

Avoiding common pitfalls and misconceptions in extractions of the proton radius

Jan C. Bernauer^{1,*} and Michael O. Distler^{2,†}

¹Laboratory for Nuclear Science, MIT, Cambridge, Massachusetts 02139, USA

²Institut für Kernphysik, Johannes-Gutenberg-Universität Mainz, D-55128 Mainz, Germany

In a series of recent publications, different authors produce a wide range of electron radii when reanalyzing electron proton scattering data. In the light of the proton radius puzzle, this is a most unfortunate situation. However, we find flaws in most analyses that result in radii around 0.84 fm. In this paper, we explain our reasoning and try to illustrate the most common pitfalls.

PACS numbers: 14.20.Dh, 13.40.-f, 31.30.jr

I. INTRODUCTION

The term “proton radius puzzle” paraphrases the disagreement between muonic hydrogen Lamb shift experiments (0.8409(4) fm) [1, 2] and both atomic and scattering experiments using electrons, summarized in the CODATA value of 0.8751(61) fm [3]. The extraction of the proton radius from scattering data is a treacherous business. In the discussion about the proton radius puzzle, many pitfalls we and others succumbed to became obvious. This paper is meant as an illustrated guide of these.

The paper is divided in two main sections: in the first section, we discuss missteps and misconceptions in general terms. The second section discusses the flaws in the analysis of some recent papers.

II. COMMENTS ON COMMON MISTAKES AND MISCONCEPTIONS

In the following sections we discuss common mistakes that are somewhat specific for the extraction of the proton radius from cross section data (II A - II F). Starting with section II G we talk about general properties of estimators which are relevant whenever a given quantity is calculated based on observed data.

A. A polynomial fit is *not* a Taylor expansion around 0, and the convergence is *not* limited by cuts in the time-like region.

A polynomial in normal form, i.e., of the form

$$\text{poly}(x, \vec{p}) = p_0 + p_1 \cdot x + p_2 \cdot x^2 + \dots$$

looks identical to a Taylor expansion around 0:

$$\text{taylor}[f](x) = f(0) + \frac{1}{1!} \frac{df}{dx} \Big|_0 \cdot x + \frac{1}{2!} \frac{d^2f}{dx^2} \Big|_0 \cdot x^2 + \dots$$

However, a fit of the polynomial does not yield the Taylor expansion. This can trivially be seen just looking at the definition: a Taylor expansion of a function around a point is given by the derivatives of that function at that point. This necessitates that the function indeed has these derivatives, and the value of the function at any other point is of no consequence for the expansion. The polynomial used in a fit might look like a Taylor expansion, but it is not: the coefficients of the polynomial are influenced by all data points, i.e., it depends on the functional value at many ordinate points. A fit with a polynomial written like a Taylor expansion around a different point x_0 , i.e.,

$$\text{poly}(x, \vec{p}) = p_0 + p_1 \times (x - x_0) + \dots$$

will find a different parameter vector \vec{p} , but transforming the polynomial to normal form by multiplying out the parenthesis will yield the same polynomial, independent of the choice of x_0 . It is worthwhile to note that the polynomial fit in general does not yield a Taylor expansion at all, i.e., there is no common point x_0 where the polynomial and the true function have the same value and derivatives.

Indeed, according to the Weierstrass theorem, any function continuous in an interval can be approximated to arbitrary precision and with global convergence (over the interval) by a polynomial. This alone does not guarantee that the first derivative is also approximated well, the requirement for an accurate extraction of the radius. However, it is trivial to show that this is true if the function is continuously differentiable.

In contrast to the theorem of Weierstrass, which concerns itself with convergence of the maximum error, i.e., norm $\|\dots\|_\infty$, the typical fit in the least squares sense minimizes according to norm $\|\dots\|_2$, a fit-technical necessity (the error function needs to be continuous close to the optimal point) which also lends it itself to the treatment of data with errors.

The prevalence of the notion that a polynomial fit is somehow related to a Taylor expansion is striking [4–9]. We want to present here an example: to this end, we generated G_E values following the standard dipole, i.e., a dipole with a parameter of 0.71 (GeV/c)², at the Q^2 points of the Mainz data set. These values are then fit with a 10th order polynomial. The data points are error-free, but we weight the points according to the uncertainty present in the Mainz data set. In Fig. 1, the

* bernauer@mit.edu

† distler@uni-mainz.de

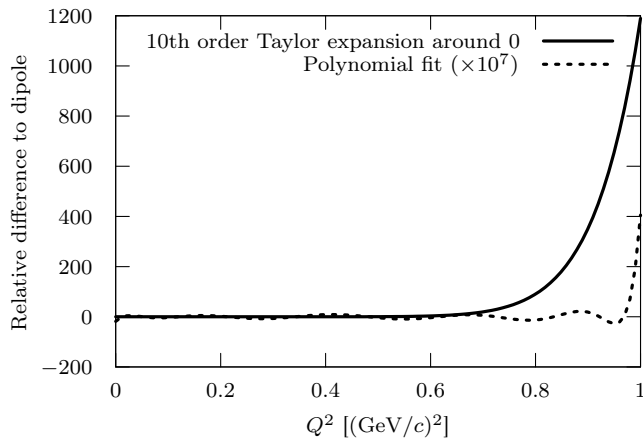


FIG. 1. Relative difference of the polynomial fit and of a Taylor expansion to the dipole, as a function of Q^2 . Please note that the difference of the polynomial is scaled up by a factor of 10 million, i.e. the difference is less than 40 ppm for the whole displayed range.

difference of the polynomial fit and of a Taylor expansions around $0 (\text{GeV}/c)^2$ truncated to 10th order. The standard dipole has a pole at $Q^2 = -0.71 (\text{GeV}/c)^2$, therefore a Taylor expansion around 0 is limited in its convergence to a radius of $0.71 (\text{GeV}/c)^2$. As expected, the Taylor expansion diverges strongly from the dipole close to $0.71 (\text{GeV}/c)^2$. In contrast, the polynomial fit does not diverge from the dipole by more than 40 ppm between 0 and 1 $(\text{GeV}/c)^2$. Indeed, the polynomial fit approximates the dipole better than the Taylor expansion for all Q^2 above $0.15 (\text{GeV}/c)^2$.

Many authors [4–7] argue that a polynomial fit is limited in its convergence to $Q^2 < 4m_\pi^2$ because of a pole at $Q^2 = -4m_\pi^2$, i.e., in the time-like region, and limit their fits to the region below, even for non-polynomial fits. As shown, this reasoning is wrong.

Of course, a Taylor expansion around a Q_0^2 more centered in the Q^2 interval one is interested in would perform better. One might be led to believe that the fit might relate to a Taylor expansion not around 0, but around a $Q_0^2 \neq 0$, an effective, weight-averaged center of gravity of the data points (indeed, the authors held this believe briefly). But this is not true in general, as can be shown for this example. From the coefficients found in the fit, one can calculate, order by order, which possible Q_0^2 these belong to. In the example case at hand, one each order, one finds 12 possible Q_0^2 , however, none of them are common to all orders, as is illustrated in Fig. 2. Therefore, the best fit polynomial is not a truncated Taylor expansion of the dipole function around any (one) point.

