step in the process” [Descartes 1927, p47]. When a reader or author does not check the gap it is a different type of gap, an untraversed gap [Fallis, 2003].

A gap is untraversed when the omissions are not checked for correctness by the mathematician [Fallis, 2003]. This is often the case when reading papers. When it is not necessary to fully understand a proof often times one glosses over it in order to understand the result better. It is considered acceptable to then assume the result true and use it.

Untraversed gaps affect the Cartesian story from a personal viewpoint. When such a gap is left one does not “trace the whole chain of intermediate conclusions” [Descartes, 1927, p62] as one needs to do for the Cartesian story. Nonetheless this is an often times accepted gap. When one takes the perspective of the

whole of mathematics, this type of gap can be consistent with the Cartesian story [Fallis, 2003]. If at least one mathematician traversed the gap then the Cartesian story is salvaged. The real problem comes with the last (sub-)type of gap. Universally untraversed gaps are those gaps in a proof that are untraversed by all mathematicians [Fallis, 2003]. This includes the author, reader and referees (more about refereeing in Section 2.4). These do exist as “mathematicians do not spend a lot of time working through tedious details when they have no doubt that the details can be filled in.” [Fallis, 2003, p60]. As we will see in Section 2.4 referees of journals will not always check these details, [Andersen, 2020; Geist et al., 2010] thus leave them universally untraversed.

Whereas the Cartesian story can fend off the threats posed by enthymematic, inferential and untraversed gaps, the existence of universally untraversed gaps turns out to be a deathblow to the Cartesian story. In that case, no mathematician has a clear view on the continuous sequence of propositions when such a gap is included. Accepting such a gap thus does not comply with the Cartesian story. In practice the proof might be accepted as “a mathematician who leaves a universally untraversed gap is convinced that he could fill in the details of an argument even though he has not done so.” [Fallis, 2003, p62] The mathematician might also be wrong by leaving an inferential gap untraversed.

Fallis shows that such gaps do exist in accepted research.⁴ Probable cases are those like *folk theorems* which “everyone quotes ... but no one can seem to cite a complete and accurate proof” of [Schattschneider, 1981, p141]. Also steps in a proof where the same reasoning is repeated for a similar, but not exactly symmetrical, object. Those could easily be dismissed as *probably the same* and thus remain universally untraversed. Another possibility comes from the following claim:

“She [Raman, 2002] claimed that some mathematicians will be fully convinced of an assertion if they construct an informal argument in support of that assertion and see a method for mapping this argument into a formal proof. These mathematicians do not need to actually carry out this translation to obtain complete conviction.”
[Weber, 2008, p450]

This too could give rise to small parts of published proofs to be left as a universally untraversed gap. Thus clearly, by the multiplicity of reasonable options, universally untraversed gaps are sometimes used and even published⁵.

In this section we have seen that the Cartesian story is unable to account for all the gaps, mostly those of the universally untraversed kind. One possible next step is to try to salvage the Cartesian story. In order to do this one needs to be able to fill these gaps. This could possibly be done through the idea of formalizability discussed in the next section. One could also accept that the Cartesian story does not tell the whole story. The whole story could be a more

⁴A known case is Gauss’ first proof of the fundamental theorem of algebra discussed in Chapter 4.

⁵This claim is supported in Section 2.4.

social endeavour. In Section 2.4 we will discuss the peer review process and the implications this might have on what proofs are deemed proper.

2.3 Formalizability

Every day mathematical proofs are informal when compared to proofs written in some type of logical language. Natural language explains steps, examples are used to clarify thought processes and heaps of background knowledge is assumed on the part of the reader. This is in stark difference to the algorithmic writing of logic. These two types of proving seem worlds apart.

The theory of formalizability assumes there is a connection between the informal and formal proofs. This connection could take many forms. One is by expanding the proof until every step is logical. Another by an informal proof indicating a certain formal proof. One might even say that informal proofs are formal ones in disguise. No matter the form or way, the goal of this connection is to fill the gap between practice and the fundamental philosophical theories. If this works, then logical soundness could explain why some proofs can be accepted even though they do not comply to the Cartesian story.

An example of such a theory is the *derivation indicator view* of mathematical practice coined by Azzouni [2004]. This view claims that “a proof of a theorem reveals a derivation of that theorem.” [Azzouni, 2004, p85] A derivation is a completely formal proof in a certain algorithmic system, say some type of logic. Thus when one takes a certain relevant logical system (like axiomatic set theory in first order logic), then the informal proof indicates a derivation of the proof in that system.

Many topics of discussion arise around such formalizability theories. Discussions like whether all can be formalized (1), if important information is lost during the process of formalizing (2) and if there is a unique derivation or formal proof for each informal proof (3) as discussed in, but not exclusively, [Azzouni, 2004; Cellucci, 2008; Larvor, 2012; Tanswell, 2015; De Toffoli, 2020]. These points of discussion will be expanded on concisely in the next few paragraphs. After we will discuss the implications for mathematical theory and practice. All to better understand the line between proper and wrong proofs. But we will first explain the discussions.

There is no certainty that a theory of formalizability will hold. Such a connection would at least need to explain a few *desiderata* as Tanswell [2015] summarizes⁶.

1. *Rigour*: How informal proofs can be deemed rigorous through the connection with formal proofs.
2. *Correctness*: How informal proofs can be detected as correct or wrong by viewing the formal proof.

⁶Tanswell specifically raises these *desiderata* in a reaction to the derivation indicator view but these can be generally posed for every formalizability theory.

3. *Content*: How the structure of the formal proof is determined by the content of the informal proof.
4. *Techniques*: Explain how seemingly inherently informal arguments can give rise to a formal proof.

These *desiderata* are central topics around which the discussions can be fueled. For instance take *correctness*. This would require the connection between the formal and informal to be strong enough that reviewing a formal proof could show an informal proof to be correct. A problem arises when there is not a single unique formal proof within a certain algorithmic system that is connected to the informal proof. Which proof would one look at for the *correctness*? What would one do with informal proofs connected to both a correct and an incorrect proof? Azzouni seems to need uniqueness for his derivation indicator view according to Tanswell [2015]. But then new difficulties arise like, how one would pick such a unique option when an informal proof left a gap which can be traversed in multiple ways.

A different example around *content* is one of *essentially informal arguments* [Larvor, 2012]. Are there informal arguments that are not formalizable as one would lose a part of the content of the proof? This could be because some arguments are *agent-dependent* as Tanswell [2015] reasons.

A last example is the question of what to do with proofs that use reasoning such as *similarly one can see that* or *by symmetry the same is true for*. These *techniques* which often give rise to untraversed gaps can give many a logician headaches when needed to formalize as Azzouni [2004] and Tanswell [2015] discuss.

There clearly are many currently unsolved discussions around the idea of formalizability. When these can be consolidated formalizability might give a test for properness. Above only the theoretical problems of formalizability are discussed. But mathematics is a practice and proofs are deemed proper by mathematicians in practice. The question is whether this could explain mathematical practice, i.e. do mathematicians use the connection of their informal proof to a formal one to judge whether a proof is proper? Could it explain why universally untraversed gaps are left in practice?

Mathematicians often do not rigorously think of the formal connection their informal proof has in practice. Sometimes this is because mathematicians do not know enough logic to do so, and often times it would just be an insensible investment of time [Tanswell, 2015]. There are projects of formalizing known proofs like the proof of the mutilated checkerboard problem.⁷ These projects take a lot of time and yield many different formal proofs. But these are the exception to the rule, generally no formal proofs are made to formalize informal ones.

Accepting an informal proof thus does not directly rely on checking logical soundness of their formal version in practice. Mathematicians would otherwise

⁷For explanation of this problem, see [Tanswell, 2015].

need a formal proof for every informal proof they write and have it computer-checkable. There are programs that can check a certain style of proofs written in mathematics. Before these programs can be widely used first the technical problems with formalizability would need to be tackled. If these are not tackled then no true general method to formalize can exist. As this is not the case now, and has not been in mathematical history before, it is safe to say mathematicians do not directly use the formalizability to understand whether a proof is proper.

An indirect way of thinking about formalizability as a ground for accepting proofs is that mathematicians intuitively understand whether a proof is formalizable [Larvor, 2012]. This indirect path has various questions to answer as how this would come to be. This conception would be intuitive as silhouetted by this quote of Weber [2008, p450] concerning a remark of one of his research participants.

“Further, she claimed that some mathematicians will be fully convinced of an assertion if they construct an informal argument in support of that assertion and see a method for mapping this argument into a formal proof. These mathematicians do not need to actually carry out this translation to obtain complete conviction [Raman, 2002].”

Even though this intuition view of formalizability is intuitive, it would still have to defend against Occam’s raiser⁸ as it does not seem to provide explanatory power. It is then still unclear to how mathematicians intuitively know a proof to be formalizable and thus proper. This view thus does not help us explain why certain proofs are accepted even though they do not fulfill the Cartesian story.

In conclusion it is clear that, currently, formalizability does not provide the clarity to understand when a proof is proper in practice. Even when mathematical proof would be theoretically formalizable, the practice of informal proofs would persist as “proofs are not only a means to certainty, but also a means to understanding. Behind each substantial formal proof there lies an idea... it will not do to bury the idea under the formalism.” [MacLane, 1986, p378]. Thus formalizability does not explain why proofs can be deemed proper without appealing to the Cartesian story.

That begs the question how mathematicians might decide on what proofs are proper in current practice. For publishing research this is social in nature through the peer review process. The next section will discuss this practice and what it can tell about when a proof is deemed proper.

2.4 Peer review process

In mathematical practice published proofs play a massive role. Mostly those in journals as “in mathematics ... conferences and their proceeding volumes play

⁸This is the idea that one should try to reduce the assumptions to a minimum for the goal at hand.

a subordinate role.” [Geist et al., 2010]. These are at the basis of widespread communication between mathematicians. Whether a paper gets published depends on the peer review process executed by referees. This process could shed light on the social nature of the notion of proper proofs.

Mathematicians use published theorems in their research, both with or without checking the proof. As Geist et al. [2010, p159-160] put it

“We know a substantial number of mathematicians who want to understand all proofs that form a part of their papers ... but we also know that many mathematicians ... accept results from the published literature as black boxes in their own research.”

Even those that do check the proofs often times focus on understanding the gist of it instead of checking it line by line [Weber, 2008].

This puts a heavy emphasis on good peer review. Wrong proofs could get published because of bad peer review, those wrong proofs that seem proper but have an inferential gap. Bad peer review could endorse something incorrect. The results ‘proven’ could then be used as a black box and start a cascade of untrue theorem, this is obviously problematic. To dissuade this possibility mathematics needs to be able to rely on its peer review practices.

In theory a referee could check a proof meticulously to get as good of a judgement as possible. The referee could comply to the Cartesian story in order to judge on a proof. In practice this is not always done. This can be seen by the fact that universally untraversed gaps exist as Fallis [2003] explains.

In order to understand the peer review process we need to understand how a referee judges a paper on validity. Andersen [2020] acknowledges two types of validation, named type 1 and type 2 validation accordingly.

Type 1 validation, or *validation by comparison* is a validation method where the referee checks (often larger steps) of the proof by holding them out to a background of knowledge gained prior.

Type 2 validation, or *line-by-line validation* is the act of checking the steps of the proof line by line to find a reasoning why they can be deemed proper.

With type 1 validation the referee uses his insights and understanding of the subject to make an educated guess. This heavily relies on plausibility of the argument rather than it being correct in detail. Type 2 validation uses more detail and precision driven checking in order to deem an argument proper and is in line with the Cartesian story. In both [Geist et al., 2010] and [Weber, 2008] these validations come to light but aren’t specifically named as such.

Which is used depends on the referee and the context. Type 2 validation is often employed when arguments look “suspicious” [Andersen, 2020] when checking the part by type 1 validation. When parts of the suspicious argument seem trivial or reasonable type 1 validation is sometimes used for those again [Andersen, 2020]. The type of validation used is also influenced by the author of the proof. When it is a well known author with a good reputation referees are

more inclined to find arguments reasonable [Geist et al., 2010]. Referees are sometimes even told by editors that it isn't their job to check the correctness of the proof, that is the author's job, "although the referee should be reasonably convinced" [Auslander, 2008, p65]. Thus what type of validation used is highly variable.

The refereeing process is sensitive to gaps because of the use of type 1 validation. When an author leaves a step untraversed and the referee only uses type 1 validation, then there is an universally untraversed gap left in the proof. Inferential gaps could also be made and overlooked.⁹

Another thing that should be mentioned is that there are changes in what is deemed 'precise enough' throughout history. View the proof of the existence of an infinite number of prime numbers given by Euclid for example. He wrote his proof by giving an example in such a way that one can intuitively generalise the method. This was considered proper in his time. In modern mathematics this would be solely an example, modern versions formulate the method of finding a new prime out of an arbitrary finite set generally. Euclid's version would not be accepted as a proof.¹⁰ This shows that opinions on which proofs can be accepted can change. Thus there is a social aspect to this.

