

RESEARCH

Open Access



Using historical data to facilitate clinical prevention trials in Alzheimer disease? An analysis of longitudinal MCI (mild cognitive impairment) data sets

Manfred Berres^{1,2*} , Andreas U. Monsch² and René Spiegel²

Abstract

Background: The Placebo Group Simulation Approach (PGSA) aims at partially replacing randomized placebo-controlled trials (RPCTs), making use of data from historical control groups in order to decrease the needed number of study participants exposed to lengthy placebo treatment. PGSA algorithms to create virtual control groups were originally derived from mild cognitive impairment (MCI) data of the Alzheimer's Disease Neuroimaging Initiative (ADNI) database. To produce more generalizable algorithms, we aimed to compile five different MCI databases in a heuristic manner to create a "standard control algorithm" for use in future clinical trials.

Methods: We compared data from two North American cohort studies ($n=395$ and 4328 , respectively), one company-sponsored international clinical drug trial ($n=831$) and two convenience patient samples, one from Germany ($n=726$), and one from Switzerland ($n=1558$).

Results: Despite differences between the five MCI samples regarding inclusion and exclusion criteria, their baseline demographic and cognitive performance data varied less than expected. However, the five samples differed markedly with regard to their subsequent cognitive performance and clinical development: (1) MCI patients from the drug trial did not deteriorate on verbal fluency over 3 years, whereas patients in the other samples did; (2) relatively few patients from the drug trial progressed from MCI to dementia (about 10% after 4 years), in contrast to the other four samples with progression rates over 30%.

Conclusion: Conventional MCI criteria were insufficient to allow for the creation of well-defined and internationally comparable samples of MCI patients. More recently published criteria for MCI or "MCI due to AD" are unlikely to remedy this situation. The Alzheimer scientific community needs to agree on a standard set of neuropsychological tests including appropriate selection criteria to make MCI a scientifically more useful concept. Patient data from different sources would then be comparable, and the scientific merits of algorithm-based study designs such as the PGSA could be properly assessed.

Keywords: Historical controls, MCI criteria, Clinical trial, Cohort study, Convenience sample, Meta-analysis

* Correspondence: berres@hs-koblenz.de

¹University of Applied Sciences Koblenz, Koblenz, Germany

²University Department of Geriatric Medicine FELIX PLATTER, Basel, Switzerland



© The Author(s). 2021, corrected publication 2021. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Introduction

Almost 10 years ago, our group published the PGSA (Placebo Group Simulation Approach) for debate to the Alzheimer community [1]. The proposed novel study design was intended to partially substitute for RPCTs (randomized placebo-controlled trials), i.e., clinical studies which by definition expose some of the participants to treatment with placebo. We argued that, in the case of Alzheimer's disease (AD), clinical prevention trials with pre-clinical subjects would typically last 18 months or longer—and that it was ethically problematic to put individuals with a high risk of developing dementia on an a priori inactive long-term medication. Instead of a concomitant placebo group, the PGSA introduced algorithm-based forecasts of trials subjects' expected own disease trajectories to account for the effects of baseline differences, time in the study, and the circumstances of trial participation. The original PGSA algorithms were derived from the Alzheimer's Disease Neuroimaging Initiative (ADNI) data [2] available at that time, i.e., from a then recent but, nonetheless, historical data set.

The pros and cons of using historical data in clinical trials have been discussed under a number of aspects [3]. Thus, in the case of rare diseases, it may be difficult to recruit sufficient numbers of patients for proper control groups in addition to the treatment group [4]. In the case of progressive, non-reversible and potentially fatal diseases, there are ethical issues limiting the use of inactive treatment, and it may also be difficult to obtain consent from patients for participation if one of the treatment options looks more promising than the other one. Finally, in situations where no effective treatments are available and a promising candidate treatment is to be tested, it may be difficult to convince all potential trial subjects to participate in a study which includes a placebo arm [5]. However, the absence of certain subgroups of patients in a clinical trial will lead to less representative samples and is then likely to cause bias in the results. Historical data could be considered in some of these situations as a potential substitute for a concomitant control group (for ALS see [6]).

A necessary prerequisite for using historical controls is that they are comparable to the current study population. This requirement is hardly ever fulfilled. If only demographic variables such as age, sex, and observable health status are different, adjustments for these covariates might solve the problem. If one attempts to set up a model for such adjustments, several potential historical controls need to be compared [7].

In the specific case of mild cognitive impairment (MCI), a clinical condition between normal aging and dementia (see definition in [8]), the problem of finding adequate historical data is aggravated by the fact that

MCI criteria have shifted over time and that there is no rigorous and generally accepted definition of the condition [9–11]. As a typical consequence thereof, different clinical drug trials with MCI subjects in the past have applied different inclusion and exclusion criteria [12]. In this paper, we investigate whether information from large MCI databases can be summarized in a heuristic manner such that a “standard control group” for use in future clinical trials with this population could be created. We will compare data from two cohort studies, one clinical drug trial and two convenience patient samples to investigate:

- Inclusion and exclusion criteria for MCI applied in five different patient datasets
- The selection of cognitive tests applied at study entry and at follow-up
- The homogeneity of patients' demographic and baseline data
- The homogeneity of disease progression as measured by cognitive tests and indicated by the proportion of transitions from MCI to dementia

While distributions of demographic and baseline data will be shown, the progression of the disorder is analyzed as the “effect of no treatment over time”. In a clinical study, this corresponds to the progression observed in the control group and would be contrasted to the progression observed in the treatment group. In the five datasets considered, an overall judgement was provided at each patient visit by an experienced clinician as to whether the patient had progressed from MCI to dementia. Progression rates and hazard ratios will be compared between studies.

In all the studies considered in this analysis, and in particular in the clinical drug trial [13], some patients were treated with anti-Alzheimer medication such as cholinesterase inhibitors or memantine. While these drugs are considered transiently effective in AD [14], none of them was shown to have significant and maintained effects on the progression of the disease from MCI to dementia [8, 13, 15]. For this reason, we will consider all subjects as “untreated” in our analyses.

Material and methods

Datasets used

We analyzed individual patient data from the following:

- The Alzheimer's Disease Neuroimaging Initiative (ADNI; <http://www.loni.ucla.edu/ADNI>)
- The National Alzheimer's Coordinating Centers (NACC; <https://www.alz.washington.edu>)
- The InDDEX clinical trial [13]

- The German Dementia Competence Network (CNG [16]);
- The Basel University Memory Clinic (BS-MC)

The ADNI and the NACC samples included only patients who were between 54 and 90 years old at entry. BS-MC included only patients with at least 7 years of education. In order to reduce variability, these restrictions were applied to all 5 datasets in our analyses.

The ADNI (Alzheimer's Disease Neuroimaging Initiative) study aims at investigating the prognostic value of biomarkers, in particular of MRI and PET images, to describe the progression of Alzheimer's disease from its preclinical to its symptomatic stages. It is led by the principal investigator [2] and representatives of the ADNI sites, the NIH (National Institutes of Health), the FDA (Food and Drug Administration), and contributing companies from the health industry. ADNI procedures follow a detailed protocol. Cognitive performance of the participants was assessed with the Alzheimer's Disease Assessment Scale – cognitive subscale (ADAScog; 11 items and modified version with 13 items) [17], the MMSE [18], a number of neuropsychological tests, and the Functional Assessment Questionnaire [19]. ADNI started in 2003 and by the time of our last data download on January 6, 2012 [20], the dataset contained 395 patients diagnosed with MCI at study entry. Two subjects were excluded from our analyses because they had less than 7 years of education.