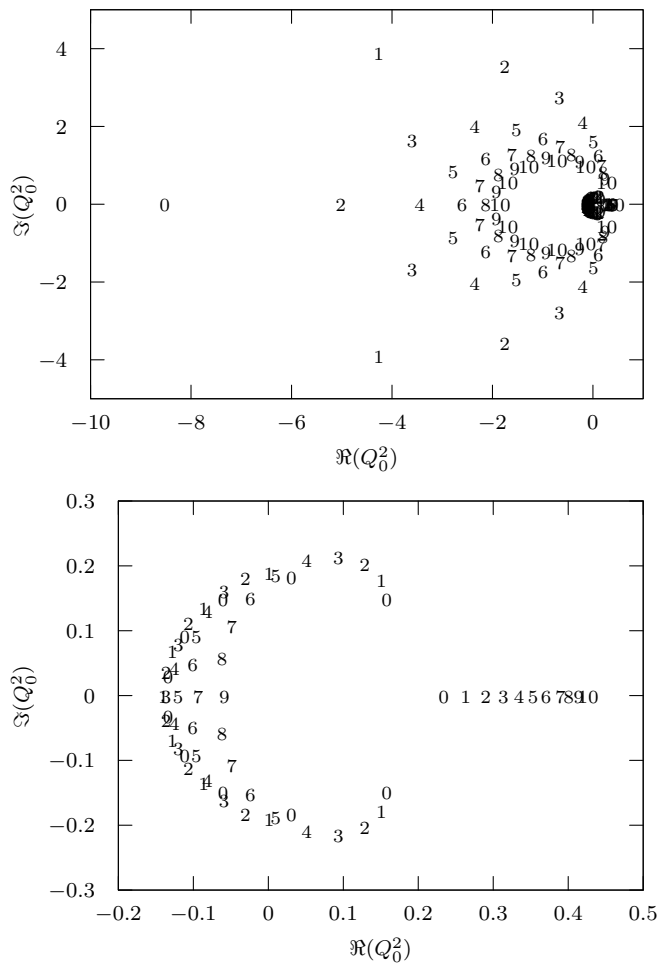


FIG. 2. A comparison of the polynomial fit coefficients with the symbolic expression for a Taylor expansion of the dipole at an arbitrary Q_0^2 order-by-order yield 12 complex-valued Q_0^2 for each order, displayed here in the complex plane and labeled by the order they stem from. None of them coincide for all orders, proving that the polynomial fit is not a Taylor expansion of the dipole at all.

B. Unconstrained fits with conformal mapping is *not* a good idea

Conformal mapping is used by some to avoid the perceived problem of the convergence radius. E.g., in [7], the authors define the function

$$z(t, t_{\text{cut}}) = \frac{\sqrt{t_{\text{cut}} - t} - \sqrt{t_{\text{cut}}}}{\sqrt{t_{\text{cut}} - t} + \sqrt{t_{\text{cut}}}}, \quad (1)$$

with $t = -Q^2$ and $t_{\text{cut}} = 4M_\pi^2$. The form factors are then expressed as a polynomial in z instead of Q^2 ,

$$G_{E/M}(t) = \sum_{k=0}^{k_{\text{max}}} a_k \cdot z(k)^k \quad (2)$$

The mapping function maps the whole positive Q^2

range into the range $[0..1]$ in a rather non-linear fashion, compressing the larger Q^2 values to a very small range in z close to 1. On the other hand, the very-low Q^2 is mapped to a comparable large range. To illustrate this point further, in the unmapped case, the fit has to “bridge” from 0 to ≈ 0.004 $(\text{GeV}/c)^2$, or about 0.4% of the range of the data. In the mapped case, it has to bridge from 0 to 0.0133, or about 2.2% of the range of the data. It follows that the flexibility of the polynomial expansion is shifted to the low- Q range, which leads to multiple problems:

- In the low- Q region, the fits are very flexible. However, the data starts at a minimal Q^2 , so that a fit can introduce arbitrary structures below the data. The extraction of the radius from the data is only meaningful if one assumes that such structures do not exist. This is warranted, as such structures typically lead to a charge density distributions with pathologically large densities at large radii [10]. Additionally, and even more relevant here, is that the analysis extracts both electric and magnetic form factor at the same time. The large flexibility of the model makes this completely unstable, as we show below. This also influences the charge radius extraction. At the lowest Q^2 , there are only measurements for one beam energy, and a Rosenbluth separation is not possible. A extraordinary flexible model for G_M in that region can “steal” from the electric form factor, affecting the extracted radius.
- The compression of the larger Q^2 to a small range of z values exacerbates a problem inherited by many polynomial-type fits: the parameters tend to get very large, but the contributions to the fit of the different orders cancel to a large extent, especially at large z . At small z , only a small difference remains which is exploited by the fit algorithm to explain the data. However, many combinations of large parameter values exist which all give similar quality of fits, but are far apart in parameter space. Care must be taken that the fit actually converges to the best minimum.

Both of these points can be somewhat addressed by constraining the parameters, as has been carried out in [9] for an older data set. On the other hand, in a fit to the Mainz data, Lee et al. [11] find a strong dependence on the cut-off in Q^2 . We believe this to be a consequence of aforementioned points.

C. A good χ^2 does *not* signal a trustworthy extraction of the radius

To rely on χ^2 to indicate a good fit is dangerous. In the original meaning, it is a test of the data quality; assuming that a) the model is correct, b) the errors are statistical and exactly known and c) the individual data points are

independent (or their correlation is at least known), it expresses how likely it is that the data are drawn from the distribution given by the model. All of these assumptions are typically violated:

- One normally does not know whether it is the correct model. Indeed, this is what one wants to test. An incorrect model, however, can produce small χ^2 values and still be wrong.
- In many experiments, especially the Mainz data set, a sufficiently large part of the errors is not driven by counting statistics but other effects. This limits the knowledge we have about the errors.
- Data have systematic errors which couple the data points. The summands in the χ^2 sum are not independent, but the correlation is unknown.

We refer to Kraus et al. [12], for an illustrative discussion.

One more caveat: the minimal sum of the weighted squares of deviations of the data from a model function should be distinguished from χ^2 and we usually call it M^2 . For the reasons given above, M^2 does not follow a χ^2 -distribution in general. However, we will adhere to the common practice and call it χ^2 in the following chapters.

D. Low-order fits are *not* a good idea

While one would hope that a linear model converges to the same value as a higher order model if the Q_{max}^2 is sufficiently small, the current state of the data clearly does not reach far enough down. We again refer to Kraus et al. [12] which discuss this at length. As an additional caveat, we want to repeat that the polynomial fit is not a Taylor expansion. In a truncated Taylor expansion, the error at the expansion point is zero, and grows from there. One expects that a lower-order expansion has a smaller radius in which the error is below a certain threshold, but the error is still zero at the expansion point. However, in a fit, this is not true. While a lower-order fit will have a bigger error, the localisation of the error is less clear. A fit will approximate the local slope of the data (i.e., at $Q^2 > 0$), not at 0.

E. Common fit algorithms do *not* always find the true minimum

In a fit, one searches for the global minimum of χ^2 , the absolute best parameters. Depending on the particular model, the χ^2 landscape can have many local minima, and many fit algorithms are prone to get stuck in one of them. In our fits, we found that both continued fraction expansion and conformal mapped polynomial type fits are especially susceptible to this problem. Except for an exhaustive search, which is prohibitively slow, there are

no generally robust algorithms available, but simulated annealing is often successful even in hard cases. In our fits, we test for this problem by fitting repeatedly with different, random start values. This can help find a better minimum in many cases, however it's impossible to prove that the found minimum is indeed the global one. We recommend to avoid models which have too many local minima; depending on the noise in the data, the true minimum might not be the global minimum using that particular data set.

