

A formally verified proof of the Central Limit Theorem

Jeremy Avigad · Johannes Hölzl · Luke
Serafin

the date of receipt and acceptance should be inserted later

Abstract We describe a proof of the Central Limit Theorem that has been formally verified in the Isabelle proof assistant. Our formalization builds upon and extends Isabelle's libraries for analysis and measure-theoretic probability. The proof of the theorem uses *characteristic functions*, which are a kind of Fourier transform, to demonstrate that, under suitable hypotheses, sums of random variables converge weakly to the standard normal distribution. We also discuss the libraries and infrastructure that supported the formalization, and reflect on some of the lessons we have learned from the effort.

Keywords interactive theorem proving, measure theory, central limit theorem

1 Introduction

If you roll a fair die many times and compute the average number of spots showing, the result is likely to be close to 3.5, and the odds that the average is far from the expected value decreases roughly as the area under the familiar bell-shaped curve. Something similar happens if the measurement is continuous rather than discrete, such as when you repeatedly toss a needle on the ground and measure the angle it makes with respect to a fixed reference line. Even if the die is not a fair die or the geometry of the needle and the ground makes some angles more likely than others, the distribution of the average still approaches the area under a bell-shaped curve centered on the expected value. The width of the bell depends on both the variance of the random measurement and the number of times it is performed. Made precise, this amounts to a statement of the Central Limit Theorem.

The Central Limit Theorem lies at the heart of modern probability. Many generalizations and variations have been studied, some of which either relax the requirement that the repeated measurements are independent of one another and

J. Avigad · L. Serafin
Carnegie Mellon University

J. Hölzl
Technische Universität München

identically distributed (cf. in particular, the results of Lyapunov and Lindberg [3]), while others provide additional information on the rate of convergence.

Here we report on a formalization of the Central Limit Theorem that was carried out in the Isabelle proof assistant. This result is noteworthy for a number of reasons. Not only is the CLT fundamental to probability theory and the study of stochastic processes, but so is the machinery developed to prove it, ranging from ordinary calculus to the properties of real distributions and characteristic functions. There is a pragmatic need to subject statistical claims made in engineering, risk analysis, and financial computation to formal verification, and our formalization along with the surrounding infrastructure can support such practical efforts.

The formalization is also a good test for Isabelle’s libraries, proof language, and automated reasoning tools. As we will make clear, the proof draws on a very broad base of facts from analysis, topology, measure theory, and probability theory, providing a useful evaluation of the robustness and completeness of the supporting libraries. Moreover, the concepts build on one another. For example, a measure is a function from a class of sets to the reals, and reasoning about convergence of measures involves reasoning about sequences of such functions. The operation of forming the characteristic function is a functional taking a measure to a function from the reals to the complex numbers, and the convergence of such functionals is used to deduce convergence of measures. The conceptual underpinnings are thus as deep as they are broad, and working with them tests Isabelle’s mechanisms for handling abstract mathematical notions.

In Section 2 we provide an overview of the Central Limit Theorem and the proof that we formalized, following the textbook presentation of Billingsley [3]. In Section 3 we describe the Isabelle proof assistant, and the parts of the library that supported our formalization. In Section 4 we describe the formal proof itself, and in Section 5 we reflect on what we have learned from the effort.

Our formalization is currently part of the Isabelle library, which can be found online at <https://isabelle.in.tum.de/>.¹ A preliminary, unpublished report on the formalization can be found on arXiv [1]. Our presentation also draws heavily on Serafin’s Carnegie Mellon MS thesis [20], which provides additional information.

Acknowledgments. We are grateful to Tobias Nipkow, Lawrence Paulson, Makarius Wenzel, and the entire Isabelle team for the ongoing development of Isabelle. We are especially grateful to Tobias for steadfast encouragement and support. We thank our two anonymous referees for a very careful reading and helpful comments. Avigad and Serafin’s work has been partially supported by NSF grant DMS-1068829, and Avigad’s work has been partially supported by AFOSR grants FA9550-12-1-0370 and FA9550-15-1-0053. Hölzl’s work has been partially supported by DFG projects Ni 491/15-1 and Ni 491/16-1.

2 Overview of the Central Limit Theorem

For our formalization we followed Billingsley’s textbook, *Probability and Measure* [3], which provides an excellent introduction to these topics. Here we provide some

¹ The probability library in particular can be found at <https://isabelle.in.tum.de/dist/library/HOL/HOL-Probability/index.html>.

historical background, briefly review the key concepts, give a precise statement of the Central Limit Theorem, and present an outline of the proof.

2.1 Historical background

In 1733, De Moivre privately circulated a proof that, as n approaches infinity, the distribution of n flips of a fair coin converges to a normal distribution. This material was later published in the 1738 second edition of his book *The Doctrine of Chances*, the first edition of which was published in 1712. That book is widely regarded as the first textbook on probability theory. De Moivre also considered the case of what we would call a biased coin, that is, an event which has value one with probability p and zero with probability $1 - p$ for some $p \in (0, 1)$. He showed that his convergence theorem continues to hold in that case.

De Moivre's result was generalized by Laplace in the period between about 1776 and 1812 to sums of random variables with various other distributions, such as the uniform distribution on an interval. Over the next three decades Laplace developed conceptual and analytical tools to extend this convergence theorem to sums of independent and identically distributed random variables with ever more general distributions, and this work culminated in his treatise *Théorie analytique des probabilités*. This included the development of the method of characteristic functions to study the convergence of sums of random variables, a move which firmly established the usefulness of analytic methods in probability theory.

Laplace's theorem later became known as the Central Limit Theorem, a designation due to Pólya, stemming from its importance both in the theory and applications of probability. In modern terms, the theorem states that the normalized sum of a sequence of independent and identically distributed random variables with finite, nonzero variance converges to a normal distribution. All of the main ingredients of the proof of the CLT are present in the work of Laplace, though of course the theorem was refined and extended as probability underwent the radical changes necessitated by its move to measure-theoretic foundations in the first half of the twentieth century.

Gauss was one of the first to recognize the importance of the normal distribution to the estimation of measurement errors. The usefulness of the normal distribution in this context is largely a consequence of the Central Limit Theorem, since errors occurring in practice are frequently the result of many independent factors which sum to an overall error in a way which can be regarded as approximated by a sum of independent and identically distributed random variables. The normal distribution also arose with surprising frequency in a wide variety of empirical contexts, from the heights of men and women to the velocities of molecules in a gas. This gave the CLT the character of a natural law, as seen in the following poetic quote from Sir Francis Galton in 1889 [8]:

I know of scarcely anything so apt to impress the imagination as the wonderful form of cosmic order expressed by the "Law of Frequency of Error." The law would have been personified by the Greeks and deified, if they had known of it. It reigns with serenity and in complete self-effacement, amidst the wildest confusion. The huger the mob, and the greater the apparent anarchy, the more perfect is its sway. It is the supreme law of Unreason.

Whenever a large sample of chaotic elements are taken in hand and marshaled in the order of their magnitude, an unsuspected and most beautiful form of regularity proves to have been latent all along.

More details on the history of the Central Limit Theorem and its proof can be found in [7].

2.2 Background from measure theory

A *measure space* (Ω, \mathcal{F}) consists of a set Ω and a σ -algebra \mathcal{F} of subsets of Ω , that is, a collection of subsets of Ω containing the empty set and closed under complements and countable unions. Think of Ω as the set of possible states of affairs, or possible outcomes of an action or experiment, and each element E of \mathcal{F} as representing the set of states or outcomes in which some *event* occurs — for example, that a card drawn is a face card, or that Spain wins the World Cup. A *probability measure* μ on this space is a function that assigns a value $\mu(E)$ in $[0, 1]$ to each event E , subject to the following conditions:

1. $\mu(\emptyset) = 0$,
2. $\mu(\Omega) = 1$, and
3. μ is countably additive: if (E_i) is any sequence of disjoint events in \mathcal{F} , then $\mu(\bigcup_i E_i) = \sum_i \mu(E_i)$.

Intuitively, $\mu(E)$ is the “probability” that E occurs.

The collection \mathcal{B} of *Borel subsets* of the real numbers is the smallest σ -algebra containing all intervals (a, b) . A *random variable* X on the measure space (Ω, \mathcal{F}) is a measurable function from (Ω, \mathcal{F}) to $(\mathbb{R}, \mathcal{B})$. Saying X is measurable means that for every Borel subset B of the real numbers, the set $\{\omega \in \Omega \mid X(\omega) \in B\}$ is in \mathcal{F} . Think of X as some real-valued measurement that one can perform on the outcome of the experiment, in which case, the measurability of X means that if we are given any probability measure μ on (Ω, \mathcal{F}) , then for any Borel set B it makes sense to talk about “the probability that X is in B .” In fact, if X is a random variable, then any measure μ on (Ω, \mathcal{F}) gives rise to a measure ν on $(\mathbb{R}, \mathcal{B})$, defined by $\nu(B) = \mu(\{\omega \in \Omega \mid X(\omega) \in B\})$. A probability measure on $(\mathbb{R}, \mathcal{B})$ is called a *real distribution*, or, more simply, a *distribution*, and the measure ν just described is called *the distribution of X* .

If X is a random variable, the *mean* or *expected value* of X with respect to a probability measure μ is $\int X d\mu$, the integral of X with respect to μ . If m is the mean, the *variance* of X is $\int (X - m)^2 d\mu$, a measure of how far, on average, we should expect X to be from its mean.

Note that passing from μ and X to its distribution ν means that instead of worrying about the probability that some abstract event occurs, we focus more concretely on the probability that some measurement on the outcome lands in some set of real numbers. In fact, many theorems of probability theory do not really depend on the abstract space (Ω, \mathcal{F}) on which X is defined, but rather the associated distribution on the real numbers. Nonetheless, it is often more intuitive and convenient to think of the real distribution as being the distribution of a random variable (and, indeed, any real distribution can be represented that way).

One way to define a real distribution is in terms of a *density*. For example, in the case where $\Omega = \{1, 2, 3, 4, 5, 6\}$, we can specify a probability on all the subsets

of Ω by specifying the probability of each of the events $\{1\}, \{2\}, \dots, \{6\}$. More generally, we can specify a distribution μ on \mathbb{R} by specifying a function f such that for every interval (a, b) , $\mu((a, b)) = \int_a^b f(x) dx$. The measure μ is then said to be the real distribution with density f . In particular, the *normal distribution* with mean m and variance σ^2 is defined to be the real distribution with density function

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}.$$

The graph of f is the bell-shaped curve centered at m . When $m = 0$ and $\sigma = 1$, the associate real distribution is called the *standard normal distribution*.

Let X_0, X_1, X_2, \dots be any sequence of independent random variables, each with the same distribution μ , mean c , and variance σ^2 . Here “independent” means that the random variables X_0, X_1, \dots are all defined on the same measure space (Ω, \mathcal{F}) , but they represent independent measurements, in the sense that for any finite sequence of events B_1, B_2, \dots, B_k and any sequence of distinct indices i_1, i_2, \dots, i_k , the probability that X_{i_j} is in B_j for each j is just the product of the individual probabilities that X_{i_j} is in B_j . For each n , let $S_n = \sum_{i < n} X_i$ be the sum of the first n random variables in the sequence. Notice that each S_n is itself a measurable function on (Ω, \mathcal{F}) (which is to say it is a random variable), and so it is natural to ask how its values are distributed. We can shift the expected value of S_n to 0 by subtracting nc , and scale the variance to 1 by dividing by $\sqrt{n\sigma^2}$. The Central Limit Theorem says that the corresponding quantity,

$$\frac{S_n - nc}{\sqrt{n\sigma^2}},$$

approaches the standard normal distribution as n approaches infinity.

All that remains to do is to make sense of the assertion that a sequence of distributions $\mu_0, \mu_1, \mu_2, \dots$ “approaches” a distribution, μ . For distributions that are defined in terms of densities, the intuition is that over time the graph of the density should look more and more like the graph of the density of the limit. For example, if you flip a coin a number of times and graph all the possible values of the average number of ones, the discrete points plotted over the possibilities $0, 1/n, 2/n, 3/n, \dots, 1$ start to look like a bell-shaped curve centered on $1/2$. The notion of *weak convergence* makes the notion of “starts to look like” precise.

If μ is any real distribution, then the function $F_\mu(x) = \mu((-\infty, x])$ is called the *cumulative distribution function* of μ . In words, for every x , $F_\mu(x)$ returns the likelihood that a real number chosen randomly according to the distribution is at most x . Clearly $F_\mu(x)$ is nondecreasing, and it is not hard to show that F_μ is right continuous, approaches 0 as x approaches $-\infty$, and approaches 1 as x approaches ∞ . Conversely, one can show that any such function is the cumulative distribution function of a unique measure. Thus there is a one-to-one correspondence between functions F satisfying the properties above and real distributions.

The notion of weak convergence can be defined in terms of the cumulative distribution function:

Definition 1 Let (μ_n) be a sequence of real distributions, and let μ be a real distribution. Then μ_n *converges weakly* to μ , written $\mu_n \Rightarrow \mu$, if $F_{\mu_n}(x)$ approaches $F_\mu(x)$ at each point x where F_μ is continuous.

To understand why we need to exclude the points of discontinuity of F_μ , consider for each n the probability measure μ_n that puts all its “weight” on $1/n$, which is to say, for any Borel set B , $\mu(B) = 1$ if and only if B contains $1/n$. Then F_{μ_n} is the function that jumps from 0 to 1 at $1/n$. Intuitively, it makes sense to say that μ_n approaches the real distribution μ that puts all its weight at 0. But for every n , $F_{\mu_n}(0) = 0$, while $F_\mu(0) = 1$, which explains why we want to exclude the point 0 from consideration. Notice that since F_μ is a monotone function, it can have at most countably many points of discontinuity, so we are excluding only countably many points.

The fact that weak convergence is a robust notion is evidenced by the fact that it has a number of equivalent characterizations, as discussed in Section 4.1 below.