The NACC (National Alzheimer's Coordinating Center) project was initiated by the National Institute on Aging /NIH. It developed a large database of standardized clinical and neuropathological research data collected from 29 Alzheimer's Disease Centers in the USA. Eight of the nine neuropsychological tests used in ADNI are also part of the NACC database [21]. We received data from the freeze of March 19, 2014. We eliminated the data from those participants that were also included in the ADNI project to avoid patients from being considered twice in our analysis. We selected MCI patients with memory impairment, with or without impairment in other domains, who were between 55 and 90 years old and had at least 7 years of education. This left 4328 MCI subjects for our analysis.

The InDDEx (Investigation of Delay of Diagnosis of AD with Exelon®) study [13] was a clinical trial sponsored by the Novartis Pharma, assessing the effect of the cholinesterase inhibitor rivastigmine on disease progression in patients with MCI. This placebo-controlled study did not show evidence of an effect of rivastigmine on either the rate of progression to dementia or the standardized Z score for a cognitive test battery. We therefore considered all patients as untreated and included them in our analysis. This conforms with the other four

patient samples where dementia-related medication was also permitted. The study applied the ADAScog, a neuropsychological battery that had only verbal fluency (animals) and the Boston Naming Test [22] in common with the battery used in ADNI and a different functional assessment scale. The InDDEx study enrolled 1018 patients randomly assigned to rivastigmine ($n=508$) and to placebo ($n=510$). After exclusions due to missing screening data, missing cognitive data, or age or education outside the admissible range, 861 subjects could be included in the present analysis.

The CNG (Competence Network Germany) study [16, 23] is a longitudinal multicenter cohort study of 14 memory clinics in Germany. It applies the ADAScog (12 subtests), six tests of the Consortium to Establish a Registry for Alzheimer's Disease – Neuropsychological Assessment Battery (CERAD-NAB) [24] and other tests. We received data of 787 patients with a diagnosis of MCI. After exclusion due to age and education restrictions, 726 were left for our analysis.

The BS-MC (Basel Memory Clinic) sample comprises data of patients referred by practicing physicians to the memory clinic of the University Hospital Basel, Switzerland, for diagnosis and treatment recommendations. Neuropsychological tests include the CERAD-NAB plus Phonemic Fluency (S-words) and Trail Making Tests A and B [25] and Digit Span Forward and Backward. Data of 2135 patients with MCI at baseline were downloaded on September 9, 2016. After application of age and education inclusion criteria data from 1558 patients were left for the analysis.

Statistical analysis

Demographics and baseline scores of frequently used cognitive tests are summarized in tables and partly in boxplots. To investigate the homogeneity of progression across the five patient samples, we performed meta-analyses for the changes from baseline of cognitive test scores. Confidence intervals in forest plots will show whether there are distinct differences between studies. Measures of heterogeneity confirm these results. Rates of transition from MCI to dementia will be shown in Kaplan-Meier curves, broken down by study, and hazard ratios for age, sex, and education will be compared in proportional hazards models.

Results

Definition of MCI

The different inclusion criteria for MCI applied in the five studies are summarized in Table 1. ADNI and CNG used the MMSE to assess cognitive status, although with different inclusion criteria: The lower limit for the MCI was 24/30 in ADNI, but 20/30 in CNG. These two studies requested a Clinical Dementia Rating (CDR) [27]

Table 1 Inclusion and exclusion criteria for the diagnosis of MCI used in five studies

	Eligibility			Inclusion		Exclusion	
	Age	MMSE	Other, medication	Functional impairment	Cognitive impairment	Depression	Other/vascular
ADNI	55–90	24–30	Stable medication, AChEIs, memantine admitted, 6 grades education or work history	No functional impairment, but many with high FAQ scores. CDR=0.5; memory ≥ 0.5	Memory complaint LogMem II, dependent on education	Geriatric Depression Scale ≥ 6	Hachinski Ischemic Score IS >5
NACC	--	--	Similar to ADNI	Essentially normal daily functions	Cognitive complaint, cognitive decline (clinician's diagnosis)	Not specified	Not specified
InDDEx	55–85	--	No AChEI in previous 2 weeks, no rivastigmine in previous 4 weeks	Cognitive symptoms (not specified); CDR=0.5	NYU delayed paragraph recall <9	HDRS >12, HDRS item 1 > 1, DSM-IV major depression	AD criteria from DSM-IV or NINCDS-ADRDA mod. Hachinski Ischemic Score >4
CNG	≥ 50	≥ 20	A broader definition of MCI was used	Complaint of cognitive deficit in daily living; minor changes were tolerated: B-ADL < 4	Decline of cog. abilities (>1 SD) in at least one neuropsychological domain	Not specified	Not specified
BS-MC	N/A	N/A	Consecutively referred patients from GPs	Essentially Winblad et al. [26] criteria; no significant functional decline	Impairment (≤ -1.28 SD; age-, education-, and gender-adjusted) in \geq one cognitive domain	Probable cause for MCI other than early AD, based on comprehensive medical exam and neuroimaging results	Not specified

score of 0.5. Cognitive complaints or symptoms (without further specification) were requested in ADNI, NACC, and CNG. ADNI used thresholds in Logical Memory II-dependent on years of education. InDDEx requested a delayed recall score in the NYU-delayed paragraph recall [28] of less than 9, CNG and BS-MC requested at least one cognitive domain below -1 SD (CNG) or -1.28 SD (BS-MC). Patients with major depression were excluded in ADNI, InDDEx, and BS-MC. It will be noted that there are wide differences between the five samples with regard to almost all inclusion and exclusion criteria. It should also be mentioned that almost all datasets contain patients who did not fulfill one or more of their own study's inclusion and/or exclusion criteria (see comments to Table 3).

Cognitive tests

A selection of tests, most of them applied in at least two studies, is listed in Table 2. Each study applied a different set of tests. The Mini Mental Status Examination (MMSE) and the Verbal Fluency Test (animals) were the only instruments used in all studies. The ADAScog [17] was applied, although with different modifications, in ADNI, InDDEx, and CNG. Eight of nine tests of a neuropsychological battery used in ADNI were applied in NACC as well, but the Auditory Verbal Learning Test was omitted. Moreover, several procedural details of the Logical Memory delayed recall from the Wechsler Memory Scale (WMS) differed significantly in ADNI and NACC (details in [20]). The Boston Naming Test [22] was performed in ADNI and NACC with 30 items, but with only 15 items in CNG and BS-MC. The Digit span forward and backward and the Trail Making Test A and

B [29] were applied in all studies except in InDDEx; however, the Trail Making Tests were conducted with different time limits (e.g., for TMTA: 150 s in NACC and CNG, 180 s in BS-MC). The Clock Drawing Test was applied in all studies except in NACC, but different scoring methods were used. The Clinical Dementia Rating scale was applied in all studies but BS-MC. Three different versions of functional assessment were used. The CERAD battery was only used in CNG and BS-MC. A few other tests were applied in only one study.