Mathematical peer review is far from foolproof in checking proofs. There are many parts of the process where errors could be made. It seems that the process as a whole acts reasonably as critical or major errors rarely occur [Andersen, 2020] and consensus on which proofs are accepted is higher than in other sciences [Geist et al., 2010]. But the possible errors left do leave a discussion to be had on the properness of published proofs, and thus the certainty of mathematical knowledge.[Geist et al., 2010]

Published proofs thus do not need to be flawless to be accepted by the mathematical community. Gaps are accepted, minor mistakes occur and published results are used as black boxes. These all indicate that, in practice, mathematics stools on a grey foundation rather than a sharp dichotomy of right and wrong. Which proofs are deemed proper is partly based upon consensus rather than exact science. Because of this it could be considered a social endeavour too.

2.5 Conclusion

In order to better understand when proofs are deemed proper three *foci* were analysed; analysing gaps, better understanding the peer review process, and understanding the (incomplete) bridge between foundations and practice. The Cartesian story introduced in Section 2.1 gives a base but does not tell the whole story as mathematical practice is sometimes inconsistent with the Cartesian story. This inconsistency cannot be easily consolidated through the idea of formalizability. A lot of problems arise and, currently, formalizability remains in the realm of theory rather than practice. We have seen that universally untraversed gaps can arise because of how peer review is practiced. The way it

⁹Davis [1972] names a few.

¹⁰Another example is that of Gauss' first proof of the Fundamental theorem of Algebra discussed in Section 4.

is practiced is in part inconsistent with the Cartesian story. It is partly based on meticulously checking steps, but partly based upon intuition and social factors too. Sometimes the correctness is even left solely to the author.

After this chapter it is clear is not an easy one to answer if a proof is proper or not. The answer may contain non-absolute nuances that require decisions based upon intuition. The answer is a social one that can change over time. Some proofs can still be fundamentally wrong without discussion. This would require a fundamental mathematical mistake or something similarly fundamental. The ideas of the Cartesian story and the intuition of formal mathematics still provide a foundation. Thus the act of judging whether a proof is proper has both social and objective aspects.

Because of this nuance in judging which proofs are wrong we would ideally study a proof that at first has been accepted into published research, but was deemed wrong later on. In this case we know for certain that the judgement on the objective and subjective aspects of the notion of wrong proof are in line.

We now have a deeper understanding of wrong proofs. Before we can commence with our case study we need some more background information, namely about proof values. The next chapter provides an explanation of values a (wrong) proof may have. Chapter 4 will explain the proofs used in the case study. Chapter 5 will provide the analysis of values on this wrong proof. The chapter thereafter will discuss the results of this analysis.

3 Values of proofs

In this thesis, I aim to investigate the possible value that wrong proofs can have. To this end, I first investigated how to make sense of the notion of a 'wrong proof'. Now I will study proof values other than validity. In the next sections, a wrong proof's values will be analysed through a case study.

In this chapter we will analyse various proof values other than validity that are described in the literature. Here *values* is meant in a broad sense; everything that could be valuable for (practitioners of) mathematics in a proof.

I submit that these values associated with proofs can be split into three categories for readability purposes: Didactic values (values considering what and how a proof might teach a reader), aesthetic values (values related to beauty and presentation) and a miscellaneous category.¹¹ I will determine which of these values are relevant for wrong proofs after discussing the values.

3.1 Didactics

I call a value didactic if it intends to teach. Mathematical proofs can be viewed through this lens. The values in this section are all centered around what and how a proof might teach its reader. Six values will be discussed here. These values concern proofs that demonstrate proof techniques, reconceive mathematical

¹¹These labels are used as a means of structuring the values. Some values might fit multiple categories.

domains, promote discovery, are explanatory, are motivated and provide visual intuition.

3.1.1 Demonstrating proof techniques

A proof can teach a multitude of things, one of which is simply showing the tricks up it's sleeve. As the methods used can be learnt and might be useful in other proofs than only the particular one. Mathematicians are “learning new proof techniques” [Weber, 2010, p32] when reading proofs, thus they “add [techniques] to their mathematical toolbox” [Morris, 2022, p4]. This is highly valuable for mathematicians.

This value can be different for each reader. The proof of the Bolzano-Weierstraß theorem where the bounded set gets split in halves consecutively could demonstrate a new technique for a student of mathematics. This proof is thus quite valuable for this student. On the other hand, for a professor of analysis this technique is canonical. Thus the professor might not value the proof as much. The proof demonstrates a proof technique for both, but only the student learns something new. Another way proof techniques could be demonstrated is with multiple proofs together that “demonstrate the power of different methodologies” [Dawson, 2006, p277]. Each proof could use a different technique, or might conceive the subject differently. Comparison between the proofs could then provide value. This could also explain why some theorems are proven over and over again, like the Pythagorean theorem [Loomis, 1968]. This demonstration could be useful as mathematicians can then learn which methods are more effective. They could also better understand the caveats of each method, thus apply those more fittingly. This could lead to more effective proving in the future.

It could be said that this method of demonstration is not a proof value, but a value of *comparison* of proofs. I argue that it could also be prescribed to a single proof. In the case that a proof reproves a known theorem it provides the new methodology for comparison. Thus the possibility for this comparison exists because of that proof. Also a proof might explicitly not use a certain technique as often done in Intuitionism where only constructive proofs are accepted.

3.1.2 Reconceive mathematical domains

Proofs can also make mathematicians “reconceive mathematical domains” [Weber, 2010, p32]. When confronted with a proof that the mathematicians in Weber’s study did not find intuitive, they would read the proof “with an eye toward understanding why the author believed the claim was true.” [Weber, 2010, p32] This would sometimes result in a different understanding of the mathematical topic at hand, thus *reconceives* it. Also when an author is publishing one interesting proof after another one mathematician said

“I try to understand how he [sic] is understanding things. I want to know how he [sic] is thinking” [Weber2010]

thus actively searching for this lesson in the proof.

For example, the complex numbers are often described by the Cartesian coordinates. When one first learns about the number i it is generally in this context. When a proof uses polar coordinates instead, it provides a different conception of the complex numbers. Because of this reconception some techniques or formulae are easier to study, e.g. the sine curve. In this case the reconception is constituted in a proof technique, the use of polar coordinates. This does not always need to be the case.

3.1.3 Promote discovery

Understanding of how a proof works might also promote discovery of new theory [Morris, 2022]. Morris bases this on [De Villiers, 1990] as he shows how proofs could lead to theorems being generalised. Furthermore proofs could lead to understanding about the theorem's boundaries and definitions, displayed in [Lakatos, 1976], thus promoting a more precise newer theorem.

For example, the solution to a real-valued quadratic equation, $ax^2+bx+c=0$ is given by $x = \frac{-b \pm \sqrt{D}}{2a}$ with $D = b^2 - 4ac$ if $D > 0$. The proof of this lemma is done with real valued equations. This lemma can be generalised to the complex numbers by making the equations in the proof complex valued.¹² The proof of the real valued lemma thus promoted the discovery of a similar statement in the complex numbers where $D > 0$ is not required.

Although promoting discovery seems to be similar to both demonstrating proof techniques and reconceiving mathematical domains, I intend to highlight a different aspect. In particular through demonstration of proof techniques and reconceiving mathematical domains the reader learns new methods of thought or techniques that might be used separately from the proof, i.e. might be used for proving other theorems. However, when proofs motivate new theorems, there is a more direct connection between the proofs and the theorems. Essentially the proof gets modified to discover new theory. This is a major difference between these values.

3.1.4 Explanatoriness

All (correct) proofs show *that* a theorem is true, but some also show *why* it is true. Proofs from the latter category are called explanatory proofs. Explanatory proofs provide mathematical understanding [Morris, 2022]. Validation alone does not provide insight into why a theorem might be true. A striking example of a non-explanatory proof that only validates a theorem is the computer proof of the four colour theorem. This proof, after reducing the possible maps to a limited amount, simply checks if each map can be coloured with four colours. This method only validates the theorem, it does not explain why the theorem is true.

¹²There are small intricacies that need attention when doing so, but the belly of the proof remains the same.

Extensive research has gone into the notion of an explanatory proof as [Morris, 2022] summarizes. No consensus on this notion has been reached. Nonetheless, mathematicians still use the notion of explanatory proof.

Morris discusses accounts of mathematical explanation by Steiner [1978], Kitcher [1989] and Lange [2014]. Steiner says that the proof should be built around a characterizing property of a central object. This characterizing property is something like the symmetry of a sum, or the prime factorisation of a number [Steiner, 1978]. The proof should then use this property to prove the theorem. Furthermore, by slightly changing this property, one gets an array of similar mathematical theorems and proofs.

Kitcher needs the proof to unify reasoning to reduce the number of ultimate facts we need to accept [Morris, 2022]. In our example of the four colour theorem, we need to accept that each of those maps can be coloured with four colours. This is a vast array of facts one needs to accept to prove the theorem. Hence, according to Kircher’s account, this proof is not explanatory.

Lange demands a salient feature that the proof needs to explain [Morris, 2022]. This salient feature is something like a symmetry or a correspondence at the centre of the theorem. In order to be explanatory, the proof needs to show where this feature originates from. For example, in the case of a symmetry this could be done by constructing a symmetry in the proof. This is different from Steiner’s account. Steiner demands that the proof *uses* the property of an object in the theorem. Lange asks that the proof explains the origin of the salient feature, the proof does not need to be built upon this feature.

3.1.5 Motivatedness

A proof can also motivate it’s steps. In doing so it provides understanding about why the proof is as it is. How this may be interpreted can differ. One view put forward by [Polya, 1949] asks for steps to be *recognisably appropriate* [Morris, 2022]. This means that the purpose of a step in the proof should not only be appropriate but that this should be understandably so for it’s reader. This is quite a powerful notion as it recognises the communication between proof and reader.

An addition to that is given by Morris [2020]. This addition asks that it is made understandable how the author conceived that step. Thus it should communicate how would one go about making such a proof. This is a stronger notion of motivation as it not only demands understanding why something is introduced, but also the thought process behind creating it. Thus when you have the thought “this works, but how on earth would you think of that?” about a step in a proof, then it is not well motivated according to Morris [2020].

An interesting example would be when a proof introduces a very specific formula out of the blue that solves a problem previously prevented. One might call it recognisably appropriate as it does precisely what it needs to do, and recognisably so. At the same time it does not communicate how the formula was conceived. Thus according to Polya it is motivated, but with Morris’ addition it is not.

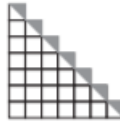
3.1.6 Visual intuitiveness

Another reason a proof can be valued is because it is intuitive. When a proof is intuitive it is easier to understand the proof. An intuitive proof presents arguments that correspond with how the reader thinks the proof should be continued. It presents a line of thought that comes naturally to the reader. This is a quite subjective and broad value thus we will look into a certain type of intuition, namely visual intuition.

Visual intuition is clearly used in mathematics in multiple ways. One is by giving a visual proof. A canonical example of this is the following.

Theorem 3.1. $\sum_{i=1}^n i = \frac{n^2}{2} + \frac{n}{2}$

Proof.



□

[Nelsen, 1993] This proof makes it visually intuitive why the equality is true.¹³ The visual suggests a line of thought that quite naturally provides a proof of the theorem.

Proofs do not need to be visual like the one above to appeal to visual intuition. It can include diagrams or mental representations of “particular geometric configurations.” [Cain, 2019, p2] This can add to the simplicity of the proof and even promote discovery [Cain, 2019].

3.2 Aesthetics

The study of aesthetics is a philosophical subject. It could be described as the philosophy of art, or a branch that studies the origin of beauty and taste. The values in this section are grouped because of their affiliation to this branch of study. It could be said that these values view proof as a form of art. The three values described below are two types of beauty, Platonist and Kantian, and the value of representation.

3.2.1 Platonist beauty

How can something so abstract as mathematics be aesthetic, or even beautiful? Giaquinto argues it can just like poetry. Poetry or a novel work the mind abstractly with it's words and bears a world to mind. Similarly mathematics

¹³This proof is also rather beautiful, more on that later on in the category of aesthetic values.

can work the mind. Also similarly to readers of poetry, mathematicians might recognise beauty too [Giaquinto, 2016].

This recognition of beauty can only be done by well enough educated mathematicians, it may not seem beautiful for any beholder. Giaquinto argues that one needs to be a *connoisseur* of mathematics, for this one needs to be open to the right perceptions and have enough background knowledge [Giaquinto, 2016].