With this background in place, we can now state the Central Limit Theorem precisely, as follows:

Theorem 1 *Let X_0, X_1, X_2, \dots be a sequence of independent random variables with mean c , strictly positive variance σ^2 , and common distribution μ . Let $S_n = X_0 + X_1 + \dots + X_{n-1}$. Then the distribution of $(S_n - nc)/\sqrt{n\sigma^2}$ converges weakly to the standard normal distribution.*

This is Theorem 27.1 in Billingsley’s book [3]. Setting $c = 0$ in the statement of the theorem does not result in loss of generality: if each X_i has mean c , we can apply Theorem 1 to the shifted sequence $(X_i - c)$ and use the result to obtain the more general statement. Our formulation of Theorem 1 in Isabelle is as follows:

```

theorem (in prob_space) central_limit_theorem:
  fixes X :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
    and  $\mu$  :: "real measure"
    and  $\sigma$  c :: real
    and S :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
  assumes X_indep: "indep_vars ( $\lambda$ i. borel) X UNIV"
    and X_integrable: " $\bigwedge$ n. integrable M (X n)"
    and X_mean: " $\bigwedge$ n. expectation (X n) = c"
    and  $\sigma$ _pos: " $\sigma > 0$ "
    and X_square_integrable: " $\bigwedge$ n. integrable M ( $\lambda$ x. (X n x)2)"
    and X_variance: " $\bigwedge$ n. variance (X n) =  $\sigma^2$ "
    and X_distrib: " $\bigwedge$ n. distr M borel (X n) =  $\mu$ "
  defines "S n x  $\equiv$   $\sum$  i < n. X i x"
  shows "weak_conv_m ( $\lambda$ n. distr M borel ( $\lambda$ x. (S n x - n * c) / sqrt (n *  $\sigma^2$ )))
    std_normal_distribution"

```

Here, M denotes the underlying probability space. We present a formal proof of the mean zero case in the appendix to this paper, and then derive the version above as a corollary.

2.3 An overview of the proof

Contemporary proofs of the Central Limit Theorem rely on the use of *characteristic functions*, a powerful method that dates back to Laplace. If μ is a real-valued distribution, its characteristic function $\varphi(t)$ is defined by

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} \mu(dx).$$

In words, $\varphi(t)$ is the integral of the function $f(x) = e^{itx}$ over the whole real line, with respect to the measure μ . Notice that for each $t \neq 0$, the function e^{itx} is periodic with period $2\pi/t$. It might be helpful to think of e^{itx} as like a sine or cosine; indeed, $e^{itx} = \cos(tx) + i\sin(tx)$. Notice that $\varphi(0)$ is equal to 1, the measure of the entire real line. The characteristic function of a real distribution μ is a Fourier transform of the measure μ , and when $t \neq 0$, $\varphi(t)$ “detects” periodicity in the way that the real distribution μ distributes its “weight” over different parts of the real line.

A key property of characteristic functions is the fact that if X_1 and X_2 are independent random variables, then the characteristic function of $X_1 + X_2$ is the product of the characteristic function of X_1 and the characteristic function of X_2 . Of course, this extends to sums with any finite number of terms, and the resulting products are often convenient to work with.

The *Lévy Uniqueness Theorem* asserts that if μ_1 and μ_2 have the same characteristic function, then $\mu_1 = \mu_2$. In other words, a measure μ can be “reconstructed” from its characteristic function, and the characteristic function of a measure determines the measure uniquely. Let (μ_n) be a sequence of distributions, where each μ_n has characteristic function φ_n , and let μ be a distribution with characteristic function φ . The *Lévy Continuity Theorem* states that μ_n converges to μ weakly if and only if $\varphi_n(t)$ converges to $\varphi(t)$ for every t .

Remember that the CLT asserts that if (X_n) is a sequence of random variable satisfying certain hypotheses, and μ_n is for each n a certain distribution defined in terms of X_1, \dots, X_n , then μ_n converges weakly to the standard normal distribution. The Lévy Continuity Theorem provides a straightforward strategy to prove the theorem: if we let φ_n denote the characteristic function of μ_n for each n , we need only show that φ_n approaches the characteristic function of the standard normal distribution pointwise.

Implementing this strategy requires two key ingredients. First, one needs to know that the characteristic function of the standard normal distribution is $\varphi(t) = e^{-t^2/2}$. Second, one needs to compute the characteristic functions of the distributions μ_n , which are defined in terms of finite sums of the independent random variables X_0, X_1, \dots , and show that they have the desired behavior. This is where the key property of characteristic functions comes into play.

Once all these components were in place, putting the pieces together was not hard. Given the continuity theorem, the characteristic function of the standard normal distribution, the result on the characteristic functions of sums of random variables, and suitable approximations to the complex exponential function, the proof of the Central Limit Theorem is quite short. In our formalization, it is only about 120 lines long, and is presented in full in the appendix.

3 Isabelle and its libraries

When we began our project, a good deal of infrastructure was already available in the Isabelle libraries, but we had to add to it substantially. The formalization thus provided a stress test, allowing us to fill in gaps in the library and ensure its practical efficacy. In this section, we will describe those features of Isabelle and its libraries that were most relevant to the formalization, and indicate some of our contributions to the latter.

3.1 The Isabelle proof assistant

The Isabelle proof assistant [17] is based on classical simple type theory [6], with variables ranging polymorphically over types, and a Hilbert choice operator (*SOME*) which returns an indeterminate element satisfying a given predicate, if there is one. Given a type α and a predicate P on α , one can introduce a new abstract type representing the elements of α satisfying P , using a *typedef* command. In addition to the references given in this section, one should consult the documentation available on the Isabelle web site for the most up-to-date information.

Isabelle is an LCF-style theorem prover. This means that stating a theorem amounts to introducing a proof goal, and one can then construct proofs by applying *tactics* that reduce that goal to other goals. Layered on top of that, the Isabelle system includes the *Isar* proof language [22], which provides a natural, declarative way of writing structured proofs. Although in some places our proofs resort to sequences of tactic applications, for the most part we relied on *Isar* to make our proofs more robust, and to make them easier to read and maintain.

One attractive feature of Isabelle is the strength of its automation. We relied extensively on built-in procedures such as its term rewriter (*simp*), its generic theorem provers (such as *auto*), and its procedure for linear arithmetic (*arith*). Occasionally we relied on Isabelle’s *sledgehammer* command [18], which invokes external theorem provers and then reconstructs the results in Isabelle.

Isabelle has two mechanisms for reasoning algebraically and generically. The first, axiomatic type classes [21], constitutes a conservative extension of the axiomatic framework. Type variables are allowed to range over types with associated functions and relations, satisfying specified axioms. Theorems can be proved generically within a type class, and then instantiated to concrete structures that have been shown to satisfy the given axioms. This is used, for example, to develop facts about arithmetic, sums, products, and orderings that are shared among various number classes, including the natural numbers, integers, rationals, reals, and complex numbers. In the Isabelle library, they are also used to associate a topological structure to a type. Thus, our reasoning about topological aspects of the reals, described in Section 3.2 below, made use of the associated type classes. There are also classes for various types of normed spaces, of which the reals, finite powers of the reals, and the complex numbers are instances. As described in Section 3.6, we made use of these structures on the real and complex numbers.

Type classes are limited by the fact that they can only be parameterized by a single type parameter. An even more serious limitation is that, in simple type theory, types cannot depend on elements of other types. For example, there is no way of using type classes to reason about \mathbb{Z}_n , the integers modulo a parameter n , as an instance of a ring. For that purpose, Isabelle has a more flexible mechanism, *locales* [2], which, however, cannot take advantage of all of the benefits of the ambient type theory. *Locales* do *not* constitute an axiomatic extension; in terms of the underlying logic, a locale is nothing more than a predicate on some data. But Isabelle provides mechanisms for reasoning “in” such a locale, that is, fixing some data and the locale assumptions, and reasoning on that basis. *Locales* can also introduce notation that implicitly depends on the locale parameters. Isabelle provides mechanisms for instantiating locales, either with fixed types, or on the fly in a proof: one shows (typically with the help of automation) that some data satisfy the locale axioms, at which point all the definitions and theorems of the

locale are made available, pre-instantiated with the relevant data and facts. For example, once we show that a relation satisfies the axioms of the `partial_order` locale, we can use notation and facts about partial orders freely for that relation, without having to repeatedly cite the fact that the relation is a partial order.

In the measure theory library, as described in Sections 3.3 and 3.5, locales are used to reason about algebras and σ -algebras, for example. Moreover, once the type of measurable spaces has been introduced, locales are used to introduce extra hypotheses, for example, the hypothesis that a measure space is finite, or is a probability space.

3.2 Topology and Limits

Isabelle's extensive library for topological spaces includes properties of open and closed sets, limits, compactness, continuity, and so on. The library is described in detail by Hölzl, Immler, and Huffman in [13]. Topological notions interact with measure-theoretic notions in various ways. For example, a real distribution is a measure on the real numbers that measures the *Borel sets*, the smallest σ -algebra containing the open sets. Continuous functions are therefore measurable. Topological notions come into the statements of many measure-theoretic theorems described below, notions including the points of continuity of a function or the boundary of a set. Proving Skorohod's theorem required showing that the set of points of continuity of an arbitrary function from reals to reals is Borel; this is done in a four-line footnote in Billingsley ([3, page 334]), and requires characterizing the set of discontinuities as a union of an intersection of open sets.

Conventional reasoning about limits was ubiquitous in our formalization. Everyday mathematics requires one to deal with expressions such as the following:

- $\lim_{x \rightarrow a} f(x) = b$
- $\lim_{n \rightarrow \infty} a_n = a$
- $\lim_{x \rightarrow \infty} f(x) = b$
- $\lim_{x \rightarrow a^-} f(x) = b$
- $\lim_{x \rightarrow a} f(x) = \infty$

Here, the source and target spaces can be any topological space, including metric spaces or the natural numbers with the order topology. One can consider limits as x approaches a value a , or ∞ , or $-\infty$. One can also restrict the allowed values for x and consider the limit as x approaches a within a set s ; saying x approaches a from the left (where x and a are real-valued, for example) is equivalent to saying that x approaches a within the interval $(-\infty, a)$. There is a similar range of variations on the output: $f(x)$ can approach a value, b , or ∞ , or $-\infty$; and it can approach the value from the left, or from the right, or within any subset of the range of f . Not only does this threaten a combinatorial explosion of definitions, but also redundancy. For example, assuming $f(x)$ and $g(x)$ converge as x approaches a , we have the identity $\lim_{x \rightarrow a} (f(x) + g(x)) = \lim_{x \rightarrow a} f(x) + \lim_{x \rightarrow a} g(x)$, but this also holds under all the variations of convergence in the source.

To handle the many instances of convergence that arose in the formalization, we used Isabelle's elegant library for dealing with limits via filters [13]. The idea is that when dealing with any notion of limit, the relevant notions of convergence in the source and the target can be represented by *filters*. A *filter* over X is a

nonempty set $\mathcal{F} \subseteq \mathcal{P}(X)$ such that if $A \subseteq B$ and $A \in \mathcal{F}$, then $B \in \mathcal{F}$, and if $A, B \in \mathcal{F}$, then $A \cap B \in \mathcal{F}$. The general notion of limit in Isabelle, `filterlim f F1 F2`, says, roughly, that the function f converges in the sense of F_2 as the input converges in the sense of F_1 . By specializing F_1 and F_2 appropriately, we obtain all the variations described in the last paragraph, and more. In addition, theorems can be proved at the appropriate level of generality. For example, we have:

```
lemma tendsto_add:
  fixes f g :: "_ => 'a::topological_monoid_add"
  assumes "(f -> a) F" and "(g -> b) F"
  shows "((λx. f x + g x) -> a + b) F"
```

Here $(f \longrightarrow x) F$ is an abbreviation for `filterlim f F (nhds x)`, where `nhds x` is the filter of topological neighborhoods of x . This avoids the need to formalize endless variations of the same theorem; we only need to instantiate the general version to the relevant filters. Details can be found in [13].

In the Isabelle library, topological facts are found in the standard `HOL` library, in files such as `Filter`, `Topological.Spaces`, and `Limits`. Topological notions for the reals, and vector spaces over the reals, can be found in `Real.Vector.Spaces`.

3.3 Measure theory and integration

Our formalization required the fundamentals of measure theory and integration, as described in any introductory textbook on the subject (including Billingsley). The fundamental development of the subject is well described in the paper “Three chapters of measure theory” [12]. Therefore we only summarize key features of this development here, and indicate some of the ways the library has changed as a result of our formalization.

Measure theory requires the use of the *extended nonnegative reals* $\overline{\mathbb{R}}_{\geq 0}$, obtained by restricting the usual reals \mathbb{R} to their non-negative part and adding the value ∞ . A σ -algebra is represented in Isabelle as a record s depending on a type α , which specifies an underlying subset of α , a set of elements *space* S , and a collection of subsets *sets* S of *space* S that contains the empty set and is closed under countable unions and complements. These assumptions are specified as a locale.

A *measure space* M extends the notion of a σ -algebra with a function *emeasure* M from subsets of α to the extended nonnegative reals, satisfying the usual axioms: the measure of the empty set is 0, and the measure of a countable disjoint union of sets in the underlying σ -algebra is equal to sum of the measures of each set in the union (which might be ∞). The underlying σ algebra corresponds to the usual notion of the collection of *measurable subsets* corresponding to the measure. In Isabelle, for any type α , the *typedef* mechanism is used to specify a new type, *measure* α , consisting of measure spaces on some subset of α .

If \mathcal{M} and \mathcal{N} are two measure spaces, a *measurable function* f from \mathcal{M} to \mathcal{N} (written $f \in \mathcal{M} \rightarrow_{\mathcal{M}} \mathcal{N}$) is a function between the underlying sets that has the property that the inverse image of any measurable subset of the codomain is a measurable subset of the domain. Note that, in fact, the property of being measurable has nothing to do with the measure; it is really a property of the function with respect to the two associated σ -algebras. Given a measurable function $f : \mathcal{M} \rightarrow \mathcal{N}$, a measure μ on \mathcal{M} gives rise to a new measure ν on \mathcal{N} , defined by $\nu(A) = \mu(f^{-1}A)$.