Demographics (Table 3)

Patients in ADNI and NACC were oldest (74.2 ± 7.5 and 74.4 ± 7.9), those in InDDEx and BS-MC were younger (70.1 ± 7.6 and 69.7 ± 9.1), and CNG (68.0 ± 7.9) comprised the youngest sample. Duration of education was highest in ADNI (15.7 ± 3.0 years) and NACC (15.2 ± 3.0 years) and lowest in InDDEx (11.8 ± 3.5). CNG (12.3 ± 2.8) and BS-MC (12.0 ± 3.1) were in between. The proportion of females was lowest in ADNI (35.7%), highest in NACC (48.8%) and InDDEx (49.0%), and intermediate in CNG (45.7%) and BS-MC (45.8%). The proportion of patients with two ApoE4 alleles was highest in ADNI (11.9%), intermediate in NACC (8.7%), InDDEx (9.6%) and BS-MC (9.4%), and lowest in CNG (5.4%).

Baseline data

In each study, the quartile range of MMSE is from 26 to 29 (Fig. 1), with several outliers in NACC, InDDEx, CNG, and BS-MC, displaying much lower values that would preclude a diagnosis of MCI. Only ADNI had no outliers, because the minimum score of 24 was an inclusion criterion. Verbal fluency is 1.4 to 1.6 points higher

Table 2 List of selected cognitive tests applied in five studies

Test	ADNI	NACC	InDDEx	CNG	BS-MC
ADAScog 11 and modified	11 & mod.		11 & mod	11	
Logical Memory II	x	x		x	
Digit Span Forward	x	x		x	x
Digit Span Backward	x	x		x	x
Category Fluency, Animals	x	x	x	x	x
Category Fluency, Vegetables	x	x			
Trail Making Test B	x	x		x	x
Boston Naming Test	x	x	x ^a	x	x
Auditory Verbal Learning Test	x				
Digit Symbol	x	x			
Trail Making Test A	x	x		x	x
Clock Drawing Test	x		x	x	x
Functional Assessment	x	x	ADCS-ADL	Bayer-ADL	
Logical Memory I	x			x	
American National Adult Reading Test	x				
Clinical Dementia Rating	x	x	x	x	
Mini Mental Status Examination	x	x	x	x	x
Phonemic fluency, S-words					x
CERAD Wordlist + intrusion errors / savings				x	x
CERAD constructional praxis				x	x
Neuropsychiatric Inventory (NPI)	x	x	x	x	
Digit cancellation task			x		

^aOnly 85 values

in InDDEx, CNG, and BS-MC compared to ADNI and NACC (Table 3, Fig. 2). InDDEx patients performed better on the ADAScog (10.0±4.7) than CNG patients (11.7±5.1), but ADNI patients (11.5±4.4) performed similar to CNG patients. In InDDEx, ADNI, and CNG, there were patients with ADAScog values typical of

dementia rather than MCI. In the Boston Naming Test [22], patients of ADNI and NACC achieved on average about 25 of 30 items, while patients of CNG and BS-MC achieved about 13.4 of 15 items. Results on the Trail Making test cannot be compared because of different time limits used.

Table 3 Descriptive statistics of demographics and baseline scores MMSE, Verbal Fluency (animals), and ADAScog (11 subtests)

	ADNI n = 395	NACC n = 4328	InDDEx n = 831	CNG n = 726	BS-MC n = 1558
Age ($\bar{x} \pm SD$)	74.2±7.5	74.4±7.9	70.1±7.6	68.0±7.9	69.7±9.1
Education ($\bar{x} \pm SD$)	15.7±3.0	15.2±3.0	11.8±3.5	12.3±2.8	12.0±3.1
Female, n (%)	141 (35.7)	2113 (48.8)	422 (49.05)	332 (45.7)	713 (45.8)
ApoE ^a	n _{ApoE} =395	n _{ApoE} =2523	n _{ApoE} =396	n _{ApoE} =577	n _{ApoE} =53
1 allele no (%)	165 (41.8)	951 (37.7)	131 (32.3)	200 (34.7)	22 (41.5)
2 alleles no (%)	47 (11.9)	219 (8.7)	39 (9.6)	200 (34.7)	5 (9.4)
MMSE ($\bar{x} \pm SD$)	27.0±1.8	27.0±2.4	27.2±2.5	27.1±2.1	27.4±2.3
Min ^b -max	23–30	2–30	16–30	17–30	14–30
Verb. Fl. ($\bar{x} \pm SD$)	15.9±4.9	16.0±5.0	17.5±5.9	17.5±5.5	17.4±5.5
Min-max	5–30	0–35	2–38	3–32	3–38
ADAScog ($\bar{x} \pm SD$)	11.5±4.4	-	10.0±4.7	11.7±5.1	-
min-max	2–27.7		1–27	0–35	

^aApoE was not determined in all patients, number of evaluations is given as n_{ApoE}, percent of ApoE evaluations shown in parentheses. ^bMMSE is below inclusion criterions 24 in ADNI (n=1) and below 20 in CNG (n=5)

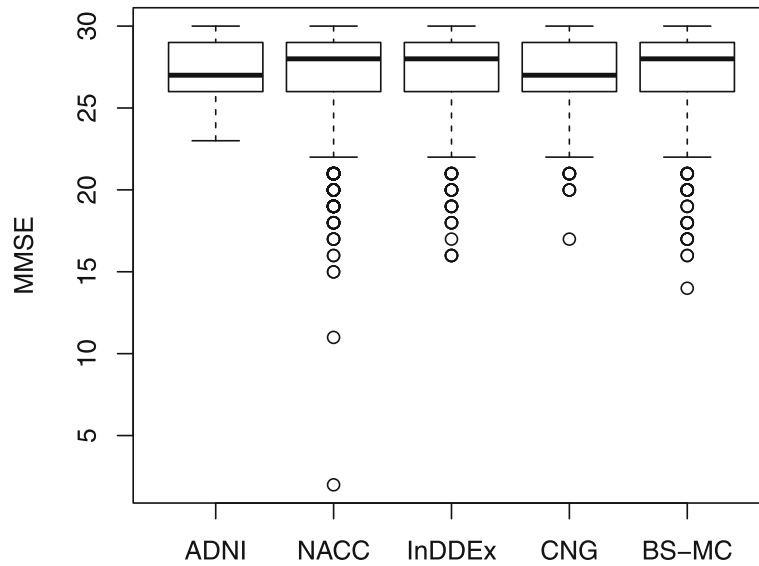


Fig. 1 Boxplot of Mini Mental Status Examination (MMSE) scores at baseline in each study

Progression of cognitive scores

Verbal fluency (animals) is—next to the MMSE—the only test score available in all five studies. From baseline to 1 year, InDDEx patients improved on average by 0.4 words, CNG patients remained stable, but ADNI, NACC, and BS-MC patients worsened by 0.6, 0.7, and 0.8 words, respectively (Fig. 3). Confidence intervals for InDDEx and the latter three studies are disjoint. The

differences between studies increased for the 2- and 3-year follow-up: InDDEx subjects still improved by 0.5 words, the latter three worsened by up to 2.1 (NACC) and 3.5 (BS-MC) words after 3 years (Fig. 3).