This necessity of background knowledge is there for multiple reasons. One is that one needs other proofs for contrast [Giaquinto, 2016; Rota, 1997]. This is such that a general sense of beauty can be developed by comparing works, or proofs in this case, against one another.

The other reason is that beauty has a connection with understanding [Giaquinto, 2016] or enlightenment [Rota, 1997]. Rota [1997] even goes to say that understanding why a proof is beautiful could be hard work as one only realises it's beauty after the proof conveys it's enlightenment to the reader.

This sense of beauty could be called Platonist, as the beauty comes from the reader touching upon the perfection of the idea [Breitenbach, 2015]. The understanding of the mathematics seems to create the beauty of the proof, as if the proof is just the means to come to an end, and that end is providing it's beauty.

3.2.2 Kantian beauty

A different notion of beauty that could be used for mathematics comes from Kant [Breitenbach, 2015]. “according to the Kantian proposal the experience of beauty in mathematics is grounded in our felt awareness of the imaginative processes that lead to mathematical knowledge.” [Breitenbach, 2015, p955] This indicates that the process to create understanding itself is beautiful, instead of the understanding that it brings.

One might imagine this similarly to Rube Goldberg machines. These are chain reaction type machines that eventually fulfil only a basic end goal that could have been reached much simpler. When one looks at such a machine, the perceived beauty is in the chain reaction. Each step fitting perfectly together with the following. The focus is on the process, rather than the end result. The end result is usually rather dull with such a machine. With Kantian beauty in mathematics a similar thing is meant. Where the progress towards understanding by the steps of the proof contain the beauty. The process, rather than the result (understanding), contains the beauty.

3.2.3 Representation

Another type of aesthetic quality associated with proofs is the way a proof is presented. This does not refer to the font in which the proof is written or other superficial aspects of proof presentation. It refers to using symmetries or appealing to visuals [Ernest, 2016]. These are techniques that are sometimes aesthetically pleasing on their own, without considering the whole of the proof.

Theorem 3.1 is an example of this. It is not just beautiful, the symmetrical visual representation could be called aesthetic.

3.3 Miscellaneous

In this section the rest of the values this thesis discusses are presented. These values do not have enough coherence to be presented under a meaningful category. Hence they are placed in a residual category. Four values are discussed below. The values presented are purity, fitting transparency, fitting simplicity and systematization.

3.3.1 Purity

Purity is a broad concept in mathematics. One can already see it arise when talking about pure opposed to applied mathematics. This type of purity is quite old, Plato even distinguished between pure and applied branches. Where applied mathematics was good enough for tradesman but pure mathematics¹⁴ was for the mathematician [Ernest, 2016].

Purity in a more local form can be found in proofs, where purity can be found as to “enforce a certain symmetry between the conceptual resources used to prove a theorem and those needed for the clarification of its content.” [Detlefsen, 2008, p192-193] Thus when a theorem is written in analytic language, it is pure when the proof uses analysis.

This is quite similar to the more global purity. It distinguishes a need to stay within a discipline of mathematics or even within a certain modality. It demands a proof to keep pure from unnecessary interference. This could aid to the accessibility of a proof as mathematicians are often specialized in (a part of) one branch of mathematics.

In their paper, Raman-Sundström and Öhman [2018] describe a sub-type of direct fit, namely coherence. This requires the proof to use the same language and conceptual ideas as the theorem. This is a different way of characterizing purity.

3.3.2 Fitting transparency

As mentioned, a very closely related notion is that of fit. [Raman-Sundström, Öhman, 2018] subdivided this notion into six different sub-types as explained by Morris [2022]. One subcategory is that of **transparency**. A proof is transparent when it is “clear, without any deus ex machina steps” [Morris, 2022, p14]. The idea of the proof should be clear to the reader. This does not mean that one should understand how each step could be conceived like with motivatedness. It rather asks the broad outlines of a proof to be transparent, even though the individual steps might be confusing at first.

¹⁴In this case number theory.

3.3.3 Fitting simplicity

Another thing that could be valuable in a proof is its fitting simplicity, to not use “a sledgehammer to crack a nut.” [Gaiquinto, 2016, p63] This idea of a proof to be fitting to a theorem runs deep in the mathematical community. This is another sub-type of direct fit according to [Raman-Sundström, Öhman, 2018]. One where the criterion of specificity is met “when the tools used by the proof are of the appropriate power.” [Morris, 2022, p14]

For example, when one wants to prove the trivial statement that $n^2 > n$ for $n > 1$. One could simply take $n^2 = n * n > n * 1 = n$. This would be a fitting proof. When one would prove by stating that $n^2 = n$ at $n = 1$, and $\frac{dn^2}{dn} = n > 1 = \frac{dn}{dn}$ for $n > 1$. This would not be of the appropriate power. It uses a relatively strong method (comparison of derivatives beyond the sole intersection) for a simple statement ($n^2 > n$).

3.3.4 Systematization

Proofs that help the systematization of mathematics are also valued. Such proofs help create a common language and grounds for mathematical theories [Morris, 2022]. Morris also explains that this is true for any proof by nature of them being proofs. As even when one uses common notation and introduces nothing new one solidifies the notations and concepts used, thus leading to systematization. For example, using epsilon-delta style definitions of continuity further solidifies this method of defining continuity as canonical mathematical practice.

3.4 Potential values in wrong proof

As mentioned, we aim to investigate the potential value of wrong proofs. In the previous subsections, we looked at a variety of values, other than validity, associated with proofs. It might be that even though these values are not the same as validity, that they are highly correlated with validity. In that case, it makes no sense to investigate these values from wrong proofs. As such, this subsection is dedicated to determining whether it, *a priori*, makes sense to investigate these values in the context of wrong proofs.

We will start with the didactic values, then continue to the aesthetic values and finally discuss the values of unity.

3.4.1 Didactics

All aforementioned didactic values are of interest in the value analysis of a wrong proof. This is as none of these values are linked to invalidity in such a way that it is fundamentally impossible to appear in a wrong proof. Let us view this per value.

If a proof is wrong, then it is still possible for the proof to demonstrate mathematical methods. This is as proof techniques could be interesting, even if

they are used wrongly. Furthermore, wrong proofs do not only consist of wrong steps.

With respect to reconceiving a mathematical domain, it might be that the proof was written with a conception of the mathematical domain in mind. This conception can still be recognised in a wrong proof. Thus it might reconceive the readers' conception of the domain.

When it comes to promoting discovery Lakatos [1976] shows in *Proofs and Refutations* that wrong proof promotes discovery. In this work it is actually fundamental to the process of discovery that wrong proofs get presented and analysed.

Can an argument be explanatory even though it is wrong? I argue that it can in an adapted way. In science education other than mathematics this is quite common to do. Simply look at a chemistry textbook where electrons are presented in layers, or in physics where Newton's laws are used as a good enough approximation even though relativity effects are not calculated. In our education many things are simplified at first, in order to create a base of intuition and understanding, to then slowly add layers of truth to create a full, comprehensible picture.

Mathematics may seem different. But as portrayed in the previous chapter, the line between correct and wrong is somewhat blurred. Thus incomplete or wrong proofs could portray a method of thought useful for learning. In effect more comprehension may be gained about the theorem at hand.

Though, in mathematics the devil is often in the details. Sometimes a wrong proof might seem very near a correct proof, even though the seemingly small fault is so fundamental that the proof cannot be consolidated in any way. In this case it is not explanatory. The reason why a proof is faulty is thus very important for this value. In the analysis of a wrong proof one has to first consider why the proof is faulty before analysing explanatoriness.

Motivatedness is not influenced by invalidity. A wrong proof can still motivate its steps similarly to a correct proof. Even for an incorrect step the proof may still recognisably present why the author thought this step appropriate.

Finally, proofs can also use visually intuitive techniques. Sometimes the wrong step coincides with the visual intuitive technique. In this case it is important whether the step is blatantly wrong, or whether its principle is correct but applied wrong. If the principle is correct the visual intuition is still valuable. If the step is fundamentally wrong, the visual intuition is thus wrong too.

Thus as we have seen, all didactic values are interesting to use in the analysis of a wrong proof.

3.4.2 Aesthetics

For the aesthetic values only Platonist beauty is too intricately linked to validity to appear in a wrong proof. The values of Kantian beauty and representation can be intelligibly analysed in a wrong proof.

As discussed in Section 3.2.1, for Platonist type beauty the proof needs to provide understanding according to Giaquinto [2016] or enlightenment according

to Rota [1997]. When a connoisseur of mathematics [Giaquinto, 2016] reads a wrong proof, he would recognise the fault. This breaks the possibility of touching upon the perfection of the idea. The proof does not provide a path to enlightenment nor understanding. Thus it is not Platonist-type beautiful.

The view of Kant on mathematical beauty does comply with the beauty of a wrong proof. This is as this highlights the process towards understanding, which can still contain beautiful steps and reasoning even though some could be incomplete or wrong.

In this case though the reason why a proof is wrong is important. If there is a blatantly wrong step that a connoisseur can recognise, it might be the case that the reasoning is considered rather ugly. When the reasoning simply makes a few too many leaps and thus is incomplete or contains small incorrect steps, the whole may still be a beautiful argument.

The way arguments are represented is not deeply dependent on validity. Representation largely depends on the argument. In what context that argument is placed does not influence if the argument itself is aesthetic. Thus if the context is an incorrect proof, the argument itself might still be aesthetic.

3.4.3 Miscellaneous

In the miscellaneous category only purity and systematization are found to fully be interesting for the analysis of a wrong proof, fitting transparency and fitting simplicity are too linked to validity to be interesting. Below the reasoning for this is presented.

A wrong proof can use the same “conceptual resources” [Detlefsen, 2008, p192] as the theorem it tries to prove. This is the requirement for a proof to be pure, thus a wrong proof can be pure.

An argument could be made that the wrong step is conceptually completely different. The wrong step is, because it is wrong, not a mathematical resource. It is rather something completely different. By this reasoning one could say such a proof might *try* to be pure, but is impure. The rest of the proof might still be pure aside from this blimp, but not the whole.

It could also be argued that such wrong step is still pure. As it might use the same conceptual resource as the theorem, but applies it wrongly. In this light such a fault is conceptually not wholly different. Thus the proof could be considered pure. This view would lead to more interesting analysis on wrong proof. Because of this I will choose this viewpoint for the analysis of a wrong proof.

For fitting transparency it is required that a proof is “without any *deus ex machina* steps” [Morris, 2022, p14]. A wrong step is certainly such a step. Because of this a wrong proof fails to be fittingly transparent.

One could say that in a wrong proof “the tools used by the proof are [not] of the appropriate power” [Morris, 2022, p14]. They do not prove the theorem, thus are not sufficiently strong. Thus the proof is not fittingly simple.

Using existing definitions, notations and concepts, or introducing new ones, is not reliant on validity of the proof as a whole. When the definitions, nota-

tions and concept are used or presented properly it does aid in systematizing mathematics. Thus a wrong proof might uphold the value of systematization similarly to a correct one.

4 The fundamental theorem of algebra

In this chapter the proofs will be introduced that will be used as case for the case study. This chapter will only present the mathematics for understanding. No interpretations of values will be given here, that will be done in the chapter hereafter.

4.1 A proof and a gap

In order to understand the values a wrong proof might have one obviously first needs a wrong proof. But such a proof in isolation would have no comparative material. The choice could be made of comparing it to proofs in general, but that still would leave a few gaps in reasoning. The theorem set out to prove might not permit certain values to come to light in the proof. For this reason a proof chosen for this task ideally is one where there is a version deemed wrong, and an adaptation of the same reasoning but with added detail that constitutes a proof deemed proper, in order to find a comparison based on most similar cases. In the published literature not many of these examples are known or to be found, even more so ones where the mathematics is not too specialised.

Gauss' first proof of the Fundamental theorem of Algebra is such a proof. Easy enough to understand, but with mistake. The mistake that later was deemed wrong is an untraversed gap at the end of the proof. This gap left could also be considered enthymematic as the reasoning Gauss presented to overcome the gap is that it "is known from higher geometry" [Struik, 2014, p121] which was not the case at the time.

Gauss himself knew he was leaving this gap as he wrote "As far as I know, nobody has raised any doubts about this. However, should someone demand it then I will undertake to give a proof that is not subject to any doubt, on some other occasion." [Struik, 2014, p121] It took many years after his 1799 dissertation before this gap was filled by [Ostrowski, 1920]. Gauss has given a different proof of the theorem later in his career, but never filled the gap that he left. There have been a couple of mathematicians that have covered the gap in published papers ([Noel, 1991], [Gersten, 1988], [Basu, Velleman, 2017] and [Cain, 2005] to name a few). The one chosen here though is [Ostrowski, 1920]. There are multiple reasons to choose Ostrowski's proof. The most straight forward one would be that he was the first to traverse the gap. This although is not the most important one. A different reason for this version of the proof is that he stays true to Gauss' method of reasoning and does not bring different fields of reasoning into play. In his own words "Insofern sind die Hilfsmittel, die wir in Folgenden benutzen, alle der Art, dass sie auch Gauss im Prinzip

bekannt gewesen sein mögen.”¹⁵ [Ostrowski, 1980] This way his proof is a good comparison to the original one and thus the conclusions based off this comparison will be grounded better.