This is sometimes called a *pushforward measure*, but in Isabelle it is denoted $\text{distr } M \ N \ f$, for reasons that are explained in Section 3.5. It is defined formally as follows:

```

definition
  distr :: "'a measure  $\Rightarrow$  'b measure  $\Rightarrow$  ('a  $\Rightarrow$  'b)  $\Rightarrow$  'b measure"
where
  "distr M N f =
    measure_of (space N) (sets N) ( $\lambda A$ . emeasure M (f -' A  $\cap$  space M))"

```

Another way to define a measure ν in terms of a measure μ on \mathcal{M} is to take a measurable function f from \mathcal{M} to $\overline{\mathbb{R}}_{\geq 0}$ and define, for every set A , $\nu(A) = \int^+ f \chi_A d\mu$. This is defined formally in the Isabelle library as follows:

```

definition
  density :: "'a measure  $\Rightarrow$  ('a  $\Rightarrow$  ennreal)  $\Rightarrow$  'a measure"
where
  "density M f =
    measure_of (space M) (sets M) ( $\lambda A$ .  $\int^+ x. f x * \text{indicator } A \ x \ \partial M$ )"

```

The function χ_A is the *characteristic function* of A , also called the *indicator function*. The integral $\int^+ f d\mu$ is the *nonnegative Lebesgue integral*, defined for functions into $\overline{\mathbb{R}}_{\geq 0}$. For measurable functions it has the expected properties: it is closed under addition, constant multiplication, and monotone convergence. It is monotone even for non-measurable functions, which simplifies certain proofs, since measurability is not always easy to prove. Because the nonnegative Lebesgue integral takes values in $\overline{\mathbb{R}}_{\geq 0}$, it is well-defined for all measurable functions, even when the integral is infinite.

The library includes a construction of the Borel sets in any topology, and the Carathéodory extension theorem. In Isabelle, the Lebesgue measure on the reals was initially constructed from the gauge integral, which is discussed in Section 3.6. After our formalization, however, the construction was replaced by the more common textbook definition as the σ -extension of the measure on finite intervals, as described in Section 3.8 below. The measure space consisting of the Lebesgue measure on the Borel subsets of the reals is denoted `lborel` in the Isabelle library.

The fundamentals of measure theory are found in the `HOL-Probability` library, including `Sigma_Algebra`, `Measure_Space`, `Caratheodory`, and `Lebesgue_Measure`.

3.4 Bochner integration

Our initial formalization of the Central Limit Theorem relied on the theory of Lebesgue integration, described in [12]. This provides a notion of integration for suitable functions $f : X \rightarrow \mathbb{R}$, where X is any space on which a measure is defined. After we completed the proof, however, the second author, Hölzl, generalized the construction to the *Bochner integral*. This provides a theory of integration for functions $f : X \rightarrow B$, where now B is any second-countable Banach space. In particular, B can be any of the spaces \mathbb{R}^n , or the complex numbers, \mathbb{C} . Our formalization made extensive use of integration of functions from \mathbb{R} to \mathbb{C} , as discussed in Section 3.6.

Similar to the Lebesgue integral, the Bochner integral approximates a function f by a sequence of simple functions s . Each simple function has a finite range, and hence its integral can be expressed by finite summation. When a function can be approximated by simple functions, its integral is the limit of the integrals

of those simple functions. Whereas approximations for Lebesgue integration are taken with respect to the pointwise order on $\overline{\mathbb{R}}_{\geq 0}$, approximations for Bochner integration are taken with respect to the L^1 -norm, defined using the Lebesgue integral by $\|f\| = \int^+ |f| d\mu$.

More formally, a function s is *simple Bochner-integrable* if s is Borel-measurable on \mathcal{M} , has a finite range $f[\mathcal{M}]$, and a support $\{x \in \mathcal{M} \mid f(x) \neq 0\}$ with finite measure. The integral of a simple Bochner function s is a finite sum over the vectors of the range of s times the measure of their support:

$$\int s \, d\mu = \sum_{y \in f[\mathcal{M}]} \mu(f^{-1}[y]) \cdot y$$

A Borel-measurable function f is Bochner-integrable if there is a sequence (s_i) of simple Bochner-integrable functions such that:

1. f is the limit of (s_i) in the L^1 norm, i.e. $\lim_{i \rightarrow \infty} \|s_i - f\| = 0$; and
2. the sequence of integrals of the functions s_i converges, i.e. $\lim_{i \rightarrow \infty} \int s_i \, d\mu$ exists.

In that case, the Bochner integral, denoted $LINT \, x|M. \, f \, x$ in Isabelle, is defined by

$$\int f \, d\mu = \lim_{i \rightarrow \infty} \int s_i \, d\mu$$

The notation $LINT$ is a holdover from Lebesgue integration, but since Bochner integration functions in similar ways, the notation is still a useful mnemonic.

From the definition it follows that each Bochner-integrable function f is Borel-measurable and has a finite L^1 -norm. In the other direction we prove that each Borel-measurable function is approximated pointwise by a sequence of simple Bochner functions. Then it follows that a function f is Bochner-integrable if and only if f is measurable and the L^1 -norm of f is finite (which is equivalent to saying that f is absolutely integrable).

```
lemma integrable_iff_bounded:
  fixes f :: "'a ⇒ 'b::{banach, second_countable_topology}"
  shows "integrable M f ⟷ f ∈ M →M borel ∧ (∫+x. norm (f x) ∂M) < ∞"
```

As one would expect of an integral, the Bochner integral respects scalar multiplication and addition. As with the Lebesgue integral, we obtain a version of the dominated convergence theorem:

```
lemma dominated_convergence:
  fixes f :: "'a ⇒ 'b::{banach, second_countable_topology}"
  and w :: "'a ⇒ real"
  assumes "f ∈ M →M borel" "∧i. s i ∈ M →M borel" "integrable M w"
  and "AE x in M. (λi. s i x) ⟶ f x"
  and "∧i. AE x in M. norm (s i x) ≤ w x"
  shows "integrable M f" and "∧i. integrable M (s i)"
  and "(λi. LINT x|M. s i x) ⟶ (LINT x|M. f x)"
```

Here the quantifier $AE \, x \, in \, M$ expresses that the subsequent statement holds for almost every element x of the measure space M , which is to say, the set of examples where it doesn't hold has measure zero. We have the monotone convergence theorem, which applies to sequences of functions taking values in the real numbers. We also obtain Fubini's theorem. These are all staples of the theory of integration, and were used throughout our formalization.

If $f : X \rightarrow B$ is any measurable function on a space X with measure μ , and S is any measurable set, one can define the integral over the set S by $\int_S f d\mu = \int f \chi_S d\mu$. Rather than introduce a new definition, we took notation for integration over sets to be an abbreviation for the definition in terms of indicators, with the notation $LINT x:S|M. f x$. But because reasoning about integrals over sets is so fundamental, we found it helpful to develop a small library to support it. For example, the following is a consequence of the dominated convergence theorem:

```
lemma integral_countable_add:
  fixes f :: "_  $\Rightarrow$  'a :: {banach, second_countable_topology}"
  assumes " $\bigwedge i::nat. A i \in sets M$ "
    and " $\bigwedge i j. i \neq j \implies A i \cap A j = \{\}$ "
    and "set_integrable M ( $\bigcup i. A i$ ) f"
  shows "LINT x:( $\bigcup i. A i$ )|M. f x = ( $\sum i. (LINT x:(A i)|M. f x)$ )"
```

The theory of Bochner integration is included in the HOL-Probability library, in `Bochner_Integration` and `Set_Integral`.

3.5 Probability

Modern probability is based on measure theory, although probabilists and statisticians tend to adopt their own distinct terminology. A *probability space* is simply a measure space in which the measure of the entire space is equal to 1. You should think of the space as the space of possible outcomes of a random event. A *random variable* on such a space is a measurable function from that space to the reals; think of it as a real number that depends on the outcome of the random event. The *expectation* of a random variable is the integral of the function over the entire space. Thus talk of probability spaces, random variables, and expectations is really talk of measure spaces, measurable functions, and integrals in disguise.

The Isabelle library defines a locale for *finite measures*, which are simply measures for which the measure of the entire space, $emeasure M (space M)$, is not infinity. For such spaces, one can work more conveniently with the associated real-valued function, $measure$, which casts the value of $emeasure$ to a real. There is also a locale for *probability measures*, which are finite measures where the measure of the entire space is equal to 1. When working with cumulative distribution functions, as described in Section 3.8, we found it convenient to define a locale for *real distributions*; a real distribution is a probability space in which the space is the set of real numbers and the measurable sets consist of exactly the Borel subsets of the reals. To capture the language of informal probability theory, the library defines all of the following abbreviations:

```
locale prob_space =
  fixes M :: "'a measure" assumes "emeasure M (space M) = 1"
begin
  abbreviation "events  $\equiv$  sets M"
  abbreviation "prob  $\equiv$  measure M"
  abbreviation "random_variable M' X  $\equiv$  X  $\in$  M  $\rightarrow_M$  M'"
  abbreviation "expectation X  $\equiv$  (LINT x|M. X x)"
  abbreviation "variance X  $\equiv$  (LINT x|M. (X x - expectation X)2)"
end
```

If X is a random variable on a measure space \mathcal{M} with measure μ , the distribution of X , as described in the previous section, has the following interpretation: it

is the measure ν on the Borel sets of \mathbb{R} such that for every A , $\nu(A)$ is the probability that X takes a value in A . Even though we think of X as depending on some underlying source of randomness, represented by \mathcal{M} , often we only care about the induced probability on the real numbers that is given by its distribution. Notice that the word “distribution” is used in probability theory in at least three distinct but related ways. In addition to the uses of the term described in this paragraph and the previous one, one also often speaks of the (cumulative) distribution function of a real distribution, as described in Section 3.8. Thus, if X is a random variable, its distribution is a real distribution, which in turn has a distribution function.

In probability theory, real distributions are often specified as densities, as described in the previous section. Thus the normal distribution with mean μ and variance σ is defined formally as follows:

```

definition
  normal_density :: "real  $\Rightarrow$  real  $\Rightarrow$  real  $\Rightarrow$  real"
where
  "normal_density  $\mu$   $\sigma$   $x$  = 1 / sqrt (2 * pi *  $\sigma^2$ ) * exp (-(x -  $\mu$ )2 / (2 *  $\sigma^2$ ))"

abbreviation
  std_normal_density :: "real  $\Rightarrow$  real"
where
  "std_normal_density  $\equiv$  normal_density 0 1"

abbreviation
  std_normal_distribution :: "real measure"
where
  "std_normal_distribution  $\equiv$  density lborel std_normal_density"

```

Various notions of independence are used in probability. Perhaps the most general is the following: suppose that for every i in some index set I , F_i is a collection of events (measurable sets) from some fixed measure space. Then the sequence $(F_i)_{i \in I}$ is said to be *independent* if for every finite subset $J \subseteq I$ and every choice of a set $A_j \in F_j$ for each j , the probability of the intersection $\bigcap_{j \in J} A_j$, i.e. the probability that all of the A_j 's occur, is the product of the individual probabilities.

```

definition (in prob_space)
  indep_sets :: "('i  $\Rightarrow$  'a set set)  $\Rightarrow$  'i set  $\Rightarrow$  bool"
where
  "indep_sets F I  $\longleftrightarrow$ 
    ( $\forall i \in I. F i \subseteq$  events)  $\wedge$ 
    ( $\forall J \subseteq I. J \neq \{\}$   $\longrightarrow$  finite J  $\longrightarrow$ 
      ( $\forall A \in (\prod i \in J. F i). \text{prob } (\bigcap j \in J. A j) = (\prod j \in J. \text{prob } (A j))$ ))"

```

If now $(A_i)_{i \in I}$ is a sequence of *events* (rather than collections of events), saying that the sequence (A_i) is independent amounts to saying that the sequence of singletons $(\{A_i\})_{i \in I}$ is an independent sequence of collections.

```

definition (in prob_space) "indep_events A I  $\longleftrightarrow$  indep_sets ( $\lambda i. \{A i\}$ ) I"

```

Finally, if $(X_i)_{i \in I}$ is a sequence of random variables with inputs in one measure space, \mathcal{M} , and values in another space, \mathcal{M}' (typically, but not necessarily, the reals), saying that the sequence X_i is independent amounts to saying that the sequence of collections of measurable sets

$$(\{X_i^{-1}(A) \mid A \text{ is a measurable subset of } \mathcal{M}'\})_{i \in I}.$$

is independent.

```

definition (in prob_space)
  indep_vars :: "('i ⇒ 'b measure) ⇒ ('i ⇒ 'a ⇒ 'b) ⇒ 'i set ⇒ bool"
where
  "indep_vars M' X I ←→
    (∀ i ∈ I. random_variable (M' i) (X i)) ∧
    indep_sets (λi. { X i -' A ∩ space M | A. A ∈ sets (M' i) }) I"

```

In probabilistic terms, this means that given any finite $J \subseteq I$ and any finite sequence A_{j_1}, \dots, A_{j_n} of events, the probability that each X_{j_u} is in A_{j_u} is just the product of the individual probabilities. Of course, we can say that any *two* events, or random variables, or collections of events, are independent by taking I to be any two-element type, such as the Booleans. Isabelle’s library defines the binary notions as well, and develops basic properties of independent sets, events, and random variables.

In the Isabelle 2016 distribution, these developments are in the `HOL-Probability` library, including the files `Probability_Measure`, `Independent_Family`, `Convolution`, and `Distributions`.

3.6 Real analysis and complex-valued functions

Isabelle has an extensive library for real multivariate analysis, which is again well-described in [13]. In Isabelle, the reals are instantiated as a complete ordered field, and as a conditionally complete lattice, which means that nonempty bounded sets have sups and infs. The library also includes definitions of transcendental functions like the sine, cosine, and exponential functions. In fact, the exponential function is defined generically for any Banach space, including the complex numbers. Of course, we have the relation $e^{ix} = \cos x + i \sin x$ for real x .