ADAScog worsened considerably in ADNI, slightly in CNG, but improved slightly in InDDEx. Boston Naming Test scores show similar decline in ADNI, NACC, and BS-MC (in relation to the number of items), while the

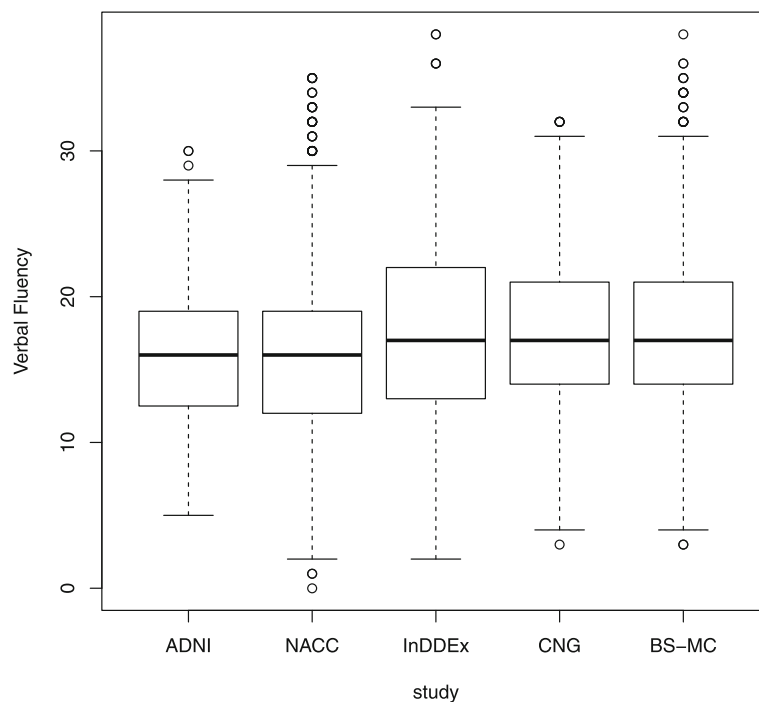


Fig. 2 Boxplot of Verbal Fluency scores (animals) at baseline in each study

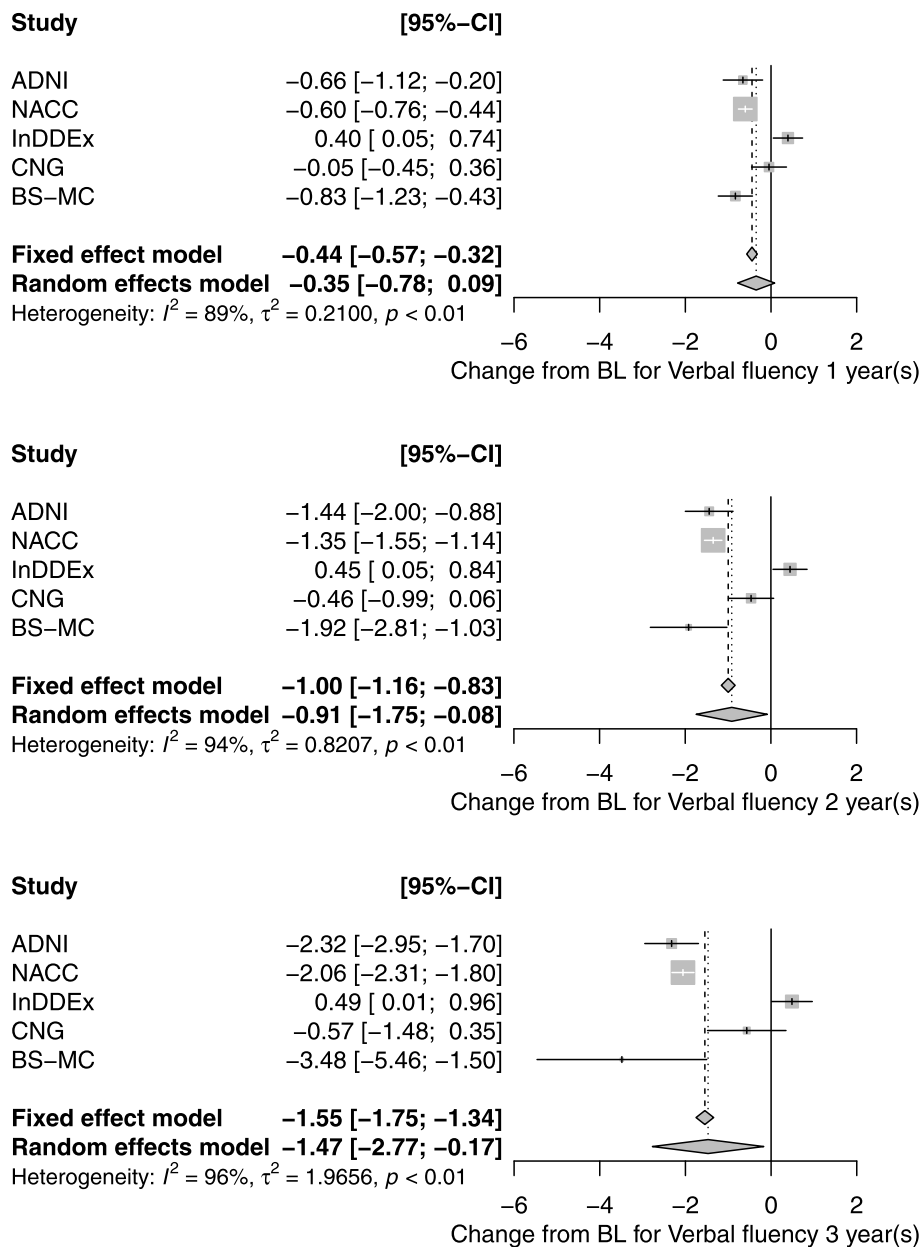


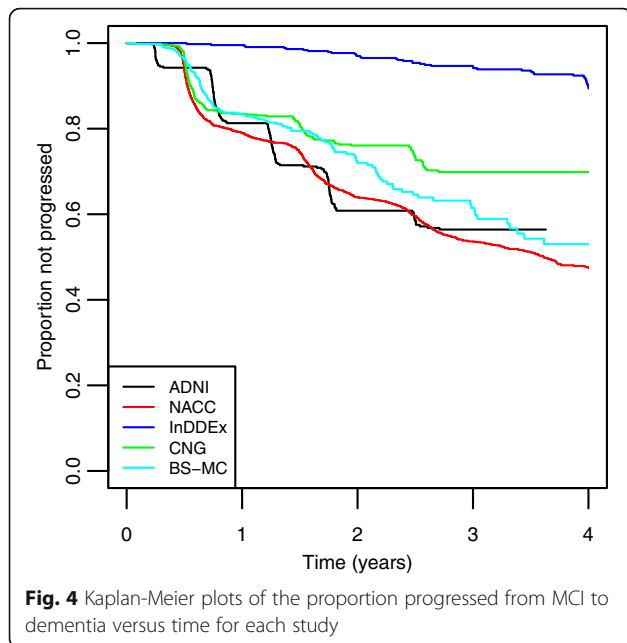
Fig. 3 Forest plots for the change of Verbal Fluency (animals) from baseline to 1, 2, and 3 years. Mean changes and 95% confidence intervals for each study and for the overall effect in the fixed effects and the random effects model are given. τ^2 is the between-study variance, I^2 measures heterogeneity (between study variance over total variance), p value for the test of heterogeneity. Graphs show study specific means and confidence intervals for each study as gray squares and lines and for overall effects as diamonds. Size of squares represents precision of individual treatment estimates

decline in CNG is very small. The CERAD-word list (delayed recall) worsened considerably in BS-MC, but remained unchanged in CNG.