In the following section Gauss’ proof will be explained, the section after Ostrowski’s adaptation of this proof will be given. In the next chapter the values of these proofs will be analysed to give a first case study of what values wrong proofs could hold.

4.2 Gauss’ first proof

In this section I will explain Gauss’ proof of the Fundamental theorem of Algebra. For this I have used [struik, 2014] a lot as it contains an English translation of the proof Gauss gave. After giving the credit here it will no longer be mentioned in the rest of this section as that would be superfluous and would reduce readability.

4.2.1 Preparatory lemmas

The equations in this theorem and the corresponding lemmas are all considered in $\mathbf{R}[x]$ except when explicitly stated otherwise. This is kind of peculiar when you consider the modern formulation of the Fundamental Theorem of Algebra (FTA) which is a statement about the complex polynomials. The formulation that Gauss used will be given later on and there will be a quick explanation as to why this is equivalent to the modern formulation.

Gauss starts with a arbitrary polynomial where $A, B, \dots, M \in \mathbf{R}$

$$\mathbf{X} = x^m + Ax^{m-1} + Bx^{m-2} + \dots + Lx + M.$$

Then he begins with two subsequent lemmas. The sin and cos in these lemmas may come out of the blue. These are more understandable when one looks at the polar notation of complex numbers, $r(\cos \phi + i \sin \phi)$ ¹⁶. Later I will clarify why specifically these lemmas are of interest. For now I will present the two lemmas that Gauss specifically states. Namely:

Lemma 4.1. Given a $m \in \mathbf{Z}$ positive. Then

$$\sin(\phi)x^m - \sin(m\phi)r^{m-1}x + \sin((m-1)\phi)r^m$$

is divisible by $x^2 - 2 \cos(\phi)rx + r^2$.

Proof. Gauss doesn’t give a proof for this other than mentioning “The proof is given by direct division.” For the reader the lemma can be more easily verified

¹⁵Translation: The tools that we use following, are all of the type that Gauss may also have known in principle.

¹⁶Which Gauss explicitly does not use.

by multiplication of factors when given the two factors resulting in the division: $x^2 - 2 \cos(\phi)rx + r^2$ and

$$\sum_{i=1}^{m-1} \sin(i\phi)r^{i-1}x^{m-i-\phi}$$

for $m > 2$. For $m = 2$ the other factor is $\sin \phi$ and for $m = 1$ the resulting factor is 0 [dawson, 2015]. The multiplication of factors is not very insightful and thus I will not include it here just as Gauss did. \square

Lemma 4.2. The polynomial \mathbf{X} will be divisible by $x^2 - 2 \cos(\phi)rm + r^2$ or by $x - r \cos \phi$ if $r \sin \phi = 0$. If there exist $r, \phi \in \mathbf{R}$ such that

$$r^m \cos(m\phi) + Ar^{m-1} \cos((m-1)\phi) + \dots + Lr \cos(\phi) + M = 0 \quad (1)$$

$$r^m \sin(m\phi) + Ar^{m-1} \sin((m-1)\phi) + \dots + Lr \sin(\phi) = 0 \quad (2)$$

To quickly explain where Equations (1) and (2) originate one can consider substituting $x = r(\cos \phi + i \sin \phi)$ in \mathbf{X} . Then together with the rule derived from de Moivre's formula that $(\cos \phi + i \sin \phi)^m = \cos(m\phi) + i \sin(m\phi)$ it follows that these equations are the real and imaginary parts of \mathbf{X} respectively.

Gauss does not (want) to use imaginaries himself. I included them here solely for our understanding of the formulas. This second lemma is previously given in Euler's *Introductio in analysin infinitorum* with the use of imaginaries as Gauss mentions in his dissertation. Gauss then continues to prove the lemma without the use of imaginary numbers through the following method;

Proof. Consider the case where $r \sin \phi = 0$. If $r = 0$ then 0 is a root of \mathbf{X} as $M = 0$ by (1). Else $\cos \phi = \pm 1$ and it follows that $\cos(m\phi) = \cos(\phi)^m = \pm 1$ in this case. Because of this $r \cos \phi$ is a root of \mathbf{X} since substituting it for x gives (1). Thus $x - r \cos \phi$ divides \mathbf{X} .

Now consider the case where $r \sin \phi \neq 0$ and look at the following addition;

$$\begin{array}{r r r r} \sin(\phi)rx^m & - & \sin(m\phi)r^m x & + \sin((m-1)\phi)r^{m+1} \\ A \sin(\phi)rx^{m-1} & - & A \sin((m-1)\phi)r^{m-1}x & + A \sin((m-2)\phi)r^m \\ B \sin(\phi)rx^{m-2} & - & B \sin((m-2)\phi)r^{m-2}x & + B \sin((m-3)\phi)r^{m-1} \\ \vdots & - & \vdots & + \vdots \\ K \sin(\phi)rx^2 & - & L \sin(2\phi)r^2 x & + K \sin(\phi)r^3 \\ L \sin(\phi)rx & - & L \sin(\phi)rx & + 0 \\ + M \sin(\phi)r & - & 0 & + M \sin(-\phi)r \\ \hline \sin(\phi)r\mathbf{X} & - & 0 & + 0 \end{array}$$

Each row in the addition is of the desired form to use Lemma 4.1. This implies that each row, and therefore the whole summation, is divisible by $x^2 - 2 \cos(\phi)rm + r^2$. The middle and rightmost column are both zero by the

assumptions of the lemma. In the leftmost column one can recognise \mathbf{X} multiplied by $\sin(\phi)r$ as shown after the addition.

As the each row in the sum is divisible by $x^2 - 2\cos(\phi)rm + r^2$ even though $\sin(\phi)r$ is not, it follows that \mathbf{X} is divisible by $x^2 - 2\cos(\phi)rm + r^2$ as stated by the lemma. \square

4.2.2 The proof

In the previous subsection we have seen two lemmas Gauss used in his proof that he specifically labeled. Before we continue I will clarify as to why these lemmas are of interest for the proof, though firstly I will first state the main theorem.

Theorem 4.3 (The Fundamental theorem of Algebra, FTA).

Any polynomial $\mathbf{X} \in \mathbf{R}[x]$ with real coefficients can be factored into linear and quadratic factors [cain, 2005].

This is written in a somewhat different form than the modern use mostly states. The modern version states that *every polynomial with complex coefficients can be factored into linear factors*. These two formulations are equivalent by a short reasoning, namely: Given a polynomial $g \in \mathbf{C}[x]$ then g multiplied by its complex conjugate \bar{g} yields a real valued polynomial $f = g\bar{g}$ which can be factored by the theorem. This shows an easy setup to proving the equivalence of the two [cain, 2005].

The Lemmas 4.1, 4.2 are of interest for the theorem because the lemmas provide a linear or quadratic factor of \mathbf{X} if the conditions of the lemmas are satisfied. Gauss attempts to do so by finding the intersection of $U = 0$ and $V = 0$. He does this by showing that such an intersection lies within a circle around zero. After he attempts to show that within this circle an intersection is inevitable, it is within this part that there lies a gap. Let us follow Gauss in his reasoning for this proof:

Proof. Given a polynomial \mathbf{X} as before. Gauss now takes the time to expand on a real valued field with polar coordinates. At this time we would simply refer to the polar coordinates of the complex field, but as Gauss was trying to show it can be done without imaginaries he has to take this extra step. I will not expand further on this field as the formulations are quite archaic and un insightful for the modern reader.

The effect is the same as the original polynomial gets split into two separate equations:

$$\begin{aligned} r^m \cos(m\phi) + Ar^{m-1} \cos((m-1)\phi) + \dots + Lr \cos(\phi) + M &= U \\ r^m \sin(m\phi) + Ar^{m-1} \sin((m-1)\phi) + \dots + Lr \sin(\phi) &= T, \end{aligned}$$

which we call U and T respectively. Here $\mathbf{X} = U + iT$ in our modern notation. This notation would immediately tell us that if we find a coordinate with U and T both 0, then \mathbf{X} would automatically be 0 at this point and thus have

a (complex) linear component, surpassing the need of Lemma 4.2. As Gauss proves the theorem without the complex numbers the relation of U and T to \mathbf{X} is given by the lemma instead. Thus we now have to, by reasoning through Lemma 4.2, find the intersection of the curves $U = 0$ and $T = 0$ in order to prove the theorem. To find an intersection we first need to understand the trajectories of the curves.

For $T = 0$ there are $2m$ trajectories that tend toward infinity. This is because the curve behaves asymptotically to the straight lines originating in the origin with angles $\frac{k\pi}{m}$ with $0 \leq k < 2m$. This asymptotic behaviour originates in the fact that \mathbf{X} is dominated by the factor x^m as x tends to infinity. Rewriting it into polar coordinates again gives a dominating factor of $r^m \sin(m\phi)$ for T . As ϕ tends towards any of the angles $\frac{k\pi}{m}$, $r^m \sin(m\phi)$ tends to zero. This explains the asymptotic behaviour of $T = 0$. Specifically interesting is a trajectory along the x-axis where $\phi = 0$ and r ranges from $(-\infty, \infty)$. (This joins the trajectories with $\phi = 0$ and $\phi = \pi$) [dawson, 2015]. We will later see how this trajectory is relevant in Gauss' argument.

For $U = 0$ one can follow a similar reasoning to find $2m$ different trajectories tending towards infinity because of the dominating factor $r^m \cos(m\phi)$. These are constituted along the lines with angles $\frac{(2k-1)\pi}{2m}$ for $0 < k \leq 2m$ [dawson, 2015].

Next we introduce a circle around the origin. This circle has to be large enough such that each trajectory of $U = 0$ and $T = 0$ is close enough to their asymptote. By 'close enough' one can choose for each trajectory to be nearer to it's own asymptote than to a different one on the circle. Gauss goes into detail as to how big the radius of the circle needs to be for this to happen. This is quite a tedious technical argument that the interested reader can find in [struik2014], Ostrowski's formulation of this argument is given in the next section as this is a bit more concise but contains the same ideas and methods.

By viewing the asymptotic behavior we can now realise that the trajectories alternately intersect the circle. This is seen by the fact that $\frac{k\pi}{m}$ and $\frac{(2k-1)\pi}{2m}$ alternate as k increases. Each zero of U on the circle lies in-between two zeros of T and vice versa.

Now Gauss formulates the following conclusion; "if a branch of an algebraic curve enters into a limited space, it necessarily has to leave it again." [Struik, 2014, p121] Gauss does not prove this fact but merely states that "this is known from higher geometry" [Struik, 2014, p121]. We will later see in the formalisation of Ostrowski that this is indeed the case but this took until 1920 to prove. This formalisation was also not nearly as trivial as Gauss made it look in his 1799 dissertation.

Next Gauss makes an argument that the curves must necessarily intersect within the circle. I will formulate this as a separate lemma, namely Lemma 4.4. When this lemma is proven the theorem can be considered proven as the necessary conditions for Lemma 4.2 are satisfied, apart from the previously mentioned gap of course. \square

Lemma 4.4. Given two algebraic curves $U = 0$, $T = 0$, these intersect if branches of these curves intersect a circle in an alternating fashion and each continuous branch of $U = 0$ and $T = 0$ that enters the circle necessarily has to leave it again through a different intersection.

Proof. We will first assume that the curves do not intersect, and later will show that this leads to contradiction. Now assume there are $2m$ intersections with the circle. We can assume m an even number as each branch into the circle has to leave it again for each curve $U = 0$ and $T = 0$ by assumption. Number these intersections in a clockwise fashion such that the intersection clockwise of k is $(k + 1) \bmod(2m)$.

All even intersections belong to the same algebraic curve, let this be $U = 0$ without loss of generality. The same is true for all odd intersections with the curve $T = 0$. This is because of the alternating intersections on the circle.

The branch through intersection 0 has to leave the circle at an even intersection k_0 . Then intersection 1 has to leave the circle at an odd intersection k_1 smaller than k_0 . When $k_1 > k_0$ then the trajectories intersect and this leads to contradiction.

In the same fashion every branch starting at intersection j will exit the circle at an intersection $k_j < k_{j-1}$. This will result in a branch B entering the circle at intersection l and exiting it at $l + 2$. Because of this the branch through $l + 1$ will have to cross the branch B in order to leave the circle through a different intersection. Thus leading to contradiction and proving the lemma. □

At the start of the proof for this lemma Gauss starts with a specific branch rather than a random one. He showed that $U = 0$ always has a straight branch through $\phi = 0$ and $\phi = \pi$. This is not necessary to complete the argument as shown above, though it may provide more structure to the reader.