Isabelle’s general notion of the derivative is the *Fréchet derivative*, which makes sense for functions f between any two Banach spaces. As with limits, the notion of Fréchet derivative supports multiple modes of convergence; the expression `(f has_derivative D) F` means that the function f has the bounded linear functional D as derivative “at” the filter F . In practice, F is usually the filter expressing that D is the derivative at a point x , or that D is the derivative at a point x when we restrict attention to a subset S of the source. The more familiar notion of the scalar derivative for functions from the reals to reals (or, more generally, from one normed field to another) is derived from the Fréchet derivative as a special case. So is the notion of a vector derivative for functions from \mathbb{R} to \mathbb{R}^n .

The characteristic function of a measure is a function from the reals, \mathbb{R} , to the complex numbers, \mathbb{C} . The theory of such functions is much simpler than the theory of functions from \mathbb{C} to \mathbb{C} , which is the subject of complex analysis. One can view a function $f : \mathbb{R} \rightarrow \mathbb{C}$ as essentially two functions from \mathbb{R} to \mathbb{R} , f^{re} and f^{im} , the first returning the real part and the second returning the imaginary part of the output. Integrals and derivatives of such functions can be understood in terms of the integrals and derivatives of these two parts.

In fact, for differentiation, we did not have to define a new notion of derivative: if we view the complex numbers as a two-dimensional real Banach space, the derivative we need is nothing more than the Fréchet derivative.

For integration, the story is more involved. Isabelle’s library now has two forms of the integral. The multivariate analysis library generally relies on the gauge

integral, which is defined for functions from \mathbb{R}^n to \mathbb{R} . When we consider the reals, \mathbb{R} , with the usual Lebesgue measure, the Bochner integral and the gauge integral agree on finite intervals, but otherwise the gauge integral is slightly more general: for a function $f : \mathbb{R} \rightarrow \mathbb{R}$ to be Bochner-integrable, both the positive and negative parts of f have to have a finite Bochner integral, whereas the gauge integral can accommodate some functions whose positive and negative parts cancel each other out in a suitable fashion. Nonetheless, for the vast majority of applications, the Bochner integral is quite sufficient. Since our formalization required integration with general measures and spaces in addition to the usual integration over \mathbb{R}^n , we used the Bochner integral throughout.

With the Bochner integral, as with the Fréchet derivative, integrating functions taking values in \mathbb{C} is no different from integrating functions taking values in \mathbb{R} . Indeed, this was the primary motivation for generalizing from the Lebesgue integral to the Bochner integral.

3.7 Calculus

Our formalization required extensive use of calculus at an undergraduate level, including integration by parts, Taylor series approximations, changes of variable, and so on. For example, the calculation of moments of the normal distribution required the following estimate on the complex exponential:

$$\left| e^{ix} - \sum_{k=0}^n \frac{(ix)^k}{k!} \right| \leq \min \left(\frac{|x|^{n+1}}{(n+1)!}, \frac{2x^n}{n!} \right).$$

We followed Billingsley [3, Section 26] in obtaining this using an inductive argument and integration by parts. Notice that this involves reasoning about functions from the real to complex numbers; as explained in the previous section, the relevant properties generally follow from the corresponding properties for real-valued functions, upon splitting functions to the real and imaginary parts. In addition, we have the general inequality $\|\int_A f d\mu\| \leq \int_A \|f\| d\mu$, which allows us to bound the modulus of a complex integral by bounding the real-valued integral of the norm. This is an instance of a more general fact about the Bochner integral:

```
lemma integral_norm_bound:
  fixes f :: "'a ⇒ 'b :: {banach, second_countable_topology}"
  shows "integrable M f ⇒ norm (LINT x|M. f x) ≤ (LINT x|M. norm (f x))"
```

Textbook results from calculus involve integrals $\int_a^b f(x) dx$ over the interval (a, b) . These can be viewed as ordinary integrals over the set (a, b) , with the following two caveats:

- Textbooks allow a to be $-\infty$ and allow b to be ∞ , which is to say, a and b should be taken to be extended real numbers (i.e. the reals extended with $\pm\infty$).
- It is convenient to adopt the convention that if $b < a$, then

$$\int_a^b f(x) dx = - \int_b^a f(x) dx.$$

We thus defined a notion of “interval integral” along these lines, together with supporting the notation $LBINT x=a..b. f x$. We could then state the first fundamental theorem of calculus in the following form, for finite intervals:

```

lemma interval_integral_FTC_finite:
  fixes f F :: "real  $\Rightarrow$  'a::euclidean_space" and a b :: real
  assumes f: "continuous_on {min a b..max a b} f"
    and F: " $\bigwedge x. \min a b \leq x \implies x \leq \max a b \implies$ 
      (F has_vector_derivative (f x)) (at x within {min a b..max a b})"
  shows "(LBINT x=a..b. f x) = F b - F a"

```

The following version, for arbitrary intervals, makes sense when the limits are infinite:

```

lemma interval_integral_FTC_integrable:
  fixes f F :: "real  $\Rightarrow$  'a::euclidean_space" and a b :: ereal
  assumes "a < b"
    and " $\bigwedge x. a < ereal x \implies ereal x < b \implies$ 
      (F has_vector_derivative f x) (at x)"
    and " $\bigwedge x. a < ereal x \implies ereal x < b \implies$  continuous (at x) f"
    and "set_integrable lborel (einterval a b) f"
    and "(F  $\circ$  real_of_ereal)  $\longrightarrow$  A) (at_right a)"
    and "(F  $\circ$  real_of_ereal)  $\longrightarrow$  B) (at_left b)"
  shows "(LBINT x=a..b. f x) = B - A"

```

Similarly, we could state the second fundamental theorem of calculus, where the variable bound to the integral can be before or after the fixed endpoint:

```

lemma interval_integral_FTC2:
  fixes a b c x :: real and f :: "real  $\Rightarrow$  'a::euclidean_space"
  assumes "a  $\leq$  c" "c  $\leq$  b" "continuous_on {a..b} f" "a  $\leq$  x" "x  $\leq$  b"
  shows "(( $\lambda u. LBINT y=c..u. f y$ ) has_vector_derivative (f x))
    (at x within {a..b})"

```

The use of such an integral was a mixed blessing. It simplified many of our theorems and proofs, but at the expense of introducing yet another notion of integral, which required another library of supporting facts, as well as, at times, translations to and from the other notions of integral.

Many textbook integration arguments require a change of variable, sometimes known as “integration by substitution.” It was not hard to prove that if a function g from \mathbb{R} to \mathbb{R} has a continuous derivative (and hence is continuous itself) on a closed interval $[a, b]$, and f is continuous on the image of $[a, b]$ under g , then $\int_a^b f(g(x))g'(x) dx = \int_{g(a)}^{g(b)} f(x) dx$.

```

lemma interval_integral_substitution_finite:
  fixes a b :: real and f :: "real  $\Rightarrow$  'a::euclidean_space"
  assumes "a  $\leq$  b" and " $\bigwedge x. a \leq x \implies x \leq b \implies$ 
      (g has_real_derivative (g' x)) (at x within {a..b})"
    and "continuous_on (g ' {a..b}) f" "continuous_on {a..b} g'"
  shows "LBINT x=a..b. g' x *R f (g x) = LBINT y=g a..g b. f y"

```

Manuel Eberl later generalized this to arbitrary Borel measurable functions f , but with the added hypothesis that g' is nonnegative on $[a, b]$. However, we also needed a version of the theorem for intervals with potentially infinite endpoints. This requires using either the monotone convergence theorem or the dominated convergence theorem to pass from finite interval approximations to the full interval. In fact, we proved two versions. The following one requires showing independently that both $f(x)$ and $f(g(x))g'(x)$ are integrable over the relevant intervals:

```

lemma interval_integral_substitution_integrable:
  fixes f :: "real  $\Rightarrow$  'a::euclidean_space" and a b A B :: ereal
  assumes "a < b"

```

```

and " $\bigwedge x. a < \text{ereal } x \implies \text{ereal } x < b \implies \text{DERIV } g \ x \ :> g' \ x$ "
and " $\bigwedge x. a < \text{ereal } x \implies \text{ereal } x < b \implies \text{continuous (at (g x)) } f$ "
and " $\bigwedge x. a < \text{ereal } x \implies \text{ereal } x < b \implies \text{continuous (at x) } g'$ "
and " $\bigwedge x. a \leq \text{ereal } x \implies \text{ereal } x \leq b \implies 0 \leq g' \ x$ "
and " $((\text{ereal} \circ g \circ \text{real\_of\_ereal}) \longrightarrow A) \text{ (at\_right a)}$ "
and " $((\text{ereal} \circ g \circ \text{real\_of\_ereal}) \longrightarrow B) \text{ (at\_left b)}$ "
and " $\text{set\_integrable lborel (einterval a b) } (\lambda x. g' \ x \ *_R \ f \ (g \ x))$ "
and " $\text{set\_integrable lborel (einterval A B) } (\lambda x. f \ x)$ "
shows " $(\text{LBINT } x=A..B. f \ x) = (\text{LBINT } x=a..b. g' \ x \ *_R \ f \ (g \ x))$ "

```

Another version assumes instead that f is nonnegative, and concludes that f is therefore integrable.

As an example where various uses of these components came together, consider the *sine integral function*. The function $\sin x/x$ is undefined at 0, but it can be made continuous at 0 by giving it the value 1 there. The resulting function is called *sinc*. The sine integral function is (confusingly) defined to be the indefinite integral of the sinc function, starting at 0:

$$\text{Si}(t) = \int_0^t \text{sinc } x \, dx.$$

The proof of the Lévy inversion formula uses the fact that

$$\lim_{t \rightarrow \infty} \text{Si}(t) = \frac{\pi}{2}.$$

A textbook proof (sketched in [3, Example 18.4]) runs as follows. By the fundamental theorem of calculus, we can verify that

$$\int_0^t e^{-ux} \sin x \, dx = \frac{1}{1+u^2} [1 - e^{-ut} (u \sin t + \cos t)]$$

by taking the derivative of both sides. Calculating, we can also show that

$$\int_0^t \left(\int_0^\infty |e^{-ux} \sin x| \, du \right) dx = \int_0^t x^{-1} |\sin x| \, dx \leq t.$$

The fact that the double-integral on the left is finite means that Fubini's theorem may be used to change the order of integration of $e^{-ux} \sin x$ over $(0, t) \times (0, \infty)$. So we have

$$\begin{aligned} \int_0^t \frac{\sin x}{x} \, dx &= \int_0^t \sin x \left(\int_0^\infty e^{-ux} \, du \right) dx \\ &= \int_0^\infty \left(\int_0^t e^{-ux} \sin x \, dx \right) du \\ &= \int_0^\infty \frac{du}{1+u^2} - \int_0^\infty \frac{e^{-ut}}{1+u^2} (u \sin t + \cos t) \, du. \end{aligned}$$

Substituting $u = \tan x$ in the first term yields

$$\int_0^\infty \frac{du}{1+u^2} = \int_0^{\pi/2} \frac{1}{1+\tan^2 x} (1+\tan^2 x) \, dx = \pi/2,$$

and the change of variable $v = ut$ can be used to show that the second integral converges to 0 as $t \rightarrow \infty$. Hence

$$\lim_{t \rightarrow \infty} \text{Si}(t) = \lim_{t \rightarrow \infty} \int_0^t \frac{\sin x}{x} dx = \frac{\pi}{2},$$

as required.

Proving this result required a tremendous amount of formal machinery: not only suitable forms of substitution, but also Fubini's theorem, the fundamental theorem of calculus, integration by parts, integral comparisons, properties of limits, and properties of the tangent function. It also required a lot of work, establishing that the relevant functions were continuous, integrable, and so on. It was somewhat demoralizing that a small calculus exercise required so much effort, but it is a good illustration of the infrastructure that is needed to carry out the kinds of calculus computations that come up routinely in engineering, modeling, and the sciences. We faced a similar calculus exercise in computing the moments of the normal distribution, as described in Section 4.2.

In the Isabelle 2016 distribution, the formulation of the fundamental theorem of calculus that we used and the substitution theorems described above are in `Interval_Integral`. Eberl's generalization is in `Lebesgue_Integral_Substitution`. The calculation concerning Si is in the file `Sinc_Integral`.

3.8 Distribution functions and the Lebesgue-Stieltjes measure

Every measure on \mathbb{R} gives rise to the real-valued function which, at each input x , returns the amount of "mass" below that argument:

Definition 2 Let μ be a finite measure on \mathbb{R} . The *cumulative distribution function* F_μ is defined by $F_\mu(x) = \mu(-\infty, x]$.

The cumulative distribution function (or *cdf*) is sometimes also called, more simply, the *distribution function* of the measure. In Isabelle, the definition is rendered as follows:

```

definition
  cdf :: "real measure  $\Rightarrow$  real  $\Rightarrow$  real"
where
  "cdf M  $\equiv$   $\lambda x$ . measure M  $\{..x\}$ "

```

It is not hard to see that the distribution function F_μ of a finite Borel measure μ is nondecreasing and right-continuous, and satisfies $\lim_{x \rightarrow -\infty} F_\mu(x) = 0$.

```

lemma (in finite_borel_measure) cdf_nondecreasing:
  "x  $\leq$  y  $\implies$  cdf M x  $\leq$  cdf M y"

```

```

lemma (in finite_borel_measure) cdf_is_right_cont:
  "continuous (at_right a) (cdf M)"

```

```

lemma (in finite_borel_measure) cdf_lim_at_bot:
  "(cdf M  $\longrightarrow$  0) at_bot"

```

Conversely, it turns out that any function with these properties is the distribution of a Borel measure on \mathbb{R} . The requisite measure μ is constructed by defining $\mu(a, b] = F(b) - F(a)$ and extending this to the Borel σ -algebra using the Carathéodory extension theorem. To that end, we defined an operation, *interval_measure*, that generates a measure from a nondecreasing, right-continuous function. To use the Carathéodory extension theorem, the key property that needs to be verified is that if a half-open interval $(a, b]$ is written as a disjoint union of countably many intervals $(a_i, b_i]$, then $b - a = \sum_i (b_i - a_i)$. This is trickier than it sounds. For example, the interval $(0, 1]$ can be written as a countable union of intervals $(1/2^{i+1}, 1/2^i]$, and any one of *those* intervals could similarly be replaced by a countable union. It is not hard to show that the infinite sum $\sum_i (b_i - a_i]$ is bounded by $b - a$. In the other direction, one picks a small ε , enlarges each interval $(a_i, b_i]$ to a slightly larger interval $(a_i - \varepsilon/2^i, b_i + \varepsilon/2^i)$, argues that the union of the enlargements covers the closed interval $[a, b]$, and then appeals to the compactness of $[a, b]$. The measure associated to a right-continuous, nondecreasing function in this way is called the *Lebesgue-Stieltjes measure*. When the function $F(x)$ is the identity function, we obtain the Lebesgue-Borel measure *lborel*, and, in fact, this now serves as the definition of Lebesgue-Borel measure in the Isabelle library.