Transition from MCI to dementia

Time to transition from MCI to dementia was usually ascertained at scheduled visits. Sometimes, however, patients were examined in an extra visit and then classified as

being demented. Kaplan-Meier curves for the time to transition from MCI to dementia (mostly AD) are shown in Fig. 4. They confirm that the InDDEx patients were in a particularly stable state. Continuation with an MCI diagnosis was considerably less frequent in the CNG sample, while ADNI and NACC patients were the fastest to progress. Transition rates after 3 years are between 5.9% (InDDEx) and 46.4% (NACC), both with a standard error of 1.1%.

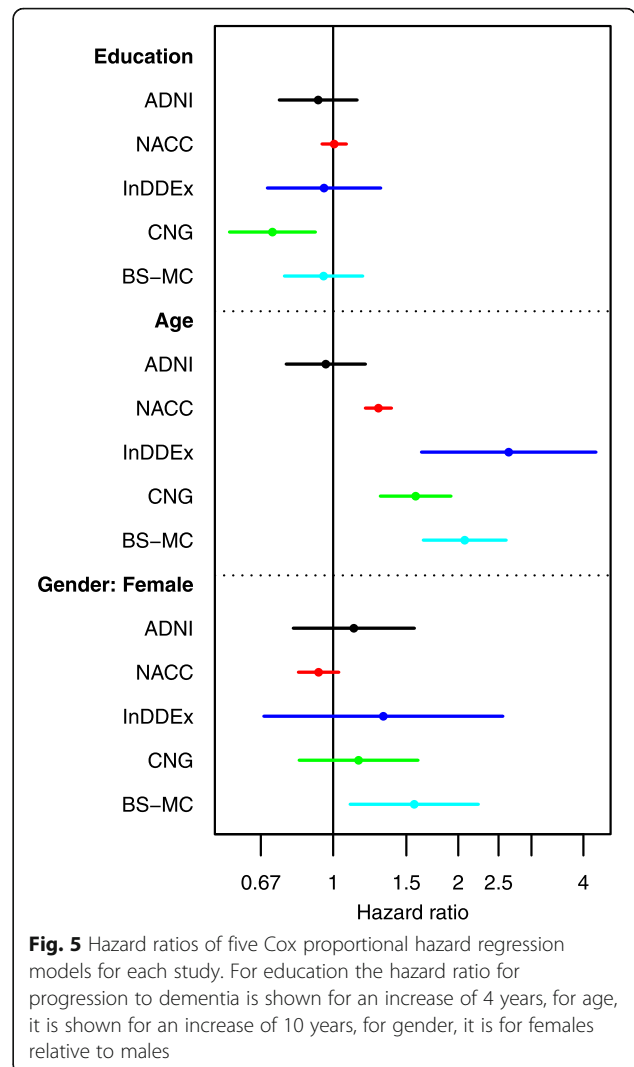


Cox regression models with covariates years of education, age, and sex were estimated for each study. Hazard ratios and 95% confidence intervals are shown in Fig. 5. The hazard ratio for education was 0.71 (for 4 years) in CNG and close to 1 in the other studies. The hazard ratio for age was close to 1 in ADNI and distinctly positive in the other studies. The hazard ratio for females was 1.56 in BS-MC and close to 1 in the other studies.

Discussion and conclusion

The PGSA (Placebo Group Simulation Approach) was submitted “for debate” to the Alzheimer community [1]. The novel study design was intended primarily to resolve an ethical problem of prevention trials involving aged subjects at risk of shifting from a pre-symptomatic into a dementia stage: the long-term use of placebo typical of RCTs (randomized controlled trials). Accordingly, instead of using a concomitant placebo group for comparison with a hopefully effective novel treatment, the PGSA applies mathematical algorithms to forecast the expected outcomes of pre-symptomatic AD patients from their baseline data and to compare those with the outcomes on experimental treatments. The PGSA was deemed to “have an advantage over the use of historical controls in futility designs in that it is based on the patient’s own observed clinical features.” [30].

Our published algorithms were derived from the ADNI database (<http://www.loni.ucla.edu/ADNI>) [31] that contained anamnestic, biological, neuroimaging, and neuropsychological findings from 397 North American patients with a diagnosis of MCI. Our analyses highlighted the strong impact of neuropsychological



performance data recorded at baseline, in addition to information from subjects’ history such as age, sex, and education, on MCI disease trajectories in the following years. A first attempt at validation of the PGSA algorithms using data from an independent MCI database (NACC; <https://www.alz.washington.edu>) confirmed the importance of neuropsychological performance data recorded at baseline to forecast cognitive decline in MCI [20]. However, we also noted that there was some slight over-estimation of cognitive decline when the ADNI-based PGSA algorithms were applied to the NACC MCI dataset for a follow-up of more than 2 years. This observation led to the question as to whether the published PGSA algorithms could be confidently applied to other longitudinal MCI data.

The current analysis comprised three MCI databases in addition to the ones from ADNI and NACC. One of these three originated from an RCT with a cholinesterase inhibitor sponsored by a drug company [13] and two

from clinical case series: one from a network of memory clinics in Germany [23], the other one from a single memory clinic in Basel, Switzerland [32]. The five databases contained information on 395 (ADNI) up to 4328 (NACC) individuals. Although all patients had a diagnosis of MCI, inspection of Table 1 shows that inclusion and exclusion as well as other eligibility criteria varied considerably between the five samples. Wide differences also existed between the five databases with regard to the neuropsychological scales and instruments used at baseline and at follow-up to document the changes in patients' cognitive performance (Table 2). Only the MMSE and the Verbal Fluency Test with "animals" were applied throughout.

Despite the differences in inclusion and exclusion criteria, the demographic and even more so the baseline cognitive performance data of the five patient samples were not as variable as one might expect. While patients' mean ages varied between 68.0 (CNG) and 74.4 (NACC) years and the ADNI sample contained a lower percentage of female participants than the other four, mean MMSE scores at baseline varied only minimally (Table 3, Fig. 1), and the same was seen with regard to the ADAS scores in the three studies where this scale was used. Nevertheless, judging from their MMSE and/or ADAS-cog scores, a number of patients in the NACC, InDDEx, CNG, and BS-MC samples should be classified as being demented rather than having MCI. Interestingly, the participants of InDDEx, CNG, and BS-MC had slightly better performance on the Verbal Fluency Test than those of ADNI and NACC (Table 3, Fig. 2), although the latter had benefited from more years of education on the average.

In contrast to the similarity of their cognitive performance data noted at baseline, the five samples differed markedly with regard to their subsequent cognitive performance and clinical development. As seen in Fig. 3, scores on the Verbal Fluency Test behaved differently in the InDDEx than in the other four groups: with increasing study duration performance deteriorated continuously in the ADNI, NACC, CNG, and BS-MC samples, but were slightly improved from baseline after 1 year and then stayed stable in the InDDEx group. An even greater disparity is seen with regard to the number of transitions from MCI to dementia (Fig. 4): Whereas some 90% of the InDDEx patients did not progress to dementia within the 4 years of observation, the respective percentages were around 70 for CNG, around 60 for ADNI and BS-MC and somewhat above 50 for NACC. What could be an explanation of these big differences?