4.3 Ostrowski's adaptation

In this section the proof that Ostrowski gave will be explained. This proof is an adaptation of Gauss' proof published in 1799. The goal of this adaptation is to fill the gap left by Gauss in his proof. For this proof some modern notation is used, e.g. imaginaries, but it will be limited to what Gauss at that time could have known as Ostrowski mentions himself [Ostrowski, 1920]. The use of imaginaries is mostly to clarify for the modern reader and is only used mathematically to avoid Lemma 4.2. Throughout this part most information will originate from [Ostrowski, 1920], after mentioning it here it will be left out to improve readability.

Ostrowski split his paper into six parts, and the same parts will be used here for clarity. The sixth part is left out as Ostrowski proves extra insights into the structures of the mathematical object there which are not necessary for the main proof. The fact that these insights can be gained is important for the analysis in the next chapter, the very nature of the proof of these insights are not.

The first two parts will contain analysis of the structure of $T = \Re(\mathbf{X})$ and $U = \Im(\mathbf{X})$ outside of, and on a large enough circle for a random $\mathbf{X} \in \mathbf{C}[z]$. These are mostly modern versions of the arguments Gauss also made. In parts three until five the structure inside the circle is analysed. Here it is also proven that each branch of the curves necessarily has to leave the circle again in a continuous fashion, thus filling Gauss' gap.

4.3.1 Part one

Let $\mathbf{X}(z) = z^m + Az^{m-1} + \dots + C$ just as in Gauss' proof but then in the complex plane (thus using z instead of x). Also let $\mathbf{X}(z) = T(r, \varphi) + iU(r, \varphi)$ in polar coordinates. Given these definitions Ostrowski starts with the following lemma.

Lemma 4.5. Divide the plane in $4m$ parts by using $\omega = \frac{\pi}{4m}$ to define the borders $\varphi = -\omega$, $\varphi = \omega$, $\varphi = 3\omega$, Name these parts (1), (2), ..., (4m) accordingly. Then there exists a radius R such that

- (i) $\forall z : |z| \geq R : \mathbf{X}(z) \neq 0$
- (ii) For the φ intervals of (\bar{R}) for all k in $[0, \dots, m-1]$ of an arbitrary circle (\bar{R}) with radius $\bar{R} \geq R$ with the following statements are true for T and U ;
 - (ii.i) On $(4k+1)$, T and $\frac{\partial U}{\partial \varphi}$ are positive and U is first negative and then positive. This transition happens only once because $\frac{\partial U}{\partial \varphi} > 0$.
 - (ii.ii) On $(4k+2)$, U is positive, $\frac{\partial T}{\partial \varphi}$ is negative and T is first positive and then negative. This transition also happens only once.
 - (ii.iii) On $(4k+3)$, T and $\frac{\partial U}{\partial \varphi}$ are negative and U is first positive and then negative. This transition happens only once.
 - (ii.iv) On $(4k+4)$, U is negative, $\frac{\partial T}{\partial \varphi}$ is positive and T is first negative then positive. This transition happens only once.

Gauss found the same constant R by a similar argument but with different notations in his dissertation. I will include the proof of Ostrowski even though I excluded Gauss' proof. The reason for this is that this version is easy to follow for the modern reader opposed to Gauss' version, even though they are in essence the same.

Proof. Let $S = |A| + |B| + \dots + |C|$ be the sum of the absolute values of the coefficients of \mathbf{X} . Then $R = \sqrt{2}S + 1$ has the right properties. For an arbitrary $\bar{R} \geq R$:

$$|T - \bar{R}^m \cos(n\varphi)| \leq |A|\bar{R}^{m-1} + \dots + |C| \leq \bar{R}^{m-1}S < \sqrt{\frac{1}{2}}\bar{R}^m$$

$$|\frac{\partial U}{\partial \varphi} - m\bar{R}^m \cos(m\varphi)| \leq (n-1)|A|\bar{R}^{m-1} + \dots < m\bar{R}^{m-1}S < m\sqrt{\frac{1}{2}}\bar{R}^m$$

These inequalities are easy to understand given

$$\begin{aligned} r^m \cos(m\varphi) + Ar^{m-1} \cos((m-1)\varphi) + \dots + C &= U \\ r^m \sin(m\varphi) + Ar^{m-1} \sin((m-1)\varphi) + \dots &= T \end{aligned}$$

These inequalities give rise to the following observation. $|T - \pm\sqrt{\frac{1}{2}}\bar{R}^n| < \sqrt{\frac{1}{2}}\bar{R}^n$ if $|\cos(m\varphi)| \leq \sqrt{\frac{1}{2}}$ and thus T has to have the same symbol as $\cos(m\varphi)$ in these areas. Same for $\frac{\partial U}{\partial \varphi}$ having the same symbol as $\cos(m\varphi)$. In the uneven intervals of φ this is the case. With T and $\frac{\partial U}{\partial \varphi}$ being positive at $(4k+1)$ and negative at $(4k+3)$.

Through an analogous inequality one can find that the sign of $\sin(m\varphi)$ is the same as the sign of U and opposite of the sign of $\frac{\partial T}{\partial \varphi}$ when $|\sin(m\varphi)| > \sqrt{\frac{1}{2}}$. Thus concluding the argument for the lemma. \square

This lemma gives us the existence of the desired circle and the behaviour of T and U on this circle. Namely the alternating zeros of T and U as Gauss proved before. Also we have some structure of the derivatives with respect to φ . These derivatives show the property that on the zeros of U on one side U is positive, and on the other side negative. This is also true for T . These facts will be important for the proof later on.

4.3.2 Part two

Most arguments presented will be about U . T is mostly left out of the picture. The reason for this is that the analysis of T is symmetrical in nature and writing everything double is superfluous.

This part is about the behaviour of $U = 0$ outside of the circle. The important conclusions are summed up in the following lemma. This lemma was mentioned by Gauss and proven in a shorthand way. Ostrowski will elaborate on it for longer and prove an additional fact.

Lemma 4.6. Outside the circle $U = 0$ consists of $2n$ branches in the uneven intervals. These branches are continuous curves which intersect every circle with radius $\bar{R} \geq R$ in a certain (uneven) interval of φ only once.

Proof. This whole argument will be in polar coordinates (r, φ) . Take an arbitrary uneven interval. Let this be (1) without loss of generality. By Lemma 4.5 there is only one zero of U on this φ -interval for a circle with radius $\bar{R} \geq R$.

We now take interest in the zero (w, l) of U with $w = \bar{R}$ in interval (1). Then we want to prove that $\forall \varepsilon \exists \delta : (w' \leq R) \wedge (w - \delta < w' < w + \delta) \Rightarrow |l - l(w')| < \varepsilon$. Where $(w', l(w'))$ is the zero of the circle with radius w' in interval (1). This would prove the lemma because all the zeros within interval (1) are shown to be part of one continuous curve.

Now the fact that U is positive on one and negative on the other side of its zero comes into play. Because of this fact there exists a sufficiently small $\varepsilon' < \varepsilon$

such that

$$U(w, l - \varepsilon') < 0, \quad U(w, l + \varepsilon') > 0$$

and both coördinates are within interval (1).

Seeing that U is continuous in r and φ there subsequently exists a δ such that

$$U(w', l - \varepsilon') < 0, \quad U(w', l + \varepsilon') > 0$$

when $|w - w'| < \delta$. Because this gives a positive on one side and a negative on the other, together with U being continuous with respect to φ give a zero $(w', l(w'))$ with $|l(w') - l| < \varepsilon$, thus proving the lemma. \square

4.3.3 Part three

This section is written, differently from the previous section, in Cartesian coordinates (x, y) . This part onward Ostrowski deviates clearly from Gauss' proof. First an important concept will be introduced, namely multiple factors.

Definition 4.7. A polynomial f in variables x_1, \dots, x_n has a multiple factor h when there exists a non-constant polynomial h and g such that $f = h^k g$ for a $k \geq 2$.

This definition is important as parts of the proof only work on polynomials without multiple factors. Ostrowski did not give a definition of multiple factors, the one added here is derived from how Ostrowski uses multiple factors in his proof.

For our purposes we can divide the possible multiple factors out of U to create a polynomial $U^{(1)}$. This does not interfere with the shape of $U = 0$ as only multiple zeros are divided out of the polynomial. Divide the parts shaped $x - x'$ for constants x' out to remain with a polynomial $U^{(0)}$, $U^{(0)} = U^{(1)}$ if these factors don't exist. $U^{(0)} = 0$ thus only differs from $U = 0$ in structure by vertical trajectories. Ostrowski mentions that it is possible to show that U does not contain multiple factors or factors of the shape $x - x'$ but it is not necessary for the argument to show this. Additionally this way the argument is more general and also could be applied to any algebraic or simple analytic curves, not only polynomials [Ostrowski1920].

What Ostrowski proves in this section is that $U = 0$ can be split into a finite number of what he calls elementary curves. These elementary curves are continuous pieces of curves with no gaps. He does so by showing that $U^{(0)}$ can be split into elementary curves and adding the possible finite number of pieces of $x - x'$ after.

To do this collect all the values of x where both $U^{(0)}$ and $\frac{\partial U^{(0)}}{\partial y}$ or both $U^{(0)}$ and $\frac{\partial U^{(0)}}{\partial x}$ are zero at the same value (x, y) . There are only a finite number of these because of the following algebraic lemma:

Lemma 4.8. A polynomial in two variables has multiple factors if it has an infinite number of coördinates at which both the polynomial and it's derivative with respect to one of the variables are zero.

Ostrowski mentions that this is easy to prove with Euclid's algorithm and a different small lemma. He does not prove this fact.

The lemma does immediately show that only finite of the above mentioned points exist, else $U^{(0)}$ has a multiple factor, which is by definition not the case.

Then also collect the values of x at which $U^{(0)} = 0$ intersects the circle (a finite amount as proven in part one) and the values x' where the parts $x - x'$ could be divided out of $U^{(1)}$ (finite because $U^{(1)}$ is of finite degree). Now we have a finite collection of values of the before mentioned and those x -values where $U^{(0)}$ and $\frac{\partial U^{(0)}}{\partial y}$ are both zero or $U^{(0)}$ and $\frac{\partial U^{(0)}}{\partial x}$ are both zero. Order all these collected x -values from leftmost to rightmost in size $x_1 < x_2 < \dots < x_k$. The interior of the circle (R) can now be split up into a finite number of slices by $x = x_1, x = x_2, \dots, x = x_k$.

The argument that follows next is one about a slice picked at random. We will use the one between $x = x_1$ and $x = x_2$ as Ostrowski did as this can be done without loss of generality. Lastly in this section we will prove the following lemma:

Lemma 4.9. Every vertical at $x \in (x_1, x_2)$ has n zeros of $U^{(0)}$ within the circle. If $U^{(0)}$ has n zeros $\eta_1, \eta_2, \dots, \eta_n$ on a certain $x = \bar{x}$ with $x_1 < \bar{x} < x_2$ within the circle (R).

Proof. The proof starts by first proving that for a \bar{x} that has n zeros within the circle that for all $\varepsilon > 0$ there exists a $\delta > 0$ such that all x with $\bar{x} - \delta < x < \bar{x} + \delta$ have m zeros and these zeros are within a radius of 2ε of the zeros $\eta_1, \eta_2, \dots, \eta_n$ of \bar{x} . Then after he goes to show that an x with $x_1 < x < x_2$ that has a different number of zeros doesn't exist because the supremum (or infimum) of this set of zeros cannot exist. Let us start with the first section of the proof.

Pick a $\delta' < \varepsilon$ small enough such that

$$U^{(0)}(\bar{x}, \eta_i - \delta') < 0, U^{(0)}(\bar{x}, \eta_i + \delta') > 0 \quad \forall i \in [n].$$

Because $U^{(0)}$ is continuous in x there exists a $\delta < \varepsilon$ such that when $|x - \bar{x}| < \delta$

$$U^{(0)}(x, \eta_i - \delta') < 0, U^{(0)}(x, \eta_i + \delta') > 0 \quad \forall i \in [n].$$

This delta needs to be small enough such that $\bar{x} - \delta$ and $\bar{x} + \delta$ stay within the interval (x_1, x_2) . Then for every such an x there is a zero η'_i within $\eta_i - \delta$ and $\eta_i + \delta$ of which the distance $|(x, \eta'_i) - (\bar{x}, \eta_i)| < \delta + \delta' < 2\varepsilon$.