Another indication for this type of problem is the dependency of χ^2 on the fit order N . For any group of models G_N , where the images in function space,

$$\mathfrak{J}(G_N) = \{G_N(Q^2, a_0, a_1 \dots a_N), \forall a_k \in \mathbb{R}\}, \quad (3)$$

fulfill the relation

$$\mathfrak{J}(G_N) \subseteq \mathfrak{J}(G_{N+1}), \quad (4)$$

the χ^2 achieved by the models must monotonically decrease as a function of N :

$$\chi_{N+1}^2 \leq \chi_N^2 \quad (5)$$

Before we implemented the randomized start value approach from above, fits of polynomial models violated this condition when the number of parameters was excessively large.

F. Rescaling the errors in the Mainz data set does *not* allow for bad fits to be correct

In the Mainz analysis [13, 14], we use the χ^2 of our best model to determine the size of point-to-point errors on top of the counting statistics errors. This might overestimate the errors in two ways; the data also contains systematic errors, and even the best model might have systematic differences from the true model. On the other hand, the model might overfit the data, giving a slight underestimate of the errors. In total, we believe the errors to be accurate to $< 10\%$. Many take this as a license to scale the errors up if their fit produces a too large χ^2 . Doing so, however, would not change the relative ordering of the fits; the better fitting models still are better, and an explanation for the worse fit of their model must be given.

G. A statistics test can *not* tell which model is the true model

When fitting data where the true model shape is unknown, as is the case for form factors, we must resort to flexible models like polynomials or splines. The crucial question is now how flexible the model actually has to be—one has to balance between minimizing bias and possible overfitting. One is tempted to try to deduce from

the data how much flexibility is needed, and indeed we do the same. However, one has to be very careful: typical statistical tests, like the F-test, are used to identify a model that best fits the data. It can not prove that the simpler or the more complicated model is true, nor that the parameters it extracts are unbiased. For an example, see Section III E.

Additionally, one has generally an interpretation problem: in the standard F-test, the zero hypothesis H_0 , which one tries to disprove, is: the simpler model is correct. The rest of the method now assumes H_0 to be correct, tests whether the data conforms to that and based on this rejects or accepts H_0 . To this end, one defines a false rejection threshold, i.e., one finds a threshold for the test function so that one would falsely reject H_0 even if it's true with a small probability. However, this is decidedly not related to the probability that H_0 is actually correct, because one does not know how often the test would accept/reject H_0 if H_0 is actually false.

The falsehood of the approach can be illustrated differently: taking a large data set, one finds that a given complexity is advocated by these methods. Reducing the data set, for example by a Q^2 cut-off, will require a simpler model. However, in truth, only one (or none) of these models can be true, invalidating the theoretical basis of the test.

For nested problems, in general, the coefficient of a lower order changes when higher orders are fitted. While the data might not be good enough to prove that these higher orders are required, they might still be there, and neglecting them in the fit leads to a bias. For polynomials, one can find a basis orthogonal in respect to the data, for example via the Forsythe method. Then, indeed, one can use a statistical criteria to select the number of basis functions without affecting the extraction of quantities related to the lower order coefficients. Unfortunately, the radius, i.e., the linear term, appears in all orders except for the constant term, so that this approach does not help for the problem at hand.

For purely polynomial fits, however, it is easy to see that any hypothesis which truncates the order must be wrong: a polynomial will always go to $\pm\infty$ for $Q^2 \rightarrow \infty$, but we know that the form factors approach 0. This means that any statistical approach which assumes any truncated hypothesis to be true is built on sand. As a consequence, the radius extracted with a truncated polynomial will always have a bias from that truncation. This does not mean all hope is lost, as this error gets smaller if one includes higher orders, a consequence of the theorem of Weierstrass.

H. An estimator is *not* guaranteed to be consistent and unbiased

It is necessary to review what it means in statistical terms to indirectly “measure” a quantity like the charge radius given a set of data, e.g., cross section data. We

will stick to the frequentist interpretation of statistics laid out in [15, 16] where probability is interpreted as the frequency of the outcome of a repeatable experiment. None of the following insights are new or original but can be found in many text books on statistics. Most of the time we are only paraphrasing.

An indirect measurement translates to an estimate of a parameter. An estimator \hat{a} is a function of the data used to estimate the value of the parameter a . Therefore the estimator \hat{a} is treated like a random variable. As there is no general rule on how to construct the estimator, one chooses a function with optimal properties. Important properties are consistency and unbiasedness which relate the estimator \hat{a} and the *true* value of the parameter, a_0 .

An estimator is called consistent if the estimator \hat{a} is equal to a_0 in the limit of an infinite sample size:

$$\lim_{n \rightarrow \infty} \hat{a} = a_0$$

The bias b of an estimator is the difference between the expected value of \hat{a} and the true value of the parameter:

$$b = E[\hat{a}] - a_0$$

Commonly used methods to construct such an estimator are the least squares method or the more general maximum likelihood method. However, there are many more possible methods to construct an estimator. Also, it can not be implied that the method of least squares results in a consistent and bias-free estimator, not even in the simplest cases.

For example, given N data points x_i , where $i = 1, 2, \dots, N$ and we assume the data points are drawn from a Gaussian distribution. Using the maximum likelihood method one gets the estimators for the mean and the variance:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (6)$$

$$s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (7)$$

However, the maximum likelihood estimator for the sample variance in equation (7) is biased. In this special case a small change leads to the well known, bias-free estimator of the true sample variance:

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (8)$$

Under very controlled circumstances, linear models, knowledge of the true model function, known statistical errors, one can rely on asymptotic properties. For all other cases a simulation with pseudo data has to be performed in order to check the consistency and the unbiasedness of the estimator used.

I. The robustness of an estimator is *not* self-evident

The robustness of an estimator describes the insensitivity of the estimator in the face of false data and false assumptions. In the case of the proton radius extraction we want that our estimate is not unduly affected by systematic errors in the data or the specific functional form of the form factors that we use. Also the precise value of the Q^2 cut-off or small changes in the cut parameter (conformal mapping) should not affect the estimation.

To illustrate the importance of the robustness criterion we can examine the estimation of the centre of an unknown, symmetric distribution. As shown in many textbooks (e.g. [15]) the well known sample mean is only optimal if the distribution is normal. For the double exponential distribution the optimal estimator is the median and if the distribution is unknown one should use the trimmed mean where the highest and lowest values of the sample are removed and the sample mean is calculated from the remaining 46% of the observations.

This demonstrates that the robustness of an estimator is not at all self evident, not even in the simplest cases. The properties of an estimator have to be studied carefully.

J. An estimator is *not* necessarily efficient

Recall that the estimate \hat{a} of a parameter a itself is a random variable. We have discussed the bias of an estimator in Section II H. Now, we focus on the efficiency. In statistics, the efficiency is about the variance of an estimator. An efficient estimator has the optimal (minimal) variance. Again there are very simple textbook examples where the standard procedure does not provide the most efficient estimator.

Consider the mean of a sample: if the underlying distribution is the uniform distribution, the use of eq. 6 will not give you the most efficient estimate of the sample mean. However, the midrange

$$\bar{x} = \frac{\hat{x} + \check{x}}{2} \quad (9)$$

which is the mean of the two extreme values within the sample has the minimal variance. The arithmetic mean, which is the most efficient estimate of the sample mean for the normal distribution, does poorly for the uniform distribution. The variance of the estimator scales with $1/n$ where n is the sample size. When using *midrange* on a sample drawn from a uniform distribution the variance is proportional to $1/n^2$.