First, one should note that the difference in transition rates between the ADNI and the NACC participants is small. This is not a surprise given that both MCI patient samples were collected in North America and that the

ADNI sample partly constituted a selection from the large NACC data collection. Thus, it is likely that both the MCI inclusion/exclusion and the transition criteria applied to the ADNI and the NACC data were similar. Transition rates for the BS-MC sample were also close to the ones seen in ADNI and NACC, suggesting that the inclusion/exclusion criteria for MCI and the criteria to diagnose transition to dementia were applied similarly in the Basel and in the North American centers. More difficult to understand are the lower transition rates in CNG and, particularly, in InDDEx. The latter dataset differed from the other four in that InDDEx was not a case series from one or more memory clinics but a clinical drug trial sponsored by a pharmaceutical company and carried out in 12 different countries in Europe, South Africa, South America, and in the USA. While one cannot exclude that this geographic variety (or perhaps some investigators' desire to include as many patients as possible) compromised proper selection of MCI patients, it is of note that other, although shorter, company-sponsored drug studies carried out in the same decade as InDDEx also differed markedly with regard to transition rates from MCI to dementia. Thus, Petersen et al. [8] reported annual transition rates of 16% in a 3-year study with donepezil and vitamin E, and relatively high percentages of transitions were also seen in two separate 2-year studies with galantamine [15]. In contrast, an RCT with a selective COX-2 inhibitor noted rather low annual transition rates: 6.4% on rofecoxib, 4.5% on placebo [33]—although these rates were still higher than the one reported in [13].

Whatever the reasons for the varying transition rates in these mostly company-sponsored trials are, it is obvious that the MCI criteria available some 20 years ago were not sufficiently precise to allow selection of clinically homogeneous and internationally comparable MCI patient groups. Ward et al. [34] and Han et al. [35] noted that—owing to the fuzzy boundaries between normal, MCI and dementia—all available estimations of the incidence and prevalence of MCI varied widely. Edmonds et al. [36] and Stephan et al. [37] stressed that samples of non-specified MCI cases will include individuals with different brain pathologies, leading to widely different clinical trajectories. As a consequence thereof, the predictive power of MCI diagnoses is poor, allowing, e.g., up to 59% of patients with MCI to revert to normal within up to 17 years after a first diagnosis [38]. The problem to delineate MCI properly was noted early on [9, 39], and efforts were made to ameliorate the situation (e.g., [26]). Nevertheless, ambiguity remained: Although an international expert group [10] stressed the importance of determining whether there is objective evidence of cognitive decline and recommended cognitive testing for quantitatively assessing the degree of cognitive

impairment for a diagnosis of MCI, these authors also emphasized that normative ranges of neuropsychological tests (typically those listed in Table 2) “are guidelines and not cutoff scores.” ([10] p. 272). In contrast to this position, it is our opinion that the scientific community needs to agree on age- and education-adjusted cutoff scores in order to make MCI a scientifically useful concept and to ensure that study results from different sources will be comparable. It simply cannot be that individuals with MMSE scores of less than 20 or ADAScog scores of more than 20 be included in so-called MCI patient samples (Table 3).

Another critical issue is the use of specific neuropsychological tests: as noted in the current analyses, only two tests were part of all five studies considered, making comparison between patient samples virtually impossible. Moreover, for several tests, different versions and scoring systems exist and are being used. This is another impediment when one tries to compare results between studies. In the USA, a series of Uniform Data Sets (UDS) have been proposed in [40] to improve this unsatisfactory situation. A separate issue refers to the underlying disorder of MCI and subsequent dementia. While neuropsychological test performance with an internationally agreed set of tests and cutoffs will allow for the determination of cognitive deficits, specific biomarker requirements would also allow to describe pathophysiologically more homogenous groups of patients with MCI, e.g., MCI due to AD. As of today, and as shown in our analyses, the heterogeneity at all these levels is way too big to allow meaningful conclusions with regard to the scientific merits of algorithm-based approaches such as the PGSA.

The scientific rationale of our earlier studies [1, 15] was to make use of well-defined historical data in clinical treatment or prevention trials. Inclusion of historical information for comparison with current treatment data has been discussed since more than 40 years ([41] and later references in [7]). Suggestions range from performing a single arm study which is compared to the historical control, over integrating historical control data with new controls—up to avoiding any use of historical data completely. Different options are available to integrate historical data [3]: (1) pooling them with new controls, (2) testing for differences between historical and new controls and pool them only if no differences are detected, (3) down-weighting the historical data by power priors dependent on the discrepancy between observed and historical control data, (4) choose a prior distribution for the means of the historical and the new control groups and apply a Bayesian model (see [3] for details on these methods), and (5) perform a random effects meta-analysis of the historical controls and down-weight their sample size according to the between-study variation [7].

In any case, data pooling is only permissible if the historical controls are exactly equivalent to the new control. This may hold for a set of pharmaceutical studies with basically the same protocol and equivalent patient populations. In less narrowly defined circumstances, this is rarely the case, due to different populations and more or less different inclusion and exclusion criteria. Down-weighting of historical data takes heterogeneity into account. It decreases the weight of the historical data from the total number of patients to a smaller number which is called the *prior effective sample size* [7]. For the five studies considered in the current analysis, the 2144 patients for whom data for change of verbal fluency after 3 years were available, would be down-weighted to 10 patients, according to formulas [7]. However, all methods of integrating historical controls are only valid under the assumption of exchangeability, i.e., if no systematic differences exist between the control groups [7, 42]. In view of the distinctly apart confidence intervals in Fig. 3, exchangeability cannot be assumed for the patient samples considered here. This makes all attempts to integrate these data in new studies futile. Insofar, our current failure of forecasting algorithms is consistent with theoretical considerations on using prior information.

Limitations

Drop-outs are a common problem in long-term clinical studies. In the five studies considered, the rate of drop-outs after 3 years was quite different: 70% in NACC, 37% in ADNI, 25% in InDDEx, 77% in CNG, and 98% in BS-MC. While studies with a stricter visit regimen (in our case InDDEx and ADNI) have lower drop-out, convenience samples such as CNG and BS-MC tend to have very high drop-out rates. This factor may cause bias and could make studies less comparable. For example, one might assume that convenience samples contain a larger number of frail patients than samples in controlled studies, that frail patients drop-out with higher probability, and that, as a consequence, the results of cognitive tests would worsen less in convenience samples. However, our results do not support this hypothesis: Patients of the InDDEx study (with the lowest dropout rate) improved their cognitive results after 2 and 3 years, whereas patients in BS-MC showed the most pronounced decrease (cf. Fig. 3). CNG and NACC, the two studies that collected data from Alzheimer's coordinating centers and/or memory clinics, also showed distinctly different changes (Fig. 3). Transition rates in the InDDEx sample were very low, and the transition rate of the CNG sample was between InDDEx and the other studies (cf. Fig. 4). This conforms to the cognitive test results, but, again, does not support the hypothesis concerning bias due to drop-out.

Transitions to dementia were in most studies ascertained at more or less strictly scheduled visits. We nevertheless applied Kaplan-Maier and Cox Regression analysis assuming continuous time. Figure 4 shows the pattern of event times for strict visit schedules of ADNI and InDDEX and less clearly for NACC and CNG. The true curves would interpolate between the top right bends of the curves in Fig. 5, but this would not change the interpretation of the transition rates.