To show that there's one and only one zero close to each η_i for each x one needs to realise that $\frac{\partial U^{(0)}}{\partial y}$ is non-zero in all these η_i 's. This is because the zeros where this derivative is zero were part of the $\{x_1, x_2, \dots, x_k\}$. Now take $\varepsilon > 0$ small enough such that $\frac{\partial U^{(0)}}{\partial y}$ is non-zero in the 2ε radius of each η_i . Then by Rolle's theorem (if f is uniformly continuous, $a \neq b$ and $f(a) = f(b)$, there exists a $c \in (a, b)$ such that $f'(c) = 0$) it is easy to see that every x value can mostly have one root of $U^{(0)}$ within this 2ε radius. Else there would be a point where $\frac{\partial U^{(0)}}{\partial y}$ is zero.

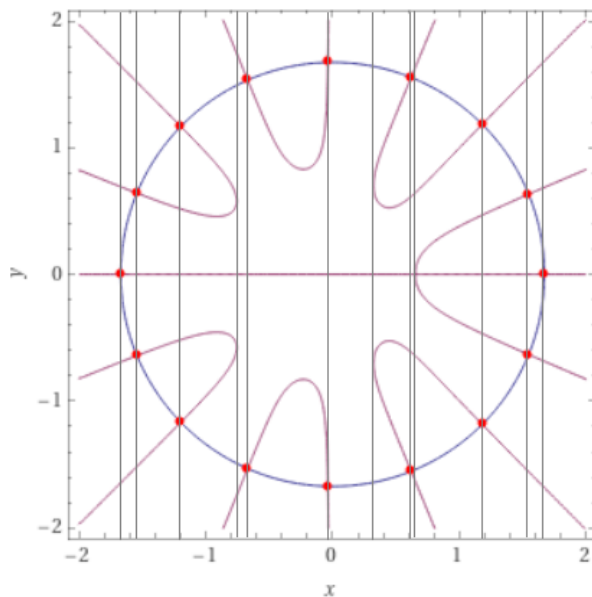


Figure 1: $U = 0$ for the function $\mathbf{X} = z^8 + 0.2z^7 - 0.1z^6 - 0.3z^5 - 0.1z^3 + 0.2z^2 - 0.3z + 0.1$) with $x = x_1, x = x_2, \dots, x = x_k$ plotted. Example function taken from [Basu, Velleman, 2016, p692]

For this to work in the extreme cases where η_i is very close to the circle (R), let ε be small enough such that the whole circle around η_i is within (R).

Ostrowski limits the ε and δ some more in order to exclude all possible anomalies to make the next part of the proof work. For a very detailed interest one can read [Ostrowski1920], for now only the general goals of these limitations will be explained. The additional limitations are there to make sure that only zeros within the circle are considered and not accidentally those outside (R) and that δ is small enough such that all zeros are within 2ε radius for close enough x values. This is mostly that in the neighbourhood of \bar{x} all zeros are behaving predictably. With these alterations we are ready for the next part, which is proving that all x in between x_1 and x_2 have n zeros of $U^{(0)}$.

First we will inspect the interval (x_1, \bar{x}) . The proof for the interval (\bar{x}, x_2) will be analog. Say there exist x in this interval with more or less than n zeros. Define \bar{x} be the supremum of these x -values with n_1 zeros. Then by reasoning earlier in the proof there should be a $\delta'' > 0$ such that all x with $\bar{x} < x < \bar{x} + \delta''$ have n_1 zeros. But these values of x have exactly n zeros because \bar{x} was the supremum. Thus no such \bar{x} exists and all $x \in (x_1, \bar{x})$ have exactly n zeros of $U^{(0)}$. Similarly the same is proven for the upper interval, and thus for whole of (x_1, x_2) . \square

4.3.4 Part four

In this part an effort is made towards proving the main theorem by building upon the previous found structure of zeros to show that $U = 0$ can be described by a finite number of elementary curves. To do this we first will show that all these n zeros in (x_1, x_2) make n elementary curves on $[x_1, x_2]$ and then analyse the structure of these curves and their surroundings. Again these arguments can be made for each $[x_i, x_{i+1}]$ for $1 \leq i \leq k - 1$.

Let $W_1(x)$ for $x_1 < x < x_2$ be the zero with the smallest ordinal, $W_2(x)$ the smallest after that, and so forth till all zeros are assigned $W_1(x), W_2(x), \dots, W_n(x)$. These functions are continuous because we have proven in the previous part that $\forall \varepsilon > 0 \exists \delta > 0 : |x' - x| < \delta \rightarrow |W_i(x') - W_i(x)| < \varepsilon$ only now we have introduced new notation for the zeros. From the mean value theorem we can derive that

$$\frac{\partial W_i(x)}{\partial x} = -\frac{\frac{\partial U^{(0)}}{\partial x}}{\frac{\partial U^{(0)}}{\partial y}}.$$

This derivative is monotonous because both the upper part and the lower part of the fraction are monotonous on $U = 0$ by definition of the choice of boundaries $x = x_1, x = x_2$. The function $W_i(x)$ is also bounded by the maximum and minimum of the circle (R) on (x_1, x_2) and thus we can use the theorem of Bolzano-Weierstraß to find that $W_i(x)$ has limits on both $x = x_1$ and $x = x_2$ and conclude because of the continuity that these are unique. Conclusively there are n of these curves on $[x_1, x_2]$ which Ostrowski calls elementary curves.

$U^{(0)} = 0$ thus falls apart into a finite amount of elementary curves as each section $[x_i, x_{i+1}]$ has a finite amount of these elementary curves and there are a finite amount of sections. $U^{(1)} = 0$ can be described the same way when adding the verticals $x - x_1$ and cutting these into (a finite amount of) pieces at each point of $W_i(x_1)$. Because $U = 0$ and $U^{(1)} = 0$ describe the same curve as stated earlier this too can be formed by a finite number of elementary curves.

In order to understand what happens at the borders of these sections we will analyse the behaviour of $U^{(0)}$ around these elementary curves. In between two subsequent curves $W_i(x)$ and $W_{i+1}(x)$ the function $U^{(0)}$ will have a constant sign. This is because for the sign to switch there need to be zeros, which only occur on these elementary curves. On every $W_i(x)$ within the section the derivatives $\frac{\partial U^{(0)}}{\partial y}$ are non-zero, thus on each elementary curve the sign of $U^{(0)}$ swaps. We will use this to understand what happens at the edges of the sections. When the elementary curve reaches the edge of its interval and the circle (R) simultaneously then only one curve will end in that limit point. One can assume this as all we have previously proven is true for every circle with radius $\bar{R} \geq R$, and thus we can pick a radius a slight bit larger if it would not be the case. This new radius should satisfy this statement because of the structure of $U = 0$ outside of the circle discussed in parts one and two.

For points on the edges $x = x_i$ which do not lie on the circle we will see that an even amount of elementary curves end at each such a point. We will see this by the following reasoning. Again this reasoning is for $x = x_2$ but can be made

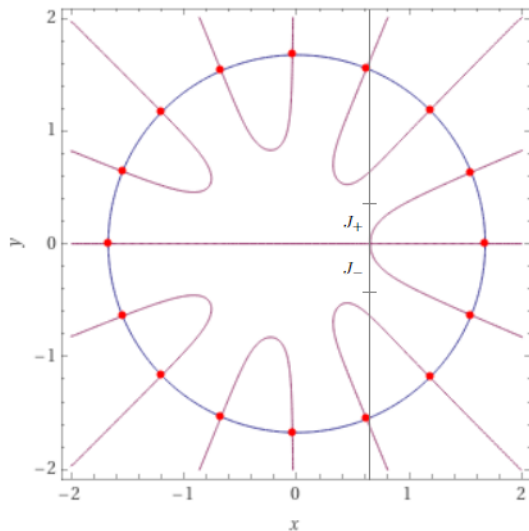


Figure 2: $U = 0$ for the function $\mathbf{X} = z^8 + 0.2z^7 - 0.1z^6 - 0.3z^5 - 0.1z^3 + 0.2z^2 - 0.3z + 0.1$) with example J_+ and J_- plotted. Example function taken from [Basu, Velleman, 2016, p692]

for every x_i . What I will write here is a slightly adapted version of the proof Ostrowski gave as Ostrowski seems to assume that at least two curves enter the point, which is more than one can assume from previous analysis. The version given here surpasses that problem and only assumes one elementary curve at first.

Let (x_2, w) be a point on $x = x_2$ not on the circle (R) on which at least one elementary curve ends. Say $W_i(x), \dots, W_{i+k}(x)$ enter this point from the right and $W'_j, \dots, W'_{j+k'}$ from the left. Now we will look at the following intervals on $x = x_2$, namely $J_+ = (w, w + \varepsilon)$ and $J_- = (w - \varepsilon, w)$ for a $\varepsilon > 0$ small enough such that no other zero of $U^{(0)}$ is within this interval. Then $U^{(0)}$ has to have the same sign on J_+ as the area above $W_{i+k}(x)$ as this vertical is non-zero in this part. The same for J_- and the area underneath $W_i(x)$. In this case the number of curves on one side of the vertical would be even and the other side would be uneven if an uneven amount of elementary curves would end on (x_2, w) . On the even side (including the possibility that there are no elementary curves on that side) J_+ and J_- will have the same symbol as there are an even number of switches of symbol. On the uneven side though there will be an uneven amount of switches. This would indicate that J_+ and J_- have the opposite symbol. This cannot be true at the same time, thus there should be an even number of elementary curves of $U^{(0)}$ ending on (x_2, w) . A vertical that isn't in $U^{(0)} = 0$ but is in $U = 0$ would consist of two elementary curves, namely one above and one below (x_2, w) thus still maintaining the even number.

We have just seen that $U = 0$ can be described by a finite number of ele-

mentary curves and have realised that at each ending of such a curve one of two things happen, namely exiting the circle or there is a different curve that also ends at the same point.

4.3.5 Part five

In this part the main theorem will be proven. This will be done by first grouping the elementary curves into continuous curves that enter and exit the circle (R) and then concluding that an intersection exists of $U = 0$ and $T = 0$.

We will first start by defining a direction on the elementary curves. The direction goes along each elementary curve such that $U^{(0)}$ is positive on the right of the curve and negative on the left. This is well defined as each curve is a transition between such a positive and negative area. We will start a larger continuous curve at the curves outside (R) that go to infinity and end on the circle. Then from there on out the curve will continue with the rightmost curve it can at each intersection of elementary curves. In this manner every elementary curve is at most used only once. This is because if a curve, h was used twice, there would be are two curves entering into the start of the curve h for which the curve is the rightmost, which would contradict the structure of these intersections we found at part four. Theoretically there could be elementary curves not used at all as this is not disproved. These would not interfere with the rest of the proof thus can safely be ignored.

Continuing this logic we will find m curves of $U = 0$ entering the circle (R) at a certain interval $(4k + 1)$ and exiting at an interval $(4k' + 3)$. We will call these curves $[k]$ respectively. In a similar manner m curves for $T = 0$ can be found. At each interval $(4k + 1)$ T is positive while at every interval $(4k' + 3)$ T is negative. Thus $T = 0$ will intersect the curve $[k]$ at least once. And thus Gauss' gap is filled and the theorem is proven.

Additionally under the assumption that every curve $[k]$ only touches a different curve of $U = 0$ and does not intersect, then one is able to prove that m roots of f exist as each curve intersects $T = 0$ uniquely.

In Ostrowski's paper he adds a section six where he proves certain additional structures of the U and T functions that could make future analysis easier. These are proven quickly with the structure already found to prove the fundamental theorem of algebra, but go beyond the initial theorem. In this paper these proofs will not be explained as only the fact that these expansions are possible will be important, not the exact nature of these additionally proven structures.

5 Value analysis

In this chapter a value analysis will be given of the proofs of the fundamental theorem of algebra presented in the previous chapter. The values discussed in Chapter 3 are the basis of this analysis.

In order to keep this analysis structured it will be done one value grouping after the other. Starting with the Didactic values, then Aesthetic values, and finally the Miscellaneous category. In Chapter 3 it is described that some values are not meaningful for wrong proofs. These will be left out of this chapter. Let us start with the didactic values

5.1 Didactic values

5.1.1 Demonstrating proof techniques

Two things need to be acknowledged before we go into depth on some proof techniques used. The first is that each proof uses a vast array of small proof techniques. Both proofs are somewhat lengthy, and thus required many separate proof steps. It is not insightful to completely list all the techniques used, thus four salient techniques will be discussed. Two by Gauss and two by Ostrowski.

The second acknowledgement is the time difference between Gauss and Ostrowski. This difference is more than a hundred years in which mathematical practice changed quite a bit. This makes that Gauss sometimes uses more archaic proof techniques compared to Ostrowski. On the other hand Ostrowski did mention in his paper that he only uses techniques that Gauss could have known [Ostrowski, 1920]. It should be mentioned, though, that to Ostrowski these techniques might be more common while to Gauss they might have been very new and unknown.