Again, a simulation with pseudo-data will help to evaluate the variance of the estimator that is used.

III. COMMENTS ON RECENT PAPERS

A. Failed fits

In the recent paper [7], the authors use the conformal mapping approach to fit the recent high precision form factor data from Mainz [13, 14], claiming a 3 sigma reduction in the proton radius puzzle. We believe that this finding is in error on multiple accounts: the fit function is, as is, not suited to analyze the data, their fitting program does not converge to the minimal solution, and their statistical approach is flawed. Additionally, the comparison with the Mainz fits is not on equal footing.

We tried to replicate the approach followed by Lorenz et al. in [7]. Our results however differ significantly from the ones reported there. The nature of the differences mainly point to a failure of the fitting algorithm used in [7] to reliably find the true minimum. Trying to reproduce Fig. 1 of [7], we find a completely different χ^2 evolution, namely significantly lower values for smaller k_{\max} , even with a naive implementation of the fitting routine.

The original paper is not clear on whether a_0 is set to 1 or fitted. Since fitting it would constitute a renormalization, we set $a_0 = 1$. In any case, this limits the flexibility of our model, that is, a fit including a_0 as a free parameter will produce an even smaller χ^2 .

For larger k_{\max} , we have to employ the more advanced fitting algorithm described in Section II E and find consistently lower numbers than what was reported in [7]. Since they use a polynomial fit, the χ^2 have to follow eq. 5, which is violated for $k_{\max} = 13$ or 14 (worse fit than $k_{\max} = 12$).

We also find a rather strong dependence on the pion mass used in the mapping and we therefore report both results (see Section III on the robustness of an estimator). Figure 3 shows a comparison of our results and the ones from [7]. In Section III, we emphasized the importance of the robustness of a model that is used to extract a parameter like the charge radius from a set of data. The strong dependence on the pion mass clearly violates that criterion of a good estimator.

For $k_{\max} \geq 10$, the best solution found by the fit sometimes produces a non-physical, that is, imaginary, magnetic radius. We therefore also keep the best solution with a real magnetic radius, which has a slightly larger χ^2 , still below the values found in [7]. Both curves are shown in Fig. 3. We only show results for valid radii in Fig. 4, which shows the dependence of the extracted radii and reached χ^2 on k_{\max} .

We find the typical “knee” in χ^2 around $k_{\max} = 6$, much smaller than $k_{\max} = 9$, found in [7]. Compared to the fits in [13, 14], the knee is softer, with visible reduction in χ^2 beyond the knee. We interpret this as a sign that the fit is already overfitting the low- Q^2 region, but still can make use of the added flexibility at larger Q^2 , where the mapping function compresses the range.

In contrast to [7], we do not observe any stable plateau

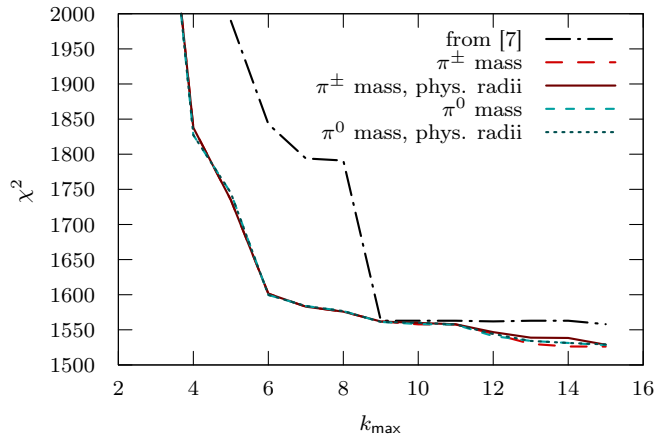


FIG. 3. Comparison of the achieved χ^2 in [7] and by us as a function of k_{\max} . Using the same data and fit function, we find substantially smaller χ^2 , with the characteristic knee at 6 instead of 9.

of the radii. From the properties of the fit function, this is somewhat expected. We can only speculate over the exact nature of what caused the plateau in [7].

Following the procedure of the Mainz analysis [13, 14], we make use of two criteria to find suitable parameter numbers. The lower bound is given by the position of the knee, while the upper bound is found by looking for a plateau in both charge and magnetic radii. The rationale behind this is easy to understand: the knee signals that the model has enough flexibility to follow the underlying shape of the data. With less flexibility, the fit has a common-mode offset from the data, leading to a large increase in χ^2 . With more flexibility, the fit starts to follow local, statistical fluctuations, which only reduce χ^2 slightly. With further flexibility, the fit gets unstable, which can be seen in the radii. Of course, these rules are not rigorous, but constitute a good guide line for the selection.

As shown in Fig. 4, there is clearly no plateau in the magnetic radius. We would therefore not accept the model at all.

However, it is interesting to note that, ignoring the magnetic radius for a moment and focusing on the charge radius, the fit extracts values in the range from 0.866 to 0.876 fm for $k_{\max} = 6 \dots 8$, slightly lower, but in good agreement with our reported results.

B. Low order polynomial fits to low- Q data

Motivated by the perceived connection of polynomial fits to Taylor expansions and their radius of convergence (see Sections II A and II D), Griffioen et al. [5] fit first and second order polynomials to the data up to $Q^2 = 0.2$ (GeV/c)² and report radii close to 0.84 fm.

To illustrate the problems of these fits, we generate two groups of pseudo-data. The first groups are gener-

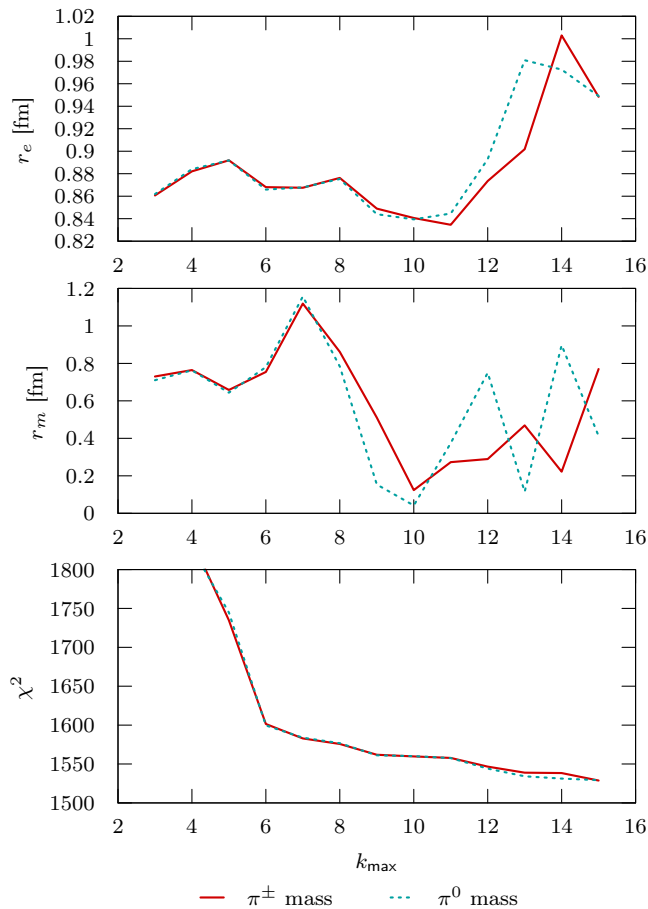


FIG. 4. Extracted radii and achieved χ^2 as a function of k_{\max} , using a polynomial fit and conformal mapping. The magnetic radius does not show a stable region for k_{\max} above the knee and we would therefore reject the model all together.

ated from the 10th order polynomial fit from [13, 14], corresponding to a radius of 0.8855 fm, the other from a 10th order polynomial fit to the data of [13, 14], with the radius forced to 0.841 fm. For each group, we simulate 2000 repetitions of the Mainz experiment, generating 2000 data sets. These pseudo-data set, and the real data set, can now be analyzed in various ways and one can compare the behavior of the fits to the real data and to the pseudo-data sets. For the real data set, we selected the normalization using the polynomial fit (see explanation in [14]), and use the (fixed) 10th order polynomial fit for G_M together with the to-be-optimized model for G_E to fit on the cross section level.