Final remark

A basic requirement of the PGSA and similar approaches is the availability of uniformly collected, high-quality data of the respective population, in our case patients with MCI. Our analysis of five differently assembled MCI datasets shows that this requirement is not met, and access to other, more recently collected databases have proven to be somewhat cumbersome. As pointed out by one of the reviewers of this paper, the availability of data in the AD field is low and this “is really problematic. In such a dreadful disease, data should be made available to any researcher in an effort to produce harmonized and robust modeling of the disease....We desperately need an open science approach in our field in which both academia and the industry would participate by sharing anonymized data. The ADNI model and thousands of publications it allowed is an indication of how much that is needed.” We definitely agree with this statement.

Acknowledgements

Data collection and sharing for this project was funded by the ADNI (National Institutes of Health Grant U01 AG024904). The ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: Abbott, AstraZeneca AB, Bayer Schering Pharma AG, Bristol-Myers Squibb, Eisai Global Clinical Development, Elan Corporation, Genentech, GE Healthcare, GlaxoSmithKline, Innogenetics, Johnson and Johnson, Eli Lilly and Co., Medpace, Inc., Merck and Co., Inc., Novartis AG, Pfizer Inc., F. Hoffman-La Roche, Schering-Plough, Synarc, Inc., as well as nonprofit partners the Alzheimer's Association and the Alzheimer's Drug Discovery Foundation, with participation from the US Food and Drug Administration. Private-sector contributions to the ADNI are facilitated by the Foundation for the National Institutes of Health. The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's Disease Cooperative Study at the University of California, San Diego. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of California, Los Angeles. This research was also supported by NIH grants P30 AG010129, K01 AG030514, and the Dana Foundation. The support provided by Michael Wagner and Steffen Wolfsgruber from the Competence Network Germany (CNG), by Walter Kukull and Sarah Monsell from the National Alzheimer Competence Center (NACC), and by Peter Quarg from Novartis, Basel, at an earlier stage of our project is gratefully acknowledged.

Authors' contributions

M.B. performed the data analysis and contributed the “Material and Methods,” “Results,” and “Limitations.” R.S. is the “spiritual father” of the PGSA and contributed the “Introduction” and most of the “Discussion,” A.U.M. added to the “Abstract” and the “Discussion.” The authors approved the final manuscript.

Funding

This study was supported by grants from the Alzheimer's Association Switzerland and the Alzheimer Forum Switzerland. Open Access funding enabled and organized by Projekt DEAL.

Availability of data and materials

The data that support the findings of this study are available from Alzheimer's Disease Neuroimaging Initiative (ADNI; <http://www.loni.ucla.edu/ADNI>), National Alzheimer's Coordinating Centers (NACC; <https://www.alz.washington.edu>), Novartis AG, Basel, Competence Network Germany [23], and University Department of Geriatric Medicine FELIX-PLATTER (A.U.M), but restrictions apply to the availability of these data, which were used under license for the current study, and so are not publicly available. Data are however available from the authors upon reasonable request and with permission of third parties mentioned before.

Declarations

Ethics approval and consent to participate

Ethics approval for using the data of the University Department of Geriatric Medicine FELIX-PLATTER has been obtained from the Ethics Committee (Ethikkommission Nordwest- und Zentralschweiz). For the other datasets, ethics approval is a prerequisite of these studies.

Consent for publication

This is not applicable beyond ethics approval.

Competing interests

The authors declare that they have no competing interests.

Received: 12 February 2021 Accepted: 19 April 2021

Published online: 07 May 2021

References

- Spiegel R, Berres M, Miserez AR, Monsch AU. For debate: substituting placebo controls in long-term Alzheimer's prevention trials. *Alz Res Ther*. 2011;3(2):9–20. <https://doi.org/10.1186/alzrt68>.
- Weiner MW, Veitch DP, Aisen PS, Beckett LA, Cairns NJ, Green RC, et al. The Alzheimer's Disease Neuroimaging Initiative: a review of papers published since its inception. *Alzheimer's Dement*. 2011;7:1–67.
- Viele K, Berry S, Neuenschwander B, Amzal B, Chen F, Enas N, et al. Use of historical control data for assessing treatment effects in clinical trials. *Pharm Stat*. 2014;13(1):41–54. <https://doi.org/10.1002/pst.1589>.
- Miller RG, Moore DH, Forsheve DA, Katz JS, Barohn RJ, Valan M, et al. Phase II screening trial of lithium carbonate in amyotrophic lateral sclerosis. *Neurology*. 2011;77(10):973–9. <https://doi.org/10.1212/WNL.0b013e31822dc7a5>.
- Grill JD, Karlawish J. Addressing the challenges to successful recruitment and retention in Alzheimer's disease clinical trials. *Alzheimer's Res Ther*. 2010;2(6):34–44. <https://doi.org/10.1186/alzrt58>.
- Cudkovic ME, Katz J, Moore DH, O'Neill G, Glass JD, et al. Toward more efficient clinical trials for amyotrophic lateral sclerosis. *Amyotroph Lateral Scler*. 2010;11(3):259–65. <https://doi.org/10.3109/17482960903358865>.
- Neuenschwander B, Capkun-Niggli G, Branson M, Spiegelhalter D. Summarizing historical information on controls in clinical trials. *Clin Trials*. 2010;7(1):5–18. <https://doi.org/10.1177/1740774509356002>.
- Petersen RC, Thomas RG, Grundman M, Bennett D, Doody R, Ferris S, et al. Vitamin E and donepezil for the treatment of mild cognitive impairment. *New Engl J Med*. 2005;352(23):2379–88. <https://doi.org/10.1056/NEJMoa050151>.
- Visser PJ, Scheltens P, Verhey FRJ. Do MCI criteria in drug trials accurately identify subjects with pre-dementia Alzheimer's disease? *J Neurol Neurosurg Psychiatry*. 2005;76(10):1348–54. <https://doi.org/10.1136/jnnp.2004.047720>.
- Albert MS, DeKosky ST, Dickson D, Dubois B, Feldman HH. The diagnosis of mild cognitive impairment due to Alzheimer's disease: recommendations from the National Institute on Aging-Alzheimer's Association workgroups on diagnostic guidelines for Alzheimer's disease. *Alzheimer's Dement*. 2011;7(3):270–9. <https://doi.org/10.1016/j.jalz.2011.03.008>.
- Morris JC. Revised criteria for mild cognitive impairment may compromise the diagnosis of Alzheimer disease dementia. *Arch Neurol*. 2012;69(6):700–8. <https://doi.org/10.1001/archneurol.2011.3152>.