The first technique arises because Gauss explicitly does not use imaginary numbers. These numbers were known in Gauss' time, but he chose to not use these. The first two lemmas that Gauss proves require quite an extensive algebraic argument because of this. Recall that these two lemmas together prove that for a random polynomial $\mathbf{X} = x^m + Ax^{m-1} + Bx^{m-2} + \dots + Lx + M$ when

$$U = r^m \cos(m\phi) + Ar^{m-1} \cos((m-1)\phi) + \dots + Lr \cos(\phi) + M = 0 \quad (3)$$

$$T = r^m \sin(m\phi) + Ar^{m-1} \sin((m-1)\phi) + \dots + Lr \sin(\phi) = 0 \quad (4)$$

then \mathbf{X} has a linear or quadratic component. The arguments that Gauss uses are algebraic calculations. This method provides explicit functions of the components. On the other hand, these two lemmas can more easily be proven by imaginary numbers, as Ostrowski does. This method does not provide explicit functions of the divisors of \mathbf{X} though. Because Gauss did not use the argument through imaginaries, it could be said that his proof provides a show of the power of algebraic reasoning without imaginaries in his first two lemmas through this comparison.

The second technique Gauss shows is possible without imaginary numbers is the use of polar coordinates. Polar coordinates are often associated with imaginary numbers as this coordinate system is very easily described with these numbers. Gauss on the other hand constructed a polar coordinate system without the use of imaginary numbers such that he could use it in his proof. This method was quite extensive as he had to describe theoretically how such a plane

might look. When one uses complex coordinates it arises more naturally from describing numbers by $e^{r+\phi i}$. Gauss made his method work nonetheless. This is thus quite an alternative proof technique he demonstrates through his proof.

Recall that Gauss left a gap by assuming that every algebraic curve that enters a circle, necessarily has to leave it again. The two proof techniques used by Ostrowski that we will discuss here are used to cross this gap. The first is that he cuts the curve $U = 0$ into elementary pieces inside the circle. He does so by first simplifying the polynomial. Afterwards he lists all the vertical lines in the circle through which the curves might behave differently. By reconstructing the curves between these lines he only has understandable, elementary, pieces left. Because of this the curve is better understood. This is quite a powerful technique that could be used in different proofs. Maybe not explicitly in the same shape, but the idea of reducing an object into understandable pieces though reconstructing it between possible anomalies is a powerful one.

The second technique is when Ostrowski continues to glue these elementary pieces back together. Eventually he finalizes the argument with the newfound continuous curves. He uses that the curves are lines in which the (part of the original) polynomial is zero. He had proven that through all these nullines (except possibly on some earlier mentioned vertical lines) the derivative is non zero. This means that on each side of the curve there is a different sign. He then views the domains created by the curves, rather than the curves. His argument is based on that a domain with positive sign has to have a border with a domain with a negative sign. Thus when an elementary curve ends at a point, an even amount of curves have to begin. If this amount is uneven then there is no border between a positive and a negative domain. Though this he is able to prove that the elementary curves piece together. He recognises that this technique of using domains is a strong technique. In his introduction he writes that Gauss used this method in his fourth proof of the Fundamental theorem of Algebra and that after it has become a more standard technique with quite some proving prowess [Ostrowski, 1920]. This proof is one that demonstrates this prowess because the argument is strongly reliant on this technique to work.

5.1.2 Reconceive mathematical domains

The way Ostrowski glues the curves together corresponds with a different conception of curves. In this conception the curves are boundaries of domains. Thus he changes the perspective from studying curves to studying domains instead.

Another interesting conception in Ostrowski's proof results from that Ostrowski uses imaginary numbers but bases his proof on Gauss' proof. Gauss' proof reduces the search for roots of $\mathbf{X} = x^m + Ax^{m-1} + Bx^{m-2} + \dots + Lx + M$ to the intersection of two affiliated functions $U = 0$ and $T = 0$ (given above in Section 5.1.1) in the polar plane. In Gauss' proof the connection to \mathbf{X} is given through his first two lemmas. In Ostrowski's proof, these two functions could also be called $Im(\mathbf{X}) = 0$ and $Re(\mathbf{X})$ respectively. Often when one wants to find the roots of a complex polynomial f , one searches when $Re(f) + Im(f) = 0$ as addition of two real-valued functions on the Cartesian plane. Though his proof,

though, Ostrowski searches for the intersection of these two curves on the polar plane instead of the addition in the Cartesian plane. This forces the reader to think about these curves different spatially.

The last different conception of curves discussed is nested in Gauss' argument. At the end of his proof he uses geometrical arguments rather than analytic. For example, he describes how a curve of $U = 0$ gets *closed in* by a curve of $V = 0$ when he is proving that the two curves intersect through *reductio ad absurdum*. Thus the curves are geometrical objects instead of analytic functions in this argument. To do so in the proof of an analytic theorem is interesting. Ostrowski does not approve of this method. He writes that using geometrical arguments in a analytic proof could be called improvident [Ostrowski, 1920]. Thus a debate could be had whether this change of conception is valuable. It provides a different conception at the very least.

5.1.3 Promote discovery

Gauss directly promoted the discovery of Ostrowski's proof. He also wrote another proof of the Fundamental theorem of Algebra on the basis of this one (his fourth proof of this theorem) [Cain, 2005]. Gauss' proof actually promoted the discovery of more proofs as multiple authors tried to fill the gap or were inspired by the proof [Gersten, Stallings, 1988; Martin et al., 2007; Basu et al., 2017].

Ostrowski's method promoted discovery of several new small statements directly based upon the methods of his proof. He has proven these in the sixth part in his paper. These are all based upon the information gained though his proof, rather than the theorem.

5.1.4 Explanatoriness

Both proofs are not very explanatory by all three accounts given by [Morris, 2022]. To recall Steiner [1978] asks that a proof is built around a characterizing property of a central object. In our proofs the only central object is the polynomial X . But not one property of this object is the centre in each one of the proofs. In Ostrowski's proof he swaps between the rotational symmetry in the limiting behaviour outside of the circle, to cutting up the curves into elementary pieces inside. In Gauss' proof first the lemmas are done with algebraic calculations, following that the same rotational symmetry in limiting behaviour is used, finishing with geometrical properties in the final arguments. Thus both do not use a characterizing property of the central object.

Kitcher [1989] asks that the proof reduces the number of ultimate facts one needs to accept. In Ostrowski's proof the cutting into elementary curves inside the circle is done very mechanically with many exclusions such that no issues can arise. It is based upon many small facts and choices made, opposed to a small amount. Gauss' proof asks us to accept a very large ultimate fact, namely that "if a branch of an algebraic curve enters into a limited space, it necessarily has to leave it again." [Struik, 2014, p121]

According to Lange’s [2014] account a theorem needs a salient feature that the proof explains. The theorem states that a random complex valued polynomial can be factored into linear factors (in the modern characterization of the theorem). This factorization is then the salient feature. Gauss’ last argument depends on *reductio ad absurdum*. Lange [2014] describes in his introduction (based upon Nicole and Arnauld’s [1717] words) that this technique necessarily is unexplanatory. When one looks at his account this is understandable. It does not show where the salient feature originates, but merely shows that it is inevitable. Thus Gauss’ proof is not explanatory through this account.

Ostrowski’s proof views the curves as edges of domains. At the edge of the circle he shows that $U = 0$ and $T = 0$ alternately intersect the circle. Also he shows that the derivative of U has alternating signs at the intersections of T . These facts that he proves in part one are later on incredibly important for his argument of why his curves intersect. Thus for the proof to be explanatory these facts need to be explained well. He proves these facts by algebraic computations. It is as if they are an “accident of algebra” [Lange, 2014]. This validates that these signs alternate, it does not explain why they do.¹⁷ As this is on the base of the reasoning for the intersection¹⁸ and thus the linear factors, the proof does not explain the salient feature. Thus it is not explanatory by Lange’s account.

5.1.5 Motivatedness

Both proofs have well motivated parts and parts that are not well motivated. They are not consistent with motivating their steps. I will highlight some important parts. Gauss’ proof starts off very unmotivated. The first lemma is introduced without explanation why. It only gets context through the second lemma. For both these lemmas Gauss does not explain why he believed these are appropriate. Also the proofs of these lemmas are done by presenting algebraic conclusions that are unmotivated.

Later on in the proof Gauss explains why he finds his gap to be appropriate. He explains that an algebraic curve does not disappear nor infinitely wind into a single point. Thus, he motivates why he found the gap reasonably appropriate.

Opposed to Gauss, Ostrowski motivates the connection between \mathbf{X} and U and T well. The step with imaginaries is recognisably appropriate and canonical. It can thus be understood how it was conceived. Thus it is strongly motivated. It should be said that such a step is so canonical for modern mathematics, that it hardly requires motivation.

The change of coordinate system from polar to Cartesian was not well motivated in Ostrowski’s proof. This was simply done in part three of his proof. After this he starts cutting the circle into slices without explanation. This slicing only gets understandable after he constructed his elementary curves through

¹⁷This is similarly to why the first proof of Zeitz’s Biased Coin given by Lange [2014] is not explanatory.

¹⁸(after it was proven for the curves that after entering the circle, in fact, “it necessarily has to leave it again.” [Struik, 2014, p121])

it. How he knew that this specific slicing of the circle would result in elementary enough curves is not presented. Thus this part was not well motivated.

Lastly I want to discuss his use of the domains to find the intersection. In his introduction he explains why he used the technique. In Gauss' forth proof of the Fundamental theorem of Algebra he uses the technique of domains to find the intersection [Ostrowski, 1920]. Ostrowski based his use on this as he explains himself. Also when he used the technique in the end he already introduced the domains earlier and the proof gets finalized immediately because of it. Thus it is recognisably appropriate and well explained how the author conceived the step. Thus this step was very well motivated.

5.1.6 Visual intuitiveness

Both Gauss' gap as his final argument are made visually intuitive. He describes the crossing of the gap with geometrical terms and talks about the curves entering and exiting as if they are in motion. Moreover this gap seems to have arisen from this visual intuition. It seemed so intuitive that every algebraic curve that enters a confined space necessarily has to leave it again that he accepted it as a fact. Ostrowski proves through his theorem that this is indeed the case for the curves in the theorem, and even mentions that this could be generalised further to a larger set of curves. Thus this intuition is right. By using it as a proof step he skips along a vast amount of mathematical reasoning necessary to make it a valid step. Thus the step is too large but the intuition is correct. The intuition could thus still be called valuable.

Gauss' final argument is a visual one. He is closing off possibilities until one intersection of the curve $U = 0$ or $T = 0$ with the circle is enclosed by a curve of $T = 0$ or $U = 0$ respectably. This argument is quite intuitive after introducing the numbering of the intersections and showing the first step. After that first step the iterative process becomes clear to the reader, thus the rest of the proof follows naturally. Thus this argument is visually intuitive.

Ostrowski explicitly strays away from visual language. He explains that it is generally ill advised to use geometrical methods in an analytic proof.¹⁹ In his proof he stays true to this statement and does not use visual language or explanations.

5.2 Aesthetic values

5.2.1 Kantian beauty

It is very hard to judge about the Kantian beauty of a proof. This notion is subjective and has no hard definitions or requirements it should uphold. It "is grounded in our felt awareness of the imaginative processes that lead to mathematical knowledge." [Breitenbach, 2015, p955] I argue that both proofs are not very Kantian beautiful according to this. Ostrowski's proof is very technical and does not leave much room for imagination. It mostly contains dry

¹⁹He states this in a reaction to Gauss' use of geometrical language to pass the gap.

statements with lots of nuances necessary to make it work. This can specifically be seen when he tries to make the elementary curves. This part has nearly a page of limitations and alterations that are introduced such that the argument works. Because of this it feels forced rather than imaginative. Thus I would not call Ostrowski's proof Kantian beautiful.

Gauss' introduction of polar coordinates without the use of imaginaries also feels forced. With the use of imaginary numbers this coordinate system is very natural. It takes Gauss quite some work to introduce this coordinate system without imaginary numbers. Also the connection between \mathbf{X} and $U = 0$ and $T = 0$ costs elaborate work, while it could be done so much more elegantly as Ostrowski shows. These constructions necessarily interrupt the general argumentation which leads away from the imaginative processes that the general argumentation might produce. The only part I might call beautiful (in the Kantian interpretation) is the final argument of Gauss after the gap. He shows through an visually intuitive (as earlier discussed) method with little clutter how the theorem arises from the beforehand proven assumptions and gap. This process could cause a "felt awareness of the imaginative processes" [Breitenbach, 2015, p955]. Aside from this small piece, as a whole Gauss' proof is not Kantian beautiful.