The results for the first order fits are shown in Fig. 5. The strong bias (see Section IIH) in the fits to pseudo-data is obvious, even for very small cut-offs. At 0.02 (GeV/c)^2 , we find an average bias of more than 0.04 fm, yielding essentially the small muonic radius of 0.84 fm despite having a true radius of 0.8855 fm. It follows that results from linear fits are unreliable; assuming that our polynomial fit is indeed an accurate representa-

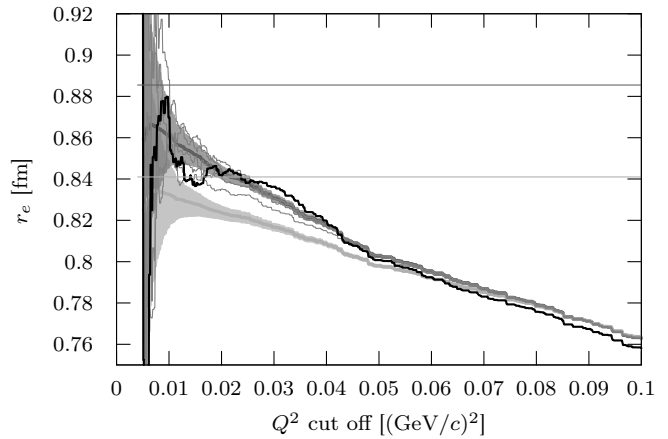


FIG. 5. Extracted radii from linear fits to pseudo and real data, as a function of the Q^2 cut-off. Black curve: fits to data; grey thick curves: average extracted radius to pseudo data (darker, upper curve: pseudo-data with large radius; lighter, lower curve: pseudo-data with small radius); bands around these curves are one-sigma point-wise error bands; dark grey thin curves: fits to the first five pseudo-data sets with large radius.

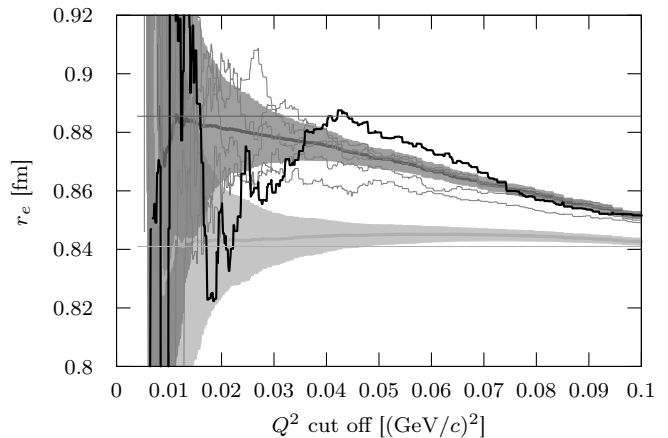


FIG. 6. Extracted radii from quadratic fits to pseudo and real data, as a function of the Q^2 cut-off. Curves as in Fig. 5.

tion of reality, the bias observed for pseudo-data explains the small radius found by Griffioen et al. [5].

It is striking how similar the fit to data is compared to the fit to the pseudo-data with large radius. It is worthwhile to note that the real data may very well have systematic errors affecting small groups of data points, not reflected in the generation of pseudo data here.

A second order fit does somewhat better, as shown in Fig. 6. The overall picture is somewhat similar to the first order fit. However, on average, the quadratic fit should have a much smaller bias. Nevertheless, the errorband which is a measure of the variance is much bigger compared to the linear fits (see Section IIJ). The fits to data show a dip around the cut-off used by Griffioen et al.

but recover and come back to higher values, until the bias lowers the extracted value again. Comparing the cut-off dependence of the fits to data to that of fits to individual sets in the pseudo-data, one can see similar swings, albeit maybe somewhat less pronounced. This might be simply the result of a statistical fluctuations, or of a local problem in the data around 0.02 (GeV/c)^2 . Both possibilities will hopefully be addressed with future data. A low order fit to small data sets will statistically be more sensitive to such perturbations: first, problems at the highest accepted Q^2 will affect the highest order most, and the effect is diminished on the first order term, more so if more data are fitted with higher-order functions. Second, even assuming that the probability that a data point is affected by such systematic effects is constant (the probability is likely smaller for higher Q^2 data, as corrections, e.g., due to backgrounds, are smaller), multiple systematic effects in the larger data sets will partially cancel, so their relative influence is likely proportional to $1/\sqrt{N}$.

For third order fits, Griffioen et al. propose to expand the form factor as

$$G_E(Q^2) = 1 - \frac{1}{6}R_E^2 Q^2 + \frac{b_2}{120}R_E^4 Q^4 - \frac{b_3}{5040}R_E^6 Q^6.$$

The coefficients are given by models where form factor and charge distributed can be expressed in terms of elementary functions with one parameter R_E and the expected values $\langle r^n \rangle$ are simple multiples of R_E^n . We have put the relevant formulas in the appendix A. However, the authors of [5] limited their analysis to three models, i.e., exponential, Gaussian and box shaped charge distribution and they did not investigate the bias (Section IIH) and the robustness (Section III) of their ansatz. We will show that this is a severe shortcoming that completely invalidates their conclusion.

The exponential or dipole model is of course an obvious choice for the proton. The other two form factor models have a smaller kurtosis than the dipole and would be suitable for light and heavy nuclei, respectively. Therefore we look at two more models: Yukawa I and II. Both are more “peaked” than the dipole model and the later, a simple pole, has been used to fit the pion form factor.

In order to evaluate the bias and the robustness of the five models we generated pseudo data equally spaced in the momentum transfer range $Q^2 = (0.004 \dots 0.02) \text{ (GeV/c)}^2$ with a constant standard deviation of 0.5%, 201 data points in total. The result of this analysis is shown in Tab. III B. With a few exceptions, any mismatch between assumed functional form and actual functional form leads to large biases. We conclude that, as long as one does not regularly win the lottery, one should not guess the functional shape.

Input model	$b_{2/3}$ of fit function according to				
	Dipole	Gauss	Box	Y. I	Y. II
Dipole	0(4)	-5(4)	-9(4)	22(5)	5(4)
Gauss	5(4)	0(4)	-3(4)	28(5)	11(4)
Box	9(4)	3(4)	0(4)	31(5)	14(4)
Yukawa I	-21(4)	-26(4)	-29(4)	-1(5)	-16(5)
Yukawa II	-5(4)	-11(4)	-14(4)	16(5)	0(5)
P×D	-8(4)	-14(4)	-17(4)	13(5)	-3(4)
Spline	-6(4)	-11(4)	-14(4)	16(4)	-1(4)

TABLE I. Bias and standard deviation in attometer. A mismatch between assumed functional form and actual functional form can lead to significant biases.