12. Petersen RC, Thomas RG, Aisen PS, Mohs RC, Carrillo MC, Albert MS, et al. Randomized controlled trials in mild cognitive impairment. Sources of variability. *Neurology*. 2017;88(18):1751–8. <https://doi.org/10.1212/WNL.0000000000003907>.
13. Feldman HH, Ferris S, Winblad B, Sfikas N, Mancione L, He Y, et al. Effect of rivastigmine on delay to diagnosis of Alzheimer's disease from mild cognitive impairment: the InDDEx study. *Lancet Neurol*. 2007;6(6):501–12. [https://doi.org/10.1016/S1474-4422\(07\)70109-6](https://doi.org/10.1016/S1474-4422(07)70109-6).
14. Rodda JWZ. Ten years of cholinesterase inhibitors. Editorial. *Int J Geriatr Psychiat*. 2009;24(5):437–42. <https://doi.org/10.1002/gps.2165>.
15. Winblad B, Gauthier S, Scinto L, Feldman H, Wilcock GK, Truyen L, et al. Safety and efficacy of galantamine in subjects with mild cognitive impairment. *Neurology*. 2008;70(22):2024–35. <https://doi.org/10.1212/01.wnl.0000303815.69777.26>.
16. Kornhuber J, Schmidtke K, Frölich L, Perneczky R, Wolf S, et al. Early and differential diagnosis of dementia and mild cognitive impairment. *Dement Geriatr Cogn Disord*. 2009;27(5):404–17. <https://doi.org/10.1159/000210388>.
17. Rosen WG, Mohs RC, Davis KL. A new rating scale for Alzheimer's disease. *Am J Psychiatry*. 1984;141:1356–64.
18. Folstein MF, Folstein SE, McHugh PR. "Mini-mental state". A practical method for grading the cognitive state of patients for the clinician. *J Psychiatr Res*. 1975;12(3):189–98. [https://doi.org/10.1016/0022-3956\(75\)90026-6](https://doi.org/10.1016/0022-3956(75)90026-6).
19. Pfeffer RI, Kurosaki TT, Harrah CH Jr, Chance JM, Filos S. Measurement of functional activities in older adults in the community. *J Gerontol*. 1982;37(3):323–9. <https://doi.org/10.1093/geronj/37.3.323>.
20. Berres M, Kukull WA, Miserez RA, Monsch AU, Monsell SE, Spiegel R, et al. A novel study paradigm for long-term prevention trials in Alzheimer disease: the Placebo Group Simulation Approach (PGSA). Application to MCI data from the NACC database. *J Prevent Alzheimer's Dis*. 2014;1:99–109.
21. Weintraub S, Salmon DS, Mercaldo N, Ferris S, Graff-Radford NR, Chui H, et al. The Alzheimer's Disease Centers' Uniform Data Set (UDS) The Neuropsychologic Test Battery. *Alzheimer Dis Assoc Disord*. 2009;23(2):91–101. <https://doi.org/10.1097/WAD.0b013e318191c7dd>.
22. Kaplan E, Goodglass H, Weintraub S. The Boston naming test. Philadelphia: Lea & Febiger; 1983.
23. Wolfsgruber S, Wagner M, Schmidtke K, Frölich L, Kurz A, Schulz S, et al. Memory concerns, memory performance and risk of dementia in patients with mild cognitive impairment. *Plos One*. 2014;9(7):e100812. <https://doi.org/10.1371/journal.pone.0100812>.
24. Morris JC, Heyman A, Mohs RC, Hughes JP, van Belle G, Fillenbaum G, et al. The Consortium to Establish a Registry for Alzheimer's Disease (CERAD): I. Clinical and neuropsychological assessment of Alzheimer's disease. *Neurology*. 1989;39(9):1159–65. <https://doi.org/10.1212/wnl.39.9.1159>.
25. Schmid NS, Ehrensperger MM, Berres M, Beck IR, Monsch AU. The Extension of the German CERAD Neuropsychological Assessment Battery with tests assessing subcortical, executive and frontal functions improves accuracy in dementia diagnosis. *Dement Geriatr Cogn Dis Extra*. 2014;4(2):322–34. <https://doi.org/10.1159/000357774>.
26. Winblad B, Palmer K, Kivipelto M, Jelic V, Fratiglioni L, et al. Mild cognitive impairment – beyond controversies, towards a consensus: report of the International Working Group on Mild Cognitive Impairment. *J Intern Med*. 2004;256:240–6.
27. Morris JC. The Clinical Dementia Rating (CDR): current version and scoring rules. *Neurology*. 1993;43(11):2412–4. <https://doi.org/10.1212/WNL.43.11.2412-a>.
28. Kluger A, Ferris SH, Golomb J, Mittelman MS, Reisberg B. Neuropsychological prediction of decline to dementia in nondemented elderly. *J Geriatr Psychiatr Neurol*. 1999;12(4):168–79. <https://doi.org/10.1177/089198879901200402>.
29. Army Individual Test Battery. Manual of directions and scoring. Washington, DC: War Department, Adjutant General's Office; 1944.
30. Cummings J, Gould H, Zhong K. Advances in designs for Alzheimer's disease clinical trials. *Amer J Neurodeger Dis*. 2012;1:205–16.
31. Alzheimer's Disease Neuroimaging Initiative. 2017. <http://www.loni.ucla.edu/ADNI>. Accessed 09-14-2020.
32. Monsch AUKRW. Specific care program for the older adults: memory clinics. *Eur Geriatr Med*. 2010;1:28–31.
33. Thal LJ, Ferris SH, Kirby L, Block GA, Lines CR, et al. A randomized double-blind, study of rofecoxib in patients with mild cognitive impairment. *Neuropharmacology*. 2005;30:1204–15.
34. Ward A, Arrighi HM, Shannon M, Cedarbaum JM. Mild cognitive impairment: disparity of incidence and prevalence estimates. *Alzheimer's & Dement*. 2012;8(1):14–21. <https://doi.org/10.1016/j.jalz.2011.01.002>.
35. Han JW, So Y, Kim TH, Lee DY, Ryu S-H, Kim SY, et al. Prevalence rates of dementia and mild cognitive impairment are affected by the diagnostic parameter changes for neurocognitive disorders in the DSM-5 in a Korean population. *Dement Geriatr Cogn Disord*. 2017;43:193–203.
36. Edmonds EC, Delano-Wood L, Clark LR, Aj J, Nation DA, CR MD, et al. Susceptibility of conventional criteria for mild cognitive impairment to false positive diagnostic errors. *Alzheimer's Dement*. 2015;11(4):415–22. <https://doi.org/10.1016/j.jalz.2014.03.005>.
37. Stephan BCM, Matthews FE, Hunter S, Savva GM, Bond J, IG MK, et al. Neuropathologic profile of mild cognitive impairment from a population perspective. *Alzheimer Dis Assoc Disord*. 2012;26(3):205–12. <https://doi.org/10.1097/WAD.0b013e31822fc24d>.
38. Malek-Ahmadi M. Reversion form mild cognitive impairment to normal cognition – a meta analysis. *Alzheimer Dis Assoc Disord*. 2016;30(4):324–30. <https://doi.org/10.1097/WAD.0000000000000145>.
39. Visser PJ, Kester A, Jolles J, Verhey F. Ten-year risk of dementia in subjects with mild cognitive impairment. *Neurology*. 2006;67(7):1201–7. <https://doi.org/10.1212/01.wnl.0000238517.59286.c5>.
40. Besser L, Kukull W, Knopman DS, Chui H, Galasko D, Weintraub S, et al. Version 3 of the National Alzheimer's Coordinating Center's Uniform Data Set. *Alzheimer Dis Assoc Disord*. 2018;32:351–8. <https://doi.org/10.1097/WAD.0000000000000279>.
41. Pocock S. Combination of randomized and historical controls in clinical trials. *J Chron Dis*. 1976;29(3):175–88. [https://doi.org/10.1016/0021-9681\(76\)90044-8](https://doi.org/10.1016/0021-9681(76)90044-8).
42. Schmidli H, Gsteiger S, Roychoudhury S, O'Hagan A, Spiegelhalter D, Neuenschwander B. Robust meta-analytic-predictive priors in clinical trials with historical control information. *Biometrics*. 2014;70(4):1023–32. <https://doi.org/10.1111/biom.12242>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://www.biomedcentral.com/submissions)