In the case of a probability measure, we have the additional property that $\lim_{x \rightarrow \infty} F_\mu(x) = 1$:

lemma (in *real_distribution*) *cdf_lim_at_top_prob*: "(cdf M \longrightarrow 1) at_top"

Conversely, any function F satisfying all four properties is a probability measure:

lemma *real_distribution_interval_measure*:
fixes $F :: \text{"real} \Rightarrow \text{real}"$
assumes "mono F" " \wedge a. continuous (at_right a) F"
and "(F \longrightarrow 0) at_bot" "(F \longrightarrow 1) at_top"
shows "real_distribution (interval_measure F)"

Recall that *real_distribution* is the name of the locale for probability measures on the Borel subsets of the reals. So, for any function F satisfying the four properties above, *interval_measure F* is the measure whose cdf is exactly F . The use of the word “the” is justified by the fact that the association is unique, in the sense that if two real distributions have the same cumulative distribution function, then they are equal:

lemma *cdf_unique*:
fixes $M1\ M2$
assumes "real_distribution M1" **and** "real_distribution M2"
and "cdf M1 = cdf M2"
shows "M1 = M2"

Thus one can pass freely between talk of measures on \mathbb{R} and of their distribution functions, a key fact in the proof of the CLT.

In the Isabelle 2016 distribution, the construction of the Lebesgue measure on the reals as a Lebesgue-Stieltjes measure is in the theory *Lebesgue_Measure*, and the correspondence between measures and their distribution functions is developed in *Distributions*.

3.9 Automation

To improve automation, Isabelle’s multivariate analysis library provides a large set of introduction rules, to establish things like openness or closedness of sets

or continuity of functions. Continuity is nicely reduced by compositionality; if we know that two functions are continuous, their composition is again continuous. Applying this as a rule requires matching terms of the form $f(g\ x)$ where both f and g are variables. But this is often not the right choice. The straightforward way to express that functions like multiplication and \ln are continuous is to write $\text{continuous_on } (\mathbb{R} \times \mathbb{R}) (\lambda(x, y). x * y)$ and $\text{continuous_on } (0, \infty) \ln$. The composition rule is then not sufficient to prove continuity of $\lambda x. \ln(1 + x * x)$, because it does not accommodate binary operations like $+$ and $*$. In addition, the composition rule does not allow for the fact that the domain of \ln has to be restricted to the positive reals.

A simple solution to these two problems is to state continuity rules *precomposed* with arbitrary continuous functions. For example, we can state the following rules for arbitrary f and g :

$$\frac{\text{continuous_on } A\ f \quad \forall x \in A. 0 < f\ x}{\text{continuous_on } A (\lambda x. \ln (f\ x))}$$

$$\frac{\text{continuous_on } A\ f \quad \text{continuous_on } A\ g}{\text{continuous_on } A (\lambda x. f\ x + g\ x)}$$

Now, to prove $\lambda x. \ln(1 + x * x)$, we just apply rules like these. We ultimately end up with the goal $\forall x. 0 < 1 + x * x$, which is proved by the simplifier. This idea goes back to a paper by Gottlieb [9], which describes an implementation in PVS.

Rules for establishing openness and closedness of sets are not as important, but nonetheless helpful. Besides the usual rules for intersections and unions, we also have rules working on logical connectives and relations in set-comprehension. For example, the set $\{x \mid f\ x < g\ x\}$ is open whenever f and g are continuous functions into real numbers.

Isabelle's automation for measurability uses precomposed rules in a similar way. A difference is that measurability is also integrated as a special-purpose simplification procedure (in Isabelle terminology, a *simproc*), called *measurable*. To use the measurability prover, the user needs to annotate all the relevant measurability assumptions with the `[measurable]` attribute. This measurability prover then tries to massage all added assumptions into the right form, and proves measurability statements by applying them as introduction rules. The massaging also includes destructions of certain compositions, e.g. the assumption that $\lambda x. (f\ x, g\ x)$ is $X \times Y$ -measurable is replaced by the fact that f is X -measurable and g is Y -measurable. Such destructions are important for higher-order proof steps like induction. As a special case it also allows us to decompose subterms with a countable range, since the measurability of $f(g\ x)\ x$ can be reduced to the measurability of g and the measurability of $f\ c\ x$ for all c in the range of g . It is also important to add measurability rules for logical connectives, including quantifiers over countable sets. As a result, predicates can also be proved measurable, and therefore expressions that depend on case distinctions.

An example of the power of such rule sets is given by the proof that the predicate “ f is continuous at x ” is measurable in x for a function f on metric spaces. We can express the continuity of f at x in the following way:

$$\forall i > 0. \exists j > 0. \forall y\ z. d(x, y) < \frac{1}{j} \wedge d(x, z) < \frac{1}{j} \implies d(f(y), f(z)) \leq \frac{1}{i}.$$

Then the proof that this is measurable is a straightforward application of rules, as follows: (1) the quantifiers over i and j are countable, hence measurable; (2) we get into a closed set by eliminating the quantifiers over y and z ; (3) for the implication, the right-hand side is constant, hence closed; and (4) the left-hand side is open, as it is a strict inequality between two continuous functions.

4 The proof of the Central Limit Theorem

4.1 Weak convergence

Recall from Section 2.2 that if (μ_n) is a sequence of real distributions and μ is a real distribution, then (μ_n) converges weakly to μ , written $\mu_n \Rightarrow \mu$, if $F_{\mu_n}(x)$ approaches $F_\mu(x)$ at each point x where F_μ is continuous. In Isabelle, this is expressed by the following two definitions:

```
definition
  weak_conv :: "(nat  $\Rightarrow$  (real  $\Rightarrow$  real))  $\Rightarrow$  (real  $\Rightarrow$  real)  $\Rightarrow$  bool"
where
  "weak_conv F_seq F  $\equiv$ 
    $\forall x.$  continuous (at x) F  $\longrightarrow$  ( $\lambda n.$  F_seq n x)  $\longrightarrow$  F x"
```

```
definition
  weak_conv_m :: "(nat  $\Rightarrow$  real measure)  $\Rightarrow$  real measure  $\Rightarrow$  bool"
where
  "weak_conv_m M_seq M  $\equiv$  weak_conv ( $\lambda n.$  cdf (M_seq n)) (cdf M)"
```

In words, a sequence of functions $(F_n)_{n \in \mathbb{N}}$ converges weakly to F if $(F_n(x))_{n \in \mathbb{N}}$ converges to $F(x)$ for each point x where F is continuous, and the sequence of measures (μ_n) converges weakly to μ if the corresponding cumulative distribution functions converge weakly.

That the notion of weak convergence is robust is supported by the fact that there are a number of equivalent characterizations. The following theorem is sometimes known as the *Portmanteau Theorem*:

Theorem 2 *The following are equivalent:*

1. $\mu_n \Rightarrow \mu$.
2. $\int f d\mu_n$ approaches $\int f d\mu$ for every bounded function f that is continuous almost everywhere.
3. $\int f d\mu_n$ approaches $\int f d\mu$ for every bounded, continuous function f .
4. If A is any Borel set, ∂A denotes the topological boundary of A , and $\mu(\partial A) = 0$, then $\mu_n(A)$ approaches $\mu(A)$.

The theorem is interesting in that it combines measure-theoretic notions (measures and the integral) with topological notions (continuity and topological boundaries). The proof from 1 to 2 uses Skorohod's theorem. This states that if (μ_n) is a sequence of real distributions that converges to a real distribution μ , there is a sequence (Y_n) of random variables and another random variable Y , all defined on a common probability space, such that each Y_n has distribution μ_n , Y has distribution μ , and Y_n converges to Y pointwise. In other words, Skorohod's theorem tells us that (μ_n) and μ can be represented in a particularly nice way.

```

theorem Skorohod:
  fixes  $\mu$  :: "nat  $\Rightarrow$  real measure" and  $M$  :: "real measure"
  assumes " $\bigwedge n$ . real_distribution ( $\mu$  n)" "real_distribution  $M$ "
  and "weak_conv_m  $\mu$   $M$ "
  shows
    " $\exists$  ( $\Omega$  :: real measure) ( $Y$ _seq :: nat  $\Rightarrow$  real  $\Rightarrow$  real) ( $Y$  :: real  $\Rightarrow$  real).
     prob_space  $\Omega$   $\wedge$ 
     ( $\forall n$ .  $Y$ _seq n  $\in$   $\Omega$   $\rightarrow_M$  borel)  $\wedge$ 
     ( $\forall n$ . distr  $\Omega$  borel ( $Y$ _seq n) =  $\mu$  n)  $\wedge$ 
      $Y \in \Omega \rightarrow_M$  lborel  $\wedge$ 
     distr  $\Omega$  borel  $Y$  =  $M$   $\wedge$ 
     ( $\forall x \in$  space  $\Omega$ . ( $\lambda n$ .  $Y$ _seq n x)  $\longrightarrow$   $Y$  x)"

```

Proving Skorohod's theorem formally presented a number of technical challenges. One was that we needed to choose a continuity point of an arbitrary probability measure in an arbitrary open interval, that is, a real number x in an open interval I such that the measure of $\{x\}$ is zero. To that end, we showed that the number of atoms of a measure (that is, points x such that $\{x\}$ has strictly positive measure) is countable:

```

lemma countable_atoms:
  "finite_borel_measure  $M \implies$  countable  $\{x$ . measure  $M$   $\{x\} > 0\}$ "

```

The result then follows from the fact that any open interval in the reals is uncountable.

The implication from 2 to 3 is immediate. Notice that 1 is equivalent to saying that for every point x of continuity of the measure μ , $\int \chi_{(-\infty, x]} d\mu_n$ approaches $\int \chi_{(-\infty, x]} d\mu$, where $\chi_{(-\infty, x]}$ is the characteristic function of the interval $(-\infty, x]$. The implication from 3 to 1 is obtained by approximating this characteristic function by continuous step functions whenever x is a point of continuity. The implication from 4 to 1 is easy, noticing that $(-\infty, x]$ is a set of the specified type, whenever x is a point of continuity of μ . To complete the proof of the theorem, it is enough to prove that 2 implies 4. This implication is also not hard, once we show that the characteristic function χ_A is bounded and continuous at any point not on the boundary.

The results discussed in this section are found in the theory `Weak_Convergence`.

4.2 Characteristic functions

Recall that the *characteristic function* φ of a probability measure μ on the real line is defined by

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} \mu(dx).$$

If X is a random variable, the characteristic function of X is defined to be the characteristic function of its distribution. In our formalization, the characteristic function of a measure is defined as follows:

```

definition
  char :: "real measure  $\Rightarrow$  real  $\Rightarrow$  complex"
  where
    "char  $M$  t = LINT x|M. iexp (t * x)"

```

The characteristic function of a random variable x defined on a measure space M is then written $\text{char } (\text{distr } M \text{ borel } X)$, since $\text{distr } M \text{ borel } X$ denotes the distribution of x with respect to the usual Borel measure on the real numbers.

The characteristic function φ of a measure is continuous, and satisfies $\varphi(0) = 1$ and $|\varphi(t)| \leq 1$ for every t :

lemma (in `real_distribution`) `continuous_char`: "*continuous (at t) (char M)*"

lemma (in `real_distribution`) `char_zero`: "*char M 0 = 1*"

lemma (in `real_distribution`) `cmod_char_le_1`: "*norm (char M t) ≤ 1*"

As noted above, a key property of characteristic functions is this: if X_1 and X_2 are independent random variables, the characteristic function of $X_1 + X_2$ is the product of the individual characteristic functions. Because we used the Bochner integral, which allows us to integrate complex-valued functions directly, our final proof of this fact is even simpler than the one in Billingsley [3]. The calculation runs as follows:

$$\begin{aligned} \varphi_{X_1+X_2}(t) &= \int e^{it(X_1+X_2)} dM \\ &= \int e^{itX_1} e^{itX_2} dM \\ &= \left(\int e^{itX_1} dM \right) \left(\int e^{itX_2} dM \right) \\ &= \varphi_{X_1}(t) \varphi_{X_2}(t). \end{aligned}$$

Here, X_1 and X_2 are really functions over the underlying probability space M , and the third equation follows from the independence of X_1 and X_2 . We reproduce our formal proof in full:

```

lemma (in prob_space) char_distr_sum:
  assumes "indep_var borel X1 borel X2"
  shows "char (distr M borel (λω. X1 ω + X2 ω)) t =
    char (distr M borel X1) t * char (distr M borel X2) t"
proof -
  have [measurable]: "random_variable borel X1" "random_variable borel X2"
    using assms by (auto dest: indep_var_rv1 indep_var_rv2)

  have "char (distr M borel (λω. X1 ω + X2 ω)) t =
    (LINT x|M. iexp (t * (X1 x + X2 x)))"
    by (simp add: char_def integral_distr)
  also have "... = (LINT x|M. iexp (t * (X1 x)) * iexp (t * (X2 x)))"
    by (simp add: field_simps exp_add)
  also have "... =
    (LINT x|M. iexp (t * (X1 x))) * (LINT x|M. iexp (t * (X2 x)))"
    by (auto intro!: indep_var_compose[unfolded comp_def, OF assms]
      integrable_iexp indep_var_lebesgue_integral)
  also have "... = char (distr M borel X1) t * char (distr M borel X2) t"
    by (simp add: char_def integral_distr)
  finally show ?thesis .
qed

```

By induction, we have that for any finite set A and any sequence $(X_i)_{i \in A}$ of mutually independent random variables,

$$\varphi_{\sum_{i \in A} X_i}(t) = \prod_{i \in A} \varphi_{X_i}(t).$$

The formal proof is as follows:

```

lemma (in prob_space) char_distr_setsum:
  "indep_vars ( $\lambda i$ . borel) X A  $\implies$ 
  char (distr M borel ( $\lambda \omega$ .  $\sum_{i \in A} X i \omega$ )) t =
  ( $\prod_{i \in A}$ . char (distr M borel (X i)) t)"
proof (induct A rule: infinite_finite_induct)
  case (insert x F) then show ?case
  using indep_vars_subset[of " $\lambda$ _. borel" X "insert x F" F]
  by (auto simp add: char_distr_sum indep_vars_sum)
qed (simp_all add: char_def integral_distr prob_space del: distr_const)

```

(We do not require finiteness of A : by definition, if A is infinite, the sum over A is 0 and the product over A is 1, and the equation still holds.)