C. G_E/G_M ratio and continued fraction expansion

In the second part of [5], the authors extract G_E and G_M from the whole Mainz data set using

$$\mu_p G_E/G_M = 1 - Q^2/Q_0^2, \quad (10)$$

with $Q_0^2 = 8 \text{ (GeV/c)}^2$. The form of eq. 10 is motivated by measurements of the form factor ratio using polarization. We believe that this approach is dangerous and wrong on multiple accounts:

- It is a well known fact that the form factor ratio extracted from polarized measurements is different from the one extracted from unpolarized experiments. The most likely explanation is the neglect of two-photon exchange, which affects mainly the unpolarized measurements. However, so far, this is only a conjecture. The (unpolarized) Mainz dataset has no full two-photon exchange corrections applied. It is questionable to extract the form factors assuming a ratio from polarized data.
- The linear fall-off describes the gross behaviour of the ratio in polarized data, but the world data set is certainly not good enough to see structures beyond that, especially below 1 (GeV/c)^2 , where the current polarized data is somewhat in disagreement with each other.

The authors then fit their extracted G_E using a continued fraction expansion

$$G_E(Q^2) = \frac{p_1}{1 + \frac{p_2 Q^2}{1 + \frac{p_3 Q^2}{1 + \dots}}}. \quad (11)$$

with 4 parameters, they achieve a $\chi^2/\text{d.o.f.}$ of 1.6 and claim that the data are well-fit on average in all regions of Q^2 . We can not follow this logic: a $\chi^2/\text{d.o.f.}$ of 1.6 is excessively high. Using our standard approach used in the Mainz analysis, we fit a continued fraction expansion of both G_E and G_M . The fits proved difficult, with

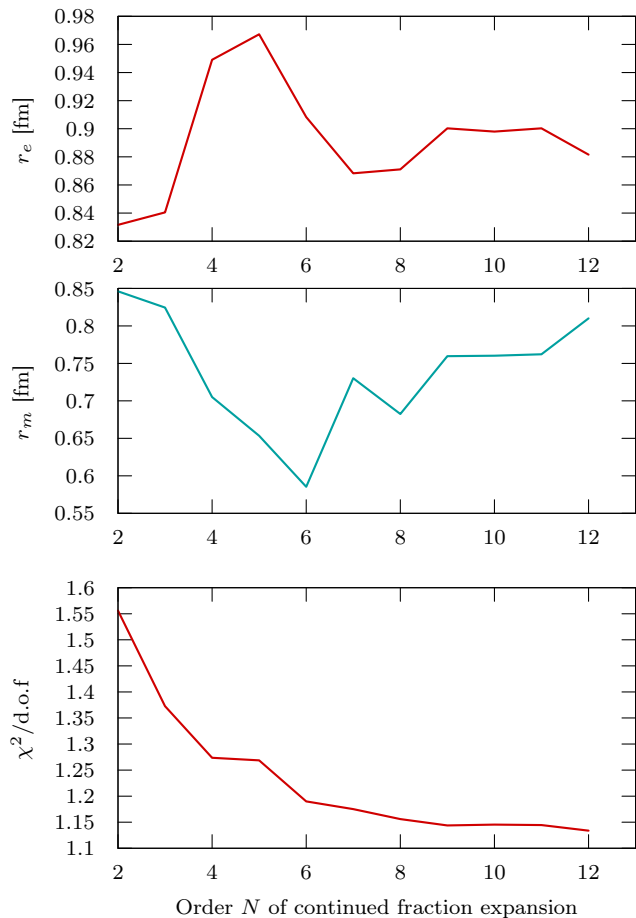


FIG. 7. Extracted radii and achieved $\chi^2/\text{d.o.f.}$ as a function of N using continued fraction expansions with N parameters for G_E and G_M . In contrast to other models, the knee is very soft. For fits with more than 3 parameters, $r_E > 0.868$ fm.

many local minima, and we can not rule out that better solutions exist. Nevertheless, the results, shown in Fig. 7, are interesting. At order 4, we already achieve a red. χ^2 substantially better than 1.6. Order 5 is only marginally better—we suspect a better solution exists, but our fit fails to find it, despite randomizing the starting values (see Section II E). At higher orders, we again see a substantial gain. Around 9, the red. χ^2 is comparable to the best models of our earlier analysis, with a somewhat larger $r_e \approx 0.899$ fm. For fits with more than 3 parameters, our extracted radius is always larger than 0.868 fm. It is unclear whether the difference to [5] is explained alone by the different extraction method. It is possible that their fitting algorithm falls victim to the adverse conditions of the fit too. Comparing their result for a double dipole (red. $\chi^2 = 1.6$) and our (1.29), the former seems likely. N.B.: we can not quite follow their remark about smoothly and monotonically falling fit functions. All our fit functions are smooth and monotonic, and achieve χ^2 around 1.15.

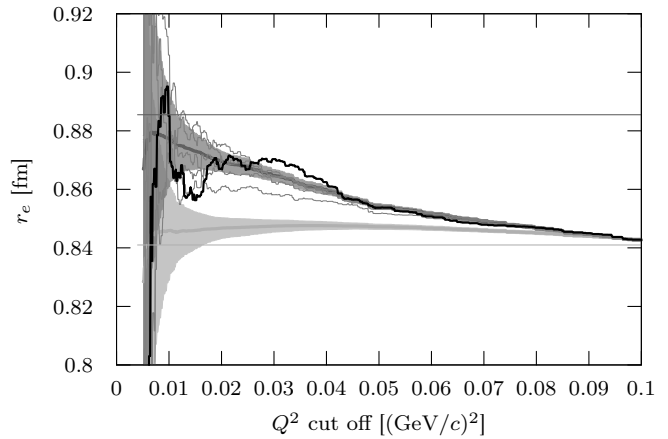


FIG. 8. Dipole fit to (pseudo)-data, same nomenclature as Fig. 5. At a cut-off of $0.1 (\text{GeV}/c)^2$, a dipole fit extracts almost identical values when fit to the two pseudo-data samples, in agreement with the value extracted from data. However, for lower cut-off values, the radii extracted from data replicate the behavior of the pseudo-data sample with $r_e = 0.8855$ fm, and does not follow the one with a small radius.

D. Dipole fit to low- Q data

In [6], Horbatsch and Hessels compare a conformal mapping polynomial fit with a dipole fit, for a range of Q^2 cut-offs and orders. Their z-expansion fits exhibit indications of the problems described in Sections II B, but generally reproduce the large radius, in agreement with our findings and in stark contrast to [7]. Their dipole fit, for data up to $0.1 (\text{GeV}/c)^2$, yields a value of $0.842(2)$ fm. While this might puzzle the reader, this is completely expected: the dipole model is known to have a strong bias, as already demonstrated in [14] for the whole data set. We repeat the procedure described in Section III B, fitting a dipole model. The results are shown in Fig. 8. At $0.1 (\text{GeV}/c)^2$, the extracted radii are identical, and no decision can be made. At lower cut-offs, the data clearly prefer the pseudo-data sample with a large radius. It is worthwhile to note that the dipole fit to the pseudo data sample with the large radius has a negative bias larger than the expected statistical error for cut-offs larger than $0.01 (\text{GeV}/c)^2$. A reliable extraction of the radius can therefore not be expected.