5.2.2 Representation

In his second lemma, Gauss places the equations under each other in a table-like fashion. Through this the terms can be added vertically. This visualisation is quite aesthetic because its spatial structure shows the, previously hidden, structure of the terms elegantly. Apart from this both proofs do not have specifically valuable representations in their arguments. They use symmetries in their arguments, specifically the symmetry between U and T . But they mostly use this symmetry to make their arguments shorter. They don't use it to highlight aspects not seen before. As an example, the alternating intersections of the circle by U and T are shown through computational measures in both proofs. These computations use symmetries to make it shorter, but not for insight. Thus this symmetry does not present representational value.

5.3 Miscellaneous

5.3.1 Purity

Ostrowski's proof is pure. The theorem is analytical in nature as it provides a statement about the real algebra. He mentions in his introduction that he believes an analytical proof should only use analytical methods. He keeps his proof purely analytical accordingly.

Gauss' proof is not pure. He starts off pure, his first lemmas and analysis of the curves outside of the circle are analytical in nature. Only at the gap and afterwards does he start using geometrical language and (unproven) theory. These geometrical methods do not contain a certain symmetry of conceptual

resources to the analytic theorem as Detlefsen [2008] requires a pure proof to have. Because of this his proof is unpure.

5.3.2 Systematization

Ostrowski's proof adds to the systematization of certain techniques. One in particular is the domain view of curves. He mentions how this technique has been used often before and has proven fruitful. Then he goes to use it again in his proof, aiding in the standardization of this technique.

In his proof he does use an undefined term though. He uses the term multiple factors ("mehrfachen Factor" [Ostrowski, 1920, p6]) that he does not define. The definition of this presented in part three of the proof in this thesis is a deduced definition from how he uses the term in the proof. This use of a term without defining it properly does not aid systematization.

Gauss' proof also adds to the systematization of certain analytical techniques. He does go against using imaginary numbers. These numbers were somewhat new concepts while he wrote his dissertation. Explicitly not using them goes against the systematization of imaginary numbers, a concept which later on proved very useful and important. Thus his proof is less valuable in this regard.

6 Discussion

This thesis tries to answer the question *What is the value of a wrong proof?* through a case study. The previous chapter has provided the value analysis for this study. In this chapter we will discuss the results generated in the analysis and provide suggestions for further study based on these results.

Results In the previous chapter we have seen that Gauss' proof of the Fundamental theorem of Algebra provides multiple proof values. We saw that Gauss' proof displays a variety of didactic values, that it demonstrates proof techniques, it reconceives mathematical domains, it promotes discovery, it is motivated and it is visually intuitive. Only the value of explanatoriness out of this category was not found in his proof. In the aesthetic value category his proof displayed one technique that is valuable in representation. The proof did not display the other aesthetic value, Kantian beauty. In the miscellaneous category only the value of systematization was fulfilled to some degree. But for this value he also made a choice that went against possible further value of systematization in his proof. The value of purity was not found in Gauss' proof.

Ostrowski's proof was not much different from Gauss' proof in terms of values found. Only the value purity could be found in Ostrowski's proof and not in Gauss' proof. On the other hand the value of visual intuitiveness was found in Gauss' proof opposed to Ostrowski's. The rest of the values found were in both proofs alike.

Implications The lack of large differences between the values of Ostrowski's and Gauss' proofs means that the values analysed stayed consistent even when the invalidity in the proof was solved by Ostrowski. This is interesting as it shows that the wrong proof is thus comparable to the right version when one reads it for these values.

That the proof does not prove something new is not a reason to directly disregard it. As discussed mathematicians value reproving enough to publish these results. This wrong proof seems to present certain mathematical values just as reproving could do, also without proving a new theorem. This displays that wrong proofs can be analysed just as proofs that reprove can to some extend.

The results also suggest that when wrong proofs get disregarded because of their invalidity, mathematical value can be lost. In this case specifically, quite some didactic values would be lost when it wouldn't have been presented because of the gap. For example, some newfound proofs of the fundamental theorem of algebra would not have been found as they were inspired by Gauss' proof.

Limitations There are some limitations the results of this study. One cannot derive general claims based on this thesis. This study is based upon only one case. This is thus not a broad enough base for general statements. Below are two reasons why this case might be too specific for general statements.

One reason is that these proofs might not contain some values because of the proof setup. These values then might be found in different wrong proofs. It could be that the method of proving used did not coincide with, for example, a beautiful one (according to the Kantian view of beauty). A different wrong proof could maybe be beautiful. This study thus cannot provide clarity whether generally some values can or cannot be found in wrong proofs.

Secondly, the proofs are of a fundamental theorem. It might be that the results would not occur in a very specialized proof where, for example, proof techniques are not as generalizable and thus less valuable to most readers.

Future research Subsequent study is needed to establish more general results. No previous research into the value of wrong proofs has been done. Thus this thesis presents a start to this discussion. It can be seen as a preliminary study into the value of wrong proofs. It seems to be a promising start as the proof analysed here did present quite some proof values.

One way the discussion could be further expanded is by analysis of a larger number of cases. This could lead to more generalizable results. This could then make a better case what values wrong proofs generally uphold.

Another part that could further be added to this discussion is by questioning practicing mathematicians about their views of wrong proofs. This could add to our understanding whether wrong proofs are undervalued or not.

Especially the didactic values could be interesting for further study. Researching what didactic values in wrong proof might mean for education could

be a fruitful line of research. The prevalence of didactic values in the wrong proof of Gauss suggests that wrong proofs might be useful for education.

Conclusion To conclude, we have seen that an incorrect proof can uphold some proof values normally found in correct proofs. Further research is required to conclusively say whether wrong proofs generally have (some of) these values. At the very least this thesis has shown that there is a case to be made for the value of wrong proofs. Specifically this could be in the values listed as didactic values in this thesis. Mathematicians often read proofs to learn, it seems that wrong proofs could teach them too.

Bibliography

- Andersen, Line Edslev (2020). “Acceptable gaps in mathematical proofs”. In: *Synthese* 197.1, pp. 233–247.
- Arana, Andrew (2015). “On the Depth of Szemerédi’s Theorem”. In: *Philosophia Mathematica* 23.2, pp. 163–176.
- Auslander, Joseph (2008). “On the roles of proof in mathematics”. In: *Proof and other dilemmas: Mathematics and philosophy*, pp. 61–77.
- Azzouni, Jody et al. (1994). *Metaphysical myths, mathematical practice: the ontology and epistemology of the exact sciences*. Cambridge University Press.
- Azzouni, Jody (2004). “The derivation-indicator view of mathematical practice”. In: *Philosophia Mathematica* 12.2, pp. 81–106.
- Basu, Soham and Daniel J Velleman (2017). “On Gauss’s first proof of the fundamental theorem of algebra”. In: *The American Mathematical Monthly* 124.8, pp. 688–694.
- Breitenbach, Angela (2015). “Beauty in Proofs: Kant on Aesthetics in Mathematics”. In: *European Journal of Philosophy* 23.4, pp. 955–977.
- Cain, Alan J (2019). “Visual thinking and simplicity of proof”. In: *Philosophical Transactions of the Royal Society A* 377.2140.
- Cain, Harel et al. (2005). “Cf Gauss’s proofs of the fundamental theorem of algebra”. In: *Citeseer*.
- Celluci, Carlo (2008). “Why proof? What is a proof?” In: *Deduction, computation, experiment*. Springer, pp. 1–27.
- Davis, Philip J (1972). “Fidelity in mathematical discourse: Is one and one really two?” In: *The American Mathematical Monthly* 79.3, pp. 252–263.
- Dawson, John W (2006). “Why do mathematicians re-prove theorems?” In: *Philosophia Mathematica* 14.3, pp. 269–286.
- Dawson Jr, John W (2015). *Why prove it again?: alternative proofs in mathematical practice*. Birkhäuser.
- De Toffoli, Silvia (2020). “Reconciling Rigor and Intuition”. In: *Erkenntnis*, pp. 1–20.
- De Villiers, Michael D (1990). “The role and function of proof in mathematics”. In: *Pythagoras* 24, pp. 17–24.
- Descartes, René and Ralph Monroe Eaton (1927). *Selections*. Charles Scribner’s Sons, New York.
- Detlefsen, Michael (2008). “Purity as an ideal of proof”. In: *The Philosophy of Mathematical Practice*, pp. 179–197.
- Ernest, Paul (2016). “Mathematics and values”. In: *Mathematical cultures*. Springer, pp. 189–214.
- Fallis, Don (2003). “Intentional gaps in mathematical proofs”. In: *Synthese* 134.1/2, pp. 45–69.
- Geist, Christian, Bart Van Kerkhove, and Benedikt Löwe (2010). “Peer Review and Knowledge by Testimony in Mathematics”. In: *Philosophy of Mathematics: Sociological Aspects and Mathematical Practice*. College Publications, pp. 155–178.

- Gersten, Steve M and John R Stallings (1988). “On Gauss’s first proof of the fundamental theorem of algebra”. In: *Proceedings of the American Mathematical Society* 103.1, pp. 331–332.
- Giaquinto, Marcus (2016). “Mathematical proofs: The beautiful and the explanatory”. In: *Journal of Humanistic Mathematics* 6.1, pp. 52–72.
- Hales, Thomas C (2008). “Formal proof”. In: *Notices of the AMS* 55.11, pp. 1370–1380.
- Kitcher, Philip (1989). “Explanatory unification and the causal structure of the world”. In: *Scientific Explanation: Minnesota Studies in the Philosophy of Science, Vol. 13*.
- Lakatos, Imre (1976). *Proofs and Refutations: The Logic of Mathematical Discovery*. Cambridge University Press.
- Lange, Marc (2014). “Aspects of mathematical explanation: Symmetry, unity, and salience”. In: *Philosophical Review* 123.4, pp. 485–531.
- Larvor, Brendan (2012). “How to think about informal proofs”. In: *Synthese* 187.2, pp. 715–730.
- Loomis, Elisha Scott (1968). *The pythagorean proposition*. National Council of Teachers of Mathematics.
- MacLane, Saunders (1986). *Mathematics form and function*. Springer Science & Business Media.
- Martin, Jeremy L, David Savitt, and Ted Singer (2007). “Harmonic algebraic curves and noncrossing partitions”. In: *Discrete & Computational Geometry* 37.2, pp. 267–286.
- Morris (2022). “The Values of Mathematical Proofs”. In: *Sriraman B. (eds) Handbook of the History and Philosophy of Mathematical Practice*.
- Morris, Rebecca Lea (2020). “Motivated proofs: What they are, why they matter and how to write them”. In: *The Review of Symbolic Logic* 13.1, pp. 23–46.
- Nelsen, RB (2001). “Proof Without Words II-Exercises in Visual Thinking”. In: *The Mathematics Association of America*.
- Nicole, Pierre and Arnauld Antoine (1717). *Logic; or the Art of Thinking*. William Taylor, London.
- Noel, Linda Hand (1991). “Fundamental Theorem of Algebra: A Survey of History and Proofs”. PhD thesis. Oklahoma State University.
- Ostrowski, Alexander (1933). “Über den ersten und vierten Gaußsschen Beweis des Fundamental-Satzes der Algebra”. In: *Carl Freiderich Gauss Werke Band X Abt. 2*.
- Poincaré, Henri and Francis Maitland (2003). *Science and method*. Courier Corporation.
- Polya, George (1949). “With, or without, motivation?” In: *The American Mathematical Monthly* 56.10, pp. 684–691.
- Raman, Manya Janaky (2002). *Proof and justification in collegiate calculus*. University of California, Berkeley.
- Raman-Sundström, Manya (2016). “The notion of fit as a mathematical value”. In: *Mathematical Cultures*. Springer, pp. 271–285.
- Raman-Sundström, Manya and Lars-Daniel Öhman (2018). “Mathematical fit: A case study”. In: *Philosophia mathematica* 26.2, pp. 184–210.

- Rav, Yehuda (2008). “The axiomatic method in theory and in practice”. In: *Logique et Analyse*, pp. 125–147.
- Rota, Gian-Carlo (1997). “The phenomenology of mathematical beauty”. In: *Synthese* 111.2, pp. 171–182.
- Schattschneider, Doris (1981). “In praise of amateurs”. In: *The Mathematical Gardner*. Springer, pp. 140–166.
- Steiner, Mark (1978). “Mathematical explanation”. In: *Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition* 34.2, pp. 135–151.
- Struik, Dirk Jan (2014). *A source book in mathematics, 1200-1800*. Vol. 445. Princeton University Press.
- Tanswell, Fenner (2015). “A problem with the dependence of informal proofs on formal proofs”. In: *Philosophia Mathematica* 23.3, pp. 295–310.
- Weber, Keith (2008). “How mathematicians determine if an argument is a valid proof”. In: *Journal for research in mathematics education* 39.4, pp. 431–459.
- (2010). “Proofs that develop insight”. In: *For the learning of mathematics* 30.1, pp. 32–36.