We also needed explicit approximations to the characteristic functions of a random variable, obtained using the calculation described at the beginning of Section 3.7. One of the results we used is as follows:

```

lemma (in prob_space) char_approx3':
  fixes  $\mu$  :: "real measure" and X
  assumes "random_variable borel X"
  and "integrable M X" "integrable M ( $\lambda x$ . (X x)^2)" "expectation X = 0"
  and "variance X =  $\sigma^2$ "
  and " $\mu = \text{distr M borel X}$ "
  shows "cmod (char  $\mu$  t - (1 - t^2 *  $\sigma^2$  / 2))  $\leq$ 
  (t^2 / 6) * expectation ( $\lambda x$ . min (6 * (X x)^2) (|t| * |X x|^3))"

```

Finally, we needed to compute the characteristic function φ of the standard normal distribution, and show $\varphi(t) = e^{-t^2/2}$. Establishing this fact took more work than we thought it would. Many textbook proofs of this invoke facts from complex analysis that were unavailable to us. Billingsley [3, page 344] sketches an elementary proof, which required calculating the moments and absolute moments of the standard normal distribution. This is where the calculations of $\int_{-\infty}^{\infty} x^k e^{-x^2/t} dx$, mentioned in Section 3.7, were needed. Specifically, we have for even k ,

```

lemma std_normal_moment_even:
  "has_bochner_integral lborel ( $\lambda x$ . std_normal_density x * x ^ (2 * k))
  (fact (2 * k) / (2^k * fact k))"

```

and for odd k ,

```

lemma std_normal_moment_odd:
  "has_bochner_integral lborel ( $\lambda x$ . std_normal_density x * x^(2 * k + 1)) 0"

```

A prior calculation by Sudeep Kanav covered the cases $k = 0, 1$, which provide the base cases for an inductive proof. Filling in the details involved carrying out careful computations with integrals and power series approximations to e^x .

In the Isabelle 2016 distribution, characteristic functions are defined in the theory `Characteristic_Functions`, and the properties cited above are proved there. The calculation of the moments of the normal distribution is found in the theory `Distributions`.

4.3 Lévy Inversion and Uniqueness

In Fourier analysis, an “inversion theorem” says that a function can be recovered from its Fourier transform, under suitable hypotheses and in a suitable sense. Along

those lines, the Lévy Inversion and Uniqueness Theorems say that a measure can be recovered from its characteristic function.

More precisely, the Lévy Inversion Theorem states the following:

Theorem 3 *Let μ be a probability measure, and φ be the characteristic function of μ . If a and b are continuity points of μ and $a < b$, then*

$$\mu(a, b] = \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt.$$

The proof is a long and subtle calculation. Let $I(T)$ denote the expression after the limit. Expanding the definition of $\varphi(t)$ and appealing to Fubini's theorem to switch the order of the two integrals, we obtain

$$I(T) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-T}^T \frac{e^{it(x-a)} - e^{it(x-b)}}{it} dt \mu(dx).$$

The idea is that as T approaches ∞ , the inner integral approaches a step function which jumps from 0 to 1 at a and then back down to 0 at b . This is shown by expanding the complex exponential in terms of sin and cos, using properties of the sine integral, and manipulating integrals and limits.

It is not hard to show that fixing the values of a measure on intervals $(a, b]$ as above is enough to determine the measure on all Borel sets. Thus the Inversion Theorem has the following result, known as the Uniqueness Theorem, as an important corollary:

Theorem 4 *If μ_1 and μ_2 are probability measures and $\varphi_{\mu_1} = \varphi_{\mu_2}$, then $\mu_1 = \mu_2$.*

In our formalization, this is expressed simply as follows:

```

theorem Levy_uniqueness:
  fixes M1 M2 :: "real measure"
  assumes "real_distribution M1" "real_distribution M2"
  and "char M1 = char M2"
  shows "M1 = M2"

```

4.4 The Lévy Continuity Theorem

Let (μ_n) be a sequence of distributions, where each μ_n has characteristic function φ_n , and let μ be a distribution with characteristic function φ . The *Lévy Continuity Theorem* states that μ_n converges to μ weakly if and only if $\varphi_n(t)$ converges to $\varphi(t)$ for every t . In our formalization, it is expressed as follows:

```

theorem Levy_continuity:
  fixes M :: "nat  $\Rightarrow$  real measure" and M' :: "real measure"
  assumes " $\bigwedge n$ . real_distribution (M n)"
  and "real_distribution M'"
  and " $\bigwedge t$ . ( $\lambda n$ . char (M n) t)  $\longrightarrow$  char M' t"
  shows "weak_conv_m M M'"

```

Proving the “only if” direction is easy, using the Portmanteau Theorem of Section 4.1, since e^{itx} is bounded and continuous. In fact, in our formalization, it has a one-line proof:

```

theorem Levy_continuity1:
  "( $\bigwedge n$ . real_distribution (M n))  $\implies$  real_distribution M'  $\implies$ 
   weak_conv_m M M'  $\implies$ 
   ( $\lambda n$ . char (M n) t)  $\longrightarrow$  char M' t"
unfolding char_def by (rule weak_conv_imp_integral_bdd_continuous_conv) auto

```

The other direction is a lot harder. Here is an outline of the proof:

1. Use a compactness argument to show that every subsequence (μ_{n_k}) of (μ_n) has a weakly convergent subsequence.
2. Suppose, for the sake of contradiction, (μ_n) does not converge weakly to μ . Then there is a subsequence (μ_{n_k}) such that no subsequence of *that* can converge weakly to μ .
3. By 1, this particular sequence (μ_{n_k}) converges weakly to some measure, ν .
4. By the “only if” direction, already proved, φ_{n_k} converges pointwise to the characteristic function of ν .
5. Since, by hypothesis, $\varphi_n(t)$ converges to $\varphi(t)$ for every t , the characteristic function of ν must be φ .
6. By the Uniqueness Theorem, this implies that $\nu = \mu$, contrary to the choice of (μ_{n_k}) in 2.

The necessary compactness principle is a consequence of the *Helly Selection Theorem*, which we now describe. Although the proof of this theorem and its consequence take up only a page-and-a-half in Billingsley’s textbook, these were among the most subtle components of our formalization. The theorem states the following:

Theorem 5 *Let $(f_n)_{n \in \mathbb{N}}$ be a uniformly bounded sequence of nondecreasing, right-continuous functions. Then there are a subsequence $(f_{n_k})_{k \in \mathbb{N}}$ and a nondecreasing, right-continuous function F such that $\lim_k f_{n_k}(x) = F(x)$ at continuity points of F .*

In our formalization, this is expressed as follows:

```

theorem Helly_selection:
  fixes f :: "nat  $\Rightarrow$  real  $\Rightarrow$  real"
  assumes " $\bigwedge n$  x. continuous (at_right x) (f n)"
  and " $\bigwedge n$ . mono (f n)"
  and " $\bigwedge n$  x. |f n x|  $\leq$  M"
  shows " $\exists s$ . subseq s  $\wedge$ 
   ( $\exists F$ . ( $\forall x$ . continuous (at_right x) F)  $\wedge$  mono F  $\wedge$  ( $\forall x$ . |F x|  $\leq$  M)  $\wedge$ 
   ( $\forall x$ . isCont F x  $\longrightarrow$  ( $\lambda n$ . f (s n) x)  $\longrightarrow$  F x))"
```

Saying that (f_{n_k}) is a subsequence of (f_n) really means that the map $k \mapsto n_k$ is strictly increasing. The statement is represented formally by explicitly asserting the existence of the strictly increasing function $s : \mathbb{N} \rightarrow \mathbb{N}$ which returns, for each k , the value n_k . The proof involves a diagonalization argument: for each rational r , we thin the sequence to guarantee convergence at r , and then take a “diagonal limit” to construct the required subsequence and limit. To that end, we used a general framework for such diagonalization arguments, provided by Fabian Immler.

To describe the relevant corollary, we need to introduce a definition. A sequence (μ_n) of real measures is said to be *tight* if, for every $\varepsilon > 0$, there is a finite interval $(a, b]$ such that $\mu_n(a, b] > 1 - \varepsilon$ for all n . Roughly, a sequence of probability measures is tight if no mass “escapes to infinity”; the sequence (μ_n) , where μ_n is a unit mass at n , is an example of a sequence that is *not* tight. Helly’s theorem can be used to show that if (μ_n) is a tight sequence of measures, then for every subsequence (μ_{n_k}) there is a further subsequence $(\mu_{n_{k(j)}})$ and a probability measure μ such that $(\mu_{n_{k(j)}})$ converges weakly to μ as j approaches infinity.

theorem `tight_imp_convergent_subsubsequence`:
assumes μ : "tight μ " "subseq s "
shows " $\exists r M$. subseq $r \wedge$ real_distribution $M \wedge$ weak_conv_m $(\mu \circ s \circ r) M$ "

In the Isabelle 2016 distribution, the Helly Selection Theorem and its corollary are proved in `Helly_Selection`, and the Lévy Continuity Theorem is proved in `Levy`. Immler's general framework for diagonal arguments can be found in the theory `Diagonal_Sequence` in `HOL-Library`.

4.5 The Central Limit Theorem

Proving the Central Limit Theorem is now just a matter of putting the pieces together. Let (X_n) be a sequence of random variables, all of which have the same distribution μ and finite variance $\sigma^2 > 0$. Without loss of generality (subtracting a common offset) we can assume that each X_n has mean 0. Let

$$S'_n = \sum_{i < n} X_i / \sqrt{n\sigma^2}$$

be the normalized sums. Our goal is to show that the distributions of S'_n converge weakly to the standard normal distribution.

For each n , let φ_n be the characteristic function of S'_n . By the Lévy continuity theorem, it suffices to show that φ_n approaches the characteristic function of the standard normal distribution pointwise. In other words, we need to show that for every t , $\varphi_n(t)$ approaches $e^{-t^2/2}$.

Since each X_i has the same distribution, all the X_i 's have the same characteristic function; call it ψ . By the key property of characteristic functions, the characteristic function of the sum S'_n is the product of the characteristic functions of the components, so

$$\begin{aligned} \varphi_n(t) &= \prod_{j=1}^n \int e^{itX_j / \sqrt{n\sigma^2}} d\mu \\ &= \prod_{j=1}^n \psi(t / \sqrt{n\sigma^2}) \\ &= (\psi(t / \sqrt{n\sigma^2}))^n. \end{aligned}$$

Now some of the explicit calculations described in Section 3.7 can be used to show that $\psi(t)$ is well approximated by

$$1 + it \int X d\mu + \frac{t^2}{2} \int X^2 d\mu,$$

which is equal to $1 - t^2\sigma^2/2$, since we are assuming X has mean 0 and variance σ^2 . Plugging $t/\sqrt{n\sigma^2}$ in for t , we obtain an approximation to $\psi(t/\sqrt{n\sigma^2})$, and substituting that in the expression for $\varphi_n(t)$, we see that $\varphi_n(t)$ is approximated by $(1 - \frac{t^2}{2n})^n$. This last expression approaches $e^{-t^2/2}$ as t approaches infinity, as required.

The Central Limit Theorem is found in the file `Central_Limit_Theorem`. The formal version of the proof we have just sketched is given in its entirety in the appendix. (As above, we derive the mean zero case first, and then derive Theorem 1 as a corollary.)

5 Reflections

We are by no means the first to formalize substantial portions of analysis in an interactive theorem prover. Both HOL4 and HOL Light have extensive theories of multivariate real analysis, and HOL Light has a substantial theory of complex analysis as well [11]. The real analysis library in HOL Light played an important part in the Flyspeck project formalization of Thomas Hales’ proof of the Kepler conjecture [10], and Isabelle’s real analysis library has also been used to formalize properties of dynamical systems [14]. Substantial portions of measure theory and measure-theoretic probability have been formalized in HOL4 [16, 19]. The C-CoRN and Coquelicot projects [15, 4] provide a libraries for real analysis based on dependent type theory for the Coq proof assistant. It would take us too far afield to discuss all this work and compare all the other approaches to ours, so, instead, we will focus on our own formalization efforts and try to convey some of the lessons we learned. We also refer the reader to an article by Boldo, Lelay, and Melquiond [5], which provides a thoughtful and thorough survey of approaches to formalizing real analysis.

5.1 Dealing with partial functions

In the logical framework of type theory, where every function is assumed to be total, one often has to deal with partial functions, such as limits, derivatives, integrals, and so on. One common way of proceeding is to represent partial functions as relations. For example, in Isabelle one can write $f \text{ sums } l$ to indicate that the finite partial sums $\sum_{i < n} f(i)$ converge to l as n approaches infinity. Another option is to make the function in question total by assigning an arbitrary value at inputs where it would otherwise be undefined, and use a predicate to pick out the “real” values. For example, $\text{summable } f$ is defined to mean that there exists an l such that $f \text{ sums } l$, and $\text{suminf } f$ specifies that value of l , if it exists, and 0 otherwise. It is not hard to see that the expression $f \text{ sums } l$ is then equivalent to the conjunction $\text{summable } f \wedge \text{suminf } f = l$. The expression $\text{suminf } f$ allows us to refer to the value of the infinite sum, as we do when we write $\sum_{i=0}^{\infty} f(i)$ in ordinary mathematics. In particular, this expression can occur in a more complicated expression; for example, we can write $\text{suminf } f + \text{suminf } g$. But one typically also wants to know that $\text{summable } f$ holds, since otherwise the value of $\text{suminf } f$ may be meaningless.