E. Statistical methods to decide order

In [4], the authors use the F-test to decide which order of polynomials to use. Besides the points addressed in Section II G, the statistical interpretation is flawed on a very basic level: they reference a critical value of 4.3 for $\text{CL}=95\%$, which is the critical value for the rejection of H_0 at this level, i.e., with an F-test value higher than 4.3, one should reject the simpler model, with a 5% probability that the rejection is wrong. They however claim that

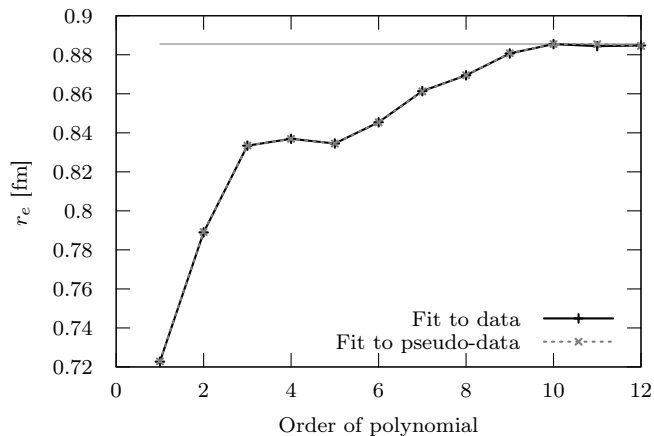


FIG. 9. Dependency of the extracted radius from polynomial models on the order, for the full Mainz data set. Fits to pseudo data track the behavior of the fits to real data, and can only recover the input radius (indicated by horizontal line) at high orders. Both F-test and AIC reject models with orders < 9 .

their value below this threshold rejects H1, i.e., the more complex model, at 95% CL. This inversion can of course not be done, and indeed, no confidence level can be given for this type of error easily, because the nominator in the F-test does not follow a standard Fisher-Snedecor distribution anymore, and because one is not restricted to just one higher order.

Nevertheless, with the pseudo data groups above, we can easily test what their flawed method would produce:

Comparing first and second order fits, the F-test would prefer (i.e., not rule out at CL=95%) the linear model up to 0.015 (large r_e group) and 0.02 (GeV/c)² (small r_e), respectively. At these Q^2 , the bias of the linear model is 0.03 fm and 0.02 fm, as can be seen in Fig. 5.

For second and third order fits, the cross-over is around 0.045 (GeV/c)² for the large radius pseudo data group, and above 0.1 for the smaller radius group (outside our simulation range), albeit with a certain fraction of the individual data sets hitting the threshold around 0.09. Comparing to Fig. 6, the method seems to work in this case for the smaller radius group—the second and third order coefficients of the input model are significantly smaller than for the large radius input model. For

the larger radius pseudo data group, however, the fit to pseudo data has a bias of 0.012 fm and 0.045 fm, and the fit to data up this point produces a *large* radius.

Let us now look at the behavior of fits of different order to the full data set. For this test, we will again use our 10th order polynomial as basis for the generation of pseudo data. We then fit polynomial models with different orders both to the real data and to the pseudo data. Instead of relying on a fixed G_M fit as we did for the low-Q fits, we fit both G_E and G_M at the same time, repeating our approach of [13, 14]. The results are shown in Fig. 9. The fits to the pseudo data replicate almost exactly the behavior of the fits to data, with also similar behavior for the F-test. It is interesting to see that lower orders linger around the muonic radius. However, this is a good example how to use statistical tests properly: The F-test rejects order 9 and below in favour of order 10, with a false rejection probability of $< 5\%$. The Akaike information criterion accepts order 9, and has a minimum at order 10. It follows that the H0 hypothesis, the lower order models are correct, is rejected by the data.

We further want to note that Table III of [4] is not consistent with its description. The listed values for χ^2 and χ^2/ν indicate that $\nu = N - j$, instead of $\nu = N - j - 1$ given in the description—and even assuming that, there seems to be a rounding error.

CONCLUSION

In summary, we inspected several recent refits of the Mainz data set which result in small radii and found flaws of various kinds in all of them. While a reanalysis of the data can not rule out faulty data—which would invalidate any extraction—we believe that the solution of the puzzle can not be found in the fit procedure. We urge anybody in the business to test their method using pseudo data generated from the Mainz fits.

ACKNOWLEDGEMENT

We thank Jörg Friedrich, Kees de Jager and Thomas Walcher for helpful discussions.

[1] R. Pohl *et al.*, Nature **466**, 213 (2010).
 [2] A. Antognini *et al.*, Science **339**, 417 (2013).
 [3] P. J. Mohr, D. B. Newell, and B. N. Taylor, (2015), arXiv:1507.07956.
 [4] D. W. Higinbotham, A. A. Kabir, V. Lin, D. Meekins, B. Norum, and B. Sawatzky, Phys. Rev. **C93**, 055207 (2016), arXiv:1510.01293 [nucl-ex].
 [5] K. Griffioen, C. Carlson, and S. Maddox, (2015), arXiv:1509.06676 [nucl-ex].

[6] M. Horbatsch and E. A. Hessels, Phys. Rev. **C93**, 015204 (2016), arXiv:1509.05644 [nucl-ex].
 [7] I. Lorenz and U.-G. Meißner, Physics Letters B **737**, 57 (2014).
 [8] G. Paz, *Particles and fields. Proceedings, Meeting of the Division of the American Physical Society, DPF 2011, Providence, USA, August 9-13, 2011*, AIP Conf. Proc. **1441**, 146 (2012), arXiv:1109.5708 [hep-ph].
 [9] R. J. Hill and G. Paz, Phys. Rev. D **82**, 113005 (2010).

- [10] I. Sick, *Progress in Particle and Nuclear Physics* **67**, 473 (2012), from Quarks and Gluons to Hadrons and Nuclei-International Workshop on Nuclear Physics, 33rd Course.
- [11] G. Lee, J. R. Arrington, and R. J. Hill, *Phys. Rev. D* **92**, 013013 (2015).
- [12] E. Kraus, K. E. Mesick, A. White, R. Gilman, and S. Strauch, *Phys. Rev. C* **90**, 045206 (2014).
- [13] J. C. Bernauer, P. Achenbach, C. Ayerbe Gayoso, R. Böhm, D. Bosnar, L. Debenjak, M. O. Distler, L. Doria, A. Esser, H. Fonvieille, J. M. Friedrich, J. Friedrich, M. Gómez Rodríguez de la Paz, M. Makek, H. Merkel, D. G. Middleton, U. Müller, L. Nungesser, J. Pochodzalla, M. Potokar, S. Sánchez Majos, B. S. Schlimme, S. Širca, T. Walcher, and M. Weinriefer (A1 Collaboration), *Phys. Rev. Lett.* **105**, 242001 (2010).
- [14] J. C. Bernauer, M. O. Distler, J. Friedrich, T. Walcher, P. Achenbach, C. Ayerbe Gayoso, R. Böhm, D. Bosnar, L. Debenjak, L. Doria, A. Esser, H. Fonvieille, M. Gómez Rodríguez de la Paz, J. M. Friedrich, M. Makek, H. Merkel, D. G. Middleton, U. Müller, L. Nungesser, J. Pochodzalla, M. Potokar, S. Sánchez Majos, B. S. Schlimme, S. Širca, and M. Weinriefer (A1 Collaboration), *Phys. Rev. C* **90**, 015206 (2014).
- [15] F. James, *Statistical methods in experimental physics* (World Scientific, 2006).
- [16] K. Olive and P. D. Group, *Chinese Physics C* **38**, 090001 (2014).
- [17] M. O. Distler, J. C. Bernauer, and T. Walcher, *Phys. Lett.* **B696**, 343 (2011), arXiv:1011.1861 [nucl-th].

Appendix A: Form factors and charge distribution of selected models

In the Breit frame electric (and magnetic) form factors can be associated with the charge (and magnetic current) density distributions through a Fourier transformation:

$$G(q) = 4\pi \int_0^\infty r^2 \rho(r) \sin\left(\frac{qr}{\hbar c}\right) \frac{\hbar c}{qr} dr \quad (\text{A1})$$

$$\rho(r) = \frac{4\pi}{(2\pi \hbar c)^3} \int_0^\infty q^2 G(q) \sin\left(\frac{qr}{\hbar c}\right) \frac{\hbar c}{qr} dq \quad (\text{A2})$$

This implies that the electric form factor can be expanded in terms of Q^2 where the coefficients are the multiples of the expected values of r^{2n} of the charge distribution:

$$\begin{aligned} G(Q^2) &= \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} \langle r^{2n} \rangle Q^{2n} \\ &= 1 - \frac{\langle r^2 \rangle}{6} Q^2 + \frac{\langle r^4 \rangle}{120} Q^4 - \frac{\langle r^6 \rangle}{5040} Q^6 + \dots \end{aligned} \quad (\text{A3})$$

The Zemach moments of the nuclear charge distributions (see [17] and references therein) are defined by

$$\langle r^n \rangle_{(2)} = \int d^3r r^n \rho_{(2)}(r) \quad (\text{A4})$$

where $\rho_{(2)}(r)$ is the convolution of the charge distribution

$$\rho_{(2)}(r) = \int d^3r_2 \rho(|\vec{r} - \vec{r}_2|) \rho(r_2). \quad (\text{A5})$$

The first and the third Zemach moment can also be expressed in momentum space:

$$\langle r \rangle_{(2)} = -\frac{4}{\pi} \int_0^\infty \frac{dQ}{Q^2} (G_E^2(Q^2) - 1) \quad (\text{A6})$$

$$\langle r^3 \rangle_{(2)} = \frac{48}{\pi} \int_0^\infty \frac{dQ}{Q^4} \left(G_E^2(Q^2) - 1 + \frac{Q^2}{3} \langle r^2 \rangle \right).$$

In the following sections we will give the form factors, the density distributions and their expected values of r^4 and r^6 for selected models as a function of $R = \sqrt{\langle r^2 \rangle}$. The first and the third Zemach moment and Zemach's convoluted density are shown as well. The latter is not available in closed form for the Yukawa I model.

1. Exponential (dipole) model

$$\begin{aligned} G(q) &= \left(1 + \frac{1}{12} \left(\frac{qR}{\hbar c} \right)^2 \right)^{-2} \\ \rho(r) &= \frac{3\sqrt{3}}{\pi R^3} \exp\left[-2\sqrt{3}\frac{r}{R}\right] \\ \rho_{(2)}(r) &= \frac{3\sqrt{3}}{8\pi R^5} \left(4r^2 + 2\sqrt{3}rR + R^2 \right) \\ &\quad \times \exp\left[-2\sqrt{3}\frac{r}{R}\right] \\ \langle r^4 \rangle &= \frac{5}{2} R^4 \\ \langle r^6 \rangle &= \frac{35}{3} R^6 \\ \langle r \rangle_{(2)} &= \frac{35}{16\sqrt{3}} R \\ \langle r^3 \rangle_{(2)} &= \frac{35\sqrt{3}}{16} R^3 \end{aligned} \quad (\text{A7})$$

2. Gaussian

$$\begin{aligned} G(q) &= \exp\left[-\frac{1}{6} \left(\frac{qR}{\hbar c} \right)^2\right] \\ \rho(r) &= \left(\sqrt{\frac{3}{2\pi}} \frac{1}{R} \right)^3 \exp\left[-\frac{3}{2} \frac{r^2}{R^2}\right] \\ \rho_{(2)}(r) &= \left(\sqrt{\frac{3}{\pi}} \frac{1}{2R} \right)^3 \exp\left[-\frac{3}{4} \frac{r^2}{R^2}\right] \\ \langle r^4 \rangle &= \frac{5}{3} R^4 \\ \langle r^6 \rangle &= \frac{35}{9} R^6 \end{aligned}$$

$$\begin{aligned}\langle r \rangle_{(2)} &= \frac{4}{\sqrt{3\pi}} R \\ \langle r^3 \rangle_{(2)} &= \frac{32}{3\sqrt{3\pi}} R^3\end{aligned}\quad (\text{A8})$$

3. Uniform

$$\begin{aligned}G(q) &= \left(\frac{3}{5} \frac{\hbar c}{qR}\right)^2 \left(-5 \cos \left[\sqrt{\frac{5}{3}} \frac{qR}{\hbar c}\right] \right. \\ &\quad \left. + \sqrt{15} \frac{\hbar c}{qR} \sin \left[\sqrt{\frac{5}{3}} \frac{qR}{\hbar c}\right] \right) \\ \rho(r) &= \frac{3}{4\pi R^3} \left(\frac{3}{5}\right)^{3/2} \Theta \left[\sqrt{\frac{3}{5}} R - r\right] \\ \rho_{(2)}(r) &= \frac{27}{8000\pi R^6} \Theta \left[2\sqrt{\frac{3}{5}} R - r\right] \\ &\quad \times \left(3r^3 - 60r R^2 + 80\sqrt{\frac{5}{3}} R^3\right) \\ \langle r^4 \rangle &= \frac{25}{21} R^4 \\ \langle r^6 \rangle &= \frac{125}{81} R^6 \\ \langle r \rangle_{(2)} &= \frac{12}{7} \sqrt{\frac{3}{5}} R \\ \langle r^3 \rangle_{(2)} &= \frac{160}{63} \sqrt{\frac{5}{3}} R^3\end{aligned}\quad (\text{A9})$$

4. Yukawa I

$$\begin{aligned}G(q) &= \sqrt{2} \frac{\hbar c}{qR} \arctan \left(\sqrt{\frac{1}{2}} \frac{qR}{\hbar c}\right) \\ \rho(r) &= \frac{1}{2\sqrt{2}\pi r^2 R} \exp \left[-\sqrt{2} \frac{r}{R}\right] \\ \langle r^4 \rangle &= 6 R^4 \\ \langle r^6 \rangle &= 90 R^6 \\ \langle r \rangle_{(2)} &= \frac{1}{3} \sqrt{2} (1 + 2 \log[2]) R \\ \langle r^3 \rangle_{(2)} &= \frac{3}{5} \sqrt{2} (3 + 4 \log[2]) R^3\end{aligned}\quad (\text{A10})$$

5. Yukawa II

$$\begin{aligned}G(q) &= \left(1 + \frac{1}{6} \left(\frac{qR}{\hbar c}\right)^2\right)^{-1} \\ \rho(r) &= \frac{3}{2\pi r R^2} \exp \left[-\sqrt{6} \frac{r}{R}\right] \\ \rho_{(2)}(r) &= \frac{3}{2\pi} \sqrt{\frac{3}{2}} \frac{1}{R^3} \exp \left[-\sqrt{6} \frac{r}{R}\right] \\ \langle r^4 \rangle &= \frac{10}{3} R^4 \\ \langle r^6 \rangle &= \frac{70}{3} R^6 \\ \langle r \rangle_{(2)} &= \sqrt{\frac{3}{2}} R \\ \langle r^3 \rangle_{(2)} &= 5\sqrt{\frac{2}{3}} R^3\end{aligned}\quad (\text{A11})$$