In Isabelle, many partial functions related to analysis, such as limits and derivatives, are represented in this way, with a relation, a predicate asserting the existence of a value, and a function that returns an arbitrary value when the predicate fails. An exception is the notion of a measure on a measure space: while the expression $s \in \text{sets } M$ expresses that s is a measurable set for the measure M , in which case, $\text{measure } M s$ is the measure of s , there is no relation between a set and its measure.

In the measure theory library, however, integration is handled in the usual way: we have $\text{has_bochner_integral } M f x$ to express that f has Bochner integral x with respect to the measure space M , $\text{integrable } M f$ to say that f is integrable, and notation $\text{integral}^L M f$ for the value of the integral, when it exists (the superscripted L is a holdover from Lebesgue measure). The library tends to favor the latter rep-

representations, however, with one theorem asserting the value of an integral, and another theorem asserting integrability. For example, we have:

lemma *integral_add*:
 $"integrable\ M\ f \implies integrable\ M\ g \implies$
 $(LINT\ x|M.\ f\ x + g\ x) = (LINT\ x|M.\ f\ x) + (LINT\ x|M.\ g\ x)"$

lemma *integrable_add*:
 $"integrable\ M\ f \implies integrable\ M\ g \implies integrable\ M\ (\lambda x.\ f\ x + g\ x)"$

This sometimes got us into trouble. A couple of times, we used theorems in the library, only to realize that the accompanying integrability assertions were missing; we then had to revise the library to provide these additional assertions. We often made the same mistake in our own developments, and in proofs we often found that we had to carry out parallel calculations: after calculating an integral, we had to go back and prove that the expression we began with was in fact integrable. In some cases, Isabelle’s automation could dispel integrability claims for us, but typically in those cases the calculations could also be carried out automatically. On the other hand, writing formulas with integrals rather than the *has_integral* predicate makes them look much more like the formulas one finds in an ordinary mathematical textbook.

We do not know the ideal solution to the problem. In a dependent type theory with propositions as types, one could require *integrable* *f* as a “precondition” — a hidden argument — to an expression *integral* *f*. Coq’s constructive C-CoRN library uses such an approach [15], but the Coquelicot project uses a classical axiomatization of the real numbers to totalize limits [4]. In simple type theory, it seems that one has to choose between using a relational version or using a function together with a definedness predicate. In the latter case, one has to take care to keep the two pieces of information close together.

5.2 Strategies for limit proofs

It is not always obvious how to carry out limit proofs at the right level of formal abstraction. With measure theory, it is often advantageous to adopt an order-theoretic point of view: instead of proving that a function approaches a certain point in a ε -environment, it is sometimes preferable to do this separately for an upper and lower bound. This is contrary to what is done in Billingsley [3], where many proofs, like that of the Helly section theorem, are performed by choosing ε ’s. That approach works when the domain is the set of real numbers or at least a metric space. But we often needed to use the extended real number structures $\overline{\mathbb{R}}$ or $\overline{\mathbb{R}}_{\geq 0}$, for two reasons: (1) we reason about measures which are not necessarily finite, and (2) we reason about \liminf ’s and \limsup ’s, which are defined on complete lattices but not on the real numbers. Working with ε ’s typically requires us to compute differences, which is difficult on $\overline{\mathbb{R}}$ or $\overline{\mathbb{R}}_{\geq 0}$, where subtraction and addition are not as well behaved as they are on \mathbb{R} . For example, if we use metric limits to prove a property on a neighborhood of $f\ x$, we may obtain a $\varepsilon > 0$, but we do not immediately have the property $f\ x < f\ x + \varepsilon$, since this fails for $f\ x = \infty$. Using the order topology on extended real number structures, however, we still have limits and filters, but instead of obtaining a ε -neighborhood of a specific point, we get upper and lower bounds on values for which the property still holds. In the

previous example, using an order-theoretic limit we instead get an upper bound y with $f x < y$, and this moreover implies $f x < \infty$.

Even using ε -proofs in a metric-space setting can be formally inconvenient. Many textbook proofs adopt a style whereby various ε 's are obtained in the course of the argument, for example, as diameters of neighborhoods in open sets or neighborhoods in the range of continuous functions. The properties and all these ε 's are combined by computing a minimal ε value and proving it correct. Textbooks often elide such details by taking an ε that is “sufficiently small.”

Formally, it is often better to avoid the uses of ε values entirely, showing instead that the required properties give rise to a set in the relevant filter. For example, suppose that the functions f and g_i , for $i < n$, are continuous at x , and that $1 < f x$ and $f y \leq g_i y$ around x . Suppose further that we want to obtain a neighborhood of x where f is above 1 and below all the g_i s. One approach is to obtain ε and ε_i 's where $1 < f y$ where $f y \leq g_i y$ in the corresponding balls around x , and then to compute the minimum of these values, with the special case where $n = 0$. But working with these ε 's is completely auxiliary to our original goal. Instead, we can easily show that $\{y \mid 1 < f y\} \cap \bigcap_{i < n} \{y \mid f y \leq g_i y\}$ is in the neighborhood filter at x , simply using the fact that a filter is closed under finite intersections.

5.3 Strategies for integrals

Measure theory gives us two different integrals on measure spaces, the nonnegative Lebesgue integral and the Bochner integral, as described in Section 3.4. The distinction is clear: the Lebesgue integral only requires a measurable function, and handles functions into the nonnegative extended reals $\overline{\mathbb{R}}_{\geq 0}$, while the Bochner integral requires an integrable function, but handles functions into arbitrary second-countable Banach spaces. There are two important advantages to using the Lebesgue integral: (1) measurability is compositional, supporting different measurable spaces, while this is not the case for integrability, and (2) the extended nonnegative reals include ∞ , and so no integrability condition is needed for the integral to be closed under addition and constant multiplication. The property *integrable.iff.bounded*, which states that a function is integrable if and only if it is bounded, provides a key way to prove integrability. Similarly, a function is integrable if it is measurable and has integrable upper and lower bounds.

A small trick that results in more convenient proof rules for the Bochner integral is to fix the value of *integral* μ f to 0 for a non-integrable function f . In interactive theorem proving, this is a common trick to totalize a function taking values in a numeric domain. Exploiting this fact avoids some auxiliary integrability rules. For the constant multiplication rules (i.e. multiplication with a real or complex value, scalar multiplication and the inner vector product) we get a rule of the form

$$\text{integral } \mu \ (\lambda x. c * f x) = c * \text{integral } \mu f$$

without any assumptions. Also, integrability is invariant under the transformations *distr* and *density*:

$$\text{integral } (\text{distr } \mu \nu f) g = \text{integral } \mu (g \circ f)$$

with the assumption that f and g are measurable. Similarly, we allow the affine transformation under the integral:

```
integral lborel f = |c| * integral lborel ( $\lambda x. f (t + c * x)$ )
```

with the only assumption $c \neq 0$.

Of course, we still have the problems mentioned in Subsection 5.1, i.e. we need to prove integrability separately. (Without a proof of integrability, the results may not mean what we think they mean, and they are generally unusable.) But in many analytical proofs, integrability is often proved separately anyhow, with a proof that may have little to do with the calculation of the integral value.

5.4 Cleanup and length

It is impossible to give a meaningful estimate of the time involved in the formalization, as the work was carried out intermittently over a long period of time, and includes time spent learning to use Isabelle by the third author, Serafin, who was an undergraduate student at Carnegie Mellon at the time. When we first obtained a proof of the CLT, we reported that our repository contained about 13,000 lines [1]. This included all the general infrastructure and additions to Isabelle’s libraries as well as the core parts of the proof. Since then, we have cleaned up and refactored most of our proof scripts, and many of them have been shortened considerably. Subsequently, Hölzl implemented the Bochner integral, eliminating the need for a separate notion of integration for functions from the reals to the complex numbers. In addition, many of the supporting theorems and facts have been moved to other parts of the Isabelle libraries.

The proof of the CLT is now part of the Isabelle distribution. One interesting observation is that, in terms of the number of lines, the majority of effort went into developing the background and general infrastructure. For example, some of our longest files involve general facts about integration:

```
Bochner_Integral:  3,066 lines
Set_Integral:      602 lines
Interval_Integral: 1,123 lines
```

In addition, the construction of Lebesgue-Stieltjes measure, described in Section 3.8, is found in `Lebesgue_Measure` and requires about 270 lines. Some of the key background for the formalization is contained in the following files:

```
Distribution_Functions: 259 lines
Weak_Convergence:      422 lines
Sinc_Integral:          403 lines
```

Also, general facts about the standard normal distribution take about about 380 lines in `Distributions`. The core development of characteristic functions and their properties, and the proof of the CLT, is found in the following files:

```
Characteristic_Functions: 554 lines
Helly_Selection:          298 lines
Levy:                     542 lines
Central_Limit_Theorem:    144 lines
```

Once all the background information was in place, many of our proofs followed those in Billingsley quite closely. This allows for some direct comparisons. The

increases in length are most dramatic in technical proofs where there are one-step arguments in the text that are indeed straightforward to verify, but nonetheless require long and tedious arguments. This includes verifying straightforward continuity claims, filling in implicit limit arguments, finding explicit choices of ε sufficiently small to make a proof go through, and so on. Thus the Helly selection theorem, only 10 lines in Billingsley's text, is 127 lines in our formalization. Billingsley derives two corollaries from that, each with a proof of 9 lines; our formal versions are 102 and 18 lines, respectively. The Lévy inversion theorem runs only 15 lines in Billingsley, and about 170 lines in our formalization. Billingsley observes that the uniqueness theorem follows from the inversion theorem with four lines of proof, which translates to 74 lines in our formal version.

5.5 Future directions

The version of the Central Limit Theorem we proved is not the most general version that is presented in Billingsley's book. With some more calculational effort one could formalize the Lindeberg central limit theorem, which relaxes the requirement that the random variables that are summed be identically distributed; we only need to assume that they do not deviate too much in distribution, as made precise by the *Lindeberg condition* [3, p. 359]. Even the condition that the variables being summed are independent can be weakened to a condition of weak dependence, as outlined in [3, p. 363]. Other generalizations include the CLT for random vectors [3, p. 385], and various versions of the CLT for martingales [3, pp. 475–478]. There are many additional refinements and generalizations of the Central Limit Theorem in the mathematical literature.

Supporting automation can always be improved, and it was at times frustrating that automated tools would get stuck on seemingly trivial matters like determining whether an instance of zero should be interpreted as a real or a nonnegative extended real. As we remarked in Section 3.7, carrying out ordinary calculations with integrals was often the most painful part of the formalization. It would be especially useful to have better automated support for such calculations, that either implement features of computer algebra systems in a proof-producing framework, or reconstruct formal proofs of such results from suitable certificates.

References

1. Jeremy Avigad, Johannes Hölzl, and Luke Serafin. A formally verified proof of the central limit theorem (preliminary announcement). 2014.
2. Clemens Ballarin. Interpretation of Locales in Isabelle: Theories and Proof Contexts. In J.M. Borwein and W.M. Farmer, editors, *Mathematical Knowledge Management 2006*, Lecture Notes in Artificial Intelligence, pages 31–43. Springer, 2006.
3. Patrick Billingsley. *Probability and measure*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons Inc., New York, third edition, 1995. A Wiley-Interscience Publication.
4. Sylvie Boldo, Catherine Lelay, and Guillaume Melquiond. Improving real analysis in coq: A user-friendly approach to integrals and derivatives. In Chris Hawblitzel and Dale Miller, editors, *Certified Programs and Proofs - Second International Conference, CPP 2012, Kyoto, Japan, December 13-15, 2012. Proceedings*, volume 7679 of *Lecture Notes in Computer Science*, pages 289–304. Springer, 2012.

5. Sylvie Boldo, Catherine Lelay, and Guillaume Melquiond. Formalization of real analysis: a survey of proof assistants and libraries. *Mathematical Structures in Computer Science*, 26(7):1196–1233, 2016.
6. Alonzo Church. A formulation of the simple theory of types. *J. Symbolic Logic*, 5:56–68, 1940.
7. Hans Fischer. *A History of the Central Limit Theorem: From Classical to Modern Probability Theory*. Springer, New York, 2011.
8. Francis Galton. *Natural Inheritance*. Macmillan, London, 1889.
9. Hanne Gottlieb. Transcendental functions and continuity checking in PVS. In *Theorem Proving in Higher-Order Logics (TPHOLs) 2000*, pages 197–214. Springer, 2000.
10. Thomas Hales, Mark Adams, Gertrud Bauer, Dat Tat Dang, John Harrison, Truong Le Hoang, Cezary Kaliszyk, Victor Magron, Sean McLaughlin, Thang Tat Nguyen, Truong Quang Nguyen, Tobias Nipkow, Steven Obua, Joseph Pleso, Jason Rute, Alexey Solovyev, An Hoai Thi Ta, Trung Nam Tran, Diep Thi Trieu, Josef Urban, Ky Khac Vu, and Roland Zumkeller. A formal proof of the Kepler conjecture. <http://arxiv.org/abs/1501.02155>.
11. John Harrison. Formalizing basic complex analysis. In R. Matuszewski and A. Zalewska, editors, *From Insight to Proof: Festschrift in Honour of Andrzej Trybulec*, volume 10(23) of *Studies in Logic, Grammar and Rhetoric*, pages 151–165. University of Białystok, 2007.
12. Johannes Hölzl and Armin Heller. Three Chapters of Measure Theory in Isabelle/HOL. In Marko C. J. D. van Eekelen, Herman Geuvers, Julien Schmaltz, and Freek Wiedijk, editors, *Interactive Theorem Proving (ITP) 2011*, volume 6898 of *Lecture Notes in Computer Science*, pages 135–151. Springer, 2011.
13. Johannes Hölzl, Fabian Immler, and Brian Huffman. Type classes and filters for mathematical analysis in Isabelle/HOL. In Sandrine Blazy, Christine Paulin-Mohring, and David Pichardie, editors, *Interactive Theorem Proving*, volume 7998 of *Lecture Notes in Computer Science*, pages 279–294. Springer Berlin Heidelberg, 2013.
14. Fabian Immler and Christoph Traut. The flow of odes. In Jasmin Christian Blanchette and Stephan Merz, editors, *Interactive Theorem Proving - 7th International Conference, ITP 2016, Nancy, France, August 22-25, 2016, Proceedings*, volume 9807 of *Lecture Notes in Computer Science*, pages 184–199. Springer, 2016.
15. Robbert Krebbers and Bas Spitters. Type classes for efficient exact real arithmetic in coq. *Logical Methods in Computer Science*, 9(1), 2011.
16. Tarek Mhamdi, Osman Hasan, and Sofïène Tahar. Formalization of entropy measures in HOL. In Marko C. J. D. van Eekelen, Herman Geuvers, Julien Schmaltz, and Freek Wiedijk, editors, *Interactive Theorem Proving - Second International Conference, ITP 2011, Berg en Dal, The Netherlands, August 22-25, 2011. Proceedings*, volume 6898 of *Lecture Notes in Computer Science*, pages 233–248. Springer, 2011.
17. Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL. A proof assistant for higher-order logic*, volume 2283 of *Lecture Notes in Computer Science*. Springer Verlag, Berlin, 2002.
18. Lawrence C. Paulson. Three years of experience with sledgehammer, a practical link between automatic and interactive theorem provers. In Renate A. Schmidt, Stephan Schulz, and Boris Konev, editors, *Proceedings of the 2nd Workshop on Practical Aspects of Automated Reasoning, PAAR-2010, Edinburgh, Scotland, UK, July 14, 2010*, volume 9 of *EPiC Series*, pages 1–10. EasyChair, 2010.
19. Muhammad Qasim, Osman Hasan, Maïssa Elleuch, and Sofïène Tahar. Formalization of normal random variables in HOL. In Michael Kohlhase, Moa Johansson, Bruce R. Miller, Leonardo de Moura, and Frank Wm. Tompa, editors, *Intelligent Computer Mathematics - 9th International Conference, CICM 2016, Białystok, Poland, July 25-29, 2016, Proceedings*, volume 9791 of *Lecture Notes in Computer Science*, pages 44–59. Springer, 2016.
20. Luke Serafin. A formally verified proof of the Central Limit Theorem. Master’s thesis, Carnegie Mellon University, 2015.
21. Markus Wenzel. Type Classes and Overloading in Higher-Order Logic. In E. Gunter and A. Felty, editors, *Proceedings of the 10th International Conference on Theorem Proving in Higher Order Logics (TPHOLs’97)*, pages 307–322, Murray Hill, New Jersey, 1997.
22. Markus Wenzel. *Isabelle/Isar—a versatile environment for human-readable formal proof documents*. PhD thesis, Institut für Informatik, Technische Universität München, 2002.

Appendix

```

theorem (in prob_space) central_limit_theorem_zero_mean:
  fixes X :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
    and  $\mu$  :: "real measure"
    and  $\sigma$  :: real
    and S :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
  assumes X_indep: "indep_vars ( $\lambda$ i. borel) X UNIV"
    and X_integrable: " $\bigwedge$ n. integrable M (X n)"
    and X_mean_0: " $\bigwedge$ n. expectation (X n) = 0"
    and  $\sigma$ _pos: " $\sigma > 0$ "
    and X_square_integrable: " $\bigwedge$ n. integrable M ( $\lambda$ x. (X n x)2)"
    and X_variance: " $\bigwedge$ n. variance (X n) =  $\sigma^2$ "
    and X_distrib: " $\bigwedge$ n. distr M borel (X n) =  $\mu$ "
  defines "S n  $\equiv$   $\lambda$ x.  $\sum$  i<n. X i x"
  shows "weak_conv_m ( $\lambda$ n. distr M borel ( $\lambda$ x. S n x / sqrt (n *  $\sigma^2$ )))
    std_normal_distribution"
proof -
  let ?S' = " $\lambda$ n x. S n x / sqrt (real n *  $\sigma^2$ )" and ?m = " $\lambda$ x. min (6 * x2)"
  define  $\varphi$  where " $\varphi$  n = char (distr M borel (?S' n))" for n
  define  $\psi$  where " $\psi$  n t = char  $\mu$  (t / sqrt ( $\sigma^2 * n$ ))" for n t

  have X_rv [simp, measurable]: " $\bigwedge$ n. random_variable borel (X n)"
  using X_indep unfolding indep_vars_def2 by simp
  interpret  $\mu$ : real_distribution  $\mu$ 
  by (subst X_distrib [symmetric, of 0], rule real_distribution_distr, simp)

  have  $\mu$ _integrable [simp]: "integrable  $\mu$  ( $\lambda$ x. x)"
  and  $\mu$ _mean_integrable [simp]: " $\mu$ .expectation ( $\lambda$ x. x) = 0"
  and  $\mu$ _square_integrable [simp]: "integrable  $\mu$  ( $\lambda$ x. x2)"
  and  $\mu$ _variance [simp]: " $\mu$ .expectation ( $\lambda$ x. x2) =  $\sigma^2$ "
  using assms by (simp_all add: X_distrib [symmetric, of 0]
    integrable_distr_eq integral_distr)

  let ?I = " $\lambda$ n t. LINT x |  $\mu$ . ?m x (|t / sqrt ( $\sigma^2 * n$ )| * |x| ^ 3)"
  have main: " $\forall_F$  n in sequentially.
    cmod ( $\varphi$  n t - (1 + (-t2) / 2) / nn)  $\leq$  t2 / (6 *  $\sigma^2$ ) * ?I n t"
    for t
  proof (rule eventually_sequentiallyI)
    fix n :: nat
    assume "n  $\geq$  nat (ceiling (t2 / 4))"
    hence n: "n  $\geq$  t2 / 4" by (subst nat_ceiling_le_eq [symmetric])
    let ?t = "t / sqrt ( $\sigma^2 * n$ )"

    define  $\psi'$  where " $\psi'$  n i = char (distr M borel
      ( $\lambda$ x. X i x / sqrt ( $\sigma^2 * n$ )))" for n i
    have *: " $\bigwedge$ n i t.  $\psi'$  n i t =  $\psi$  n t"
      unfolding  $\psi$ _def  $\psi'$ _def char_def
      by (subst X_distrib [symmetric]) (auto simp: integral_distr)

    have " $\varphi$  n t = char (distr M borel
      ( $\lambda$ x.  $\sum$  i<n. X i x / sqrt ( $\sigma^2 * n$ ))) t"
      by (auto simp:  $\varphi$ _def S_def sum_divide_distrib ac_simps)
    also have "... = ( $\prod$  i < n.  $\psi'$  n i t)"
      unfolding  $\psi'$ _def
      apply (rule char_distr_sum)
      apply (rule indep_vars_compose2[where X=X])
      apply (rule indep_vars_subset)
      apply (rule X_indep)
      apply auto
  end

```

```

done
also have "... = ( $\psi$  n t)n"
  by (auto simp add: * prod_constant)
finally have  $\varphi_{eq}$ : " $\varphi$  n t = ( $\psi$  n t)n" .

have "norm ( $\psi$  n t - (1 -  $t^2 * \sigma^2 / 2$ ))  $\leq t^2 / 6 * ?I$  n t"
  unfolding  $\psi_{def}$  by (rule  $\mu.char\_approx3$ , auto)
also have " $t^2 * \sigma^2 = t^2 / n$ "
  using  $\sigma_{pos}$  by (simp add: power_divide)
also have " $t^2 / n / 2 = (t^2 / 2) / n$ "
  by simp
finally have **: "norm ( $\psi$  n t - (1 +  $(-t^2) / 2 / n$ ))  $\leq t^2 / 6 * ?I$  n t"
  by simp

have "norm ( $\varphi$  n t - (complex_of_real (1 +  $(-t^2) / 2 / n$ ))n)  $\leq n * norm (\psi$  n t - (complex_of_real (1 +  $(-t^2) / 2 / n$ )))"
  using n unfolding  $\varphi_{eq}$   $\psi_{def}$ 
  by (auto intro!: norm_power_diff  $\mu.cmod\_char\_le_1$  abs_leI
    simp del: of_real_diff simp: of_real_diff[symmetric] divide_le_eq)
also have "...  $\leq n * (t^2 / 6 * ?I$  n t)"
  by (rule mult_left_mono [OF **], simp)
also have "... =  $(t^2 / (6 * \sigma^2)) * ?I$  n t)"
  using  $\sigma_{pos}$  by (simp add: field_simps min_absorb2)
finally show "norm ( $\varphi$  n t - (1 +  $(-t^2) / 2 / n$ )n)  $\leq (t^2 / (6 * \sigma^2)) * ?I$  n t)"
  by simp
qed

show ?thesis
proof (rule levy_continuity)
  fix t
  have " $\bigwedge x. (\lambda n. ?m x (|t| * |x| ^ 3 / \sqrt{\sigma^2 * \text{real } n})) \longrightarrow 0$ "
    using  $\sigma_{pos}$ 
    by (auto simp: real_sqrt_mult min_absorb2
      intro!: tendsto_min[THEN tendsto_eq_rhs]
        sqrt_at_top[THEN filterlim_compose]
        filterlim_tendsto_pos_mult_at_top
        filterlim_at_top_imp_at_infinity
        tendsto_divide_0
        filterlim_real_sequentially)
  then have " $(\lambda n. ?I$  n t)  $\longrightarrow (LINT x | \mu. 0)$ "
    by (intro integral_dominated_convergence [where w = " $\lambda x. 6 * x^2$ "]) auto
  then have *: " $(\lambda n. t^2 / (6 * \sigma^2)) * ?I$  n t  $\longrightarrow 0$ "
    by (simp only: integral_zero tendsto_mult_right_zero)

  have " $(\lambda n. \text{complex\_of\_real } ((1 + (-t^2) / 2) / n)^n) \longrightarrow \text{complex\_of\_real } (\exp (-t^2 / 2))$ "
    by (rule isCont_tendsto_compose [OF _ tendsto_exp_limit_sequentially])
    auto
  then have " $(\lambda n. \varphi$  n t)  $\longrightarrow \text{complex\_of\_real } (\exp (-t^2 / 2))$ "
    by (rule Lim_transform) (rule Lim_null_comparison [OF main *])
  then show " $(\lambda n. \text{char } (\text{distr } M \text{ borel } (?S' n)) t) \longrightarrow \text{char\_std\_normal\_distribution } t$ "
    by (simp add:  $\varphi_{def}$  char_std_normal_distribution)
qed (auto intro!: real_dist_normal_dist simp: S_def)
qed

theorem (in prob_space) central_limit_theorem:
  fixes X :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"

```

```

    and  $\mu$  :: "real measure"
    and  $c$   $\sigma$  :: real
    and  $S$  :: "nat  $\Rightarrow$  'a  $\Rightarrow$  real"
  assumes  $X_{indep}$ : "indep_vars ( $\lambda i$ . borel)  $X$  UNIV"
    and  $X_{integrable}$ : " $\bigwedge n$ . integrable  $M$  ( $X$   $n$ )"
    and  $X_{mean}$ : " $\bigwedge n$ . expectation ( $X$   $n$ ) =  $c$ "
    and  $\sigma_{pos}$ : " $\sigma > 0$ "
    and  $X_{square\_integrable}$ : " $\bigwedge n$ . integrable  $M$  ( $\lambda x$ . ( $X$   $n$   $x$ )2)"
    and  $X_{variance}$ : " $\bigwedge n$ . variance ( $X$   $n$ ) =  $\sigma^2$ "
    and  $X_{distrib}$ : " $\bigwedge n$ . distr  $M$  borel ( $X$   $n$ ) =  $\mu$ "
  defines " $S$   $n$   $x \equiv \sum_{i < n} X$   $i$   $x$ "
  shows "weak_conv_m ( $\lambda n$ . distr  $M$  borel ( $\lambda x$ . ( $S$   $n$   $x$  -  $n$  *  $c$ ) / sqrt ( $n$ * $\sigma^2$ )))
    std_normal_distribution"
proof -
  have "weak_conv_m
    ( $\lambda n$ . distr  $M$  borel ( $\lambda x$ . ( $\sum_{i < n} X$   $i$   $x$  -  $c$ ) / sqrt ( $n$  *  $\sigma^2$ )))
    std_normal_distribution"
  proof (intro central_limit_theorem_zero_mean)
    show "indep_vars ( $\lambda i$ . borel) ( $\lambda i$   $x$ .  $X$   $i$   $x$  -  $c$ ) UNIV"
      using  $X_{indep}$  by (rule indep_vars_compose2) auto
    show "integrable  $M$  ( $\lambda x$ .  $X$   $n$   $x$  -  $c$ )"
      "expectation ( $\lambda x$ .  $X$   $n$   $x$  -  $c$ ) = 0" for  $n$ 
      using  $X_{integrable}$   $X_{mean}$  by (auto simp: prob_space)
    show " $\sigma > 0$ " "integrable  $M$  ( $\lambda x$ . ( $X$   $n$   $x$  -  $c$ )2)"
      "variance ( $\lambda x$ .  $X$   $n$   $x$  -  $c$ ) =  $\sigma^2$ " for  $n$ 
      using  $\langle 0 < \sigma \rangle$   $X_{integrable}$   $X_{mean}$   $X_{square\_integrable}$   $X_{variance}$ 
      by (auto simp: prob_space power2_diff)
    show "distr  $M$  borel ( $\lambda x$ .  $X$   $n$   $x$  -  $c$ ) = distr  $\mu$  borel ( $\lambda x$ .  $x$  -  $c$ )" for  $n$ 
      unfolding  $X_{distrib}$ [of  $n$ , symmetric] using  $X_{integrable}$ 
      by (subst distr_distr) (auto simp: comp_def)
  qed
  moreover have " $(\sum_{i < n} X$   $i$   $x$  -  $c$ ) =  $S$   $n$   $x$  -  $n$  *  $c$ " for  $n$   $x$ 
    by (simp add: sum_subtractf  $S_{def}$ )
  ultimately show ?thesis
    by simp
qed

```